

## Correct-by-Construction Controller Synthesis for Safety-Critical Systems

Victor Gaßmann

Vollständiger Abdruck der von der TUM School of Computation, Information and Technology der Technischen Universität München zur Erlangung des akademischen Grades eines

**Doktors der Ingenieurwissenschaften (Dr.-Ing.)**

genehmigten Dissertation.

**Vorsitz:**

Prof. Dr. Debarghya Ghoshdastidar

**Prüfende der Dissertation:**

1. Prof. Dr.-Ing. Matthias Althoff
2. Prof. Dr. Majid Zamani,  
University of Colorado Boulder, USA

Die Dissertation wurde am 01.02.2024 bei der Technischen Universität München eingereicht und durch die TUM School of Computation, Information and Technology am 08.04.2024 angenommen.



# Foreword

This thesis summarizes my research over the past years at the Technical University of Munich. First, I would like to thank my supervisor Prof. Matthias Althoff who gave me the chance to pursue this PhD. His academic guidance and invaluable feedback has profoundly impacted the quality of my work, for which I extend my sincere appreciation and gratitude.

I also thank all of my colleagues for the great moments we had together during my time at the chair. The atmosphere was always welcoming and all the fun activities, afternoons with cake, and stimulating conversations were a nice change of pace, especially when things did not go according to plan. I would especially like to thank Lukas Schäfer and Mark Wetzlinger, who always had an open ear and provided me with many invaluable comments and suggestions; I always very much enjoyed our fruitful and honest conversations. I would also like to thank Dr. Niklas Kochdumper and Adrian Kulmburg for their helpful suggestions and valuable insight.

Last, but certainly not least, I am especially grateful to my family and friends. Thank you to my parents Theresia and Volker who always encouraged and supported me. And thank you to my brother David and his wife Lena who were always very inclusive, and showed sympathy and understanding when I had hard times.



# Abstract

Due to recent advances in the field of formal verification, the safety of complex autonomous systems, such as a spaceship attempting to dock or a robot collaborating with humans, may be formally verified. The verification of such reach-avoid problems, where a set of initial states is to be steered to a given target under input and state constraints, is often non-constructive since it only provides a yes-or-no answer. As a result, first synthesizing the controller and then adapting it to guarantee safety turns out to be difficult. In this thesis, we thus combine numerical optimization and set-based reachability analysis to synthesize controllers that ensure formal guarantees by construction.

The efficient optimization of these controllers requires differentiability of the considered synthesis problems. Hence, we first introduce an abstraction of these synthesis problems, for which differentiability can be shown and which then implies differentiability of all considered synthesis problems.

Since an exact formulation of these problems is often hard to solve, we need to derive tractable approximations thereof. As we compute with sets, this often requires the use of different set representations, as they are closed under different set operations. To that end, we present the efficient conversions between different set representations. Additionally, we describe set containment checks from the literature, which are required to check input and constraint satisfaction for sets of states.

To synthesize piecewise constant controllers, we start by extending existing approaches from previous work. To improve controller performance, for instance, we extend the class of control laws that can be synthesized. Further, we propose a novel synthesis algorithm based on trust regions, which iteratively synthesizes a piecewise constant controller under input and state constraints. In contrast to previous work, this synthesis algorithm does not require additional user input to achieve non-conservative results for state-constrained systems. All presented algorithms have polynomial complexity in the state dimension and we demonstrate their advantages compared to other state-of-the-art approaches using numerical experiments.

To incorporate continuous state feedback, we present a novel algorithm for the combined synthesis of a piecewise constant feedforward controller with continuous state feedback. In contrast to other state-of-the-art approaches, we do not synthesize the feedforward controller and the state feedback separately but combine the synthesis into a single optimization problem and thus avoid the introduction of additional user-provided algorithm parameters. The complexity of the proposed algorithm is polynomial in the state dimension. As numerical experiments demonstrate, it outperforms existing synthesis approaches while requiring reduced user input compared to the current state of the art.



# Zusammenfassung

Die Sicherheit komplexer autonomer Systeme – wie beispielsweise das Andocken von Raumschiffen oder die Kollaboration zwischen Mensch und Roboter – kann dank jüngster Fortschritte im Bereich der formalen Verifikation sichergestellt werden. Der Sicherheitsbeweis solcher “reach-avoid”-Probleme, bei welchen die initiale Zustandsmenge eines Systems zu einem gegebenen Zielzustand unter Einhaltung gegebener Spezifikationen gesteuert werden soll, liefert dabei häufig lediglich eine nicht-konstruktive Ja-Nein-Aussage. Somit ist die Trennung von Design und anschließender Verifikation des geregelten Systems oft schwierig. In dieser Arbeit kombinieren wir daher numerische Optimierung mit mengenbasierter Erreichbarkeitsanalyse, wodurch die so berechneten Regler die gewünschten formalen Garantien “by design” erfüllen.

Eine effiziente Optimierung dieser Regler setzt allerdings die Differenzierbarkeit des zugrundeliegenden Synthese-Problems voraus. Daher führen wir zu Beginn ein abstrahiertes Synthese-Problem ein, welches alle relevanten Optimierungsprobleme generalisiert und zeigen die Differenzierbarkeit dieser Abstrahierung, was damit direkt die Differenzierbarkeit aller abstrahierten Synthese-Probleme impliziert.

Da eine exakte Formulierung dieser Probleme oft schwierig zu lösen ist, müssen effizientere Approximationen konstruiert werden. Da wir in dieser Arbeit mit Mengen rechnen, ist dazu oft die Zuhilfenahme verschiedener Mengenrepräsentationen nötig, da unterschiedliche Repräsentationen unter verschiedenen Mengenoperationen geschlossen sind. Daher zeigen wir die effiziente Konvertierung verschiedener Mengenrepräsentationen. Zur Überprüfung von Zustands- und Eingangsbeschränkungen für Mengen beschreiben wir außerdem Methoden zur Überprüfung der Teilmengen-Eigenschaft zweier Mengen aus der Literatur.

Zur Synthese von stückweise konstanten Regelgesetzen beginnen wir zunächst mit der Erweiterung existierender Ansätze zur Regler-Synthese aus vorangegangenen Arbeiten. Um die Leistung der berechneten Regler zu verbessern, erweitern wir beispielsweise die Klasse an berechenbaren Regelgesetzen. Wir präsentieren zudem einen neuartigen Synthese-Algorithmus, welcher – basierend auf sogenannten “Trust Regions” – iterativ ein stückweise konstantes Regelgesetz berechnet, welches sowohl Zustands- als auch Eingangsspezifikationen berücksichtigt. Dabei benötigt der präsentierte Algorithmus – im Gegensatz zu vorhandenen Ansätzen – neben den gegebenen Spezifikationen keine zusätzlichen Angaben des Nutzers, um performante Regelgesetze zu berechnen. Die beschriebenen Methoden haben hierbei alle eine polynomielle Laufzeit im Bezug auf die Zustandsdimension. Wir demonstrieren die Vorteile der entwickelten Algorithmen im Vergleich zu existierenden Methoden mit Hilfe numerischer Experimente.

Zur Kombination des stückweise konstanten Regelgesetzes mit einem kontinuierlichem Feedback-Regler beschreiben wir dann einen Synthese-Algorithmus für das kombinierte

## *Zusammenfassung*

Regelgesetz. Im Gegensatz zu bekannten Ansätzen synthetisieren wir dabei den stückweise konstanten und den kontinuierlichen Feedback-Regler in einem Optimierungsproblem, wodurch wir zusätzliche Parameter vermeiden, die vom Nutzer spezifiziert werden müssten. Der beschriebene Algorithmus hat polynomielle Laufzeit im Bezug auf die Zustandsdimension, benötigt im Vergleich zu existierenden Ansätzen reduzierten Nutzereingriff, und zeigt in numerischen Experimenten eine verbesserte Regelperformance.



# Contents

<b>Foreword</b>	<b>iii</b>
<b>Abstract</b>	<b>v</b>
<b>Zusammenfassung</b>	<b>vii</b>
<b>List of Figures</b>	<b>xiii</b>
<b>List of Tables</b>	<b>xv</b>
<b>Acronyms</b>	<b>xvii</b>
<b>Notations</b>	<b>xix</b>
<b>Computing Platform and Implementation</b>	<b>xxv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 State of the Art . . . . .	2
1.2.1 Optimal Control . . . . .	2
1.2.2 Model Predictive Control . . . . .	3
1.2.3 Abstraction-Based Synthesis . . . . .	4
1.2.4 Reachability-Based Synthesis . . . . .	4
1.3 Outline . . . . .	5
<b>2 Background</b>	<b>7</b>
2.1 Standard Operations . . . . .	7
2.2 Differentials . . . . .	9
2.3 Optimization Theory . . . . .	12
2.3.1 Karush-Kuhn-Tucker Conditions . . . . .	13
2.3.2 Constraint Qualifications . . . . .	13
2.3.3 Convex Optimization . . . . .	14
2.3.4 Non-Convex Optimization . . . . .	18
2.3.5 Lagrangian Duality . . . . .	18
2.3.6 Modeling . . . . .	19
2.4 Linear-Quadratic Regulator . . . . .	21
2.5 Set Operations . . . . .	22
2.6 Set Representations . . . . .	23
2.6.1 Intervals . . . . .	23

2.6.2	Arbitrary Convex Set . . . . .	25
2.6.3	Ellipsoids . . . . .	25
2.6.4	H-Polytopes . . . . .	27
2.6.5	Zonotopes . . . . .	29
2.6.6	Polynomial Zonotopes . . . . .	31
2.7	Reachability Analysis . . . . .	34
2.7.1	Linear Time-Invariant Systems . . . . .	35
2.7.2	Nonlinear Systems using Conservative Linearization . . . . .	38
2.7.3	Nonlinear Systems using Conservative Polynomialization . . . . .	39
2.8	Reduction Techniques . . . . .	40
<b>3</b>	<b>The Abstracted Synthesis Problem</b>	<b>43</b>
3.1	Problem Statement . . . . .	43
3.2	Abstraction . . . . .	45
3.2.1	Smooth Reformulation . . . . .	47
3.2.2	Regularity of Smooth Reformulation . . . . .	48
3.3	Different Norms for Optimization . . . . .	48
3.4	Summary . . . . .	49
<b>4</b>	<b>Set Conversions and Set Containment</b>	<b>51</b>
4.1	Set Conversions . . . . .	51
4.1.1	Zonotope To Ellipsoid . . . . .	51
4.1.2	Ellipsoid to Zonotope . . . . .	59
4.1.3	H-Polytope to Parallelotope . . . . .	63
4.1.4	Polynomial Zonotope to Zonotope . . . . .	63
4.2	Set Containment . . . . .	64
4.2.1	Zonotope in Zonotope . . . . .	64
4.2.2	Zonotope in H-Polytope . . . . .	64
4.3	Summary . . . . .	65
<b>5</b>	<b>Piecewise Constant Controller Synthesis</b>	<b>67</b>
5.1	Problem Statement . . . . .	67
5.2	Generator-Space Control . . . . .	69
5.2.1	Reference Trajectory . . . . .	71
5.2.2	Controller Template . . . . .	71
5.2.3	Parameterized Reachable Set . . . . .	73
5.2.4	Controller Computation . . . . .	74
5.2.5	Closed-Loop Reachable Set . . . . .	76
5.2.6	Computational Complexity . . . . .	78
5.2.7	Discussion . . . . .	80
5.3	Polynomial Generator-Space Control . . . . .	81
5.3.1	Reference Trajectory . . . . .	82
5.3.2	Controller Template . . . . .	84
5.3.3	Parameterized Reachable Set . . . . .	86

5.3.4	Controller Computation . . . . .	88
5.3.5	Closed-Loop Reachable Set . . . . .	90
5.3.6	Computational Complexity . . . . .	90
5.3.7	Generalization of Generator-Space Control . . . . .	93
5.3.8	Experiments . . . . .	94
5.3.9	Discussion . . . . .	98
5.4	Iterative Polynomial Generator-Space Control . . . . .	100
5.4.1	Cost of the Controller . . . . .	104
5.4.2	Parameterized Reachable Set . . . . .	107
5.4.3	Trust-Region Subproblem . . . . .	109
5.4.4	Tuning of the Trust-Region Radius . . . . .	110
5.4.5	Computational Complexity . . . . .	111
5.4.6	Experiments . . . . .	112
5.4.7	Discussion . . . . .	113
5.5	Summary . . . . .	115
<b>6</b>	<b>Piecewise Constant Controller Synthesis with Continuous State Feedback</b>	<b>117</b>
6.1	Problem Statement . . . . .	117
6.2	Reachset Optimal Control . . . . .	119
6.2.1	Feedback Matrix Parameterization . . . . .	121
6.2.2	Closed-Loop Reachable Set . . . . .	122
6.2.3	Feedforward Controller Computation . . . . .	123
6.2.4	Feedback Controller Computation . . . . .	123
6.2.5	Initial Guess For Feedback Synthesis . . . . .	126
6.2.6	Computational Complexity . . . . .	126
6.2.7	Discussion . . . . .	127
6.3	Iterative Polynomial Reachset Optimal Control . . . . .	129
6.3.1	Feedback Matrix Parameterization . . . . .	130
6.3.2	Closed-Loop Reachable Set . . . . .	137
6.3.3	Cost of the Controller . . . . .	137
6.3.4	Parameterized Reachable Set . . . . .	141
6.3.5	Controller Computation . . . . .	145
6.3.6	Tuning of Trust-Region Radii . . . . .	147
6.3.7	Feedback Derivative Information . . . . .	149
6.3.8	Initial Guess . . . . .	153
6.3.9	Computational Complexity . . . . .	153
6.3.10	Experiments . . . . .	155
6.3.11	Discussion . . . . .	164
6.4	Summary . . . . .	164
<b>7</b>	<b>Conclusion and Future Directions</b>	<b>167</b>
7.1	Conclusion . . . . .	167
7.2	Future Directions . . . . .	168

*Contents*

<b>A Appendix</b>	<b>171</b>
A.1 Abstracted Synthesis Problem . . . . .	171
A.1.1 Reformulation Equivalence . . . . .	171
A.1.2 Necessity of First-Order Optimality Conditions . . . . .	175
A.2 Complexity Analysis . . . . .	177
A.3 Derivatives of the LQR Feedback Matrix . . . . .	179
A.3.1 Jacobian of Feedback Matrix . . . . .	180
A.3.2 Hessian of Feedback Matrix . . . . .	181
<b>Bibliography</b>	<b>187</b>

# List of Figures

1.1	Reach-avoid problem . . . . .	2
2.1	Examples of convex and non-convex sets . . . . .	15
2.2	Examples of convex and non-convex functions . . . . .	16
2.3	Construction of an interval . . . . .	24
2.4	Construction of an ellipsoid . . . . .	26
2.5	Example of V-polytope and H-polytope . . . . .	28
2.6	Construction of a zonotope . . . . .	30
2.7	Construction of a non-convex polynomial zonotope . . . . .	32
2.8	Subset property of polynomial zonotopes . . . . .	34
2.9	Reachability computation for linear systems . . . . .	37
2.10	Reduction of 2D zonotope . . . . .	41
4.1	Distribution and covariance of zonotope vertex candidate points . . . . .	54
4.2	Lower bound on minimum zonotope norm . . . . .	57
4.3	Example of enclosing and inscribed ellipsoids for a given zonotope . . . . .	58
4.4	Uniformly sampled vs. approx. equi-distantly sampled zonotope . . . . .	60
4.5	Example of enclosing and inscribed zonotopes for a given ellipsoid . . . . .	62
5.1	Discrete-time state feedback and piecewise constant inputs . . . . .	68
5.2	Overview of the GSC approach . . . . .	70
5.3	Comparison of polynomial zonotope starting point vs. center of zonotope outer approximation . . . . .	83
5.4	Dependency-preserved part of parameterized reachable set . . . . .	88
5.5	Visualization of the Single-Track model . . . . .	96
5.6	Comparison of GSC and PGSC for the bicycle model . . . . .	97
5.7	Comparison of GSC and PGSC for the Van-der-Pol benchmark . . . . .	99
5.8	Linear vs. quadratic controller template for PGSC . . . . .	99
5.9	Visualization of the iterative optimization of the iPGSC approach . . . . .	105
5.10	Comparison of PGSC and iPGSC for the bicycle model . . . . .	113
5.11	iPGSC solver progress for the bicycle benchmark . . . . .	114
5.12	Comparison of PGSC and iPGSC for the Van-der-Pol model . . . . .	115
6.1	Overview of ROC . . . . .	120
6.2	Overview of iPROC . . . . .	131
6.3	Approximation accuracy of the recursive LQR for a 1D example . . . . .	136
6.4	1D visualization of the support function approximation for iPROC . . . . .	143
6.5	Visualization of the spacecraft coordinate system . . . . .	156

*List of Figures*

6.6	Comparison of ROC and iPROC for the space rendezvous model . . . . .	157
6.7	Final closed-loop reachable set of ROC and iPROC for the space rendezvous model . . . . .	158
6.8	iPGSC solver progress for the space rendezvous model . . . . .	159
6.9	Comparison of ROC and iPROC for the car model . . . . .	161
6.10	Final closed-loop reachable set of ROC and iPROC for the car model . . .	161
6.11	Visualization of the water tank dynamics . . . . .	162
6.12	Comparison of ROC and iPROC for the 4D tank model . . . . .	163

# List of Tables

2.1	Closedness of sets under common operations . . . . .	23
5.1	Parameters for the bicycle model . . . . .	97
5.2	Parameters for the Van-der-Pol model . . . . .	98
6.1	Parameters for the space benchmark . . . . .	157
6.2	Parameters for the car model . . . . .	160
6.3	Parameters for the water tank model . . . . .	162
6.4	Scalability comparison of iPROC and ROC for the tank benchmark . . .	163





# Acronyms

<b>ACADO</b>	automatic control and dynamic optimization <sup>1</sup> [55]
<b>AI</b>	artificial intelligence
<b>AROC</b>	automated reachset optimal control <sup>2</sup> [63]
<b>CBF</b>	control barrier function
<b>CLF</b>	control Lyapunov function
<b>CORA</b>	continuous reachability analyzer <sup>3</sup> [2]
<b>GSC</b>	generator-space control
<b>HJB</b>	Hamilton-Jacobi-Bellman
<b>iPGSC</b>	iterative polynomial generator-space control
<b>IPOPT</b>	interior point optimizer <sup>4</sup> [111]
<b>iPROC</b>	iterative polynomial reachset optimal control
<b>KKT</b>	Karush-Kuhn-Tucker
<b>LCQ</b>	Linearity constraint qualification
<b>LIDAR</b>	light detection and ranging
<b>LMI</b>	linear matrix inequality
<b>LP</b>	linear program
<b>LQR</b>	linear-quadratic regulator
<b>LTI</b>	linear time-invariant
<b>MATLAB</b>	matrix Laboratory
<b>MFCQ</b>	Mangasarian-Fromovitz constraint qualification
<b>MILP</b>	mixed-integer linear program
<b>MIQP</b>	mixed-integer quadratic program
<b>MPC</b>	model predictive control
<b>MVEE</b>	minimum-volume enclosing ellipsoid
<b>MVIE</b>	maximum-volume inscribed ellipsoid

---

<sup>1</sup><https://acado.github.io/>

<sup>2</sup><https://aroc.in.tum.de/>

<sup>3</sup><https://cora.in.tum.de/>

<sup>4</sup><https://github.com/coin-or/Ipopt>

## *Acronyms*

<b>ODE</b>	ordinary differential equation
<b>PCA</b>	principle component analysis
<b>PGSC</b>	polynomial generator-space control
<b>ROC</b>	reachset optimal control
<b>SC</b>	Slater's condition
<b>SDP</b>	semi-definite program

# Notations

In this thesis, we denote scalars and vectors with lowercase letters, matrices with uppercase letters, and sets with uppercase calligraphic letters.

## Definitions

### Number Sets

The set of real numbers, non-negative real numbers, and positive real numbers is denoted by  $\mathbb{R}$ ,  $\mathbb{R}_{\geq 0}$ , and  $\mathbb{R}_+$ . Further, the set of complex numbers is denoted by  $\mathbb{C}$ . Additionally,  $\mathbb{N}$  and  $\mathbb{N}_+$  denote the sets of natural and positive natural numbers. Lastly, the sets of symmetric, symmetric semi-positive definite, and symmetric positive definite matrices are given by  $\mathbb{S}$ ,  $\mathbb{S}_+$ , and  $\mathbb{S}_{++}$ .

### Vectors, Matrices, and Sets

The element of a matrix  $M \in \mathbb{R}^{n \times m}$  in the  $i$ -th row and  $j$ -th column for  $1 \leq i \leq n$  and  $1 \leq j \leq m$  is denoted by  $M_{ij}$ ; further, we denote with  $M_{(i,:)}$  and  $M_{(:,j)}$  the  $i$ -th row and  $j$ -th column of  $M$ , and define the shorthand  $M_{(i,:)}^T = \left(M_{(i,:)}\right)^T$  and  $M_{(:,j)}^T = \left(M_{(:,j)}\right)^T$ . For a vector  $v \in \mathbb{R}^n$ , we similarly denote with  $v_i \in \mathbb{R}$  the  $i$ -th component of  $v \in \mathbb{R}^n$  for  $1 \leq i \leq n$ . For  $a \in \mathbb{N}$  and  $b \in \mathbb{N}$  with  $a \leq b$ , the vector of all integers between  $a$  and  $b$  (both included) is denoted by  $a : b = [a, a + 1, \dots, b]$ ; for a matrix  $M \in \mathbb{R}^{n \times m}$ , we extend this notation and denote with  $M_{(:,a)} = M_{(:,a:m)}$  the matrix consisting of the  $a$ -th to the last column with  $a \leq m$ . Similarly, we denote with  $\mathcal{M}_{(a:b)}$  the projection of the set  $\mathcal{M} \subseteq \mathbb{R}^n$  onto its  $a$ -th to  $b$ -th dimension, where  $a \leq b \leq n$ . For two matrices of equal size, inequality relations between them are to be interpreted element-wise. Moreover,  $\mathbf{1}_n$  and  $I_n = [e_{(n)}^{(1)}, \dots, e_{(n)}^{(n)}]$  denote the  $n$ -dimensional all-ones vector and the  $n$ -dimensional identity matrix, where  $e_{(n)}^{(i)} \in \mathbb{R}^n$  denotes the  $i$ -th unit vector of dimension  $n$  for  $1 \leq i \leq n$ .

### Mathematical Operators

For a given vector  $v \in \mathbb{R}^n$ ,  $\text{diag}(v) \in \mathbb{R}^{n \times n}$  denotes the diagonal matrix with entries of  $v$  on its diagonal. Further,  $\text{mod}(v, k) = [\text{mod}(v_1, k), \text{mod}(v_2, k), \dots, \text{mod}(v_n, k)]$  denotes the modulo operation on all components of  $v \in \mathbb{R}^n$  where  $k \in \mathbb{R}$ . The trace of a square matrix  $M \in \mathbb{R}^{n \times n}$  is defined as  $\text{trace}(M) = \sum_{i=1}^n M_{ii}$ . For two sets

## Notations

$\mathcal{A} \in \mathbb{N}^n$  and  $\mathcal{B} \in \mathbb{N}^n$ , their set difference is given by  $\mathcal{A} \setminus \mathcal{B} = \{a \in \mathcal{A} \mid a \notin \mathcal{B}\}$ . Moreover,  $|\mathcal{S}| \in \mathbb{N}$  denotes the cardinality of a set  $\mathcal{S} \in \mathbb{N}^n$ . Further, the empty set is denoted by  $\emptyset$ . The transpose of a real matrix  $R \in \mathbb{R}^{m \times n}$  and the conjugate transpose of a complex matrix  $C \in \mathbb{C}^{m \times n}$  are denoted by  $R^T$  and  $C^H$ , respectively. The minimum and maximum entry of a given matrix  $A \in \mathbb{R}^{n \times m}$  are denoted by  $\min(A)$  and  $\max(A)$ , respectively; for  $B \in \mathbb{R}^{n \times m}$ , we further denote with  $\min(A, B) \in \mathbb{R}^{n \times m}$  (“max” analogous) the matrix that results from taking the minimum over each matrix element separately. For a finite number of vectors  $v^{(i)} \in \mathbb{R}^n$  with  $1 \leq i \leq m$ , we denote by  $\min_{1 \leq i \leq m} v^{(i)} = \left[ \min_{1 \leq i \leq m} v_1^{(i)}, \dots, \min_{1 \leq i \leq m} v_n^{(i)} \right]^T$  the finite minimization over each vector component, where the result is again collected in an  $n$ -dimensional vector (analogous for max). Further,  $A \boxtimes B \in \mathbb{R}^{mp \times nq}$  denotes the Kronecker product of matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{p \times q}$ , and  $C \boxplus D = C \boxtimes I_m + I_n \boxtimes D$  denotes the Kronecker sum for  $C \in \mathbb{R}^{n \times n}$  and  $D \in \mathbb{R}^{m \times m}$ . For a set of  $m$  points  $\{v^{(1)}, v^{(2)}, \dots, v^{(m)}\}$  with  $v^{(i)} \in \mathbb{R}^n$  and  $1 \leq i \leq m$ , we denote its convex hull, i.e., the smallest convex set containing all  $v^{(i)}$ , with  $\text{convh}\left(\{v^{(1)}, v^{(2)}, \dots, v^{(m)}\}\right) \subset \mathbb{R}^n$ . For a complex matrix  $C \in \mathbb{C}^{n \times m}$ ,  $\text{real}(C) \in \mathbb{R}^{n \times m}$  and  $\text{imag}(C)$  denote the real and imaginary part of  $C$ , respectively. For a matrix  $M \in \mathbb{S}_+^{n \times n}$ , we alternatively write  $M \succeq 0$ .

## Functions

For a function  $f : \mathbb{R} \mapsto \mathbb{R}^n$  and time  $t \in \mathbb{R}_{\geq 0}$ , its time derivative is denoted by  $\dot{f}(t) = \left[ \frac{df_1(t)}{dt}, \dots, \frac{df_n(t)}{dt} \right]^T$ . For once and twice continuously differentiable functions  $f : \mathbb{R}^n \mapsto \mathbb{R}^m$  and  $g : \mathbb{R}^n \mapsto \mathbb{R}$  with vector  $x \in \mathbb{R}^n$ , the Jacobian matrix  $J_f(x) \in \mathbb{R}^{m \times n}$  of  $f$  is defined by  $[J_f(x)]_{ij} = \frac{\partial f_i(x)}{\partial x_j}$  for  $1 \leq i \leq m$  and  $1 \leq j \leq n$ , and the Hessian matrix  $H_g(x) \in \mathbb{R}^{n \times n}$  of  $g$  is defined by  $[H_g(x)]_{ij} = \frac{\partial^2 g_i(x)}{\partial x_i \partial x_j}$  for  $1 \leq i \leq n$  and  $1 \leq j \leq n$ .

## List of Symbols

### Scalars

Symbol	Dimension	Description
$t$	$\in \mathbb{R}_{\geq 0}$	current point in time
$t_f$	$\in \mathbb{R}_+$	time horizon length
$m_s$	$\in \mathbb{N}_+$	number of discrete-time feedback steps
$m_c$	$\in \mathbb{N}_+$	number of piecewise constant control inputs per feedback step
$t_s$	$\in \mathbb{R}$	time between discrete-time feedback
$t_c$	$\in \mathbb{R}_+$	duration of piecewise constant control inputs per feedback step
$a$	$\in \mathbb{N}_+$	number of monomials for the controller template
$\kappa$	$\in \mathbb{N}_+$	controller order

Symbol	Dimension	Description
$\pi$	$\in \mathbb{N}_+$	abstraction order for reachable set computations
$o$	$\in \mathbb{N}_+$	reduction order for zonotopes and polynomial zonotopes
$\gamma$	$\in (0, 1]$	trust-region radius for piecewise constant controller parameters $P$
$\eta$	$\in (0, 1]$	trust-region radius for feedback controller parameters $Q$ and $R$
$m_d$	$\in \mathbb{N}_+$	number of steps for the time discretization of time-varying matrices
$n_x$	$\in \mathbb{N}_+$	state dimension
$n_u$	$\in \mathbb{N}_+$	input dimension
$n_w$	$\in \mathbb{N}_+$	disturbance dimension
$o_{\mathcal{X}_f}$	$\in \mathbb{N}_+$	number of halfspaces making up the final state constraints
$o_{\mathcal{X}}$	$\in \mathbb{N}_+$	number of halfspaces making up the state constraints
$o_{\mathcal{U}}$	$\in \mathbb{N}_+$	number of halfspaces making up the input constraints
$h$	$\in \mathbb{N}_+$	extended optimization horizon
$m_r$	$\in \mathbb{N}_+$	number of steps for reachability analysis
$\mu$	$\in \mathbb{R}_+$	optimality tolerance
$\sigma$	$\in \mathbb{R}_+$	constraint penalty multiplier
$\zeta$	$\in \mathbb{R}_+$	input penalty factor
$l_{\max}$	$\in \mathbb{N}_+$	maximum number of iterations for a given algorithm
$n_z$	$\in \mathbb{N}_+$	number of controller parameters for the combined controller
$\omega$	$\in \mathbb{R}_+$	exponent for the computational complexity of matrix multiplication

## Vectors

Symbol	Dimension	Description
$x(t)$	$\in \mathbb{R}^{n_x}$	state vector
$u(t)$	$\in \mathbb{R}^{n_u}$	controllable input vector
$w(t)$	$\in \mathbb{R}^{n_w}$	uncontrollable input vector (disturbance)
$x_f$	$\in \mathbb{R}^{n_x}$	target state
$x_{\text{ref}}$	$\in \mathbb{R}^{n_x}$	reference state trajectory
$u_{\text{ref}}$	$\in \mathbb{R}^{n_u}$	reference input trajectory
$\vartheta$	$\in \mathbb{R}_{\geq 0}^h$	weight vector for the objective function using the extended horizon $h$

## Matrices

Symbol	Dimension	Description
$P$	$\in [-1, 1]^{m_c m_s n_u \times a}$	controller parameter matrix for all $m_c m_s$ steps
$Q_{\text{ref}}$	$\in \mathbb{S}_+^{n_x \times n_x}$	state weighting matrix for reference trajectory

## Notations

Symbol	Dimension	Description
$R_{\text{ref}}$	$\in \mathbb{S}_+^{n_u \times n_u}$	input weighting matrix for reference trajectory
$K$	$\in \mathbb{R}^{n_u \times n_x}$	state feedback matrix
$Q$	$\in \mathbb{S}_{++}^{n_x \times n_x}$	state weighting matrix for LQR control
$R$	$\in \mathbb{S}_{++}^{n_u \times n_u}$	input weighting matrix for LQR control

## Sets

Symbol	Dimension	Description
$\mathcal{R}^{(e)}$	$\subset \mathbb{R}^{n_x}$	exact reachable set
$\mathcal{R}$	$\subset \mathbb{R}^{n_x}$	outer approximation of the reachable set
$\tilde{\mathcal{R}}$	$\subset \mathbb{R}^{n_x}$	approximation (without formal guarantees) of the reachable set
$\mathcal{X}^{(0)}$	$\subset \mathbb{R}^{n_x}$	set of initial states
$\tilde{\mathcal{X}}^{(0)}$	$\subset \mathbb{R}^{n_x}$	parallelotope outer approximation of $\mathcal{X}^{(0)}$
$\mathcal{U}$	$\subset \mathbb{R}^{n_u}$	input constraint set
$\mathcal{W}$	$\subset \mathbb{R}^{n_w}$	disturbance set
$\mathcal{X}$	$\subseteq \mathbb{R}^{n_x}$	state constraint set
$\mathcal{X}_f$	$\subseteq \mathbb{R}^{n_x}$	final state constraint set
$\tilde{\mathcal{X}}$	$\subseteq \mathbb{R}^{n_x}$	adapted state constraint set
$\tilde{\mathcal{U}}$	$\subset \mathbb{R}^{n_u}$	adapted input constraint set with $\tilde{\mathcal{U}} \subseteq \mathcal{U}$
$\tilde{\mathcal{X}}_f$	$\subseteq \mathbb{R}^{n_x}$	adapted final state constraint set
$\mathcal{M}^{(Q)}$	$\subset \mathbb{S}_{++}^{n_x \times n_x}$	bounded set of state weighting matrices for LQR control
$\mathcal{M}^{(R)}$	$\subset \mathbb{S}_{++}^{n_u \times n_u}$	bounded set of input weighting matrices for LQR control





# Computing Platform and Implementation

All numerical experiments and computations are carried out in the matrix Laboratory (MATLAB) 2022b on an Intel(R) Core(TM) i7-9700K with 16 GB of RAM. Further, we implemented all presented set operations and conversions in the continuous reachability analyzer<sup>1</sup> [2] (CORA) toolbox, and all proposed synthesis algorithms in the automated reachset optimal control<sup>2</sup> [63] (AROC) toolbox, which are both publicly available.

---

<sup>1</sup><https://cora.in.tum.de/>

<sup>2</sup><https://aroc.in.tum.de/>



# 1 Introduction

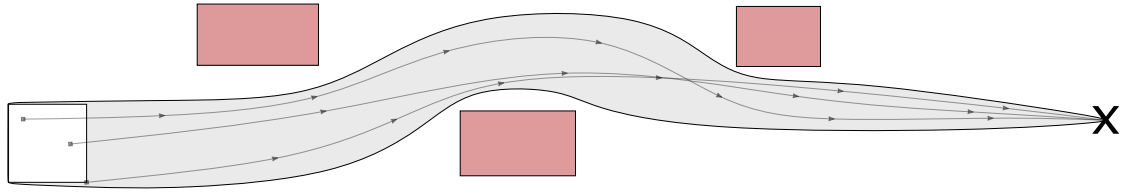
## 1.1 Motivation

In 1965, Gordon Moore, co-founder of one of today's biggest semiconductor chip manufacturers in the world, posited that the number of components on a chip of fixed size would double every year. His revised version, which stated that this number would double every two years and which has held true every since, became known as Moore's law and is de facto synonymous with the exponential rise of computational power we have seen over the last decades. Incidentally, the birth of Moore's law in 1965 coincided with the publication of the first machine learning paper. While back then, machine learning was restricted to very basic applications, today's availability of computational resources predicted by Moore has truly enabled the application of artificial intelligence (AI) to nearly every aspect of our lives, with applications ranging from smart kitchen appliances and vacuuming robots to highly sophisticated chat bots.

Among many, autonomous driving is one prominent area of research that benefited from the increasing capabilities of AI. However, while smart vacuuming robots may have some form of human interaction, guaranteeing the safety of their machine-learning-guided journey through one's home is not safety critical. In contrast, AI systems which enable fully autonomous driving are highly safety critical. However, the high complexity of such AI systems, often using millions or even billions of neurons, makes it hard or even impossible to verify them. As a result, these autonomous vehicles are required to validate their safety to within some statistical margin by extensive testing. For self-driving cars, this involves driving thousands or even millions of kilometers in varying environments. Furthermore, changes to any safety-critical system in the car can require additional testing, or even demand a complete restart of the testing campaign.

Instead of only validating a given system through extensive testing, one may formally verify it to obtain guarantees about its behavior. With the increase of digital computers in all aspects of our society, cyber-physical systems, which combine both physical and software components, are becoming increasingly popular to model these systems. Reachability analysis is one popular formal verification tool to verify cyber-physical systems: The reachable set contains all possible future states of a system that can be reached, starting from a given set of initial states under any possible disturbance realization. A state hereby collects the necessary information of a system to describe it fully. For autonomous vehicles, this means that its subsystems can be formally verified by checking intersections of unsafe states with the reachable set; changes to a subsystem then only require the re-verification of that subsystem but leaves the formal guarantees of all other subsystems untouched. Additionally, reachability analysis can account for uncertainty in

## 1 Introduction



**Figure 1.1:** Concept of a reach-avoid problem: From a given initial set of states (white), the system is steered to the target state (“x”) within a given time horizon while avoiding all obstacles (red). The grey area hereby is the reachable set, where the grey lines visualize possible trajectories for a specific initial state (grey square).

the system, such as measurement errors of radar or light detection and ranging (LIDAR) sensors.

However, reachable sets of more complex systems cannot be computed exactly. Therefore, an outer approximation of the reachable set, which is an approximation containing the exact reachable set, is computed. That said, autonomous systems, such as autonomous cars, are often first designed, i.e., a controller that executes a certain task is synthesized, and only then verified. As a result, tuning the controller manually in order to verify safety may be impractical for highly complex systems. Therefore, instead of verifying existing controllers, we focus in this thesis on combining numerical optimization and reachability analysis to obtain correct-by-construction controllers that already provide formal guarantees about the behavior of the closed-loop system.

A wide range of these tasks can be classified as reach-avoid problems, which are visualized in Fig. 1.1: For a given system with a set of initial states, we try to find a controller that, within a given time horizon, avoids all obstacles and steers the system state close to a given target state while minimizing the required control effort. In previous work, the combination of reachability analysis and optimization theory for the synthesis proved to be a powerful combination as it allows to optimize over a set of initial states instead of a single initial state. That said, the computation of reachable sets is computationally expensive when executed repeatedly. Thus, in this thesis, we propose novel synthesis approaches that avoid the computation of reachable sets as much as possible while still using their formal verification properties.

## 1.2 State of the Art

The task of synthesizing a controller for a given system has already been extensively explored in literature. We subsequently summarize the most important work.

### 1.2.1 Optimal Control

When only a single initial state needs to be controlled, many approaches in the field of optimal control have been proposed [92]. For a given system of possibly nonlinear ordinary differential equations (ODEs), optimal control approaches try to find a function for the input vector subject to these ODEs that minimizes a functional while respecting

constraints on the input and state of the system. When there are no constraints on a linear time-invariant (LTI) system and one optimizes over a quadratic cost function in both state and input, Kalman derived in his seminal paper [56] the optimal control law – now known as linear-quadratic regulator (LQR) control – using Pointryagin’s maximum principle, which provides necessary (but not sufficient) conditions for an optimal control law. For general nonlinear optimal control problems, however, it becomes necessary to use numerical methods [105]. We differentiate between two approaches: Starting from the first-order necessary conditions of the original optimal control problem given by the Hamiltonian boundary-value problem, indirect methods numerically solve this boundary-value problem for given boundary values and thus reduce the original problem to the solution of a system of nonlinear equations [92]. Due to the recent progress in optimization theory, however, indirect methods [105, 96] have become less popular. In contrast, direct methods parameterize the input and state functions using a finite number of parameters and then solve the resulting finite-dimensional optimization problem. Here, methods include direct multiple shooting [13], which subdivides the time horizon of the optimal control problem into multiple steps and subsequently parameterizes the control function, and has, e.g., been applied to optimal robot control [26]. Further, direct collocation methods [104] parameterize both the input and state vector and have, e.g., been used for the estimation of muscle forces during motion [43]. There exist many software packages, such as the automatic control and dynamic optimization<sup>1</sup> [55] (ACADO) toolbox and CasADi [10], that implement these optimal control methods while providing user-friendly interfaces in multiple programming languages.

As an extension to Pointryagin’s maximum principle, the Hamilton-Jacobi-Bellman (HJB) equations provide necessary and sufficient conditions for an optimal control law over the entire state space instead of a single trajectory [57]. Further, they allow the direct inclusion of state and input constraints into the problem. However, while there have been attempts to reduce the curse of dimensionality when trying to solve the HJB equations, solving them requires the solution of partial differential equations and thus still has worst-case exponential complexity in the number of state variables [85].

### 1.2.2 Model Predictive Control

When using a direct method for the solution of an optimal control problem for a given initial state, choosing the input to be piecewise constant as a parameterization means that the algorithm returns a sequence of piecewise constant control inputs that minimizes the given objective function for the given initial state. In model predictive control (MPC), this procedure is applied iteratively: For a given initial state, the sequence of control inputs is computed and the first input is applied to the system. After one step, the system state is again measured and the input sequence is recomputed. Since this procedure easily allows for the incorporation of state and input constraints into the optimal control problem, MPC has been used extensively in industry for the past decades [82, 93]. When the given system has inherent uncertainty, tube-based MPC guarantees

---

<sup>1</sup><https://acado.github.io/>

## 1 Introduction

that the system stays within a tube around a reference trajectory [80, 72, 79, 91, 78]. However, since MPC repeatedly solves an optimal control problem for a given state, its application requires optimization online. To alleviate some of that computation cost online, explicit MPC approaches parameterize the optimal control problem in the initial state and solve the resulting parametric optimization problem offline [1, 30], albeit at the cost of high computation times when trying to solve high-dimensional, parametric, non-convex optimization problems. With major advances in optimization theory, more recent work thus focuses again on solving the optimal control problem online, which allows the online controller synthesis for LTI systems in real-time [44, 119, 102].

### 1.2.3 Abstraction-Based Synthesis

Similarly to the HJB equations, controller synthesis using model abstraction – and particularly symbolic model abstraction – allows for the synthesis of controllers for sets of initial states instead of a single initial state. In model abstraction, the continuous or hybrid system is abstracted to a finite-state system which has a finite number of discrete states [107, 88, 89, 94, 58, 39]. For such automata, standard algorithms for finite systems can be employed to synthesize a controller for the abstracted system; the resulting controller then needs to be suitably refined such that we obtain a controller for the original system [95]. When abstracting the state space with a uniform grid, the number of abstracted finite states grows exponentially in the number of state variables. Additionally, early work made assumptions about the system class [58, 46, 12] or the stability of the system [42, 90]; recent work [118, 95] has managed to avoid such assumptions. Recently, attempts have been made to also circumvent the curse of dimensionality due to state space discretization; however, the proposed techniques either use mixed-integer linear program (MILP) with exponential worst-case complexity in the number of integer variables [116] or focus on linear systems [87].

### 1.2.4 Reachability-Based Synthesis

Reachability analysis has recently attracted increased attention in the formal verification community [9, 4, 23, 5, 21, 6]. Previous work has already combined reachability analysis with optimization to solve reach-avoid problems: The authors in [100] synthesize a continuous state feedback controller by formulating a non-convex optimization problem, where in each optimization iteration an outer approximation to the reachable set is computed, which in turn is used to compute the current value of the objective function and the constraints. In [99], the authors consider the controller synthesis for nonlinear system dynamics by first computing a piecewise constant control sequence for each vertex of the initial set. These sequences are then used in the final controller to form a piecewise constant input sequence for any initial state. However, the number of vertices of the initial set grows exponentially with the system dimension. Thus, the authors in [99] instead parameterize the controller using zonotopes and compute a parameterized reachable set based on the linearized dynamics to obtain a synthesis problem which can be formulated as a linear program (LP). This parameterization grows linearly in the

state dimension and thus can efficiently synthesize piecewise constant controllers for disturbed nonlinear systems under input and state constraints. While the parameterization of the control input by a zonotope means that input constraints can be easily verified, the parameterized reachable set is only an approximation and thus the inclusion of state constraints requires user input to achieve non-conservative results. Furthermore, while this approach can be applied iteratively to realize discrete-time state feedback, it fails to efficiently reject larger disturbances. Therefore, the authors combine the feedback synthesis from [100] and the feedforward synthesis from [99] for the synthesis of a combined piecewise constant feedforward controller with continuous state feedback [101], albeit at the cost of again computing an outer approximation to the reachable set of the current controller in each optimization iteration.

## 1.3 Outline

We propose several new algorithms for formally verified controller synthesis of constrained, disturbed nonlinear systems. Since we make extensive use of set-based computations, we further introduce required conversions between set representations in order to pose the synthesis problems efficiently.

**Chapter 2: Background** We introduce all necessary background in Chap. 2. Specifically, we start by introducing mathematical standard operations in Sec. 2.1 and the theory of differentials in Sec. 2.2. We then describe the basics of optimization theory in Sec. 2.3 and continue with a review of the well-known linear-quadratic regulator (LQR) control scheme in Sec. 2.4. Since we compute with sets, we describe in Sec. 2.5 and Sec. 2.6 all necessary set operations and set representations. In Sec. 2.7, we introduce the concept of reachable sets and give a short overview over the relevant algorithms for this thesis. In Sec. 2.8, we then briefly discuss how the representation complexity of zonotopes and polynomial zonotopes can be controlled, which is, e.g., necessary to keep reachability computations efficient.

**Chapter 3: Abstracted Synthesis Problem** In Chap. 3, we discuss the abstracted synthesis problem, i.e., an optimization problem which generalizes all major optimization problems of subsequent chapters and which we use to derive important properties. First, we introduce reach-avoid problems and define the general synthesis problem in Sec. 3.1. Because many optimization solvers find possible extrema by solving for the first-order critical points that are characterized by the Karush-Kuhn-Tucker (KKT) conditions, we then define the abstracted synthesis problem in Sec. 3.2 and derive a smooth optimization problem, for which we prove that it shares its global minimum with the original optimization problem. Lastly, we discuss the effect of different norms in optimization in Sec. 3.3.

**Chapter 4: Set Conversions and Set Containment** In order to efficiently pose and solve optimization problems of subsequent chapters, we further introduce existing and

## 1 Introduction

novel conversions between different set representations in Chap. 4. First, we introduce all necessary set conversions in Sec. 4.1. In order to check, e.g., the constraint satisfaction of input or state constraints, we then shortly describe existing approaches from the literature to check set containment in Sec. 4.2.

**Chapter 5: Piecewise Constant Controller Synthesis** We describe existing and novel algorithms for the synthesis of piecewise constant controllers, i.e., controllers which are constant in time for a given time interval, in Chap. 5. After the introduction of the problem statement in Sec. 5.1, we first review the generator-space control (GSC) approach from [99] in Sec. 5.2 and then propose the polynomial generator-space control (PGSC) approach as a novel extension of GSC in Sec. 5.3. Because both GSC and PGSC require additional user inputs when state constraints are to be enforced, we introduce iterative polynomial generator-space control (iPGSC), a novel iterative synthesis approach using trust regions, in Sec. 5.4. We demonstrate the applicability of our algorithms using numerical examples.

**Chapter 6: Piecewise Constant Controller Synthesis with Continuous State Feedback** Since piecewise constant controllers cannot quickly counteract disturbances due to their inherent discrete-time feedback, we discuss the synthesis of piecewise constant feedforward controllers with continuous state feedback in Chap. 6. We start by introducing the problem statement in Sec. 6.1 and then describe the reachset optimal control (ROC) approach from [101] in Sec. 6.2, where feedforward synthesis using GSC is combined with the synthesis of a continuous linear feedback term. Because ROC executes the feedforward and feedback synthesis sequentially, i.e., the feedforward controller is synthesized first and then kept fixed during the feedback controller optimization, we introduce iterative polynomial reachset optimal control (iPROC) in Sec. 6.3, which – for the first time – optimizes over the feedforward and feedback controller simultaneously reduces the required user input compared to ROC. The applicability of our algorithms is demonstrated with numerical benchmarks.



## 2 Background

In this chapter, we introduce all concepts that are relevant for this thesis. We start by introducing necessary definitions and well-known theorems in Sec. 2.1 and then briefly describe differentials in Sec. 2.2. In Sec. 2.3, we then give an overview of optimization theory and introduce the well-known LQR controller in Sec. 2.4. In Sec. 2.5, we define all necessary set operations and introduce all required set representations in Sec. 2.6. We continue with a short review of the relevant reachability algorithms from the literature in Sec. 2.7 and conclude with the concept of order reduction in Sec. 2.8, which is required for the efficient application of reachability analysis.

### 2.1 Standard Operations

All controllers in this thesis are synthesized using different numerical optimization techniques. Since numerical solvers will generally not reach an optimum exactly but only approximately, we measure numerical tolerances as defined next.

**Definition 2.1** (Numerical Tolerance). Given two matrices  $A \in \mathbb{R}^{n \times m}$  and  $B \in \mathbb{R}^{n \times m}$ , we define the tolerance between  $A$  and  $B$  as the minimum between the absolute tolerance and the relative tolerance, i.e.

$$\text{tol}(A, B) = \max_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}} \min \left( |A_{ij} - B_{ij}|, \lim_{\epsilon \rightarrow 0^+} \frac{|A_{ij} - B_{ij}|}{\min(|A_{ij}|, |B_{ij}|) + \epsilon} \right),$$

where “ $\lim_{\epsilon \rightarrow 0^+}$ ” denotes the limit from  $\epsilon > 0$  to 0. ■

For later convenience, we subsequently introduce the vectorization and convex combination operation.

**Definition 2.2** (Vectorization). The column-wise vectorization of a given matrix  $X \in \mathbb{R}^{n \times m}$  into a vector  $x \in \mathbb{R}^{nm}$  is given by

$$x = X_{(:,\cdot)} = \begin{bmatrix} X_{(:,1)} \\ X_{(:,2)} \\ \vdots \\ X_{(:,m)} \end{bmatrix}.$$
■

## 2 Background

For the vectorization of matrix products, it further holds that [74, Chap. 2, Sec. 4]

$$[ABC]_{(:,)} = \left( C^T \boxtimes A \right) B_{(:,)} \quad (2.1a)$$

$$= \left( C^T B^T \boxtimes I_m \right) A_{(:,)} \quad (2.1b)$$

$$= (I_q \boxtimes AB) C_{(:,)}, \quad (2.1c)$$

$$A_{(:,)}^T D_{(:,)} = \text{trace} \left( A^T D \right), \quad (2.1d)$$

for  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times o}$ ,  $C \in \mathbb{R}^{o \times q}$ , and  $D \in \mathbb{R}^{n \times m}$ .

**Definition 2.3** (Convex Combination). The convex combination of two sets  $A \subseteq \mathbb{R}^n$  and  $B \subseteq \mathbb{R}^n$  is defined by

$$\text{convc}(A, B) = \{la + (1 - l)b \mid a \in A, b \in B, l \in [0, 1]\}.$$

■

Next, we quickly review well-known theorems that will be required later.

**Theorem 2.1** (Singular Value Decomposition [54, Th. 2.6.3]). *The singular value decomposition of a matrix  $A \in \mathbb{R}^{n \times m}$  is*

$$A = U \Sigma V^H,$$

where  $U \in \mathbb{C}^{n \times n}$  and  $V \in \mathbb{C}^{m \times m}$  are unitary matrices, and  $\Sigma \in \mathbb{R}^{n \times m}$  is a rectangular matrix of zeros except for  $\Sigma_{(i,i)} = \sigma_i$ ,  $1 \leq i \leq \min(n, m)$ , with  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(n,m)} \geq 0$  being the singular values of  $A$ .

We sometimes use the shorthand  $\sigma(A)$  to refer to the vector of singular values of  $A$ .

More widely known is the eigen decomposition, which we introduce next.

**Theorem 2.2** (Eigen decomposition [54, Th. 1.3.7, adapted]). *Let  $A \in \mathbb{R}^{n \times n}$  be a diagonalizable matrix, i.e., there exists an invertible matrix  $P \in \mathbb{R}^{n \times n}$  and a diagonal matrix  $D \in \mathbb{R}^{n \times n}$  such that  $P^{-1}AP = D$ . Then we can decompose  $A$  as*

$$A = Q \Lambda Q^{-1},$$

where  $\Lambda = \text{diag} \left( \left[ \lambda_1, \lambda_2, \dots, \lambda_n \right] \right) \in \mathbb{R}^{n \times n}$ , and  $Q \in \mathbb{C}^{n \times n}$  is an invertible matrix. For  $1 \leq i \leq n$ ,  $Q_{(:,i)} \in \mathbb{R}^n$  and  $\lambda_i \in \mathbb{R}$  denote the eigenvector and corresponding eigenvalue of  $A$ , where  $AQ_{(:,i)} = \lambda_i Q_{(:,i)}$  holds.

We denote with  $\lambda(A)$  the vector of eigenvalues. Next, we introduce the square root of a positive semi-definite matrix.

**Theorem 2.3** (Matrix Square Root [54, Th. 7.2.6, adapted]). *A square root of a positive semi-definite matrix  $Q \in \mathbb{S}_+^{n \times n}$  is*

$$\sqrt{Q} = Q^{\frac{1}{2}} = U \text{diag} \left( \sqrt{\sigma(Q)} \right),$$

such that  $Q = \sqrt{Q} \sqrt{Q}$ , where  $Q = U \text{diag}(\sigma(Q)) U^T$  with singular values  $\sigma(Q) \in \mathbb{R}_{\geq 0}^n$ .

## 2.2 Differentials

Whenever matrices are involved, taking their derivative with respect to a vector, i.e., the derivative of each matrix element, is often cumbersome. Differentials offer a convenient formulation for taking derivatives that also works well for matrices. Hence, we introduce differentials and relevant properties thereof in this section – which is inspired by [74, Chap. 18] – starting with a simple example.

**Example 2.1.** Let  $f(x) = [x_1^2, x_2^2]^T$  with  $x \in \mathbb{R}^2$ . Its Jacobian matrix is given by

$$J_f(x) = \begin{bmatrix} 2x_1 & 0 \\ 0 & 2x_2 \end{bmatrix},$$

with  $J_f(x) \in \mathbb{R}^{2 \times 2}$  since each component of  $f$  is differentiated with respect to each component of  $x$ , i.e., we have  $2 \cdot 2 = 4$  derivatives. In contrast, the differential of  $f$  is

$$df = d \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix} = \begin{bmatrix} [2x_1, 0] dx \\ [0, 2x_2] dx \end{bmatrix}.$$

Immediately, it is obvious that  $df \in \mathbb{R}^2$  and thus the differential has the same dimensionality as  $f$ , whereas  $\frac{df(x)}{dx} \in \mathbb{R}^{2 \times 2}$ . Further, one can directly identify the Jacobian matrix from the above expression. ■

This identification of the Jacobian matrix from the differential is possible in general.

**Theorem 2.4** (First Identification Theorem [74, Th. 5.6, adapted]). *Let  $x \in \mathbb{R}^n$  and  $f : \mathbb{R}^n \mapsto \mathbb{R}^m$ . Then*

$$df = A(x) dx \iff J_f(x) = A(x).$$

In a nutshell, the differential  $df$  does not represent the derivative directly but rather is a linear approximation around  $f(x)$  with an infinitesimally small variation around  $f(x)$ , i.e.,  $df = A(x) dx$ , where  $dx$  is that very small deviation. The fact that the dimension of differentials stays constant makes them particularly useful when handling matrix derivatives and simplifies the following rules, which will be helpful in subsequent chapters.

**Proposition 2.1** (Collection of Differential Rules [74, Chap. 18]). *Let  $A \in \mathbb{R}^{n \times m}$  be a constant matrix,  $F \in \mathbb{R}^{n \times m}$ ,  $G \in \mathbb{R}^{n \times m}$ ,  $H \in \mathbb{R}^{m \times p}$ ,  $X \in \mathbb{R}^{n \times n}$  some non-constant functions, and  $c \in \mathbb{R}$  some scalar. We have*

$$\begin{aligned} dA &= 0, \\ d(cF) &= c dF, \\ d(F^T) &= (dF)^T, \\ d(F + G) &= dF + dG, \\ d(FH) &= (dF)H + F dH, \\ dX^{-1} &= -X^{-1} (dX) X^{-1} \text{ (for nonsingular } X \text{)}. \end{aligned}$$

## 2 Background

For a scalar field  $f : \mathbb{R}^n \mapsto \mathbb{R}$ , the Hessian matrix is well-defined and can also be computed using differentials as the following example demonstrates.

**Example 2.2.** Let  $f(x) = x_1^2 x_2$ . We have

$$\begin{aligned}
 d^2 f &= d^2 (x_1^2 x_2) \\
 &= d \left( \begin{bmatrix} 2x_1 x_2 & x_1^2 \end{bmatrix} dx \right) \\
 &= d \begin{bmatrix} 2x_1 x_2 & x_1^2 \end{bmatrix} dx + \begin{bmatrix} 2x_1 x_2 & x_1^2 \end{bmatrix} d^2 x \\
 &= \left[ d(2x_1 x_2), d(x_1^2) \right] dx \\
 &= \left[ dx^T \begin{bmatrix} 2x_2 \\ 2x_1 \end{bmatrix}, dx^T \begin{bmatrix} 2x_1 \\ 0 \end{bmatrix} \right] dx \\
 &= dx^T \begin{bmatrix} 2x_2 \\ 2x_1 \end{bmatrix} dx_1 + dx^T \begin{bmatrix} 2x_1 \\ 0 \end{bmatrix} dx_2 \\
 &= dx^T \left( \begin{bmatrix} 2x_2 \\ 2x_1 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix} + \begin{bmatrix} 2x_1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \end{bmatrix} \right) dx \\
 &= dx^T \begin{bmatrix} 2x_2 & 2x_1 \\ 2x_1 & 0 \end{bmatrix} dx.
 \end{aligned}$$

Computing the Hessian matrix of  $f(x)$ , we see that  $H_f(x) = \begin{bmatrix} 2x_2 & 2x_1 \\ 2x_1 & 0 \end{bmatrix}$ . ■

Ex. 2.2 can be generalized as follows:

**Theorem 2.5** (Second Identification Theorem [74, Th. 6.6, adapted]). *Let  $f : \mathbb{R}^n \mapsto \mathbb{R}$  with  $x \in \mathbb{R}^n$ . We have*

$$d^2 f = (dx)^T B(x) dx \iff H_f(x) = \frac{B(x) + B(x)^T}{2}.$$

To compute with differentials, we can interpret them as tensors. Let  $dA_{ij} = \frac{dA_{ij}(x)}{dx} \in \mathbb{R}^{1 \times n}$  for  $1 \leq i \leq m$  and  $1 \leq j \leq o$ , where  $A(x) \in \mathbb{R}^{m \times o}$  denotes a matrix function in  $x \in \mathbb{R}^n$  with differentiable matrix elements  $A_{ij}(x)$ . Then

$$[dA]_{ij} = dA_{ij} dx,$$

and we can alternatively collect the Jacobian matrix of each matrix element with respect to  $x$  in the first two dimensions of a 4D tensor, where the last two dimensions reflect the dimensions of  $A$ . Thus, the differential  $dA$  can be represented by the tensor  $dA \in \mathbb{R}^{1 \times n \times m \times o}$ . Similarly, let  $d^2 A_{ij} = \frac{d^2 A_{ij}(x)}{dx^2} \in \mathbb{R}^{n \times n}$ . Then

$$d^2 A_{ij} = dx^T dA_{ij} dx,$$

and thus  $d^2 A$  can – analogously to the first differential – be interpreted as the 4D tensor  $d^2 A \in \mathbb{R}^{n \times n \times m \times o}$ .

It remains to define arithmetic operations on these differential tensors. To that end, let  $B(x) \in \mathbb{R}^{m \times o}$  be a twice differentiable matrix function. We have

$$\begin{aligned} d[A + B]_{ij} &= [dA + dB]_{ij} dx, \\ d^2[A + B]_{ij} &= dx^T [d^2A + d^2B]_{ij} dx, \end{aligned}$$

and thus the tensor representation of the sum of two differentials is equivalent to the sum of the two tensor representations. Further, let  $C(x) \in \mathbb{R}^{o \times l}$  and  $D \in \mathbb{R}^{q \times m}$ . Since

$$\begin{aligned} [D dA]_{ij} &= \sum_{k=1}^m D_{ik} dA_{kj} \\ &= \left( \sum_{k=1}^m D_{ik} dA_{kj} \right) dx, \quad 1 \leq i \leq q, \quad 1 \leq j \leq o, \\ [dAC]_{ij} &= \sum_{k=1}^o dA_{ik} C_{kj} \\ &= \left( \sum_{k=1}^m dA_{ik} C_{kj} \right) dx, \quad 1 \leq i \leq m, \quad 1 \leq j \leq l, \\ [dA dC]_{ij} &= \sum_{k=1}^o dA_{ik} dC_{kj} \\ &= \sum_{k=1}^o dA_{ik} dx dC_{kj} dx \\ &= dx^T \left( \sum_{k=1}^o dA_{ik}^T dC_{kj} \right) dx, \quad 1 \leq i \leq o, \quad 1 \leq j \leq l, \end{aligned}$$

we define the following arithmetic operations for these tensors:

**Definition 2.4** (Arithmetic Operations for Differential Tensors). Let  $A(x) \in \mathbb{R}^{m \times o}$ ,  $B(x) \in \mathbb{R}^{m \times o}$ ,  $C(x) \in \mathbb{R}^{o \times l}$ ,  $D(x) \in \mathbb{R}^{q \times m}$ , and  $x \in \mathbb{R}^n$  with the corresponding differentials represented as the tensors  $dA \in \mathbb{R}^{1 \times n \times m \times o}$ ,  $dB \in \mathbb{R}^{1 \times n \times m \times o}$ ,  $dC \in \mathbb{R}^{1 \times n \times o \times l}$ , and  $dD \in \mathbb{R}^{1 \times n \times q \times m}$  (dependency on  $x$  omitted for readability). We define

$$[dA + dB]_{ij} = dA_{ij} + dB_{ij}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq o, \quad (2.2)$$

$$[D dA]_{ij} = \sum_{k=1}^m D_{ik} dA_{kj}, \quad 1 \leq i \leq q, \quad 1 \leq j \leq o, \quad (2.3)$$

$$[dAC]_{ij} = \sum_{k=1}^m dA_{ik} C_{kj}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq l, \quad (2.4)$$

$$[dAdC]_{ij} = \sum_{k=1}^o dA_{ik}^T dC_{kj}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq l. \quad (2.5)$$

■

## 2 Background

Thus, any differential expressions with sums, products or transpose operations can be directly computed by simply replacing the differentials with their differential tensor counterparts and executing the corresponding operations as defined in Def. 2.4.

### 2.3 Optimization Theory

We extensively use optimization in this thesis to find optimal controllers. Thus, we briefly introduce necessary concepts from optimization theory. We only deal with differentiable optimization problems and therefore assume that all involved functions are sufficiently smooth, which generally means that they are at least twice continuously differentiable. We start with the definition of a smooth optimization problem.

**Definition 2.5.** Let  $f : \mathbb{R}^n \mapsto \mathbb{R}$ ,  $g : \mathbb{R}^n \mapsto \mathbb{R}^m$ , and  $h : \mathbb{R}^n \mapsto \mathbb{R}^o$  be sufficiently smooth functions. A smooth optimization problem can be written as

$$p^* = \arg \min_x f(x), \quad (2.6a)$$

$$\text{s.t. } g(x) \leq 0, \quad (2.6b)$$

$$h(x) = 0, \quad (2.6c)$$

where  $x \in \mathbb{R}^n$  denotes the multivariate optimization variable,  $f(x)$  is the objective function,  $g(x)$  are the inequality constraints,  $h(x)$  denote the equality constraints, and  $p^*$  is the optimal objective value at the optimizer  $x^*$ . Further, we call the set

$$\mathcal{S} = \{x \in \mathbb{R}^n \mid g(x) \leq 0, h(x) = 0\},$$

the feasible set of (2.6). ■

In this thesis, we only consider optimization problems for which an optimum exists, i.e.,  $p^* > -\infty$ . When it is not clear that the optimal objective value is finite, limit expressions (“inf” and “sup”) need to be used instead of “min” and “max” operations. Further, we only consider minimization problems in this section, as any maximization problem can be transformed into a minimization problem via

$$\max_{x \in \mathcal{S}} f(x) = - \min_{x \in \mathcal{S}} -f(x).$$

To solve such smooth optimization problems, practically all known efficient algorithms use the Jacobian matrix and sometimes Hessian matrix information of the objective function and the constraints, respectively. Methods that solely rely on first-order derivatives are called first-order methods while second-order methods additionally use the Hessian matrices of the objective function and constraints.

The remainder of this section is structured as follows: In Sec. 2.3.1, we introduce the Karush-Kuhn-Tucker (KKT) optimal conditions for the original optimization problem in (2.6): Any optimizer of (2.6) necessarily fulfills these optimality conditions, assuming that any applicable constraint qualification – discussed next in Sec. 2.3.2 – holds. We then describe convex optimization as an important subclass in Sec. 2.3.3, followed by a discussion of non-convex optimization problems in Sec. 2.3.4. In Sec. 2.3.5, we then

introduce the dual problem, whose maximum always bounds the minimum of the original problem from below, and which is always a convex optimization problem. We conclude the section with a short review of popular reformulation techniques for optimization problems in Sec. 2.3.6.

### 2.3.1 Karush-Kuhn-Tucker Conditions

While (2.6) can be intuitively interpreted, it remains unclear how it can be solved in general. Instead of trying to solve (2.6) directly, one may solve for the optimizer using the first-order KKT optimality conditions, which we introduce in this section.

For the introduction of these optimality conditions, we first require the Lagrangian function of (2.6).

**Definition 2.6** (Lagrangian Function [14, Chap. 5, adapted]). Let  $f$ ,  $g$ , and  $h$  be defined as in Def. 2.5. The Lagrangian of (2.6) is

$$\phi(x, \mu, \lambda) = f(x) + \mu^T g(x) + \lambda^T h(x), \quad (2.7)$$

where  $\lambda \in \mathbb{R}^m$  and  $\mu \in \mathbb{R}^o$  are the Lagrange multipliers for the inequality and equality constraints, respectively. ■

Given that (2.6) is regular at a critical point  $x^*$ , i.e., it satisfies an appropriate constraint qualification (see Sec. 2.3.2), one can use the Lagrangian function to derive the first-order necessary KKT conditions for optimality of  $x^*$ , which are given by [14, Sec. 5.5]

$$\nabla_x f(x^*) + \sum_{i=1}^m \mu_i^* \nabla_x g_i(x^*) + \sum_{j=1}^o \lambda_j^* \nabla_x h_j(x^*) = 0 \text{ (stationarity)}, \quad (2.8a)$$

$$\mu^{*T} g(x^*) = 0 \text{ (complementary slackness)}, \quad (2.8b)$$

$$\mu^* \geq 0 \text{ (dual feasibility)}, \quad (2.8c)$$

$$g(x^*) \leq 0 \text{ (primal feasibility)}, \quad (2.8d)$$

$$h(x^*) = 0 \text{ (primal feasibility (cont.))}, \quad (2.8e)$$

where  $\mu^* \in \mathbb{R}_{\geq 0}^m$  and  $\lambda^* \in \mathbb{R}^o$  are the optimal Lagrange multipliers for (2.6b) and (2.6c), respectively. Intuitively, complementary slackness means that whenever the  $i$ -th constraint is inactive, i.e.,  $g_i(x^*) < 0$ ,  $\mu_i = 0$  is enforced by (2.8b) so that  $g_i(x^*)$  is treated as if it was not included in the optimization problem.

### 2.3.2 Constraint Qualifications

As described in Sec. 2.3.1, the KKT conditions are only necessary optimality conditions when regularity conditions, also called constraint qualifications, hold. In this section, we introduce an important constraint qualification for convex optimization problems and then introduce a constraint qualification which is also applicable to non-convex optimization problems that is used in a later section.

### 2.3.2.1 Slater's Condition

Arguably the most well-known constraint qualification is Slater's condition (SC).

**Proposition 2.2** (Slater's Condition [14, Sec. 5.2.3, adapted]). *Consider a convex minimization problem as defined in Def. 2.9. If there exists an interior point  $x \in \mathbb{R}^n$ , i.e.,  $h(x) = 0$  and  $g(x) < 0$ , then Slater's condition (SC) holds.*

This constraint qualification is widely used, e.g., to prove strong duality for instances of semi-definite programs (SDPs) (strong duality does not generally hold for SDPs even though they are convex) since finding an interior point within the feasible set is often easier than checking other constraint qualifications.

### 2.3.2.2 Mangasarian-Fromovitz Constraint Qualification

In this section, we provide the Mangasarian-Fromovitz constraint qualification (MFCQ) which will prove useful in a later section (see Sec. 3.2.2).

**Proposition 2.3** (Mangasarian-Fromovitz Constraint Qualification [76, Sec. 3, adapted]). *Let a smooth optimization problem be given as in (2.6) with critical point  $x^* \in \mathbb{R}^n$ , inequality constraints  $g(x^*) \leq 0$ , and equality constraints  $h(x^*) = 0$ , where  $g : \mathbb{R}^n \mapsto \mathbb{R}^m$  and  $h : \mathbb{R}^n \mapsto \mathbb{R}^o$ . Further, let  $\mathcal{A} = \{a \in \{1, \dots, m\} \mid g_a(x^*) = 0\}$  denote the set of active inequality constraints. If the gradients  $\nabla_x h_j(x^*)$  for  $1 \leq j \leq o$  of the equality constraints are linearly independent and there exists a vector  $d \in \mathbb{R}^n$  such that*

$$\begin{aligned} \nabla_x g_a(x^*)^T d &< 0, a \in \mathcal{A}, \\ \nabla_x h_j(x^*)^T d &= 0, 1 \leq j \leq o, \end{aligned}$$

*then  $x^*$  satisfies the Mangasarian-Fromovitz constraint qualification (MFCQ).*

### 2.3.3 Convex Optimization

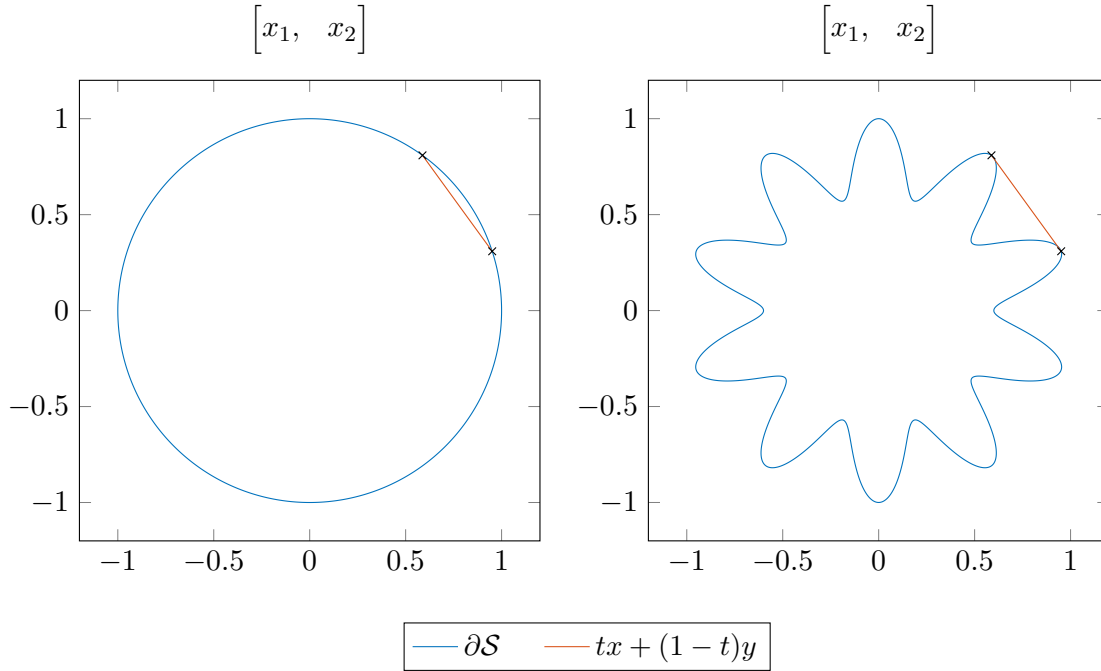
Being essential in a variety of everyday applications, convex optimization is widely adopted since, once a problem has been formulated as a convex optimization problem, one can generally (with a few edge cases as an exception) efficiently solve the problem [14, Sec. 1.3.2]. In more technical terms, convex optimization is such an important field in optimization since all stationary points of the Lagrangian are global optima under a suitable constraint qualification (see Sec. 2.3.2). Furthermore, there exist many mature numerical solvers for convex optimization problems, such as MOSEK<sup>1</sup>, Gurobi<sup>2</sup>, or SDPT3<sup>3</sup>, to only name a few. For the definition of a convex optimization problem, we first introduce convex sets and convex functions. For a thorough discussion on convex optimization, we refer the interested reader to [14].

<sup>1</sup><https://www.mosek.com/>

<sup>2</sup><https://www.gurobi.com/>

<sup>3</sup><https://blog.nus.edu.sg/mattohkc/software/sdpt3/>





**Figure 2.1:** Example of a convex set given by  $\left\{ r \begin{bmatrix} \cos \phi \\ \sin \phi \end{bmatrix} \mid 0 \leq r \leq 1, 0 \leq \phi \leq 2\pi \right\}$  and non-convex set  $\left\{ r (0.8 + 0.2 \cos(10\phi)) \begin{bmatrix} \cos \phi \\ \sin \phi \end{bmatrix} \mid 0 \leq r \leq 1, 0 \leq \phi \leq 2\pi \right\}$ . Additionally, we show the line connecting  $x = \begin{bmatrix} \cos(\frac{\pi}{10}) \\ \sin(\frac{\pi}{10}) \end{bmatrix}$  and  $y = \begin{bmatrix} \cos(\frac{3\pi}{10}) \\ \sin(\frac{3\pi}{10}) \end{bmatrix}$ .

**Definition 2.7** (Convex Set [14, Sec. 2.1.4]). A set  $\mathcal{S} \subseteq \mathbb{R}^n$  is said to be convex if and only if

$$\forall x, y \in \mathcal{S} \forall t \in [0, 1] : tx + (1-t)y \in \mathcal{S}. \quad (2.9)$$

Fig. 2.1 shows an example for a convex and a non-convex set, respectively, where the second set is not convex as it violates (2.9). ■

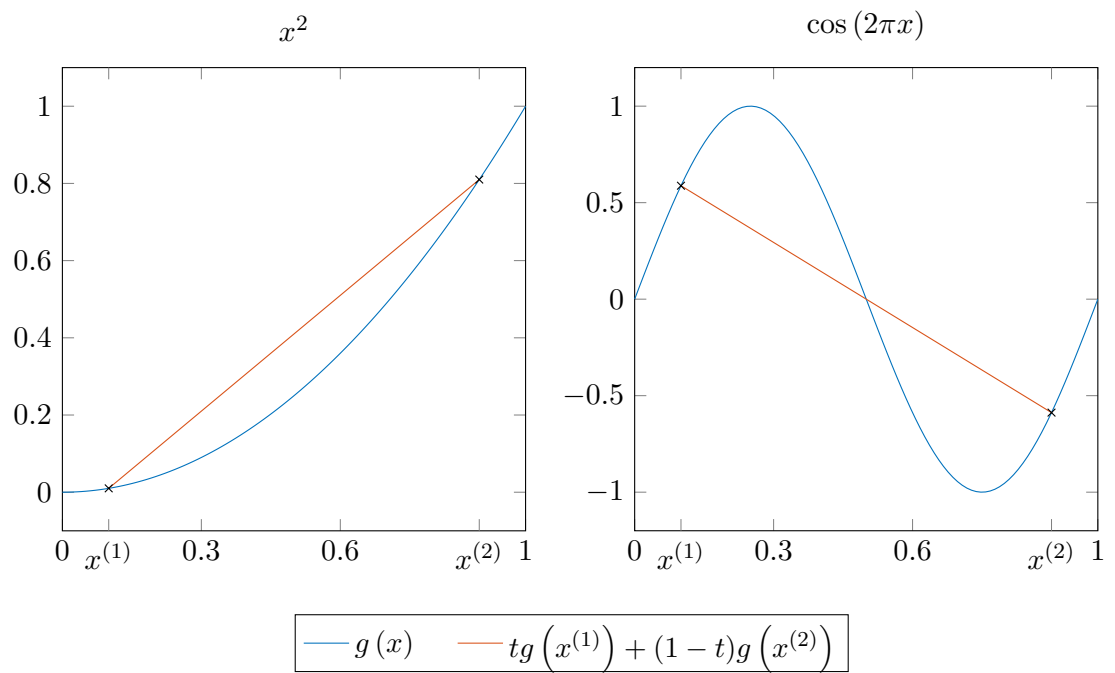
**Definition 2.8** (Convex Function [14, Sec. 3.1.1]). A function  $g : \mathcal{G} \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$  is convex if and only if

$$\forall x, y \in \mathcal{G} \forall t \in [0, 1] : g(tx + (1-t)y) \leq tg(x) + (1-t)g(y), \quad (2.10)$$

where  $\mathcal{G}$  is a convex set as introduced in Def. 2.7. ■

A function  $g$  is concave if  $-g$  is convex. Fig. 2.2 shows an example for one convex and non-convex function. Clearly, the second function is not convex as (2.10) is violated for  $x^{(1)} = 0.1$  and  $x^{(2)} = 0.9$ . Geometrically speaking, the second function is not convex since not all points on the line connecting  $x^{(1)}$  and  $x^{(2)}$  are above their respective function values (see (2.10)). We are now ready to define a convex optimization problem.

## 2 Background



**Figure 2.2:** Example of a convex function ( $g(x) = x^2$ ) and non-convex function ( $g(x) = \cos(2\pi x)$ ) evaluated over  $x \in [0, 1]$ . Additionally we show the line connecting  $(x^{(1)}, g(x^{(1)}))$  to  $(x^{(2)}, g(x^{(2)}))$ .

**Definition 2.9** (Convex Minimization Problem). We say that (2.6) is a convex minimization problem if and only if  $f$  and  $g$  are convex and  $h$  is affine. ■

Put differently, both the objective function  $f$  and the feasible set  $\mathcal{S}$  have to be convex. Next, we briefly introduce important subclasses of convex optimization problems which are relevant for this thesis. For a more detailed introduction, see, e.g., [14].

### 2.3.3.1 Linear Program

The objective and constraint functions in linear programs (LPs) are linear in the optimization variable, i.e.

$$\min_x c^T x, \quad (2.11a)$$

$$\text{s.t. } Cx \leq d, \quad (2.11b)$$

$$Ax = b, \quad (2.11c)$$

with optimization variable  $x \in \mathbb{R}^n$ , where  $c \in \mathbb{R}^n$ ,  $C \in \mathbb{R}^{m \times n}$ ,  $d \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{q \times n}$ , and  $b \in \mathbb{R}^q$ . Linear programs can be solved reliably and efficiently, even for a large number of optimization variables and constraints [14, Sec. 1.2.2], [22]. Since we assume that all optimization problems have a finite minimum, strong duality (see Prop. 2.4) always holds [25, Sec. 7.4].

### 2.3.3.2 Semi-Definite Program

With SDPs, one can formulate optimization problems of the form

$$\min_{x, X^{(i)}} c^T x + \sum_{i=1}^m \text{trace} \left( C^{(i)}, X^{(i)} \right), \quad (2.12a)$$

$$\text{s.t. } \underline{b}_j \leq l^{(j)T} x + \sum_{i=1}^N \text{trace} \left( A^{(ij)}, X^{(i)} \right) \leq \bar{b}_j, \quad 1 \leq j \leq q, \quad (2.12b)$$

$$\underline{d} \leq x \leq \bar{d}, \quad (2.12c)$$

$$x \in \mathcal{K}, \quad (2.12d)$$

$$X^{(i)} \succeq 0, \quad i \in \{1, \dots, N\}, \quad (2.12e)$$

where  $x \in \mathbb{R}^n$  contains  $n$  scalar optimization variables, and  $X^{(i)} \in \mathbb{S}_+^{n \times n}$  are  $m$  positive semi-definite optimization matrices. Further, we have  $c \in \mathbb{R}^n$ ,  $C^{(i)} \in \mathbb{S}^{n \times n}$ ,  $\underline{b} \in \mathbb{R}^q$ ,  $\bar{b} \in \mathbb{R}^q$ ,  $l^{(j)} \in \mathbb{R}^n$ ,  $A^{(ij)} \in \mathbb{S}^{n \times n}$ ,  $\underline{d} \in \mathbb{R}^n$ ,  $\bar{d} \in \mathbb{R}^n$ , and  $\mathcal{K} \subseteq \mathbb{R}^n$  with the Cartesian product of convex cones<sup>4</sup> (e.g., quadratic cone, rotated quadratic cone etc), where  $\text{trace} \left( C^{(i)}, X^{(i)} \right) = \sum_{k,l} C_{kl}^{(i)} X_{kl}^{(i)}$  denotes the Frobenius inner product ( $\square_{kl}$  denotes the element in the  $k$ -th row and  $l$ -th column).

In contrast to LPs, strong duality for SDPs additionally requires Slater's condition (see Prop. 2.2) to hold [51, Sec. 2.2]. Under that assumption, however, SDPs can be efficiently solved [110].

<sup>4</sup>for details, see, e.g.: <https://docs.mosek.com/MOSEKModelingCookbook-a4paper.pdf>

### 2.3.4 Non-Convex Optimization

While convex optimization offers great advantages, optimization problems in practice often cannot be formulated as a convex optimization problem. Thus, we introduce non-convex optimization problems in this section.

**Definition 2.10** (Non-Convex Optimization Problem). Let the optimization problem as defined in Def. 2.5 be given. We say that (2.6) is a non-convex optimization problem if it is not a convex optimization problem according to Def. 2.9. ■

In general, non-convex optimization problems cannot be solved efficiently to global optimality in general. As it turns out, even checking whether a feasible solution is a local minimum is in the NP-complete complexity class [84]. That said, when second-order methods are used, the number of function evaluations of the objective function and the constraint function required to converge to an  $\epsilon$ -critical point of the optimization problem is polynomially bounded in the inverse of the accuracy parameter  $\epsilon > 0$  [18]. A critical point hereby is a feasible point of (2.6) that is either a (global or local) minimum or a saddle point of (2.6).

### 2.3.5 Lagrangian Duality

When the optimization problem at hand cannot be reformulated as a convex optimization problem, finding a solution is substantially more difficult. However, instead of trying to find the optimum  $p^*$ , one may be satisfied with a lower bound on the original optimization problem in (2.6). In this section, we thus introduce the Lagrangian dual problem which is always a convex optimization problem and provides a lower bound on the optimum  $p^*$ .

**Definition 2.11** (Lagrangian Dual Problem [14, Chap. 5, adapted]). Let  $f$ ,  $g$ , and  $h$  be defined as in Def. 2.5, and let  $\phi$  be defined as in (2.7). Further, let  $\tilde{\mathcal{S}}$  collect any additional constraints on  $x$  that we do not wish to dualize. The dual problem of (2.6) is

$$d^* = \sup_{\lambda \geq 0, \mu} \inf_{x \in \tilde{\mathcal{S}}} \phi(x, \mu, \lambda), \quad (2.13)$$

where  $\lambda \in \mathbb{R}_{\geq 0}^m$  and  $\mu \in \mathbb{R}^o$  are the Lagrange multipliers for the inequality and equality constraints, respectively. ■

We collect relevant properties of the primal-dual relationship of (2.6) and (2.13) in the following proposition.

**Proposition 2.4** (Primal-Dual Relationship [14, Chap. 5, adapted]). *Let the primal and dual problem be given as in (2.6) and (2.13). The following properties are true:*

- *The dual problem is always a convex optimization problem, even if the primal problem is non-convex.*
- *The dual optimal value bounds the primal optimal value from below, i.e.,  $d^* \leq p^*$  (weak duality).*

- If the primal problem is convex and constraint qualifications (see Sec. 2.3.2) hold, then  $d^* = p^*$ , i.e., there is zero duality gap between the primal and the dual problem (strong duality).

### 2.3.6 Modeling

When formulating optimization problems, it often happens that this formulation cannot be directly fed to numerical optimization solvers. Examples include the objective function or constraints containing absolute value expressions or maximizations over multiple functions; then the optimization problem is no longer differentiable. Subsequently, we derive equivalent but smooth reformulations of such optimization problems.

#### 2.3.6.1 Absolute Values

Let

$$\min_x \left\{ f(x) + \sum_{i=1}^q |r_i(x)| \right\}, \quad (2.14a)$$

$$\text{s.t. } g(x) + \sum_{j=1}^p |c_j(x)| \leq 0, \quad (2.14b)$$

be given, where  $x \in \mathbb{R}^n$ ,  $f : \mathbb{R}^n \mapsto \mathbb{R}$ ,  $r : \mathbb{R}^n \mapsto \mathbb{R}^q$ ,  $g : \mathbb{R}^n \mapsto \mathbb{R}$ ,  $c : \mathbb{R}^n \mapsto \mathbb{R}^p$ , and we assume that all functions are sufficiently smooth. We only consider a single inequality constraint here as the following derivation can be easily extended to multiple inequality constraints. The non-differentiability of (2.14) due to the absolute values can be avoided by transforming (2.14) into the equivalent but differentiable optimization problem

$$\min_{x,s,v} \left\{ f(x) + \sum_{i=1}^q s_i \right\}, \quad (2.15a)$$

$$\text{s.t. } g(x) + \sum_{j=1}^p v_j \leq 0, \quad (2.15b)$$

$$\begin{bmatrix} +r(x) \\ -r(x) \end{bmatrix} \leq \begin{bmatrix} s \\ s \end{bmatrix}, \quad (2.15c)$$

$$\begin{bmatrix} +c(x) \\ -c(x) \end{bmatrix} \leq \begin{bmatrix} v \\ v \end{bmatrix}, \quad (2.15d)$$

with  $s \in \mathbb{R}^q$  and  $v \in \mathbb{R}^p$ . This transformation can be justified as follows:

Let  $\hat{x}$  denote the minimizer of (2.14) and let  $(x^*, s^*, v^*)$  denote the minimizer of (2.15). Further, we define

$$\mathcal{S} = \left\{ x \in \mathbb{R}^n \mid g(x) + \mathbf{1}_p^T |c(x)| \leq 0 \right\},$$

$$\tilde{\mathcal{S}} = \left\{ x \in \mathbb{R}^n \mid \exists v \geq |c(x)| : g(x) + \mathbf{1}_p^T v \leq 0 \right\},$$

## 2 Background

where  $\tilde{\mathcal{S}} = \mathcal{S}$  follows since choosing  $v = |c(x)|$  is always possible. We have

$$\begin{aligned} f(\hat{x}) + 1_q^T |r(\hat{x})| &= \min_{x \in \mathcal{S}} f(x) + 1_q^T |r(x)| = \min_{x \in \tilde{\mathcal{S}}} f(x) + 1_q^T |r(x)| \\ &\leq \min_{\substack{x \in \tilde{\mathcal{S}} \\ |r(x)| \leq s}} f(x) + 1_q^T s = f(x^*) + 1_q^T s^*. \end{aligned} \quad (2.16)$$

On the other hand,  $(\hat{x}, |r(\hat{x})|, |c(\hat{x})|)$  is a feasible point of (2.15), and by definition

$$f(x^*) + 1_q^T s^* \leq f(\hat{x}) + 1_q^T |r(\hat{x})|. \quad (2.17)$$

Combining (2.16) and (2.17) yields the equivalence of (2.14) and (2.15).

### 2.3.6.2 Finite Maximization

Let

$$\min_x f(x) + \max_{1 \leq i \leq q} r_i(x), \quad (2.18a)$$

$$\text{s.t. } g(x) + \max_{1 \leq j \leq p} c_j(x) \leq 0, \quad (2.18b)$$

with  $x \in \mathbb{R}^n$ ,  $f : \mathbb{R}^n \mapsto \mathbb{R}$ ,  $r : \mathbb{R}^n \mapsto \mathbb{R}^q$ ,  $g : \mathbb{R}^n \mapsto \mathbb{R}$ ,  $c : \mathbb{R}^n \mapsto \mathbb{R}^p$ , and where all functions are assumed to be sufficiently smooth. We limit ourselves to one inequality constraint, as an extension to multiple inequality constraints is straightforward. Further, we define  $\hat{x}$  as a minimizer of (2.14) and  $(x^*, s^*, v^*)$  as a minimizer of (2.15). To avoid the maximizations in (2.18), we reformulate the problem as

$$\min_{x, s, v} f(x) + s, \quad (2.19a)$$

$$\text{s.t. } g(x) + v \leq 0, \quad (2.19b)$$

$$r(x) \leq s 1_q, \quad (2.19c)$$

$$c(x) \leq v 1_p, \quad (2.19d)$$

with  $s \in \mathbb{R}$  and  $v \in \mathbb{R}$ . To show equivalence of (2.18) and (2.19), we first define

$$\mathcal{S} = \left\{ x \in \mathbb{R}^n \mid g(x) + \max_{1 \leq j \leq p} c_j(x) \leq 0 \right\},$$

$$\tilde{\mathcal{S}} = \left\{ x \in \mathbb{R}^n \mid \exists v \geq \max_{1 \leq j \leq p} c_j(x) : g(x) + v \leq 0 \right\},$$

from  $\tilde{\mathcal{S}} = \mathcal{S}$  follows since choosing  $v = \max_{1 \leq j \leq p} c_j(x)$  is always possible. Further, we have

$$\begin{aligned} f(\hat{x}) + \max_{1 \leq i \leq q} r_i(\hat{x}) &= \min_{x \in \mathcal{S}} f(x) + \max_{1 \leq i \leq q} r_i(x) = \min_{x \in \tilde{\mathcal{S}}} f(x) + \max_{1 \leq i \leq q} r_i(x) \\ &\leq \min_{\substack{x \in \tilde{\mathcal{S}}, s \\ \max_{1 \leq i \leq q} r_i(x) \leq s}} f(x) + s = f(x^*) + s^*. \end{aligned} \quad (2.20)$$

On the other hand,  $(\hat{x}, \max_{1 \leq i \leq q} r_i(\hat{x}), \max_{1 \leq j \leq p} c_j(\hat{x}))$  is a feasible point of (2.19), and thus by definition

$$f(x^*) + \max_{1 \leq i \leq q} r_i(x^*) \leq f(\hat{x}) + \max_{1 \leq i \leq q} r_i(\hat{x}), \quad (2.21)$$

which yields the equivalence of (2.18) and (2.19) by combining (2.20) and (2.21).

## 2.4 Linear-Quadratic Regulator

In this section, we introduce LQR control [56], which will be used in later sections of this thesis. Let

$$\dot{x}(t) = Ax(t) + Bu(t),$$

be a given LTI system  $(A, B)$ , where  $x(t) \in \mathbb{R}^{n_x}$  denotes the state at time  $t \in \mathbb{R}_{\geq 0}$ ,  $u(t) \in \mathbb{R}^{n_u}$  denotes the controllable input,  $A \in \mathbb{R}^{n_x \times n_x}$  is the system matrix, and  $B \in \mathbb{R}^{n_x \times n_u}$  denotes the input matrix. For the introduction of LQR control, we first require the concept of controllability.

**Corollary 2.1** (Controllability [56, adapted]). *An LTI system  $(A, B)$  as defined above, for which*

$$\text{rank} \left( \begin{bmatrix} B & AB & A^2B & \dots & A^{n_x-1}B \end{bmatrix} \right) = n_x,$$

*holds, is controllable.*

Put differently, an LTI system is controllable if and only if there exists an input trajectory such that any finite state  $x(t)$  can be reached from any initial state  $x(0)$  in finite time. Thus, assume that the LTI system  $(A, B)$  is controllable as per Corr. 2.1. The unique, optimal controller which minimizes

$$\min_{u(t)} \int_0^\infty \left( x(t)^T Qx(t) + u(t)^T Ru(t) \right) dt, \quad (2.22)$$

with  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $R \in \mathbb{S}_{++}^{n_u \times n_u}$  (weaker conditions for  $Q$  and  $R$  are possible) is

$$u(t) = Kx(t), \quad (2.23)$$

where

$$K = -R^{-1}B^T X, \quad (2.24)$$

and  $X \in \mathbb{S}_{++}^{n_x \times n_x}$  is the unique, positive-definite solution to the algebraic Riccati equation

$$A^T X + XA - XBR^{-1}B^T X + Q = 0. \quad (2.25)$$

For later convenience, we further state that solving the algebraic Riccati equation has the same computational complexity as computing the eigendecomposition of an  $n_x$ -by- $n_x$  dimensional matrix [69]; thus, computing the gain matrix  $K$  for a given system has the same computational complexity since it otherwise only involves matrix multiplications.

## 2.5 Set Operations

In this thesis, we synthesize controllers for sets of initial states instead of a single initial state. As a result, we need to perform arithmetic operations on sets, which we define next.

**Definition 2.12** (Minkowski Sum). The Minkowski sum of two sets  $\mathcal{A} \subseteq \mathbb{R}^n$  and  $\mathcal{B} \subseteq \mathbb{R}^n$  is defined by

$$\mathcal{A} \oplus \mathcal{B} = \{a + b \mid a \in \mathcal{A}, b \in \mathcal{B}\}.$$

**Definition 2.13** (Linear Map). The linear map of a set  $\mathcal{S} \subseteq \mathbb{R}^n$  with matrix  $M \in \mathbb{R}^{m \times n}$  is given by

$$M\mathcal{S} = \{Ms \mid s \in \mathcal{S}\}.$$

**Definition 2.14** (Quadratic Map [5]). The quadratic map of a set  $\mathcal{S} \subseteq \mathbb{R}^n$  with a set  $\mathcal{Q} = \{Q^{(1)}, \dots, Q^{(m)}\}$  of matrices  $Q^{(i)} \in \mathbb{R}^{n \times n}$  for  $1 \leq i \leq m$  is defined as

$$\text{sq}(\mathcal{S}, \mathcal{Q}) = \left\{x \in \mathbb{R}^m \mid x_i = s^T Q^{(i)} s, 1 \leq i \leq m, s \in \mathcal{S}\right\}.$$

In addition to operations that again produce a set, we further define the following operations on sets which return a scalar value.

**Definition 2.15** (Support Function [41, Def. 1, adapted]). The support function of a convex set  $\mathcal{S} \subseteq \mathbb{R}^n$  in direction  $l \in \mathbb{R}^n$  is defined by

$$\rho_{\mathcal{S}}(l) = \sup_{s \in \mathcal{S}} l^T s.$$

For a set  $\mathcal{S}(z)$  parameterized in some vector  $z \in \mathbb{R}^n$ , we denote its support function with  $\rho_{\mathcal{S}}(l, z)$ . For later convenience, we further define the short-hand

$$\rho_{\mathcal{S}}(L) = \left[ \rho_{\mathcal{S}}(L_{(1,:)}^T), \rho_{\mathcal{S}}(L_{(2,:)}^T), \dots, \rho_{\mathcal{S}}(L_{(m,:)}^T) \right]^T,$$

where  $L \in \mathbb{R}^{m \times n}$ .

Lastly, we later also require a measure of similarity for two sets, for which the Hausdorff distance is defined next.

**Definition 2.16** (Hausdorff Distance [113, Def. 2]). The Hausdorff distance of two non-empty sets  $\mathcal{X} \subseteq \mathbb{R}^n$  and  $\mathcal{Y} \subseteq \mathbb{R}^n$  is defined by

$$d_H(\mathcal{X}, \mathcal{Y}) = \max \left( \sup_{x \in \mathcal{X}} \inf_{y \in \mathcal{Y}} \|x - y\|_2, \sup_{y \in \mathcal{Y}} \inf_{x \in \mathcal{X}} \|x - y\|_2 \right).$$

Intuitively, we can say that the smaller the Hausdorff distance between two sets, the more similar (e.g. visually) these two sets are.



**Table 2.1:** Closedness of intervals ( $\mathcal{I}$ ), zonotopes ( $\mathcal{Z}$ ), polynomial zonotopes ( $\mathcal{PZ}$ ), H-polytopes ( $\mathcal{H}$ ), and ellipsoids ( $\mathcal{E}$ ) under common operations [6, Tab. 1]. We differentiate between: Result can be efficiently and exactly computed ( $\checkmark$ ); result cannot be exactly represented ( $\times$ ).

	$\mathcal{I}$	$\mathcal{Z}$	$\mathcal{PZ}$	$\mathcal{H}$	$\mathcal{E}$
Linear Map	$\times$	$\checkmark$	$\checkmark$	$\times$	$\checkmark$
Minkowski Sum	$\checkmark$	$\checkmark$	$\checkmark$	$\times$	$\times$
Quadratic Map	$\times$	$\times$	$\checkmark$	$\times$	$\times$

## 2.6 Set Representations

In order to compute with sets, we need to represent them appropriately. Ideally, we find a set representation that is closed under all relevant operations: A set representation is said to be closed under an operation if the result of applying said operation can again be represented in the same set representation. That said, while being closed under relevant operations is important, we also need to be able to perform these set operations efficiently. As a result, there exist many different set representations in the literature which are closed under different set operations and have varying computational effort (see Tab. 2.1). Choosing the right set representation thus heavily depends on the algorithm we wish to implement. In this section, we therefore introduce all set representations which are relevant in this thesis.

Intuitively, the subsequently introduced set representations differ in how they parameterize an element of a set. The set itself is then generated by including all states that can be parameterized in that particular way.

**Definition 2.17** (Set Generation [35, Def. 3]). We define

$$\{s(M)\}_M = \left\{s(M) \mid M \in [-1, 1]^{m \times n}\right\} = \mathcal{S},$$

where  $s : [-1, 1]^{m \times n} \rightarrow \mathbb{R}^n$  is the generating function of  $\mathcal{S} \subseteq \mathbb{R}^n$  with dependent factors  $M \in [-1, 1]^{m \times n}$ . We say that  $\mathcal{S}$  is generated by  $s(M)$  over  $M$ . ■

Thus, the generating function allows us to express any element of a set in a unified way, i.e., by only varying  $M$  within the unit hypercube.

### 2.6.1 Intervals

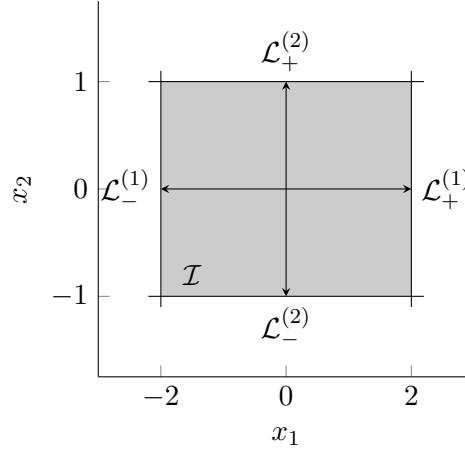
As the most fundamental of set representations, intervals are widely used due to their compact representation size.

**Definition 2.18** (Interval [4, Def. 2.4, adapted]). The interval with lower bound  $\underline{x} \in \mathbb{R}^n$  and upper bound  $\bar{x} \in \mathbb{R}^n$  is defined as

$$\mathcal{I} = [\underline{x}, \bar{x}] = \left\{ \frac{1}{2}(\bar{x} + \underline{x}) + \frac{1}{2}(\bar{x} - \underline{x})\beta \right\}_\beta = \{x \in \mathbb{R}^n \mid \underline{x} \leq x \leq \bar{x}\},$$

where inequalities are to be interpreted element-wise and  $\underline{x} \leq \bar{x}$ . ■

## 2 Background



**Figure 2.3:** Visualization of the interval  $\mathcal{I} = \left[ [-2, -1]^T, [2, 1]^T \right]$ , where  $\mathcal{L}_{\pm}^{(i)} = \left\{ x \in \mathbb{R}^2 \mid \pm e_{(2)}^{(i)T} x = \rho_{\mathcal{I}}(\pm e_{(2)}^{(i)}) \right\}$  for  $1 \leq i \leq 2$  defines the hyperplanes which bound this interval.

Fig. 2.3 visualizes a simple interval  $\mathcal{I} = \left[ \begin{bmatrix} -2 \\ -1 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right]$ .

**Minkowski sum [4, Sec. 2.6, adapted]:** Given two intervals  $\mathcal{I} = [\underline{x}, \bar{x}] \subseteq \mathbb{R}^n$  and  $\mathcal{L} = [\underline{l}, \bar{l}] \subseteq \mathbb{R}^n$ , their Minkowski sum is

$$\mathcal{I} \oplus \mathcal{L} = [\underline{x} + \underline{l}, \bar{x} + \bar{l}].$$

**Linear map:** Since intervals are axis-aligned boxes and because linear transformations generally include rotations, intervals are not closed under linear maps. That said, the tight outer approximation of a linear map  $M \in \mathbb{R}^{m \times n}$  of an interval  $\mathcal{I} \subseteq \mathbb{R}^n$  is given by

$$M\mathcal{I} \subseteq \left\{ x \in \mathbb{R}^m \mid \begin{bmatrix} I_m \\ -I_m \end{bmatrix} x \leq \rho_{M\mathcal{I}} \left( \begin{bmatrix} I_m \\ -I_m \end{bmatrix} \right) \right\} = [-\rho_{M\mathcal{I}}(-I_m), \rho_{M\mathcal{I}}(I_m)],$$

since  $\forall x \in M\mathcal{I} : \pm e_m^{(i)T} x \leq \rho_{M\mathcal{I}}(\pm e_m^{(i)})$  for  $1 \leq i \leq m$  holds by definition (see Def. 2.15). Here,  $\rho_{M\mathcal{I}}(\cdot)$  can be computed by first representing  $M\mathcal{I}$  as a zonotope and then computing the resulting support function (see Sec. 2.6.5).

**Quadratic map:** Naturally, intervals cannot be closed under quadratic maps as they are convex sets. That said, one may simply write the interval as a zonotope and apply the outer approximation of the quadratic map for zonotopes from [7, Lem. 1].

**Support function:** Since each dimension of an interval  $\mathcal{I} \subseteq \mathbb{R}^n$  is independent, it follows from Def. 2.15 that its support function in direction  $l \in \mathbb{R}^n$  is given by

$$\rho_{\mathcal{I}}(l) = \max_{x \leq \bar{x}} \sum_{k=1}^n l_k x_k = \sum_{k=1}^n \max_{x_k \leq \bar{x}_k} l_k x_k = \sum_{k=1}^n \begin{cases} l_k \bar{x}_k, & \text{sign}(x_k) \geq 0, \\ l_k \underline{x}_k, & \text{sign}(x_k) < 0 \end{cases}.$$

### 2.6.2 Arbitrary Convex Set

While support functions as defined in Def. 2.15 are often used as set operations, they can also be used to represent arbitrary convex sets:

$$\mathcal{S} = \left\{ x \in \mathbb{R}^n \mid \forall l \in \mathbb{R}^n : l^T x \leq \rho_{\mathcal{S}}(l) \right\}.$$

Fig. 2.3 visualizes the construction of a 2D interval using support functions.

**Minkowski sum [41, Prop. 2, adapted]:** The Minkowski sum between two convex sets  $\mathcal{A} \in \mathbb{R}^n$  and  $\mathcal{B} \in \mathbb{R}^n$  is given by

$$\mathcal{A} \oplus \mathcal{B} = \left\{ x \in \mathbb{R}^n \mid \forall l \in \mathbb{R}^n : l^T x \leq \rho_{\mathcal{A}}(l) + \rho_{\mathcal{B}}(l) \right\}, \quad (2.26)$$

i.e.,  $\rho_{\mathcal{A} \oplus \mathcal{B}}(l) = \rho_{\mathcal{A}}(l) + \rho_{\mathcal{B}}(l)$  for  $l \in \mathbb{R}^n$ .

**Linear map [41, Prop. 2, adapted]:** The linear map  $M \in \mathbb{R}^{m \times n}$  of a convex set  $\mathcal{A} \subseteq \mathbb{R}^n$  represented using support functions is given by

$$M\mathcal{A} = \left\{ x \in \mathbb{R}^m \mid \forall l \in \mathbb{R}^m : l^T x \leq \rho_{\mathcal{A}}(M^T l) \right\}, \quad (2.27)$$

i.e.,  $\rho_{M\mathcal{A}}(l) = \rho_{\mathcal{A}}(M^T l)$  for  $l \in \mathbb{R}^m$ .

### 2.6.3 Ellipsoids

Next, we define ellipsoids which – while being similarly compact in representation size as intervals – further allow for arbitrary rotations.

**Definition 2.19** (Non-Degenerate Ellipsoid [67, Def. A.1]). The non-degenerate ellipsoid  $\langle c, Q \rangle_E$  with center  $c \in \mathbb{R}^n$  and shape matrix  $Q \in \mathbb{S}_{++}^{n \times n}$  is given by

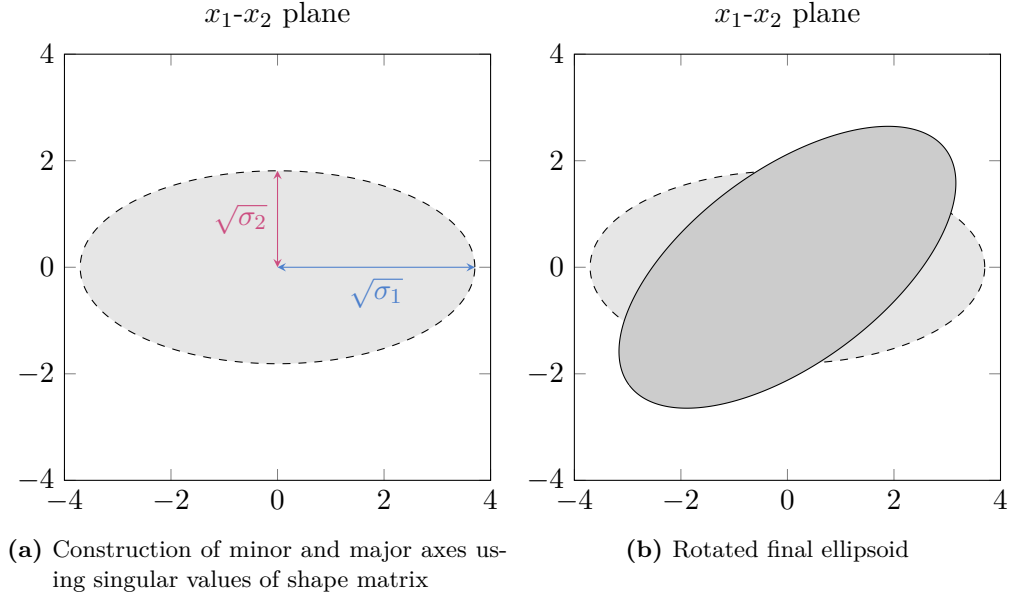
$$\langle c, Q \rangle_E = \left\{ x \in \mathbb{R}^n \mid (x - c)^T Q^{-1} (x - c) \leq 1 \right\}.$$

■

The construction of an ellipsoid is demonstrated subsequently.

**Example 2.3.** Let  $\mathcal{E} = \langle 0, Q \rangle_E = \left\langle 0, \begin{bmatrix} 10 & 5 \\ 5 & 7 \end{bmatrix} \right\rangle_E$ . Since  $Q \succ 0$ , one can construct the ellipsoid easily using the singular value decomposition, which is visualized in Fig. 2.4: With the singular value decomposition given by  $Q = U^T \text{diag}(\sigma) U$ , we first construct

## 2 Background



**Figure 2.4:** Construction of the ellipsoid from Ex. 2.3.

an axis-aligned ellipsoid  $\langle 0, \text{diag}(\sigma) \rangle_E$  which extends  $\sqrt{\sigma_1}$  in the first and  $\sqrt{\sigma_2}$  in the second unit direction. The original ellipsoid is then constructed by a simple rotation since  $U \langle 0, \text{diag}(\sigma) \rangle_E = \langle 0, U \text{diag}(\sigma) U^T \rangle_E = \mathcal{E}$  (see below for details on the linear map of an ellipsoid). ■

**Minkowski sum:** Ellipsoids are not closed under the Minkowski sum. However, for  $m$  ellipsoids  $\langle q^{(i)}, Q^{(i)} \rangle_E \subset \mathbb{R}^n$  with  $1 \leq i \leq m$ , an efficient outer approximation  $\langle -\hat{W}^{-1}\hat{b}, \hat{W}^{-1} \rangle_E$  can be found as the solution of the SDP [15, Sec. 3.7.4] with

$$\begin{aligned} \hat{W}, \hat{b} &= \arg \min_{W, b, \tau} \log \det W^{-1}, \\ \text{s.t. } & \begin{bmatrix} E^{(0)T} W E^{(0)}, & E^{(0)T} b, & 0 \\ \left( E^{(0)T} b \right)^T, & -1, & b^T \\ 0, & b, & -W \end{bmatrix} - \sum_{i=1}^m \tau_i \begin{bmatrix} \tilde{A}^{(i)}, & \tilde{b}^{(i)}, & 0 \\ \tilde{b}^{(i)T}, & c^{(i)}, & 0 \\ 0, & 0, & 0 \end{bmatrix} \preceq 0, \\ & \tau_i \geq 0, \forall i \in \{1, \dots, m\}, \\ & W \succ 0, \end{aligned}$$

where

$$\begin{aligned}
 E^{(i)} &= \begin{bmatrix} 0_{n \times ((i-1)n)}, & I_n, & 0_{n \times ((m-i)n)} \end{bmatrix}, \\
 E^{(0)} &= \sum_{i=1}^m E^{(i)}, \\
 \tilde{A}^{(i)} &= E^{(i)T} A^{(i)} E^{(i)}, \\
 \tilde{b}^{(i)} &= E^{(i)T} b^{(i)}, \\
 A^{(i)} &= Q^{(i)-1}, \\
 b^{(i)} &= -Q^{(i)-1} q^{(i)}, \\
 c^{(i)} &= q^{(i)T} Q^{(i)-1} q^{(i)}.
 \end{aligned}$$

Here,  $0_{q \times p}$  denotes the  $q$ -by- $p$ -dimensional zero matrix. For alternative approaches to the computation of the Minkowski sum, we refer the reader to [47].

**Linear map [66, Sec. II]:** For  $A \in \mathbb{R}^{m \times n}$ , the linear map of an ellipsoid  $\langle c, Q \rangle_E \subset \mathbb{R}^n$  is

$$A \langle c, Q \rangle_E = \langle Ac, AQA^T \rangle_E. \quad (2.28)$$

**Quadratic map:** Since ellipsoids are convex sets, they cannot be closed under the quadratic map.

**Support function [66, Sec. II]:** The support function of an ellipsoid  $\mathcal{E} = \langle c, Q \rangle_E \subset \mathbb{R}^n$  for a direction  $l \in \mathbb{R}^n$  is

$$\rho_{\mathcal{E}}(l) = l^T c + \sqrt{l^T Q l}. \quad (2.29)$$

Using the support function of an ellipsoid allows for a definition which includes degenerate ellipsoids.

**Definition 2.20** (Ellipsoid [67, Def. A.3]). The ellipsoid  $\mathcal{E} = \langle c, Q \rangle_E$  with center  $c \in \mathbb{R}^n$  and shape matrix  $Q \in \mathbb{S}_+^{n \times n}$  is given by

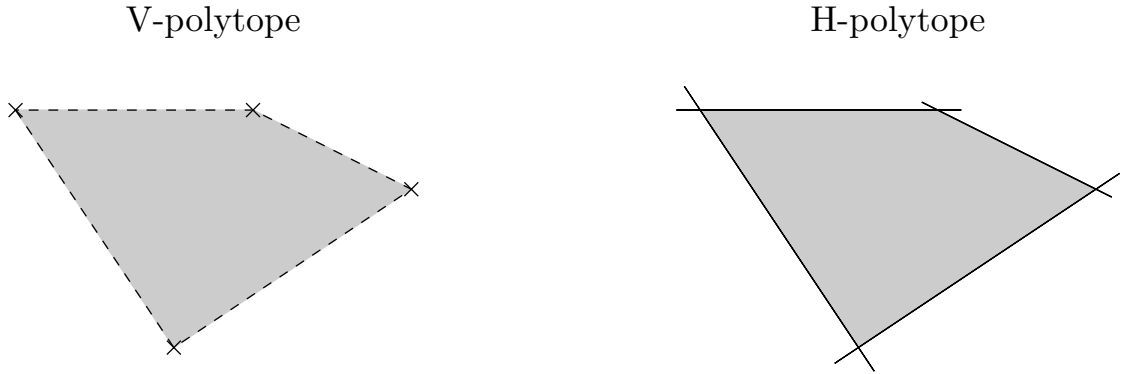
$$\mathcal{E} = \left\{ x \in \mathbb{R}^n \mid \forall l \in \mathbb{R}^n : l^T x \leq \rho_{\mathcal{E}}(l) \right\}.$$

■

## 2.6.4 H-Polytopes

Typically, one distinguishes between V-polytopes, which are represented using their vertices, and H-polytopes, which enclose a given polytope with a finite number of halfspaces (see Fig. 2.5). In this thesis, we only require convex H-polytopes as defined next.

## 2 Background



**Figure 2.5:** Visualization of a two-dimensional polytope, defined by its vertices (left) and its half-spaces (right).

**Definition 2.21** (H-Polytope [4, Def. 2.21, adapted]). A polytope  $\mathcal{H} = \langle A, b \rangle_H \subseteq \mathbb{R}^n$  with its halfspaces contained in  $A \in \mathbb{R}^{m \times n}$  and its offsets given in  $b \in \mathbb{R}^m$  is defined as

$$\mathcal{H} = \{x \in \mathbb{R}^n \mid Ax \leq b\}.$$

■

Other representations for polytopes exist, such as the Z-representation which uses polynomial zonotopes to represent polytopes [60].

**Minkowski sum:** H-polytopes are closed under the Minkowski sum since V-polytopes are closed under linear maps and both representations are equivalent [120, Th. 1.1]. However, the Minkowski sum between two H-polytopes is NP-hard since the conversion to an V-polytope requires vertex enumeration [109, Sec. 3].

That said, an efficiently computable outer approximation of the Minkowski sum  $\mathcal{A} \oplus \mathcal{B}$  for two H-polytopes  $\mathcal{A} \subseteq \mathbb{R}^n$  and  $\mathcal{B} \subseteq \mathbb{R}^n$  is given by

$$\mathcal{A} \oplus \mathcal{B} \subseteq \langle L, \rho_{\mathcal{A} \oplus \mathcal{B}}(L) \rangle_H = \langle L, \rho_{\mathcal{A}}(L) + \rho_{\mathcal{B}}(L) \rangle_H,$$

where  $L \in \mathbb{R}^{m \times n}$ : With a given choice of directions in  $L$ ,  $\rho_{\mathcal{A} \oplus \mathcal{B}}(L) = \rho_{\mathcal{A}}(L) + \rho_{\mathcal{B}}(L)$  follows from (2.26) and thus the result follows since  $\forall c \in \mathcal{A} \oplus \mathcal{B} : L_{(i,:)}c \leq \rho_{\mathcal{A} \oplus \mathcal{B}}(L_{(i,:)})$  holds due to Def. 2.15 for  $1 \leq i \leq m$  in. Thus, the outer approximation of  $\mathcal{A} \oplus \mathcal{B}$  is also tight in directions  $L_{(i,:)}$ .

**Linear map:** H-polytopes are closed under linear maps since V-polytopes are closed under linear maps and both representations are equivalent [120, Th. 1.1]. However, a linear transformation generally rotates and projects the polytope, and computing the projection of an H-polytope is NP-hard [108].

Similarly to its Minkowski sum, however, an efficiently computable outer approximation of the linear map  $M\mathcal{A}$  of an H-polytope  $\mathcal{A} \subseteq \mathbb{R}^n$  with matrix  $M \in \mathbb{R}^{q \times n}$ , which is

tight in  $m$  directions collected in  $L \in \mathbb{R}^{m \times q}$ , is given by

$$MA \subseteq \langle L, \rho_{MA}(L) \rangle_H = \langle L, \rho_A(LM) \rangle_H.$$

This follows, analogously to the Minkowski sum above, from Def. 2.15 and (2.27).

**Quadratic Map:** Since H-polytopes are convex sets, they cannot be closed under the quadratic map.

**Support function:** The support function of an H-polytope  $\mathcal{H} = \langle A, b \rangle_H$  follows directly from Def. 2.15 as

$$\rho_{\mathcal{H}}(l) = \max_{Ax \leq b} l^T x.$$

### 2.6.5 Zonotopes

Even though intervals and ellipsoids can be compactly represented, they lack representational power. Zonotopes are a popular choice for set-based computations, as they are closed under linear maps and the Minkowski sum.

**Definition 2.22** (Zonotope [40, Def. 1, adapted]). A zonotope  $\mathcal{Z} = \langle c, G \rangle_Z$  with center  $c \in \mathbb{R}^n$  and generator matrix  $G \in \mathbb{R}^{n \times m}$  consisting of generators  $G_{(:,i)}$  for  $1 \leq i \leq m$  is given by

$$\mathcal{Z} = \langle c, G \rangle_Z = \{c + G\nu\}_\nu,$$

where  $c + G\nu$  is its generating function with dependent factors  $\nu \in [-1, 1]^m$ . ■

We call an  $n$ -dimensional zonotope with  $n$  generators a parallelotope. We illustrate the construction of a zonotope in the following example.

**Example 2.4.** Let  $\mathcal{Z} = \left\langle 0, \begin{bmatrix} 1, & 3, & 1 \\ 2, & 0, & 1 \end{bmatrix} \right\rangle_Z$  be given. Fig. 2.6 shows the construction of the zonotope for all three generators: Since

$$\mathcal{Z} = \left\langle 0, \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right\rangle_Z \oplus \left\langle 0, \begin{bmatrix} 3 \\ 0 \end{bmatrix} \right\rangle_Z \oplus \left\langle 0, \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\rangle_Z,$$

the zonotope can be constructed by performing sequential Minkowski additions of each generator. ■

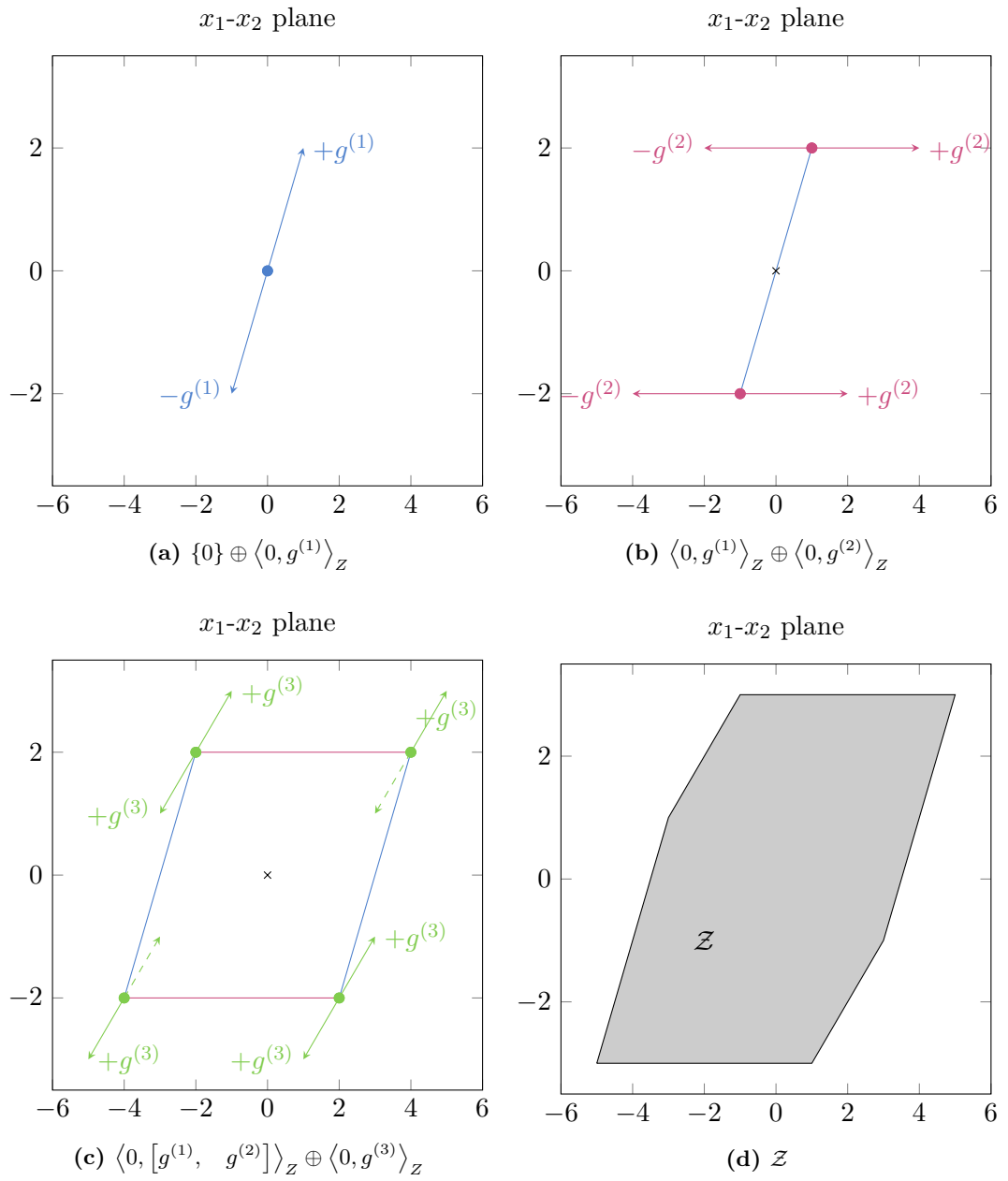
**Minkowski sum [40, Sec. 2]:** Given two zonotopes  $\langle c, G \rangle_Z \subset \mathbb{R}^n$  and  $\langle d, M \rangle_Z \subset \mathbb{R}^n$ , their Minkowski sum is

$$\langle c, G \rangle_Z \oplus \langle d, M \rangle_Z = \left\langle c + d, \begin{bmatrix} G, & M \end{bmatrix} \right\rangle_Z.$$

**Linear map [40, Sec. 2]:** The linear map  $A \in \mathbb{R}^{m \times n}$  of a zonotope  $\langle c, G \rangle_Z \subset \mathbb{R}^n$  is

$$A \langle c, G \rangle_Z = \langle Ac, AG \rangle_Z.$$

2 Background



**Figure 2.6:** Construction of the zonotope  $\mathcal{Z} = [g^{(1)}, g^{(2)}, g^{(3)}]$  from Ex. 2.4.



**Quadratic map:** Since sets described with support functions are convex by definition, they cannot be closed under quadratic maps. That said, one can apply the efficient outer approximation from [7, Lem. 1].

**Support function [45, Prop 2.2, adapted]:** The support function of a zonotope  $\mathcal{Z} = \langle c, G \rangle_Z \subset \mathbb{R}^n$  in direction  $l \in \mathbb{R}^n$  is given by

$$\rho_{\mathcal{Z}}(l) = c^T l + \left\| G^T l \right\|_1. \quad (2.30)$$

### 2.6.6 Polynomial Zonotopes

As an extension to zonotopes, polynomial zonotopes were first introduced in [5] to represent non-convex reachable sets (see Sec. 2.7.3).

**Definition 2.23** (Polynomial Zonotope [61, Def. 1, adapted]). Let  $c \in \mathbb{R}^n$  be the starting point,  $G \in \mathbb{R}^{n \times m}$  the dependent generator matrix with generators  $G_{(:,i)} \in \mathbb{R}^n$  for  $1 \leq i \leq m$ ,  $G_{\text{indep}} \in \mathbb{R}^{n \times o}$  the independent generator matrix, and  $E \in \mathbb{N}^{d \times m}$  the exponent matrix. We define the generating function of a polynomial zonotope as

$$g(\nu, \xi) = c + \sum_{i=1}^m G_{(:,i)} \nu^{E_{(:,i)}} + G_{\text{indep}} \xi,$$

with dependent factors  $\nu \in [-1, 1]^d$  and independent factors  $\xi \in [-1, 1]^o$ . The polynomial zonotope is then constructed as  $\mathcal{PZ} = \{g(\nu, \xi)\}_{\nu, \xi}$  with the shorthand  $\mathcal{PZ} = \langle c, G, G_{\text{indep}}, E \rangle_{\mathcal{PZ}}$ . ■

The following example demonstrates the construction of a simple, non-convex polynomial zonotope.

**Example 2.5.** Let

$$\mathcal{PZ} = \left\langle 0, \begin{bmatrix} 1, & 2, & 1 \\ 2, & 1, & 0 \end{bmatrix}, 0, \begin{bmatrix} 1, & 0, & 2 \\ 0, & 1, & 0 \end{bmatrix} \right\rangle_{\mathcal{PZ}} = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix} \beta_1 + \begin{bmatrix} 2 \\ 1 \end{bmatrix} \beta_2 + \begin{bmatrix} 1 \\ 0 \end{bmatrix} \beta_1^2 \right\}_{\beta},$$

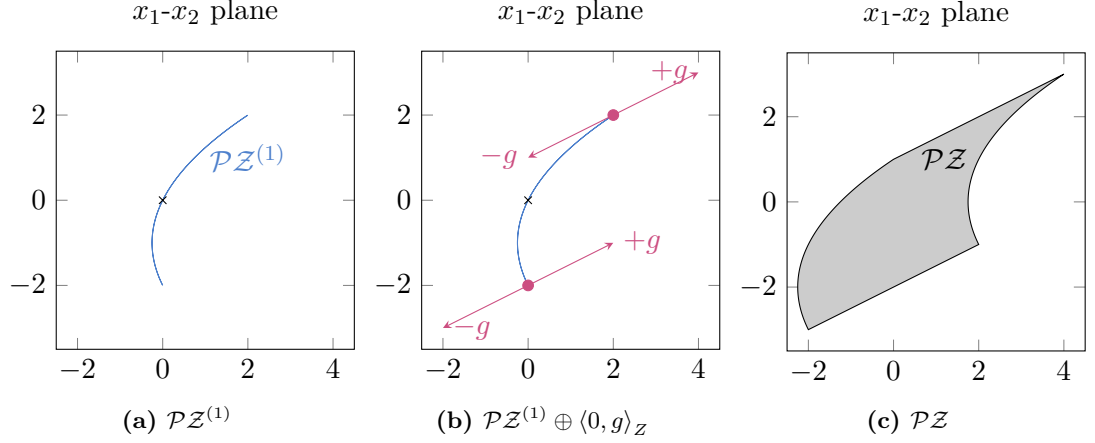
where  $\beta \in [-1, 1]^2$ . Since

$$\mathcal{PZ} = \mathcal{PZ}^{(1)} \oplus \langle 0, g \rangle_Z = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix} \beta_1 + \begin{bmatrix} 1 \\ 0 \end{bmatrix} \beta_1^2 \right\}_{\beta_1} \oplus \left\langle 0, \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right\rangle_Z,$$

the construction of  $\mathcal{PZ}$  is similar to that of a zonotope, with the difference being that the first generator is not linear. Fig. 2.7 shows its construction. ■

We introduce the following relevant operations on polynomial zonotopes (see [61] for details).

## 2 Background



**Figure 2.7:** Construction of the non-convex polynomial zonotope  $\mathcal{PZ} = \mathcal{PZ}^{(1)} \oplus \langle 0, g \rangle_Z$  from Ex. 2.5.

**Minkowski sum [61, Prop. 9, adapted]:** Given two polynomial zonotopes  $\mathcal{A} = \langle c, G, G_{\text{indep}}, E \rangle_{PZ} \subset \mathbb{R}^n$  and  $\mathcal{B} = \langle d, M, M_{\text{indep}}, O \rangle_{PZ} \subset \mathbb{R}^n$ , their Minkowski sum is given by

$$\mathcal{A} \oplus \mathcal{B} = \left\langle c + d, \begin{bmatrix} G & \\ & M \end{bmatrix}, \begin{bmatrix} G_{\text{indep}} & \\ & M_{\text{indep}} \end{bmatrix}, \text{blkdiag}(E, O) \right\rangle_{PZ}.$$

**Exact sum:** Given  $\mathcal{A} = \{a(\xi, \alpha)\}_{\xi, \alpha} \subset \mathbb{R}^n$  and  $\mathcal{B} = \{b(\xi, \beta)\}_{\xi, \beta} \subset \mathbb{R}^n$  as polynomial zonotopes, the exact sum of  $\mathcal{A}$  and  $\mathcal{B}$  is

$$\mathcal{A} \oplus_e \mathcal{B} = \{a(\xi, \alpha) + b(\xi, \beta)\}_{\xi, \alpha, \beta}.$$

Note that  $\mathcal{A} \oplus_e \mathcal{B} \subseteq \mathcal{A} \oplus \mathcal{B}$  since  $\{a(\xi, \alpha) + b(\xi, \beta)\}_{\xi, \alpha, \beta} \subseteq \{a(\xi, \alpha)\}_{\xi, \alpha} \oplus \{b(\xi, \beta)\}_{\xi, \beta}$  where  $\mathcal{A} \oplus_e \mathcal{B} = \mathcal{A} \oplus \mathcal{B}$  if and only if the generating functions  $a$  and  $b$  do not share any dependent factors. We refer the reader to [61, Prop. 10] for details.

**Linear map [61, Prop. 8, adapted]:** The linear map  $A \in \mathbb{R}^{m \times n}$  of an  $n$ -dimensional polynomial zonotope  $\langle c, G, G_{\text{rest}}, E \rangle_{PZ}$  is

$$A \langle c, G, G_{\text{rest}}, E \rangle_{PZ} = \langle Ac, AG, AG_{\text{rest}}, E \rangle_{PZ}.$$

**Quadratic map:** The quadratic map of a polynomial zonotope  $\mathcal{A} = \{a(\alpha)\}_{\alpha} \subset \mathbb{R}^n$ , following the definition in Def. 2.14, is given by

$$\text{sq}(\mathcal{A}, \mathcal{Q}) = \left\{ \begin{bmatrix} a(\alpha)^T Q^{(1)} a(\alpha) \\ \vdots \\ a(\alpha)^T Q^{(m)} a(\alpha) \end{bmatrix} \right\}_{\alpha},$$

where  $\mathcal{Q}$  collects  $m$  matrices  $Q^{(i)} \in \mathbb{R}^n$  for  $1 \leq i \leq m$ . For an efficient implementation of this quadratic map, we refer the reader to [61, Prop. 12, Prop. 13].

**Support function:** Since a polynomial zonotope in general represents a non-convex set, its support function cannot be evaluated exactly as the evaluation of the support function essentially requires range bounding of a polynomial, which is NP-hard in general [33]. We refer the reader to [61, Prop. 7] for an efficient outer approximation.

That said, using an efficient outer approximation of a polynomial zonotope to a zonotope – introduced in Sec. 4.1.4 – one can use the support function value of a zonotope (see Sec. 2.6.5) as an upper bound.

**Subset property:** Let  $\mathcal{X}^{(0)} = \{x^{(0)}(\beta)\}_\beta \subset \mathbb{R}^n$  be a polynomial zonotope and let  $\mathcal{R} = \{r(\beta)\}_\beta$  be the resulting polynomial zonotope of an algorithm ( $A$ ) that takes  $\mathcal{X}^{(0)}$  as an input and returns  $\mathcal{R}$  by applying a finite number of elementary set operations, such as Minkowski sums or linear maps. The subset property of polynomial zonotopes as described in [62] allows us to extract the solution of ( $A$ ) for a subset of  $\mathcal{X}^{(0)}$  without executing ( $A$ ) again, as the following example demonstrates.

**Example 2.6** (Subset property of polynomial zonotopes). Let  $\mathcal{X}^{(0)} = \{x^{(0)}(\beta)\}_\beta = \{\beta\}_\beta$  for  $\beta \in [-1, 1]^2$  and assume that the algorithm ( $A$ ) computes

$$\mathcal{R} = \{r(\beta)\}_\beta = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right\} \oplus \text{sq}(\mathcal{X}^{(0)}, \mathcal{Q}) = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right\} \oplus \left\{ \begin{bmatrix} 2\beta_1\beta_2 \\ \beta_2^2 \end{bmatrix} \right\}_\beta = \left\{ \begin{bmatrix} 1 + 2\beta_1\beta_2 \\ 2 + \beta_2^2 \end{bmatrix} \right\}_\beta,$$

for

$$\mathcal{Q} = \left\{ \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \right\}.$$

Now, we want to compute the output  $\check{\mathcal{R}}$  of algorithm ( $A$ ) for a subset

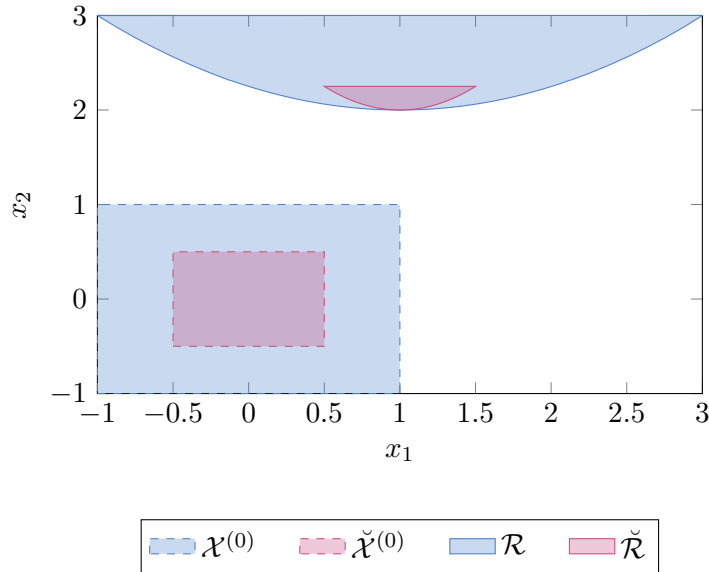
$$\check{\mathcal{X}}^{(0)} = \{x^{(0)}(b(\beta))\}_\beta = \left\{ \begin{bmatrix} b_1(\beta) \\ b_2(\beta) \end{bmatrix} \right\}_\beta = \left\{ \frac{1}{2}\beta \right\}_\beta,$$

where  $b(\beta) = \frac{1}{2}\beta$ . Since we have an analytic expression for the generating function  $r(\beta)$  of  $\mathcal{R}$ , this output is simply given by substituting  $\beta$  with  $b(\beta)$ , i.e., we have

$$\check{\mathcal{R}} = \{r(b(\beta))\}_\beta = \left\{ \begin{bmatrix} 1 + \frac{1}{2}\beta_1\beta_2 \\ 2 + \frac{1}{4}\beta_2^2 \end{bmatrix} \right\}_\beta.$$

Fig. 2.8 visualizes  $\mathcal{X}^{(0)}$  and its output  $\mathcal{R}$  as well as the subset  $\check{\mathcal{X}}^{(0)} \subseteq \mathcal{X}^{(0)}$  and the corresponding output  $\check{\mathcal{R}}$ . ■

In this example,  $\mathcal{R}$  can be computed exactly since it only involves a single quadratic map. However, an algorithm often contains many more elementary operations, which causes the representation of the resulting polynomial zonotope to grow. To reduce this complexity, reduction techniques (see Sec. 2.8) are necessary.



**Figure 2.8:** Illustration of the subset property of polynomial zonotopes using a 2D example.

## 2.7 Reachability Analysis

In order to synthesize a controller for a given task, the underlying model needs to be mathematically described so that predictions about its future evolution can be made, based on which a controller can be synthesized. Since we compute with sets of initial states instead of a single initial state, this chapter introduces reachability analysis: For a given set of initial states and a set of inputs, reachability analysis computes a set of future states that is guaranteed to contain all possible future states of a given system. Here, we only provide a short introduction to the most relevant aspects of reachability analysis; we refer the reader to [4, 59] for more details.

**Definition 2.24** (System). An  $n_x$ -dimensional system with state  $x(t) \in \mathbb{R}^{n_x}$ , controllable input  $u(x(t), t) \in \mathbb{R}^{n_u}$ , and uncontrollable disturbance  $w(t) \in \mathbb{R}^{n_w}$  is governed by the system of ordinary differential equations (ODEs)

$$\dot{x}(t) = f(x(t), u(x(t), t), w(t)), \quad (2.31)$$

with  $t \in \mathbb{R}_{\geq 0}$  being the time and where  $\xi(t, x(0), u(\cdot, \cdot), w(\cdot))$  denotes the solution of (2.31) at time  $t$  starting at the initial state  $x(0)$ . We assume that  $f$  is sufficiently smooth. ■

Without loss of generality, we assume that any input  $u$  dependent on the state has already been substituted into  $f$  and define

$$\dot{x} = f_{\text{cl}}(x, w), \quad (2.32)$$

which absorbs the fixed controller. With slight abuse of notation, the solution  $x(t)$  of the closed-loop dynamics in (2.32) is then denoted by  $\xi(t, x(0), w(\cdot))$ . The reachable set is then defined as follows:

**Definition 2.25** (Reachable Set). The reachable set at time  $t \in \mathbb{R}_{\geq 0}$  for an initial set  $\mathcal{X}^{(0)} \subset \mathbb{R}^{n_x}$  and a set of possible disturbances  $\mathcal{W} \subset \mathbb{R}^{n_w}$  is defined by

$$\mathcal{R}^{(e)}(t) = \left\{ x(t) \in \mathbb{R}^{n_x} \mid \forall \tau \in [0, t] : \dot{x}(\tau) = f_{\text{cl}}(x(\tau), w(\tau)), x(0) \in \mathcal{X}^{(0)}, w(\tau) \in \mathcal{W} \right\}.$$

■

Fig. 1.1 visualizes the main idea of reachability analysis: For a given initial state and disturbance  $w(t)$ , the reachable states are given by the trajectory of the system for these values. The reachable set is then the union of all trajectories starting in  $\mathcal{X}^{(0)}$  with any realization of the disturbance  $w(t)$ . Since in general, computing the exact reachable set  $\mathcal{R}^{(e)}$  is not tractable [68], we use an outer approximation  $\mathcal{R}$  of the reachable set for formal verification, or an approximation  $\tilde{\mathcal{R}}$  (without formal guarantees) for optimization.

We now briefly introduce the reachability algorithms relevant for this thesis from [9, 5, 4]. All algorithms presented subsequently are available in the continuous reachability analyzer<sup>5</sup> [2] (CORA), a MATLAB toolbox for reachability analysis and set computations. We compute the reachable sets for  $t \in [0, t_f]$ , where  $t_f \in \mathbb{R}_+$  is the final time and  $r = \frac{t_f}{m_r}$  is the duration of one reachability step with  $m_r \in \mathbb{N}_+$  being the number of reachability steps.

### 2.7.1 Linear Time-Invariant Systems

In this section, we introduce reachability analysis for an LTI system (see [4, Sec. 3.2] for details)

$$\dot{x}(t) = Ax(t) + w(t). \quad (2.33)$$

The well-known solution of (2.33) given  $x(0) \in \mathcal{X}^{(0)}$  is

$$x(t) = e^{At}x(0) + \int_0^t e^{A(t-\tau)}w(\tau) d\tau. \quad (2.34)$$

Interpreting (2.34) in a set-based manner yields the reachable set

$$\mathcal{R}^{(e)}(t) = \mathcal{H}^{(e)}(t) \oplus \mathcal{P}^{(e)}(t), \quad (2.35)$$

where

$$\mathcal{H}^{(e)}(t) = \left\{ e^{At}x(0) \mid x(0) \in \mathcal{X}^{(0)} \right\}, \quad (2.36)$$

$$\mathcal{P}^{(e)}(t) = \left\{ \int_0^t e^{A(t-\tau)}w(\tau) d\tau \mid \forall \tau \in [0, t] : w(\tau) \in \mathcal{W} \right\}. \quad (2.37)$$

Subsequently, we first focus on a sub interval  $t \in [0, r] \subset [0, t_f]$ . The extension to  $t \in [0, t_f]$  is then made at the end of this section.

<sup>5</sup><https://cora.in.tum.de/>

## 2 Background

**Time-point solution:** While the homogeneous solution directly follows as  $\mathcal{H}^{(e)}(r) = \mathcal{H}(r) = e^{Ar} \mathcal{X}^{(0)}$  from (2.36), only an outer approximation  $\mathcal{P}(r) \supseteq \mathcal{P}^{(e)}(r)$  of the exact particular solution  $\mathcal{P}^{(e)}(r)$  in (2.37) can be computed using [4, Th 3.1], i.e.

$$\mathcal{P}(r) = \bigoplus_{k=0}^{\mu} \left( \frac{A^k r^{k+1}}{(k+1)!} \mathcal{W} \right) \oplus \mathcal{I}_{\text{rem}}(r) r \mathcal{W}, \quad (2.38)$$

with

$$\mathcal{I}_{\text{rem}}(r) = [-1, 1]^{n_x \times n_x} \frac{(\|A\|_{\infty} r)^{\mu+1}}{(\mu+1)!} \frac{1}{1-\epsilon}, \quad \epsilon = \frac{\|A\|_{\infty} r}{\mu+2},$$

and where it is required that  $\epsilon < 1$ . As a result,  $\mu$  and the number of time steps  $m_r$  are typically chosen large enough to ensure  $\epsilon < 1$ . Plugging  $\mathcal{H}(r)$  and  $\mathcal{P}(r)$  from (2.38) into (2.35) yields an outer approximation of the time-point reachable set at  $t = r$  by

$$\mathcal{R}(r) = \mathcal{H}(r) \oplus \mathcal{P}(r) \supseteq \mathcal{R}^{(e)}(r).$$

**Time-interval solution:** In [4, Sec. 3.2.2], it is proven that  $\mathcal{P}([0, r]) = \mathcal{P}(r)$  if  $0 \in \mathcal{W}$ . Hence, we introduce  $\tilde{w} \in \mathcal{W}$  so that  $0 \in \tilde{\mathcal{W}} = \{-\tilde{w}\} \oplus \mathcal{W}$  and define  $\tilde{\mathcal{P}}(r) = \tilde{\mathcal{P}}([0, r])$  analogously to (2.38) using  $\tilde{\mathcal{W}}$ . Thus, it only remains to find an outer approximation of the time-interval homogeneous solution  $\mathcal{H}^{(e)}([0, r])$  and the time-interval solution  $\mathcal{P}_{\tilde{w}}^{(e)}([0, r])$  due to the constant disturbance  $\tilde{w}$ . If we approximate the homogeneous solution  $x^{(h)}(t)$  and the constant disturbance solution  $x_{\tilde{w}}(t) \in \mathcal{P}_{\tilde{w}}^{(e)}(t)$  for  $t \in [0, r]$  as

$$x^{(h)}(t) \approx x(0) + \frac{t}{r} \left( e^{Ar} x(0) - x(0) \right) = \left( 1 - \frac{t}{r} \right) x(0) + \frac{t}{r} e^{Ar} x(0), \quad (2.39)$$

$$x_{\tilde{w}}(t) \approx \frac{t}{r} \int_0^r e^{A\tau} \tilde{w} \, d\tau = \left( 1 - \frac{t}{r} \right) 0 + \frac{t}{r} x_{\tilde{w}}(r), \quad (2.40)$$

summation of (2.37) and (2.39) and its set-based evaluation for  $t \in [0, r]$  yields (see Def. 2.3)

$$\mathcal{H}^{(e)}([0, r]) \oplus \mathcal{P}_{\tilde{w}}^{(e)}([0, r]) \approx \text{convc} \left( \mathcal{X}^{(0)}, e^{Ar} \mathcal{X}^{(0)} \oplus \mathcal{P}_{\tilde{w}}^{(e)}(r) \right).$$

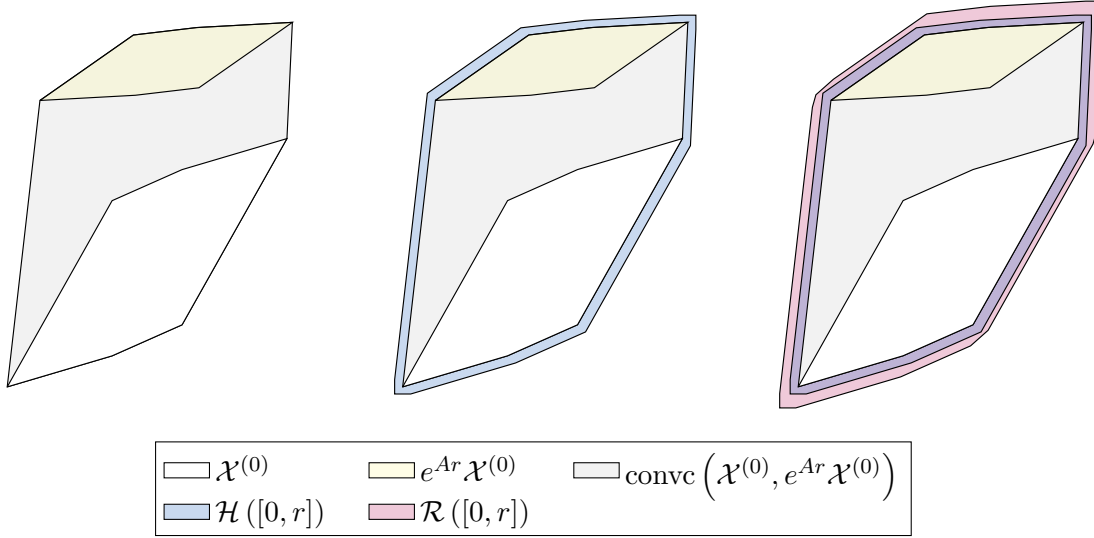
Since (2.39) and (2.40) only approximate the time interval solutions, we additionally need to compute interval matrices  $\mathcal{F}$  and  $\tilde{\mathcal{F}}$  (see [4, Sec. 3.2] for details) to account for the curvature of solutions, so that an outer approximation is given by

$$\mathcal{H}^{(e)}([0, r]) \oplus \mathcal{P}_{\tilde{w}}^{(e)}([0, r]) \subseteq \text{convc} \left( \mathcal{X}^{(0)}, e^{Ar} \mathcal{X}^{(0)} \oplus \mathcal{P}_{\tilde{w}}(r) \right) \oplus \mathcal{F} \mathcal{X}^{(0)} \oplus \tilde{\mathcal{F}} \tilde{w},$$

where  $\mathcal{P}_{\tilde{w}}(r) \supseteq \mathcal{P}_{\tilde{w}}^{(e)}(r) = \int_0^r e^{A\tau} \tilde{w} \, d\tau$  is computed using, e.g., (2.38) with  $\mathcal{W} = \{\tilde{w}\}$ , since  $\mathcal{P}_{\tilde{w}}^{(e)}(r)$  generally cannot be computed exactly. Thus, the time-interval solution is given by

$$\mathcal{R}([0, r]) = \text{convc} \left( \mathcal{X}^{(0)}, e^{Ar} \mathcal{X}^{(0)} \oplus \mathcal{P}_{\tilde{w}}(r) \right) \oplus \mathcal{F} \mathcal{X}^{(0)} \oplus \tilde{\mathcal{F}} \tilde{w} \oplus \tilde{\mathcal{P}}(r) \supseteq \mathcal{R}^{(e)}([0, r]).$$

Fig. 2.9 visualizes the computation of the reachable set, where we assume that  $0 \in \mathcal{W}$ , i.e.,  $\tilde{w} = 0$ , for illustrative purposes:



**Figure 2.9:** One step of reachability analysis for linear systems.

- (i) First, we approximate the set of homogeneous solutions  $\mathcal{H}([0, r])$  for  $t \in [0, r]$  as  $\text{convc}(\mathcal{X}^{(0)}, e^{Ar} \mathcal{X}^{(0)})$ .
- (ii) To arrive at an outer approximation of the homogeneous solution for  $t \in [0, r]$ , we additionally need to account for the curvature of trajectories, which yields  $\mathcal{H}([0, r]) = \text{convc}(\mathcal{X}^{(0)}, e^{Ar} \mathcal{X}^{(0)}) \oplus \mathcal{F}\mathcal{X}^{(0)}$ .
- (iii) As a last step, we add the effect of the disturbance, i.e., the particular time-interval solution, which finally yields the time-interval reachable set  $\mathcal{R}([0, r]) = \text{convc}(\mathcal{X}^{(0)}, e^{Ar} \mathcal{X}^{(0)}) \oplus \mathcal{F}\mathcal{X}^{(0)} \oplus \mathcal{P}(r)$ .

Following [4, Sec. 3.2.2], one can then iteratively compute the outer approximations  $\mathcal{R}([k, k+1]r)$  of the subsequent time-interval reachable sets for  $0 \leq k \leq m_r - 1$ , and we obtain  $\mathcal{R}([0, t_f]) = \bigcup_{k=0}^{m_r-1} \mathcal{R}([k, k+1]r)$ . If we use zonotopes for these iterative computations, increasing the number of reachability steps  $m_r$  increases the number of required Minkowski additions, so that the number of generators grows. To keep computations efficient, this requires reduction techniques (see Sec. 2.8).

Linear systems have gained much attention recently. Most notably, significant improvements to set algorithm parameters, such as the number of reachability steps, automatically [112]. Furthermore, polynomial-time algorithms in the state-dimension for fully automated verification of LTI systems have been proposed in [114] and can even be extended to include safe and unsafe set specifications [115].

### 2.7.2 Nonlinear Systems using Conservative Linearization

Next, we look at a reachability algorithm for nonlinear systems from [9], i.e., we again have  $\dot{x} = f_{\text{cl}}(x, w) = f_{\text{cl}}(z)$  with  $z = [x^T, w^T]^T$ . While the reachability computations are similar to Sec. 2.7.1, we linearize the nonlinear system in each reachability step for  $t \in [k, k+1]r$  with  $0 \leq k \leq m_r - 1$ , which yields

$$\dot{x}(t) = c(z^*) + A(z^*)x(t) + B(z^*)w(t) + l(k), \quad (2.41)$$

with

$$\begin{aligned} A(z^*) &= \left. \frac{\partial f_{\text{cl}}(z)}{\partial x} \right|_{z=z^*}, \\ B(z^*) &= \left. \frac{\partial f_{\text{cl}}(z)}{\partial w} \right|_{z=z^*}, \\ c(z^*) &= f_{\text{cl}}(z^*) - A(z^*)x^* - B(z^*)w^*, \end{aligned}$$

where  $l(k) \in \mathcal{L}(k)$  with  $\mathcal{L}(k)$  being the Lagrange remainder to account for linearization errors,  $z^* = [x^{*T}, w^{*T}]^T$  is the linearization point, and

$$\mathcal{L}(k) = \left\{ l \in \mathbb{R}^{n_x} \left| l_i = \frac{1}{2} (z - z^*)^T \frac{\partial^2 f_{\text{cl}i}(\chi(z, z^*, \alpha))}{\partial z^2} (z - z^*), \right. \right. \\ \left. \left. z \in \mathcal{R}([k, k+1]r) \times \mathcal{W}, \alpha \in [0, 1] \right\}, \quad (2.42)$$

with  $\chi(z, z^*, \alpha) = z^* + \alpha(z - z^*)$  and  $\mathcal{R}([k, k+1]r)$  denoting the outer approximation of the reachable set for  $t \in [k, k+1]r$  that we wish to compute. The Lagrange remainder can be evaluated using interval arithmetic [9, Sec. V]. Further, since the Lagrange remainder is constant in (2.41), the superposition principle can be applied to compute the reachable set  $\mathcal{R}_{\text{lin}}([k, k+1]r)$  of the corresponding LTI system for  $l(k) = 0$  according to Sec. 2.7.1, and the reachable set  $\mathcal{R}_{\text{err}}([k, k+1]r)$  due to the linearization error separately. However, the Lagrange remainder set computation in (2.42) requires the reachable set  $\mathcal{R}([k, k+1]r)$ , whose computation in turn requires  $\mathcal{L}(k)$ .

To break this circular dependency, we assume an initial candidate  $\bar{\mathcal{L}}(k)$  (e.g.  $\bar{\mathcal{L}}(k) = \{0\}$ ) for the Lagrange remainder and use the superposition principle to arrive at a reachable set candidate

$$\bar{\mathcal{R}}([k, k+1]r) = \mathcal{R}_{\text{lin}}([k, k+1]r) \oplus \bar{\mathcal{R}}_{\text{err}}([k, k+1]r),$$

where  $\bar{\mathcal{R}}_{\text{err}}([k, k+1]r)$  is the reachable set due to the linearization errors captured in  $\bar{\mathcal{L}}(k)$  (can be computed using (2.38)). Then, we evaluate (2.42) with  $\bar{\mathcal{R}}([k, k+1]r)$  to obtain  $\mathcal{L}(k)$  and check if  $\mathcal{L}(k) \subseteq \bar{\mathcal{L}}(k)$ . If this check fails, we enlarge  $\bar{\mathcal{L}}(k)$  and try again; otherwise,  $\mathcal{R}([k, k+1]r) = \bar{\mathcal{R}}([k, k+1]r)$  is an outer approximation of the reachable set.



### 2.7.3 Nonlinear Systems using Conservative Polynomialization

In general, reachable sets of nonlinear systems may be non-convex. Therefore, we provide a brief overview on abstracting a nonlinear system  $\dot{x} = f(x, w)$  for reachability analysis using conservative polynomialization and polynomial zonotopes as proposed in [5]. To improve readability, we focus on the first reachability step, i.e., we assume  $t \in [0, r]$  without loss of generality; the remaining steps for  $t \in [r, t_f]$  follow analogously.

We start by computing a Taylor expansion of  $f_{\text{cl}}(z)$  of order  $\pi \in \mathbb{N}_+$  around the linearization point  $z^* = [x^{*T}, w^{*T}]^T$  with  $x^* \in \mathbb{R}^{n_x}$  and  $w^* \in \mathbb{R}^{n_w}$  being linearization points, which yields

$$\dot{x} = A(z^*)x + v(z(t), z^*, w(t)) + l, \quad (2.43)$$

for  $l \in \mathcal{L}$ , and where the Jacobian of  $f$  is given by  $A(z^*)$ , the expression  $v(z(t), z^*, w(t))$  contains all remaining terms of the Taylor expansion up to (and including) order  $\pi$ , and  $\mathcal{L}$  is the Lagrange remainder evaluated over the rest term of the Taylor expansion (as in (2.42)). Using, e.g., a quadratic abstraction, the advantage of using polynomial zonotopes becomes clear since – in contrast to zonotopes – they are closed under quadratic maps (see Tab. 2.1). With  $w^\Delta(t) = w(t) - \tilde{w}$  for  $\tilde{w} \in \mathcal{W}$  chosen such that  $0 \in \mathcal{W}^\Delta = \{-\tilde{w}\} \oplus \mathcal{W}$  and  $z^\Delta(t) = z(t) - z(0)$ , we split

$$v(z(t), z^*, w(t)) = v^{(0)}(z(0), z^*, \tilde{w}) + v^\Delta(z^\Delta(t), z(0), w^\Delta(t)),$$

into  $v^{(0)}(z(0), z^*, \tilde{w})$  only dependent on known values  $z(0)$ ,  $z^*$ ,  $\tilde{w}$  at  $t = 0$ , and  $v^\Delta(z^\Delta(t), z(0), w^\Delta(t))$  for  $t \in [0, r]$ , which depends on  $z(t)$  and thus the reachable set  $\mathcal{R}([0, r])$  which we aim to compute. Therefore, (2.43) can be rewritten as

$$\dot{x}(t) = A(z^*)x(t) + v^{(0)}(z(0), z^*, \tilde{w}) + \tilde{l}, \quad (2.44)$$

for  $\tilde{l} \in \tilde{\mathcal{L}}$ , where

$$\begin{aligned} \tilde{\mathcal{L}} &= \mathcal{V}^\Delta([0, r]) \oplus \mathcal{L}, \\ \mathcal{V}^\Delta([0, r]) &= \left\{ v^\Delta(z^\Delta, z(0), w^\Delta) \mid z^\Delta \in \mathcal{R}^\Delta([0, r]) \times \mathcal{W}^\Delta, z(0) \in \mathcal{X}^{(0)} \times \{0\}, \right. \\ &\quad \left. w^\Delta \in \mathcal{W}^\Delta \right\}, \\ \mathcal{R}^\Delta([0, r]) &= \left\{ \xi(t, x(0), w(\cdot)) - x(0) \mid t \in [0, r], x(0) \in \mathcal{X}^{(0)}, w(t) \in \mathcal{W} \right\}, \end{aligned}$$

where  $\xi(t, x(0), w(\cdot))$  again denotes the state solution. Since (2.44) can be interpreted as an LTI system with uncertainty  $\tilde{l} \in \tilde{\mathcal{L}}$  and additional constant  $v^{(0)}(z(0), z^*, \tilde{w})$ , its reachable set can be computed using the methods from Sec. 2.7.1. In contrast to (2.41), however, (2.44) captures some nonlinearities in  $v^{(0)}(z(0), z^*, \tilde{w})$ .

**Remark:** By choosing the reachability step  $r$  smaller, the set  $\tilde{\mathcal{L}}$  shrinks and thus the accuracy of the approximation is increased [5], albeit at the cost of additional computational effort.

## 2.8 Reduction Techniques

When calculating reachable sets (see Sec. 2.7), reduction techniques are essential to keep computations efficient. In this thesis, we mainly use zonotopes and polynomial zonotopes for set-based computations. In this section, we thus introduce the concept of reduction using a simple zonotope example. We refer the reader to [64], [59, Sec. 2.6] for details on zonotope order reduction techniques and to [61, Sec. 2.E], [59, Sec. 3.1.7] for polynomial zonotope order reduction techniques.

For a given set  $\mathcal{S} \subset \mathbb{R}^n$  (zonotope or polynomial zonotope) with  $m$  generators, we denote by  $\mathcal{S} \downarrow_o \supseteq \mathcal{S}$  the reduced set with a smaller number of generators  $\tilde{m} \leq m$ , where  $o = \frac{\tilde{m}}{n}$  denotes the generator order of the new set. We demonstrate the concept of order reduction in the following example:

**Example 2.7** (Zonotope Order Reduction). Let

$$\mathcal{Z} = \left\langle \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \begin{bmatrix} -6, & -2, & 2, & 0, & 3, & -1 \\ 3, & 4, & 0, & -2, & 0, & -2 \end{bmatrix} \right\rangle_Z,$$

and we want to compute  $\mathcal{Z} \downarrow_2$ , i.e., the reduced zonotope should have  $2 \cdot 2 = 4$  generators. Here, we choose the order reduction method from [40]: For a given zonotope  $\mathcal{Z} = \langle c, G \rangle_Z$  of order  $o = \frac{m}{n} \in \mathbb{N}$  with  $c \in \mathbb{R}^n$  and  $G \in \mathbb{R}^{n \times m}$ , we sort all generators  $G_{(:,i)}$  for  $1 \leq i \leq m$  according to the metric  $l(g) = \|g\|_1 - \|g\|_\infty$ . To compute  $\mathcal{Z} \downarrow_q$  with  $q \leq o$ , one then chooses the  $n(q-1)$  smallest generators according to the metric  $l$ , and adds the resulting interval outer approximation back to the remaining generators.

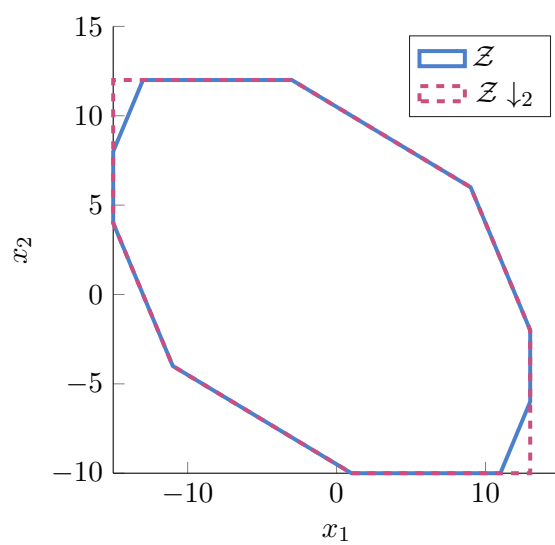
In this example, we select the 4 smallest generators of  $\mathcal{Z}$ , which correspond to the last 4 generators, i.e.,  $\begin{bmatrix} 2, & 0, & 3, & -1 \\ 0, & -2, & 0, & -2 \end{bmatrix}$ . The resulting box enclosure is then given by

$$\left\langle \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \text{diag} \left( \left| \begin{bmatrix} 2 \\ 0 \end{bmatrix} \right| + \left| \begin{bmatrix} 0 \\ -2 \end{bmatrix} \right| + \left| \begin{bmatrix} 3 \\ 0 \end{bmatrix} \right| \left| \begin{bmatrix} -1 \\ -2 \end{bmatrix} \right| \right) \right\rangle_Z = \left\langle \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 6, & 0 \\ 0, & 4 \end{bmatrix} \right\rangle_Z.$$

As a result, we have

$$\mathcal{Z} \downarrow_2 = \left\langle \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \begin{bmatrix} -6, & -2, & 6, & 0 \\ 3, & 4, & 0, & 4 \end{bmatrix} \right\rangle_Z.$$

Fig. 2.10 shows both the original and the reduced zonotope. ■



**Figure 2.10:** Simple reduction approach for a 2D zonotope.



## 3 The Abstracted Synthesis Problem

In this thesis, we synthesize controllers that steer a system from any initial state, chosen from a bounded set, close to a given target state while formally guaranteeing that no constraints are violated. We solve these so-called reach-avoid problems by combining optimization theory with reachability analysis. Here, optimization theory is generally used to formulate an optimization problem in order to synthesize the controller, and reachability analysis is used for formal verification. While subsequent chapters introduce various algorithms to synthesize a controller for these reach-avoid problems, they all share a common structure; this allows the introduction of the abstracted synthesis problem, which captures this common structure. Since practically all available efficient optimization solvers require differentiability, we derive the smooth but equivalent formulation of the abstracted problem. Furthermore, we show that the first-order necessary conditions for this smooth problem are always necessary.

This chapter is structured as follows: First, we describe the general synthesis problem we solve in this thesis in Sec. 3.1. In Sec. 3.2, we then discuss the abstracted problem, from which we derive an equivalent but smooth formulation and show that the first-order optimality conditions of this smooth problem are always necessary. Lastly, we briefly motivate the use of the 1-norm when synthesizing controllers using reachable sets in Sec. 3.3.

### 3.1 Problem Statement

All synthesis problems in this work are so-called reach-avoid problems.

**Definition 3.1** (Reach-Avoid Problem). Let the system  $\dot{x}(t) = f(x(t), u(x(t), t), w(t))$  be given as in Def. 2.24, and let  $x(0) \in \mathcal{X}^{(0)}$ ,  $u(x(t), t) \in \mathcal{U}$ ,  $x(t) \in \mathcal{X}$ ,  $x(t_f) \in \mathcal{X}_f$ , and  $w(t) \in \mathcal{W}$ , where  $\mathcal{X}^{(0)} \subset \mathbb{R}^{n_x}$  is the set of bounded initial states,  $\mathcal{U} \subset \mathbb{R}^{n_u}$  is the bounded set of admissible inputs,  $\mathcal{X} \subset \mathbb{R}^{n_x}$  denotes the set of admissible states,  $\mathcal{X}_f \subset \mathbb{R}^{n_x}$  are the admissible states at final time  $t_f \in \mathbb{R}_+$ , and  $\mathcal{W} \subset \mathbb{R}^{n_w}$  is the set of bounded, uncontrollable disturbances. We define the problem of steering the system  $f$  from a given set of initial states  $\mathcal{X}^{(0)}$  to a target state  $x_f \in \mathcal{X}_f$  within time  $t_f$  – while respecting the input and state constraints – as a reach-avoid problem. ■

### 3 The Abstracted Synthesis Problem

Fig. 1.1 visualizes this concept. In this thesis, we solve a specific instance of such a reach-avoid problem, given by

$$\hat{u}_{\text{ctrl}}(x(t), t) = \arg \min_{u_{\text{ctrl}}(x(t), t)} \max_{x(t)} \left\{ \|x(t_f) - x_f\|_1 + \zeta \int_0^{t_f} \|u_{\text{ctrl}}(x(\tau), \tau)\|_1 d\tau \right\}, \quad (3.1a)$$

$$\text{s.t. } x(t) \in \mathcal{R}(t), \quad (3.1b)$$

$$\forall t \in [0, t_f] : u_{\text{ctrl}}(x(t), t) \in \mathcal{U}, \quad (3.1c)$$

$$\mathcal{R}([0, t_f]) \subseteq \mathcal{X}, \quad (3.1d)$$

$$\mathcal{R}(t_f) \subseteq \mathcal{X}_f, \quad (3.1e)$$

where  $\mathcal{R}(t)$  denotes the (not necessarily exact) reachable set of states  $x(t)$  at time  $t$  for the controller term  $u_{\text{ctrl}}(x(t), t)$  under disturbances  $w(t) \in \mathcal{W}$ , and  $\zeta \in \mathbb{R}_+$  is a weighting factor to penalize large inputs. No assumptions about the statistical nature of  $\mathcal{W}$  are made, only that  $\mathcal{W} = \langle c_{\mathcal{W}}, G_{\mathcal{W}} \rangle_Z$  is given as a zonotope. That said, we only consider zero-centered disturbance zonotopes  $\mathcal{W} = \langle 0, G_{\mathcal{W}} \rangle_Z$  since any non-zero center  $c_{\mathcal{W}}$  can be absorbed by the dynamics  $f$ . We focus on state, final state, and input constraints over the time interval  $t \in [0, t_f]$ ; while an extension to consider different constraints over a finite number of time points or time intervals is conceptually straightforward, we omit further constraints for notational simplicity. Naturally, we not only want to satisfy constraints, but also steer the final state  $x(t_f)$  as close as possible to a target state  $x_f \in \mathbb{R}^{n_x}$  while ideally minimizing the expended control effort. This is reflected in (3.1a) by finding the controller that minimizes the maximum deviation of  $x(t_f)$  from  $x_f$  while penalizing large inputs. For computational reasons, we choose the 1-norm (see Sec. 3.3).

**Remark:** The solution to (3.1) may be used for the computation of motion primitives to plan formally verified trajectories using a maneuver automaton [52]. As an illustrative example, imagine a take-over maneuver of an autonomous vehicle (see Fig. 3.1): To perform the take-over (see Fig. 3.1a), a suitable controller has to be found. Instead of trying find a suitable controller for an uncertain initial state (Fig. 3.1b), the trajectory is split in time into left and right turn primitives for  $t \in [0, t^{(1)}]$  and  $t \in [t^{(1)}, t_f]$  (see Fig. 3.1c); here,  $x_f^{(1)}$  at  $t = t^{(1)}$  is used as an intermediate target state. When only considering static constraints, a controller for a given motion primitive and its corresponding reachable set  $\mathcal{R}(\cdot)$  can be computed offline for various different initial sets and target states. These motion primitives are then collected into a maneuver automaton which links appropriate motion primitives so that they can be concatenated online to form a wide range of trajectories while respecting dynamic constraints (see [98, Sec. 1.2] for a detailed example on a maneuver automaton). Linking motion primitives requires – among other factors – that the final reachable set is contained in the initial set of the next motion primitive (see Fig. 3.1c). While this thesis does not consider maneuver automata for trajectory planning, it motivates both the fact that we minimize

the worst-case deviation from the target state in (3.1a) and allow final state constraints in the form of  $\mathcal{X}_f$ .

## 3.2 Abstraction

We now introduce the abstraction of the synthesis problem from (3.1). To enable efficient numerical optimization, we then derive its smooth reformulation in Sec. 3.2.1 and show in Sec. 3.2.2 that the KKT conditions are always necessary. Since these properties hold for the abstracted synthesis problem, they follow for all subsequent synthesis problems as they can be generalized to this abstraction.

Thus, we define the abstracted synthesis problem as

$$\hat{z}, \hat{s} = \arg \min_{z, s} J(z, s), \quad (3.2a)$$

$$\text{s.t. } g(z) \leq s, \quad (3.2b)$$

$$\underline{z} \leq z \leq \bar{z}, \quad (3.2c)$$

$$s \geq 0, \quad (3.2d)$$

with

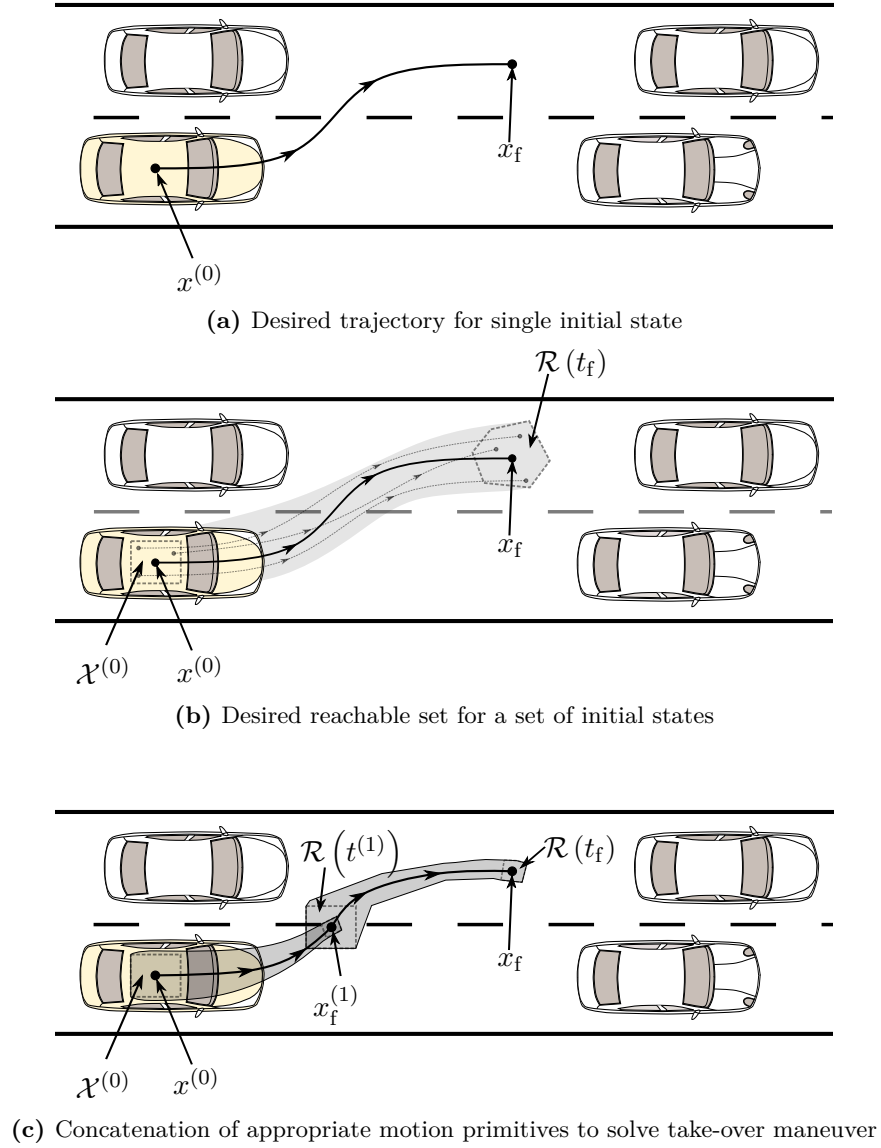
$$J(z, s) = l^T |h(z)| + 1_{n_y}^T \max(\underline{R}\rho(z), \bar{R}\rho(z)) + \sigma 1_{n_s}^T s, \quad (3.3)$$

$$g_m(z) = \max_{1 \leq k \leq n_k} \left\{ \max(F_{(m)}^{(k)}(Ah(z) + M|h(z)|)) + \max(G_{(m)}^{(k)}\rho(z)) \right\} - b_m, \quad (3.4)$$

for  $1 \leq m \leq n_s$ , where  $A \in \mathbb{R}^{n_r \times n_w}$ ,  $M \in \mathbb{R}_{\geq 0}^{n_r \times n_w}$ ,  $\underline{R} \in \mathbb{R}^{n_y \times g}$ ,  $\bar{R} \in \mathbb{R}^{n_y \times g}$ ,  $F_{(m)}^{(k)} \in \mathbb{R}^{q \times n_r}$ , and  $G_{(m)}^{(k)} \in \mathbb{R}^{q \times g}$  are matrices,  $l \in \mathbb{R}_{\geq 0}^{n_w}$  is a vector, and  $\sigma \in \mathbb{R}_+$  is a given scalar; further,  $h(z) \in \mathbb{R}^{n_w}$  and  $\rho(z) \in \mathbb{R}^g$  are sufficiently smooth functions,  $z \in \mathbb{R}^{n_z}$  and  $\hat{z} \in \mathbb{R}^{n_z}$  as well as  $s \in \mathbb{R}^{n_s}$  and  $\hat{s} \in \mathbb{R}^{n_s}$  are the optimization variables and their respective optimizers with  $\underline{z} \in \mathbb{R}^{n_z}$  and  $\bar{z} \in \mathbb{R}^{n_z}$  for  $\underline{z} < \bar{z}$  denoting lower and upper bounds on  $z$ , and  $\sigma \in \mathbb{R}_+$ . If  $\underline{z}_i = \bar{z}_i$  for  $1 \leq i \leq n_z$ ,  $z_i$  can simply be removed from the optimization problem since then  $\hat{z}_i = \underline{z}_i = \bar{z}_i$ .

We briefly motivate the structure of (3.2) here; the intuition behind the exact choice for (3.2) will become clear in subsequent chapters (see, e.g., Sec. 6.3.5): Instead of maximizing over the unknown controller function directly as in (3.1), we fix the structure of the controller by parameterizing it in the controller parameters  $z$ . We use  $h(z)$  to, e.g., collect generators and centers of reachable sets which are required for the optimization. Any additional differentiable functions not originating from reachable sets are collected in  $\rho(z)$ .

The abstracted objective function in (3.3) is obtained by outer approximating the maximization over the state in (3.1a). The constraints (3.1c) to (3.1e) of the original synthesis problem are collected in (3.2b). The maximizations in (3.2b) may be necessary if, e.g., generators from reachable sets for different time intervals are to be considered in a single constraint.



**Figure 3.1:** Illustration of the concatenation of motion primitives online (adapted from [77]).



Since the original synthesis problem contained constraints on both input and states, the optimization problem in (3.1) and therefore also its abstraction in (3.2) may not have a feasible solution. Thus, we also use the slack variable  $s \geq 0$  in (3.2d) to relax the constraints in (3.2b) so that we always find a feasible solution; by choosing  $\sigma$  large enough, any constraint violation is therefore penalized in (3.3) to prioritize no constraint violation over the minimization of the objective function for  $s = 0$ . Lastly, synthesis problems of subsequent chapters restrict the set of controller parameters  $z$  to a bounded set, which is reflected in (3.2c).

### 3.2.1 Smooth Reformulation

Many efficient numerical solvers require the optimization problem at hand to be sufficiently smooth, i.e., the objective function as well as all constraints need to be sufficiently often differentiable. Using insights from Sec. 2.3.6.1 and Sec. 2.3.6.2, the abstracted synthesis problem in (3.2) can be replaced by the smooth optimization problem

$$\min_{z,s,w,y,V,T,p} J_{\text{smooth}}(z, s, w, y), \quad (3.5a)$$

$$\text{s.t.} \quad \begin{bmatrix} +h(z) \\ -h(z) \end{bmatrix} \leq \begin{bmatrix} w \\ w \end{bmatrix}, \quad (3.5b)$$

$$\begin{bmatrix} \underline{R}\rho(z) \\ \bar{R}\rho(z) \end{bmatrix} \leq \begin{bmatrix} y \\ y \end{bmatrix}, \quad (3.5c)$$

$$F_{(m)}^{(k)}(Ah(z) + Mw) \leq 1_q V_{km}, \quad (3.5d)$$

$$G_{(m)}^{(k)}\rho(z) \leq 1_q T_{km}, \quad (3.5e)$$

$$V_{(:,m)} + T_{(:,m)} \leq 1_{n_k} p_m, \quad (3.5f)$$

$$p - b \leq s, \quad (3.5g)$$

$$\underline{z} \leq z \leq \bar{z}, \quad (3.5h)$$

$$s \geq 0, \quad (3.5i)$$

for  $1 \leq m \leq n_s$  and  $1 \leq k \leq n_k$  with

$$J_{\text{smooth}}(z, s, w, y) = l^T w + 1_{n_y}^T y + \sigma 1_{n_s}^T s, \quad (3.6)$$

where  $w \in \mathbb{R}^{n_w}$ ,  $y \in \mathbb{R}^{n_y}$ ,  $V \in \mathbb{R}^{n_k \times n_s}$ ,  $T \in \mathbb{R}^{n_k \times n_s}$  and  $p \in \mathbb{R}^{n_s}$  are auxiliary variables to resolve absolute values and maximization expressions. The equivalence of (3.2) and (3.5) is stated in the following proposition.

**Proposition 3.1** (Equivalence of the Smooth Reformulation). *The abstracted synthesis problem in (3.2) and its reformulation in (3.5) attain the same global minimum. Additionally, the objective function (3.6) at any critical point  $(\hat{z}, \hat{s}, \hat{w}, \hat{y}, \hat{V}, \hat{T}, \hat{p})$  of the reformulated problem can be written as*

$$\tilde{J}(\hat{z}) = J_{\text{smooth}}(\hat{z}, \hat{s}, \hat{w}, \hat{y}) = J(\hat{z}, \max(0, g(\hat{z}))).$$

*Proof.* See App. A.1.1. □

### 3.2.2 Regularity of Smooth Reformulation

Because the abstracted synthesis problem can be rewritten as a smooth optimization problem, efficient gradient-based solvers can be used for its solution. That said, many solvers implicitly solve optimization problems by finding an approximate solution to the corresponding first-order Karush-Kuhn-Tucker (KKT) conditions. However, the KKT conditions are only necessary conditions if the optimization problem at hand is regular at a critical point (see Sec. 2.3.1). As it turns out, this holds for the smooth reformulation.

**Proposition 3.2** (Regularity of Abstracted Synthesis Problem). *The KKT conditions are necessary for the smooth optimization problem in (3.5) at any critical point.*

*Proof.* See App. A.1.2. □

### 3.3 Different Norms for Optimization

As subsequent chapters will show, we frequently make use of norm functions inside the objective function of a controller synthesis problem. Thus, briefly motivate the choice of the 1-norm for optimization over reachable sets in this thesis.

In the following, we will make use of the well-known norm inequalities

$$\|x\|_\infty \leq \|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2 \leq n \|x\|_\infty, \quad (3.7)$$

for  $x \in \mathbb{R}^n$ . As is clear from (3.7), the 2-norm is a better approximation of the infinity norm than the 1-norm, especially if  $\exists i \in \{1, \dots, n\} \forall j \in \{1, \dots, n\} \setminus \{i\} : |x_i| \gg |x_j|$ , since then  $\|x\|_2 \approx |x_i| = \|x\|_\infty$ . However, both the Euclidean and the infinity norm have a downside when optimizing sizes of sets, as the following example demonstrates.

**Example 3.1.** Let the zonotope  $\mathcal{R}(p) = \left\langle 0, \begin{bmatrix} 0.001p + 1000 & -p + 1 \end{bmatrix} \right\rangle_Z$  for  $p \in [-1, 1]$  be given. We are interested in finding the optimal  $\hat{p}$  which minimizes the size of  $\mathcal{R}(p)$ . Since  $\mathcal{R}(p)$  is zero-centered, one way to minimize the size of  $\mathcal{R}(p)$  is to minimize the norm of its generators by solving

$$\hat{p}^* = \arg \min_{p \in [-1, 1]} \left\| \begin{bmatrix} 0.001p + 1000 \\ -p + 1 \end{bmatrix} \right\|,$$

i.e., we shrink the size of  $\mathcal{R}(p)$  by minimizing the magnitude of each generator. Using different norms, we have

$$\left\| \begin{bmatrix} 0.001p + 1000 \\ -p + 1 \end{bmatrix} \right\|_\infty = \max(|0.001p + 1000|, |-p + 1|) = |0.001p + 1000|, \quad (3.8)$$

$$\left\| \begin{bmatrix} 0.001p + 1000 \\ -p + 1 \end{bmatrix} \right\|_2 = \sqrt{(0.001p + 1000)^2 + (-p + 1)^2} \approx |0.001p + 1000|, \quad (3.9)$$

$$\left\| \begin{bmatrix} 0.001p + 1000 \\ -p + 1 \end{bmatrix} \right\|_1 = |0.001p + 1000| + |-p + 1|. \quad (3.10)$$

Since the first generator is much bigger than the second generator, a subsequent minimization over  $p$  will effectively keep the objective value for (3.8) and (3.9) constant since  $\min_{p \in [-1, 1]} |0.001p + 1000| = 999.999 \approx 1000$ . As a result,  $\hat{p} = -1$ , which does not really affect the size of the set in the first dimension, and even increases the size in the second dimension.

The  $\|\cdot\|_1$  norm in (3.10), however, values reduction in all dimensions equally, and thus reduces the size in the second dimension significantly ( $\hat{p} = 1$ ) while only marginally increasing the set size in the first dimension. ■

Since we frequently optimize over (reachable) sets in this thesis, we choose 1-norms for their optimization as motivated by Ex. 3.1.

### 3.4 Summary

In this chapter, we introduced the abstracted synthesis problem as a generalization of many synthesis problems of subsequent chapters. By deriving an equivalent but smooth formulation of this abstracted problem and proving that its Karush-Kuhn-Tucker (KKT) first-order optimality conditions are always necessary, we showed that any problem that can be written in the form of the abstracted problem can be solved by off-the-shelf optimization solvers.

We reviewed reach-avoid problems as the category of synthesis problems we solve in this thesis. In reach-avoid problems, we try to steer a given system for a set of initial states to a target state within a specified time horizon while ensuring both input and state constraints. This is reflected in the general problem formulation, where – for the controller to be synthesized – we penalize the distance of the resulting center of the reachable set to the target and furthermore minimize the size of the final reachable set while also minimizing the accumulated control effort and satisfying input and state constraints.

Since the abstracted problem is not continuously differentiable, we described a smooth formulation to allow for efficient optimization. Because many optimization solvers solve the KKT optimality conditions instead of the optimization problem, we further showed that these optimality conditions are necessary for any critical point of the smooth abstracted problem. As optimization problems for the controller synthesis in subsequent chapters can be written in the form of the abstracted problem, this ensures that these problems can be solved by a variety of off-the-shelf optimization solvers.

All synthesis problems described in this thesis use 1-norms to bound the size of the respective (reachable) sets. As a result, we motivated the choice of the 1-norm by comparing the impact of different norms on the optimization of reachable set sizes.



## 4 Set Conversions and Set Containment

Depending on the application, the choice of set representation can make a significant difference as different set representations are closed under different set operations (see Tab. 2.1) or have varying computational complexities. Thus, set conversions are essential when converting the original intractable synthesis problem into a tractable optimization problem: Examples include the conversion of a non-convex reachable set given by a polynomial zonotope to a simple zonotope to efficiently bound the size of the reachable set, or the conversion from a zonotope to an ellipsoid to avoid an increasing number of generators when many Minkowski sums are required (see Sec. 6.3). Similarly, efficiently checking the containment of these sets within each other is essential, e.g., when using a given reachable set to check state constraints.

Thus, we first introduce all necessary set conversion methods in Sec. 4.1. We then conclude the chapter with a short discussion of relevant set containment methods in Sec. 4.2.

### 4.1 Set Conversions

Since not all sets are closed under all operations, conversions between set representations can become necessary. In Sec. 4.1.1 and Sec. 4.1.2, we first introduce inner and outer approximations of zonotopes by ellipsoids and vice versa, where we refer the interested reader to [34] for their implementation in CORA. While we do not need all zonotope-ellipsoid set conversions introduced here in subsequent chapters, we include them for the sake of completeness. Then, we describe an outer approximation of an H-polytope by a parallelotope in Sec. 4.1.3, and introduce the outer approximation of a polynomial zonotope by a zonotope from [61] in Sec. 4.1.4.

To keep the presentation simple, we assume for the remainder of this section that all sets are non-degenerate, i.e., they have non-zero volumes and refer the reader to [34, Sec. 4.1.3] for a possible approach to handle non-degenerate sets.

#### 4.1.1 Zonotope To Ellipsoid

In this section, we describe approximations of a zonotope by an ellipsoid from our work in [36]: In Sec. 4.1.1.1, we derive an outer approximation of a given zonotope by an ellipsoid, and then discuss an inner approximation of a zonotope by an ellipsoid in Sec. 4.1.1.2. First, we compute a template ellipsoid to describe its shape and then scale this template appropriately using the maximum and minimum zonotope norms, which are defined next.

**Definition 4.1** (Squared Maximum Zonotope Norm). The squared maximum norm from any point on the boundary of a zonotope  $\langle c, G \rangle_Z$  to its center  $c$  is defined by

$$[\langle c, G \rangle_Z]^2 = \max_{\|\beta\|_\infty \leq 1} \beta^T G^T G \beta. \quad (4.1)$$

■

**Definition 4.2** (Squared Minimum Zonotope Norm). The squared minimum norm from any point on the boundary of a zonotope  $\langle c, G \rangle_Z$  to its center  $c$  is defined by the radius of the largest hyper-sphere  $\langle c, [\langle c, G \rangle_Z]^2 I_n \rangle_E$  still contained in  $\langle c, G \rangle_Z$ , where

$$[\langle c, G \rangle_Z]^2 = \max_{r \geq 0} r, \text{ s.t. } \rho_{\langle 0, r I_n \rangle_E}(A) \leq b, \quad (4.2)$$

and where  $\langle A, b \rangle_H$  with  $A \in \mathbb{R}^{o \times n}$  and  $b \in \mathbb{R}^o$  is the half-space representation of  $\langle 0, G \rangle_Z$ .

■

When the halfspace representation for the zonotope is given, the optimization problem in Def. 4.2 can be avoided.

**Proposition 4.1** (Squared Minimum Zonotope Norm). *Let the zonotope  $\mathcal{Z} \subset \mathbb{R}^n$  and its half-space representation  $\langle A, b \rangle_H$  with  $A \in \mathbb{R}^{o \times n}$  and  $b \in \mathbb{R}^o$  be given, where we assume without loss of generality that  $\|A_{(i,:)}\|_2 > 0$  for all  $1 \leq i \leq o$ . Then*

$$[\mathcal{Z}]^2 = \min_{1 \leq i \leq o} \frac{b_i^2}{\|A_{(i,:)}^T\|_2^2}.$$

*Proof.* Since  $\rho_{\langle 0, r I_n \rangle_E}(A_{(i,:)}^T) = \sqrt{A_{(i,:)}^T r I_n A_{(i,:)}^T} = \sqrt{r} \|A_{(i,:)}\|_2$  for  $1 \leq i \leq o$ , it follows from (4.2) that  $\sqrt{r} \|A_{(i,:)}\|_2 \leq b_i \iff r \leq \frac{b_i^2}{\|A_{(i,:)}\|_2^2}$ , which concludes the proof.  $\square$

#### 4.1.1.1 Outer Approximation

If we want to find an ellipsoid  $\langle e_{\min}, E_{\min} \rangle_E$  enclosing a zonotope  $\langle c, G \rangle_Z \subset \mathbb{R}^n$  with minimum volume (also called the Löwner-John ellipsoid), it is clear that  $e_{\min} = c$  must hold since both zonotopes and ellipsoids are centrally symmetric. The optimal shape matrix  $E_{\min}$  is then found by solving the SDP [14, Sec. 8.4.1]

$$E_{\min} = \arg \min_{E \succ 0} -\log \det E, \quad (4.3a)$$

$$\text{s.t. } (v^{(i)} - c)^T E^{-1} (v^{(i)} - c) \leq 1, \quad 1 \leq i \leq m, \quad (4.3b)$$

with  $v^{(i)} \in \mathbb{R}^n$  for  $1 \leq i \leq m$  being the vertices of  $\langle c, G \rangle_Z$ . While even more scalable approaches exist to compute this minimum-volume enclosing ellipsoid (MVEE) [106], they still require us to have the vertex representation of  $\langle c, G \rangle_Z$  available, which is intractable to compute for higher dimensions [32, Theorem 2.5]. Therefore, we subsequently derive a sub-optimal solution that does not require the vertex representation of a given zonotope.

To that end, we first describe an initial guess for the relative proportions of the sub-optimal shape matrix  $\bar{E} \in \mathbb{S}_{++}^{n \times n}$ , and then use the idea of the maximum zonotope norm (see Def. 4.1) to find a scaling factor  $\hat{r} \in \mathbb{R}_+$  to ensure  $\langle c, G \rangle_Z \subseteq \langle c, \hat{r} \bar{E} \rangle_E$ .

**Shape matrix:** In [19], a variant of Goffin’s algorithm for zonotopes is used to compute an approximation to the MVEE. While their resulting approximation is somewhat conservative, they propose the following initial guess [19, Sec. 4.1]: The zonotope  $\langle 0, G \rangle_Z$  with  $G \in \mathbb{R}^{n \times m}$  can be interpreted as the projection of the  $m$ -dimensional hypercube under  $G$ , where the smallest hyper-ball enclosing this hypercube is given by  $\langle 0, mI_m \rangle_E$ . Applying the linear map  $G$  to both hypercube and hyper-ball yields the original zonotope and the enclosing ellipsoid  $G \langle 0, mI_m \rangle_E = \langle 0, mGG^T \rangle_E$ , where the shape matrix we are looking for is then given by  $\bar{E} = GG^T$ .

This choice of  $\bar{E}$  can also be motivated by showing that the covariance matrix over all vertex candidates of  $\langle 0, G \rangle_Z$  is identical to  $\bar{E}$ : Let  $B \in \{-1, 1\}^{m \times 2^{m-1}}$ , where  $B_{i1} = 1$  and we alternate the sign of the  $i$ -th row  $B_{(i,:)}$  every  $2^{i-1}$  entries for  $1 \leq i \leq m$ . For example, this results in

$$B = \begin{bmatrix} 1, & -1, & 1, & -1 \\ 1, & 1, & -1, & -1 \\ 1, & 1, & 1, & 1 \end{bmatrix},$$

for  $m = 3$ . As a result,  $L = G \begin{bmatrix} B, & -B \end{bmatrix} \in \mathbb{R}^{n \times 2^m}$  contains all  $2^m$  possible vertex candidates of  $\langle 0, G \rangle_Z$ . Since the column-wise mean  $\sum_{k=1}^{2^m} L_{(:,k)} = 0_n$  by construction, we can use principle component analysis (PCA) to construct the covariance matrix of all  $2^m$  samples collected in  $L$ , which is given by<sup>1</sup>

$$\begin{aligned} Q_L &= \frac{1}{2^m} LL^T \\ &= \frac{2}{2^m} GBB^T G^T \\ &= GG^T, \end{aligned}$$

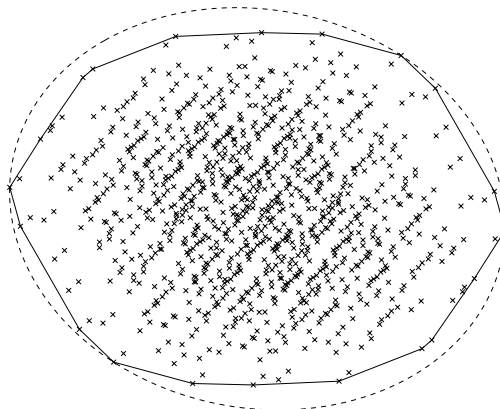
since

$$B_{(i,:)} \begin{bmatrix} B^T \end{bmatrix}_{(:,j)} = \begin{cases} 2^{m-1}, & i = j, \\ 0, & \text{otherwise} \end{cases}, \quad (4.4)$$

for  $1 \leq i \leq m$  and  $1 \leq j \leq m$  due to the construction of  $B$ . The ellipsoid  $\langle 0, Q_L \rangle_E = \langle 0, GG^T \rangle_E$  induced by the covariance matrix  $Q_L \in \mathbb{S}_+^{n \times n}$  then describes the principal components of all samples in  $L$ . Fig. 4.1 shows the distribution of these  $2^m$  samples for a random two-dimensional zonotope  $\langle c, G \rangle_Z$  with ten generators, where the covariance matrix is scaled so that  $\langle c, G \rangle_Z$  is tightly enclosed (see Th. 4.1).

**Radius:** With  $\bar{E}$  given, it only remains to find an appropriate scaling factor  $\hat{r} > 0$  to ensure that  $\langle c, G \rangle_Z \subseteq \langle c, \hat{r}\bar{E} \rangle_E$ . As it turns out, the maximum zonotope norm will be useful in finding such a scaling factor. The maximum squared zonotope norm from Def. 4.1 can then be rewritten as a mixed-integer quadratic program (MIQP).

<sup>1</sup>for details, see Sec. IV.C in <https://arxiv.org/abs/1404.1100>



**Figure 4.1:** All candidate vertex points (black x) of a random two-dimensional zonotope with ten generators (solid black) and the corresponding covariance matrix (dashed black), scaled using Th. 4.1.

**Proposition 4.2.** *Given the zonotope  $\langle \cdot, G \rangle_Z$  with generator matrix  $G \in \mathbb{R}^{n \times m}$ , we have*

$$\lceil \langle \cdot, G \rangle_Z \rceil^2 = - \min_{\beta \in \{-1,1\}^m} \beta^T M \beta + \lambda_{\max} m, \quad (4.5)$$

where  $M = - \left( G^T G - \lambda_{\max} I \right) \succeq 0$  and  $\lambda_{\max} = \max \left( \lambda \left( G^T G \right) \right)$ .

*Proof.* From  $\beta^T \beta = m$ , it follows that  $\beta^T G^T G \beta = \beta^T \left( G^T G - \lambda_{\max} I \right) \beta + \lambda_{\max} m$ . Because  $G^T G \succeq 0$ , using  $M \succeq 0$  together with Def. 4.1 provides the desired result.  $\square$

While (4.5) in Prop. 4.2 is a MIQP and thus has a worst-case complexity which is polynomial in the number of generators [86], it is often quickly solvable for smaller values of  $m$ : Finding  $\lceil \langle c, G \rangle_Z \rceil^2$  to scale the covariance matrix in Fig. 4.1 took around 0.05 s. That said, computing (4.5) exactly quickly becomes intractable for larger  $m$ . Interestingly, [31] showed that solving (4.5) was polynomial in the number of generators of  $G \in \mathbb{R}^{n \times m}$  if  $\text{rank } G = d$  was fixed. They argue that (4.5) can be solved by enumerating all extreme points of  $\langle c, G \rangle_Z$ , and that the number of these extreme points  $m_v$  can be bounded from above by  $m_v \leq \sum_{k=0}^{d-1} 2 \binom{m-1}{k}$  [31, Th. 3.1]. By using the known upper bound for the binomial coefficient  $\binom{m-1}{k} \leq \left( \frac{e(m-1)}{k} \right)^k$  [24, Eq. (C.3)], it follows for  $m \geq 1$  that

$$m_v \leq \sum_{k=0}^{d-1} 2 \binom{m-1}{k} \leq \sum_{k=0}^{d-1} 2 \left( \frac{e(m-1)}{k} \right)^k,$$

and therefore  $m_v$  is of complexity  $O \left( m^{d-1} \right)$  in  $m$  if  $d$  is fixed. However, generally  $d = n$ , and therefore this approach still scales exponentially with the number of dimensions.

Hence, we now propose an upper bound on the maximum zonotope norm, which can similarly be found in [75]. Since the proof was omitted in [75], we provide it here for completeness.



**Lemma 4.1** (Bound on Maximum Zonotope Norm [75, Eq. (3)], adapted). *Let the zonotope  $\langle \cdot, G \rangle_Z$  with  $G \in \mathbb{R}^{n \times m}$  and  $\text{rank}(G) = n$  be given. Then the SDP*

$$\hat{r} = \min_{\lambda \geq 0} 1_m^T \lambda, \quad (4.6a)$$

$$\text{s.t. } \text{diag}(\lambda) - G^T G \succeq 0, \quad (4.6b)$$

with  $\lambda \in \mathbb{R}^m$  provides an upper bound  $\hat{r} \geq \lceil \langle \cdot, G \rangle_Z \rceil^2$ .

*Proof.* The exact squared maximum zonotope norm from Def. 4.1 can be equivalently computed by solving

$$\max_{\beta_i^2 \leq 1, 1 \leq i \leq m} \beta^T G^T G \beta,$$

since  $G^T G \succeq 0$  and thus the maximum is attained for  $\beta_i^2 = 1$ . Its Lagrangian dual function is given by

$$g(\lambda) = \inf_{\beta \in \mathbb{R}^m} \left( \beta^T G^T G \beta + \sum_{i=1}^m \lambda_i (\beta_i^2 - 1) \right) = \inf_{\beta \in \mathbb{R}^m} \beta^T (\text{diag}(\lambda) - G^T G) \beta - 1_m^T \lambda,$$

where  $\lambda \in \mathbb{R}^m$ ; the subsequent maximization  $\sup_{\lambda \geq 0} g(\lambda)$  to form the Lagrangian dual problem then results in (4.6) since

$$g(\lambda) = \begin{cases} -1_m^T \lambda, & \text{diag}(\lambda) - G^T G \succeq 0, \\ -\infty, & \text{otherwise} \end{cases},$$

which concludes the proof.  $\square$

Naturally, there generally is a duality gap between (4.1) and (4.6) since the former is not a convex optimization problem. Rewriting (4.1) with  $B = \beta\beta^T$ , we obtain

$$\max_B \text{trace}(G^T G B), \quad (4.7a)$$

$$\text{s.t. } B \succeq 0, \quad (4.7b)$$

$$\text{trace}(\text{diag}(e(m)^{(i)}) B) = 1, \quad 1 \leq i \leq m, \quad (4.7c)$$

$$\text{rank}(B) = 1. \quad (4.7d)$$

If we drop the rank constraint (4.7d) and formulate its Lagrangian dual problem (see Def. 2.11), we arrive at (4.6). Since there is a strictly feasible  $B$  (e.g.,  $B = I_m$ ), strong duality holds according to SC (see Prop. 2.2). As a result, the gap between (4.1) and (4.6) can be attributed to the “loss” of (4.7d) in (4.6). A result for tightening this duality gap is provided in [75, 48]. We are now ready to state the main result of this section.

**Theorem 4.1** (Enclosing Ellipsoid). *An  $n$ -dimensional enclosing ellipsoid of the non-degenerate zonotope  $\mathcal{Z} = \langle c, G \rangle_Z$  is given by*

$$\hat{\mathcal{E}}(\mathcal{Z}) = \begin{cases} \langle c, \hat{r} \bar{E} \rangle_E, & m > n, \\ \langle c, n \bar{E} \rangle_E, & m = n \end{cases},$$

#### 4 Set Conversions and Set Containment

where  $\hat{r}$  is an upper bound on the maximum squared zonotope norm  $[\bar{E}^{-\frac{1}{2}} \langle c, G \rangle_Z]^2$  using, e.g., Lem. 4.1, and  $\bar{E} = GG^T$ .

*Proof.* Let  $m > n$ . Define  $T = \bar{E}^{-\frac{1}{2}}$  and compute  $T \langle c, \bar{E} \rangle_E = \langle Tc, I_n \rangle_E$  and  $T \langle c, G \rangle_Z$ , where  $T$  exists since  $\mathcal{Z}$  is non-degenerate. Let  $\hat{r} \geq [T \langle c, G \rangle_Z]^2$  denote an upper bound on the squared maximum zonotope norm of  $T \langle c, G \rangle_Z$ . Then  $\max_{\tilde{x} \in T \langle c, G \rangle_Z} \|\tilde{x} - Tc\|_2^2 \leq \hat{r}$  by definition,  $T \langle c, G \rangle_Z \subseteq \langle Tc, \hat{r}I_n \rangle_E$ , and applying  $T^{-1}$  to both  $T \langle c, G \rangle_Z$  and  $\langle Tc, \hat{r}I_n \rangle_E$  yields the desired result.

Now let  $m = n$ . The MVEE enclosing an  $n$ -dimensional hypercube is the hypersphere with radius  $\sqrt{n}$ . Applying  $G$ , the minimum-volume enclosing ellipsoid of  $\langle c, G \rangle_Z$  is  $\langle c, nGG^T \rangle_E$ .  $\square$

##### 4.1.1.2 Inner Approximation

Finding the maximum-volume inscribed ellipsoid (MVIE)  $\langle e_{\max}, E_{\max} \rangle_E$  for a zonotope  $\langle c, G \rangle_Z \subset \mathbb{R}^n$  with center  $c \in \mathbb{R}^n$  and generator matrix  $G \in \mathbb{R}^{n \times m}$ , where  $\langle A, b \rangle_H$  for  $A \in \mathbb{R}^{q \times n}$  and  $b \in \mathbb{R}^q$  denotes its half-space representation, can be formulated as the SDP [14, Sec. 8.4.2]

$$\hat{B} = \arg \min_{B \succ 0} -\log \det B, \quad (4.8a)$$

$$\text{s.t. } \left\| BA_{(i,:)}^T \right\|_2 + A_{(i,:)}c \leq b_i, \quad 1 \leq i \leq q, \quad (4.8b)$$

where  $E_{\max} = \hat{B}\hat{B}$  and again  $e_{\max} = c$  due to the central symmetry of ellipsoids and zonotopes. That said, the half-space representation of a zonotope is generally not available and existing algorithms produce a number of half-spaces exponential in the number of dimensions for a given zonotope [3, Th. 2].

In this section, we thus compute an approximation to  $\langle e_{\max}, E_{\max} \rangle_E$ , analogously to Sec. 4.1.1.1, but now use a lower bound on the minimum zonotope norm from Def. 4.2 to scale the shape matrix appropriately since its exact value requires the half-space representation of the zonotope (see Prop. 4.1). We again use  $\bar{E} = GG^T$ , and propose the following lemma to find a lower bound on the minimum zonotope norm for scaling.

**Lemma 4.2** (Lower Bound on Minimum Zonotope Norm). *Given is a zonotope  $\mathcal{Z} = \langle c, G \rangle_Z \subset \mathbb{R}^n$  with center  $c \in \mathbb{R}^n$  and generator matrix  $G \in \mathbb{R}^{n \times m}$ . Using the test from [19, Sec. 4.4], we can find a lower bound  $\check{r} \in \mathbb{R}_+$  on  $[\mathcal{Z}]^2$  as*

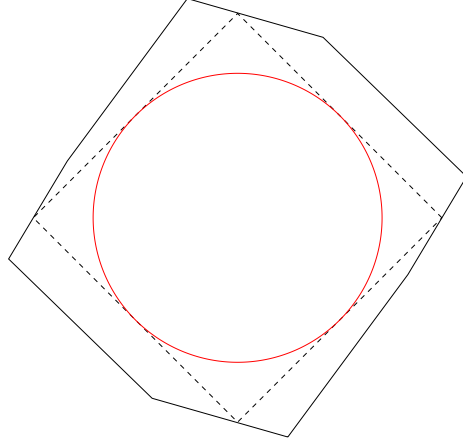
$$\check{r} = \frac{\hat{\nu}^2}{n} \leq [\mathcal{Z}]^2, \quad (4.9)$$

where  $\hat{\nu}$  is optimizer of the LP

$$\hat{\nu} = \arg \max_{\nu, B} \nu, \quad (4.10a)$$

$$\text{s.t. } \nu I_n = GB, \quad (4.10b)$$

$$B \in [-1, 1]^{m \times n}. \quad (4.10c)$$



**Figure 4.2:** Lower bound computation of the minimum zonotope norm as described in Lem. 4.2: For a given  $n$ -dimensional zonotope (black), we compute its boundary points along all  $2n$  unit directions, which define a convex polytope (dashed black). Using the MVEE of this polytope, we apply the Löwner-John property of symmetric sets (see (4.11)) to obtain an ellipsoidal inner approximation (red).

*Proof.* By solving (4.10) for  $\nu$  and due to the symmetry of  $\mathcal{Z}$ , we have

$$\mathcal{V} = \text{convh} \left( \pm \hat{\nu} e_{(n)}^{(1)}, \dots, \pm \hat{\nu} e_{(n)}^{(n)} \right) \subseteq \mathcal{Z}.$$

To inscribe a hyper-ball into  $\mathcal{V}$  and thus find a lower bound on the squared minimum zonotope norm by its radius, we can use the Löwner-John ellipsoidal approximation for symmetric sets [14, Sec. 8.4.1]

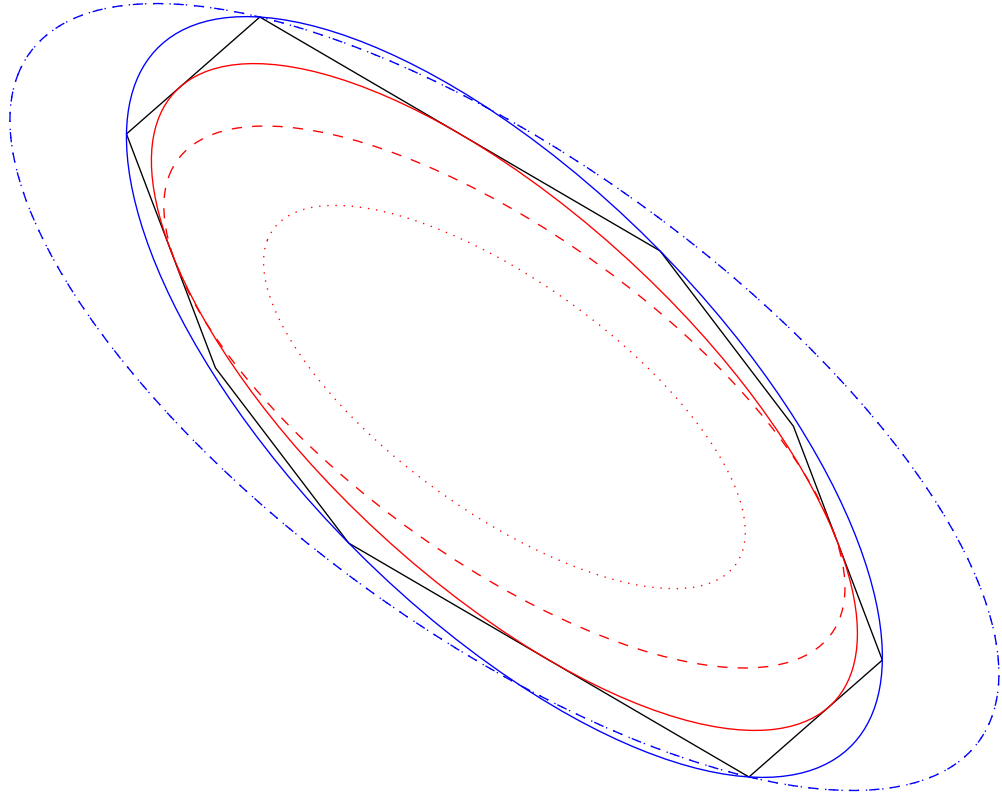
$$\frac{1}{\sqrt{n}} \langle 0, \hat{\nu}^2 I_n \rangle_E \subseteq \mathcal{V} \subseteq \langle 0, \hat{\nu}^2 I_n \rangle_E, \quad (4.11)$$

since  $\mathcal{V} = \{-v \mid v \in \mathcal{V}\}$  is symmetric by construction, where  $\langle 0, \hat{\nu}^2 I_n \rangle_E$  is the MVEE (or Löwner-John ellipsoid) of  $\mathcal{V}$ . Applying the definition of the linear map for ellipsoids (see Sec. 2.6.3) to (4.11), it follows that  $\langle 0, \frac{\hat{\nu}^2}{n} I_n \rangle_E \subseteq \mathcal{V} \subseteq \mathcal{Z}$ , which concludes the proof.  $\square$

The application of Lem. 4.2 to a 2D zonotope is visualized in Fig. 4.2. Since Lem. 4.2 essentially only requires the solution of an LP, for which algorithms with complexity polynomial in the number of optimization variables exist [22], a lower bound on the minimum zonotope norm can be efficiently computed. That said, since we cannot efficiently test whether a given hyper-ball is contained in a zonotope, we need to resort to checking whether its enclosing hyper-cube is contained, which naturally adds conservatism. We can now inscribe an ellipsoid into a zonotope.

**Theorem 4.2** (Inscribed Ellipsoid). *An  $n$ -dimensional ellipsoid inscribed into the non-degenerate zonotope  $\mathcal{Z} = \langle c, G \rangle_{\mathcal{Z}}$  with center  $c \in \mathbb{R}^n$  and generator matrix  $G \in \mathbb{R}^{n \times m}$  is given by*

$$\check{\mathcal{E}}(\mathcal{Z}) = \begin{cases} \langle c, \check{r} \bar{E} \rangle_E, & m > n, \\ \langle c, \bar{E} \rangle_E, & m = n \end{cases},$$



**Figure 4.3:** Different enclosing and inscribed ellipsoids for a given two-dimensional zonotope. Shown are the MVEE and MVIE (solid blue and red), the approximated enclosing and inscribed ellipsoid with exact zonotope norm (dashed blue and red), and the approximated enclosing and inscribed ellipsoids with bounded zonotope norm (dotted blue and red). For this example, the enclosing ellipsoid with exact and bounded zonotope norm coincide.

where  $\check{r} \leq \lfloor \bar{E}^{-\frac{1}{2}} \langle c, G \rangle_Z \rfloor^2$  can, e.g., be computed using Lem. 4.2, and  $\bar{E} = GG^T$ .

*Proof.* Let  $T = \bar{E}^{-\frac{1}{2}}$  and further let  $\check{r} \leq \lfloor T \langle c, G \rangle_Z \rfloor^2$  be a lower bound on the respective squared minimum zonotope norm. The hyper-sphere  $\langle Tc, \check{r}I_n \rangle_E$  is still contained in  $T \langle c, G \rangle_Z$  by definition of the minimum zonotope norm and thus applying the inverse transform  $T^{-1}$  yields the desired result.

Now let  $m = n$ . The MVIE inscribed into an  $n$ -dimensional hypercube is the hyper-sphere with radius 1. Applying  $G$ , the minimum-volume enclosing ellipsoid of  $\langle c, G \rangle_Z$  is  $\langle c, \bar{E} \rangle_E$ .  $\square$

Fig. 4.3 visualizes different enclosing as well as inscribed ellipsoids for a given two-dimensional zonotope.

### 4.1.2 Ellipsoid to Zonotope

In this section, we now describe approximations of a given ellipsoid by a zonotope for a given number of generators from our work in [36]: In Sec. 4.1.2.1, we derive an outer approximation of an ellipsoid with a zonotope, and then describe an inner approximation of an ellipsoid with a zonotope in Sec. 4.1.2.2. Similarly to Sec. 4.1.1, we compute these approximations by finding a template zonotope approximating the general shape of the ellipsoid, and then using the maximum and minimum zonotope norms (Defs. 4.1 and 4.2) to scale them appropriately.

To that end, we first describe how a unit hyper-ball can be approximated with a zonotope and then apply this approximation to derive inner and outer approximations for general ellipsoids. Thus, we momentarily assume that the given ellipsoid is a hyper-ball, as we can always apply the appropriate linear maps to achieve that. Intuitively, a zonotope with  $m \in \mathbb{N}_+$  generators, sampled equally distant from the surface of a hyper-sphere, might approximate a hyper-sphere increasingly closely (in the Hausdorff distance). Since computing equally distant points on a hyper-sphere is impossible for arbitrary  $m$  [117], we first approximate this distribution by sampling  $m$  uniformly distributed points on the surface of a hyper-ball [83]: Let the random variable  $Y$  be normally distributed, and denote with  $y^{(i)} \in \mathbb{R}^n$  for  $1 \leq i \leq m$  any  $m$  independent realizations of  $Y$ . Defining  $x^{(i)} = \frac{y^{(i)}}{\|y^{(i)}\|_2}$ , the vectors  $x^{(i)}$  will be uniformly distributed on the surface of the unit hyper-ball for large enough  $m$ . This approach does not necessarily result in a uniform distribution for small  $m$  in practice. In [71], it is claimed that a polynomial-time deterministic sampling approach is found that solves this distribution problem approximately. Fig. 4.4 visualizes the two variants of approximating a hyper-ball with a zonotope. However, its derivation seems to have a flaw<sup>2</sup>; thus, we use the probabilistic sampling method to obtain an approximately uniform distribution on the surface of the hyper-sphere, since it was proven that a zonotope constructed from these sampled surface points approximates a hyper-sphere increasingly closer for an increasing number of samples [73, 3., Theorem].

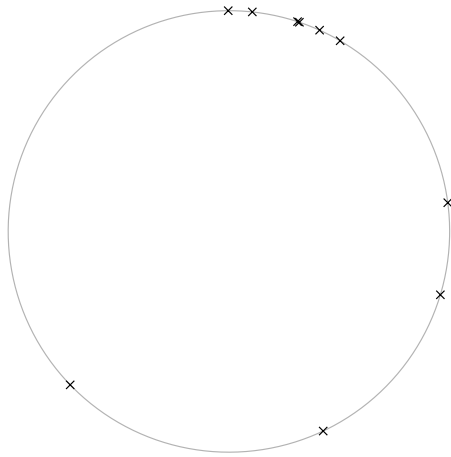
#### 4.1.2.1 Outer Approximation

With the shape of the ellipsoid given by samples of the hyper-sphere as described previously, an outer approximation of an ellipsoid by a zonotope is described next.

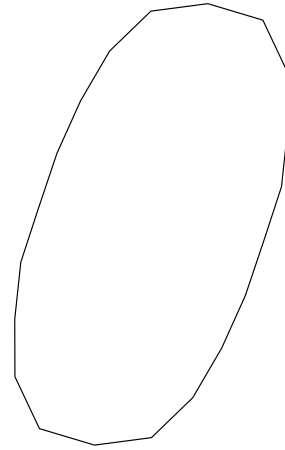
---

<sup>2</sup>I thank Adrian Kulmberg for bringing the following issue to my attention: From [71, Eq. (3.5.8)] and [71, Eq. (3.5.11)], it follows that  $\rho_H(\mathcal{N}_0(\frac{1}{2})) = 1 + \frac{1}{2^{\frac{1}{2}-1}}$ , which is undefined and thus causes a problem in [71, Lem. 3.5.23], which in turn causes a problem in [71, Th. 5.4.1]. Thus, this invalidates [36, Lem. 4] from our original paper since we used [71, Th. 5.4.1]. However, this issue has no consequence on the remainder of our paper since we can simply use the probabilistic sampling method to obtain an approximation of the hyper-sphere by a zonotope as described in [36, Sec. V].

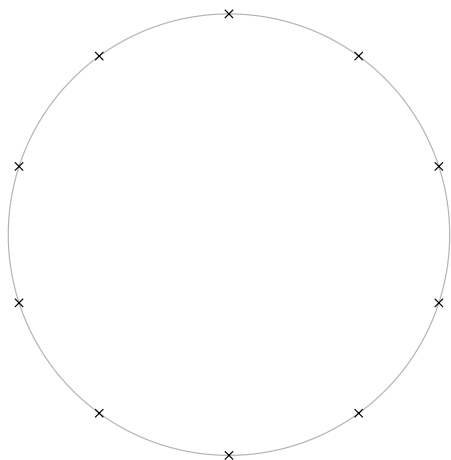
uniformly distributed points



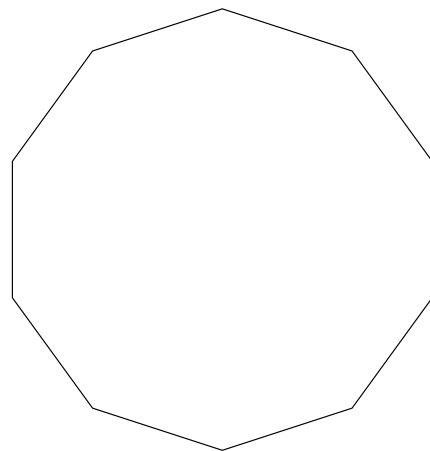
zonotope (uni. points)



approx. equidistant points



zonotope (approx. equidistant points)



**Figure 4.4:** Uniform sampling of sphere directions and approximately equi-distant sampling of sphere directions and their resulting zonotopes.

**Theorem 4.3** (Enclosing Zonotope). *An  $n$ -dimensional zonotope with  $m \geq n$  generators enclosing the ellipsoid  $\langle c, Q \rangle_E$  is given by*

$$\hat{\mathcal{Z}}(\langle c, Q \rangle_E)_m = \begin{cases} \left\langle c, \frac{1}{\sqrt{\check{r}}} \sqrt{Q} S \right\rangle_Z, & m > n, \\ \left\langle c, \sqrt{Q} \right\rangle_Z, & m = n \end{cases},$$

where  $S \in \mathbb{R}^{n \times m}$  collects the points  $S_{(:,i)}$  on the surface of the hyper-ball for  $1 \leq i \leq m$  and  $\check{r}$  is a lower bound of  $[\langle \cdot, S \rangle_Z]^2$  using Lem. 4.2.

*Proof.* Let  $m > n$ . With  $T = Q^{-\frac{1}{2}}$ , we have  $T \langle c, Q \rangle_E = \langle Tc, I_n \rangle_E$ . Further, we know that  $\langle Tc, I_n \rangle_E \subseteq \left\langle Tc, \frac{1}{\sqrt{\check{r}}} S \right\rangle_Z$  due to the definition of  $\check{r}$ . Applying the inverse transform  $T^{-1}$  then gives the desired result.

Now let  $m = n$ . A zonotope enclosing a hyper-ball  $\langle 0, I_n \rangle_E$  is given by  $\langle 0, I_n \rangle_Z$ , and therefore  $\hat{\mathcal{Z}}(\langle c, Q \rangle_E)_n = \sqrt{Q} \langle 0, I_n \rangle_Z \oplus \{c\}$  since  $\langle c, Q \rangle_E = \sqrt{Q} \langle 0, I_n \rangle_E \oplus \{c\}$ .  $\square$

#### 4.1.2.2 Inner Approximation

The inner approximation similarly follows by using the maximum zonotope norm.

**Theorem 4.4** (Inscribed Zonotope). *An  $n$ -dimensional zonotope with  $m \geq n$  generators inscribed into the ellipsoid  $\langle c, Q \rangle_E$  is given by*

$$\check{\mathcal{Z}}(\langle c, Q \rangle_E)_m = \begin{cases} \left\langle c, \frac{1}{\sqrt{\hat{r}}} \sqrt{Q} S \right\rangle_Z, & m > n, \\ \left\langle c, \frac{1}{\sqrt{\hat{r}}} \sqrt{Q} \right\rangle_Z, & m = n \end{cases},$$

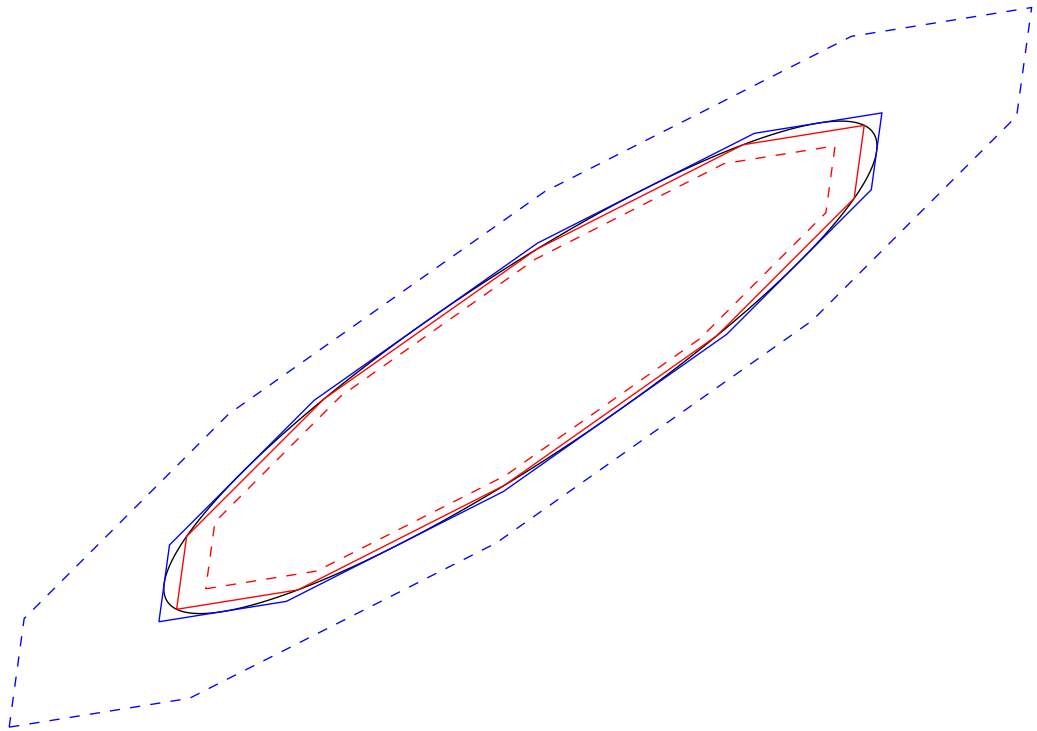
where  $S \in \mathbb{R}^{n \times m}$  collects the points  $S_{(:,i)}$  on the surface of the hyper-ball for  $1 \leq i \leq m$  and  $\hat{r}$  is an upper bound on  $[\langle \cdot, S \rangle_Z]^2$  using Lem. 4.1.

*Proof.* Let  $m > n$ . Setting  $T = Q^{-\frac{1}{2}}$ , we have  $T \langle c, Q \rangle_E = \langle Tc, I_n \rangle_E$ . Further, we know that  $\left\langle Tc, \frac{1}{\sqrt{\hat{r}}} S \right\rangle_Z \subseteq \langle Tc, I_n \rangle_E$  due to the definition of  $\hat{r}$ . Applying the inverse transform  $T^{-1}$  then gives the desired result.

Now let  $m = n$ . A zonotope inscribed into a hyper-sphere  $\langle 0, I_n \rangle_E$  is given by  $\left\langle 0, \frac{1}{\sqrt{\hat{r}}} I_n \right\rangle_Z$ , and therefore  $\check{\mathcal{Z}}(\langle c, Q \rangle_E)_n = \sqrt{Q} \left\langle 0, \frac{1}{\sqrt{\hat{r}}} I_n \right\rangle_Z \oplus \{c\}$  since it holds that  $\langle c, Q \rangle_E = \sqrt{Q} \langle 0, I_n \rangle_E \oplus \{c\}$ .  $\square$

Fig. 4.5 visualizes enclosing and inscribed zonotopes for a given two-dimensional ellipsoid.

For computation times and tightness results on the zonotope-to-ellipsoid and ellipsoid-to-zonotope conversions from Sec. 4.1.1 and Sec. 4.1.2, we refer the reader to [36, Sec. VI] as exact numbers are not relevant for the content of this thesis.



**Figure 4.5:** Different enclosing and inscribed zonotopes for a given two-dimensional ellipsoid. Shown are the approximated enclosing (see Th. 4.3) and inscribed (see Th. 4.4) zonotopes with exact zonotope norm (solid blue and red), and the approximated enclosing and inscribed zonotopes with bounded zonotope norm (dashed blue and red). All zonotopes have  $m = 10$  generators.



### 4.1.3 H-Polytope to Parallelotope

In this thesis, we specify constraints for the reach-avoid problem using H-polytopes. An outer approximation of a bounded H-polytope with a parallelotope, as introduced in our work in [35], is described next.

**Proposition 4.3** (H-Polytope to Parallelotope). *Let  $\mathcal{M} = \langle A, b \rangle_H$  with  $A \in \mathbb{R}^{o \times n}$  and  $b \in \mathbb{R}^o$  be given as a bounded polytope. Then a parallelotope outer approximation of  $\mathcal{M}$  is given by  $\mathcal{Z} = Q^{\frac{1}{2}} \langle \tilde{c}, \text{diag } \tilde{r} \rangle_Z \oplus \{c\}$ , where*

$$\begin{aligned}\tilde{c} &= \frac{1}{2} (\rho_{\tilde{\mathcal{M}}}(I) + \rho_{\tilde{\mathcal{M}}}(-I)), \\ \tilde{r} &= \frac{1}{2} (\rho_{\tilde{\mathcal{M}}}(I) - \rho_{\tilde{\mathcal{M}}}(-I)), \\ \tilde{\mathcal{M}} &= \langle AQ^{\frac{1}{2}}, b - Ac \rangle_H,\end{aligned}$$

and  $\langle c, Q \rangle_E$  is the MVIE of  $\mathcal{M}$  [14, Sec. 8.4.2].

*Proof.* We shift  $\mathcal{M}$  by  $c$  so that  $0 \in \mathcal{M} \oplus \{-c\}$ , and use the shape matrix  $Q$  to transform  $\mathcal{M}$  into roughly a hypercube. Then  $Q^{-\frac{1}{2}} (\langle c, Q \rangle_E \oplus \{-c\}) = \langle 0, I_n \rangle_E$  and it holds that  $\tilde{\mathcal{M}} = Q^{-\frac{1}{2}} (\mathcal{M} \oplus \{-c\}) = \langle AQ^{\frac{1}{2}}, b - Ac \rangle_H$ . A tight parallelotope enclosing  $\tilde{\mathcal{M}}$  is given by  $\tilde{\mathcal{Z}} = \langle \tilde{c}, \text{diag } \tilde{r} \rangle_Z$  with  $\tilde{c}$  and  $\tilde{r}$  given as above. Applying the inverse transform and shift yields  $\mathcal{Z} = Q^{\frac{1}{2}} \tilde{\mathcal{Z}} \oplus \{c\}$  which concludes the proof.  $\square$

Since we use the MVIE as a shape template for the enclosing parallelotope, the approximation quality of Prop. 4.3 depends on a shape similarity of the inscribed and the enclosed ellipsoid, which is not necessarily given. Depending on the required approximation quality, Prop. 4.3 may thus not be approximate; since Prop. 4.3 is used in this thesis for normalization purposes only, this is of no concern here.

### 4.1.4 Polynomial Zonotope to Zonotope

Reachable sets of nonlinear systems are often computed using polynomial zonotopes to represent non-convex sets (see Sec. 2.7.3). However, for operations such as the support function, only outer approximations can be computed, which often requires the solution of a (possibly non-convex) optimization problem. Alternatively, one can enclose a polynomial zonotope with a zonotope and then use its support function.

**Proposition 4.4** (Polynomial Zonotope to Zonotope [61, Prop. 5, adapted]). *A zonotope outer approximation of the polynomial zonotope  $\mathcal{PZ} = \langle c, G, G_{\text{rest}}, E \rangle_{PZ}$  with starting point  $c \in \mathbb{R}^n$ , generator matrix  $G \in \mathbb{R}^{n \times m}$ , rest matrix  $G_{\text{rest}} \in \mathbb{R}^{n \times r}$ , and exponent matrix  $E \in \mathbb{N}^{q \times m}$  is given by*

$$\hat{\mathcal{Z}}(\mathcal{PZ}) = \left\langle c + \frac{1}{2} \sum_{k \in \mathcal{I}_{\text{even}}} G_{(:,k)}, \left[ G_{(:,\mathcal{I}_{\text{odd}})}, \frac{1}{2} G_{(:,\mathcal{I}_{\text{even}})} \right], G_{\text{rest}}, I_m \right\rangle_{PZ},$$

## 4 Set Conversions and Set Containment

where

$$\begin{aligned}\mathcal{I}_{\text{even}} &= \left\{ i \in \{1, \dots, m\} \mid \text{mod} \left( E_{(:,i)}, 2 \right) = 0 \right\}, \\ \mathcal{I}_{\text{odd}} &= \{1, \dots, m\} \setminus \mathcal{I}_{\text{even}},\end{aligned}$$

and  $G_{(:,\mathcal{I}_{\text{odd}})} \in \mathbb{R}^{n \times |\mathcal{I}_{\text{odd}}|}$  and  $G_{(:,\mathcal{I}_{\text{even}})} \in \mathbb{R}^{n \times |\mathcal{I}_{\text{even}}|}$  collect all corresponding odd and even generators.

By inspection, Prop. 4.4 is of at most polynomial complexity in the state dimension (see also [61, Prop. 5]). For a demonstration of the tightness of the resulting zonotope, see, e.g., [61, Fig. 3].

## 4.2 Set Containment

In this thesis, we want to solve controller synthesis problems under input and state constraints. As a result, efficient methods to check the containment of a set within, e.g., a given constraint set, are necessary.

Thus, we first introduce a containment check from previous work for two zonotopes in Sec. 4.2.2. Since constraints in this thesis will be given as H-polytopes, we then describe the containment check of a zonotope in an H-polytope and briefly motivate the superiority of this check over the zonotope-zonotope-containment for our purposes.

### 4.2.1 Zonotope in Zonotope

Since zonotopes are a popular choice for reachable sets in general, it is only natural to first try zonotopes as constraint sets.

**Proposition 4.5** (Zonotope in Zonotope [97, Cor. 4, adapted]). *Given are two zonotopes  $\mathcal{X} = \langle \bar{x}, X \rangle_Z$  with  $\bar{x} \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times m}$ , and  $\mathcal{Y} = \langle \bar{y}, Y \rangle_Z$  with  $\bar{y} \in \mathbb{R}^n$ ,  $Y \in \mathbb{R}^{n \times q}$ . Then  $\mathcal{X} \subseteq \mathcal{Y}$  if there exists  $\Gamma \in \mathbb{R}^{q \times m}$  and  $\beta \in \mathbb{R}^q$  such that*

$$\begin{aligned}X &= Y\Gamma, \\ \bar{y} - \bar{x} &= Y\beta, \\ \left\| \begin{bmatrix} \Gamma^T & \beta^T \end{bmatrix} \right\|_{\infty} &\leq 1.\end{aligned}$$

Thus, Prop. 4.5 gives numerical conditions which can be easily checked using any linear programming solver. Note, however, that Prop. 4.5 is not sufficient and necessary, i.e., failure to find  $\Gamma$  and  $\beta$  that fulfill Prop. 4.5 does not necessarily mean that  $\mathcal{X}$  is not contained in  $\mathcal{Y}$ .

### 4.2.2 Zonotope in H-Polytope

As mentioned in Sec. 4.2.1, the zonotope-zonotope-containment check is only sufficient. When using such a containment check for controller synthesis to check constraints, this

might result in an overly conservative controller. Since we later assume that all constraint sets are represented as H-polytopes anyway, we next present an efficient way to check the containment of a zonotope within an H-polytope which is both necessary and sufficient.

**Proposition 4.6** (Zonotope in Polytope [98, Lem. 1, adapted]). *A zonotope  $\mathcal{Z} = \langle c, G \rangle_{\mathcal{Z}}$  with  $c \in \mathbb{R}^n$  and  $G \in \mathbb{R}^{n \times p}$  is contained within a given H-polytope  $\mathcal{H} = \langle A, b \rangle_{\mathcal{H}}$  with  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$  if and only if*

$$\rho_{\mathcal{Z}}(A) \leq b.$$

*Proof.* We provide an alternative to the proof in [98] here.

We have

$$\begin{aligned} \mathcal{Z} &\subseteq \mathcal{H} \\ \iff \forall l \in \mathbb{R}^n : \rho_{\mathcal{Z}}(l) &\leq \rho_{\mathcal{H}}(l) \\ \iff \rho_{\mathcal{Z}}(A) &\leq \rho_{\mathcal{H}}(A), \end{aligned}$$

where the last equivalence follows from the fact that all halfspaces defining  $\mathcal{H}$  are contained in  $A$ , which concludes the proof since  $\rho_{\mathcal{H}}(A) = b$  (see Def. 2.15).  $\square$

## 4.3 Summary

In this chapter, we both reviewed and introduced novel set conversion methods in Sec. 4.1 which are necessary for later chapters. Additionally, we also described a zonotope-in-zonotope and zonotope-in-polytope containment check from the literature in Sec. 4.2 and motivated why the zonotope-in-polytope check is superior for our purposes.

In Sec. 4.1, we described – for the first time – ellipsoid-zonotope and zonotope-ellipsoid inner and outer approximations from our work in [36]: By combining upper and lower bounds on the maximum and minimum zonotope norms with appropriate initial shapes for both ellipsoids and zonotopes, we were able to find efficient inner and outer approximations for given zonotopes and ellipsoids. Furthermore, we described a novel method for efficiently outer approximating an H-polytope by a parallelotope and review the existing zonotope outer approximation of a given polynomial zonotope from [61].

In order to check constraints for the controller synthesis in subsequent chapters, we then introduced two existing set containment checks in Sec. 4.2: First, we reviewed the zonotope-zonotope containment checks from [97]. Because the zonotope-in-zonotope check is only sufficient, we avoid possible conservatism by specifying all constraints as H-polytopes in this thesis and thus introduced a zonotope-in-H-polytope check from [101] which is both necessary and sufficient.



# 5 Piecewise Constant Controller Synthesis

In order to arrive at an efficient method to synthesize controllers for nonlinear systems, it is often necessary to restrict the form of the controller one wishes to compute. In this chapter, formally correct controller synthesis using a combination of reachability analysis and optimization is discussed, where the controllers to be synthesized are piecewise constant in time.

In Sec. 5.1, we first introduce the general problem we aim to solve, and introduce an approach to the solution of this problem using a linearization of the system dynamics from previous work [99] in Sec. 5.2. In Sec. 5.3, we then extend this approach to polynomial controllers with higher-order reachability abstractions. Finally, we propose a novel iterative solution approach in Sec. 5.4 which can handle larger input sets efficiently and does not require manual tuning of the state constraint sets.

## 5.1 Problem Statement

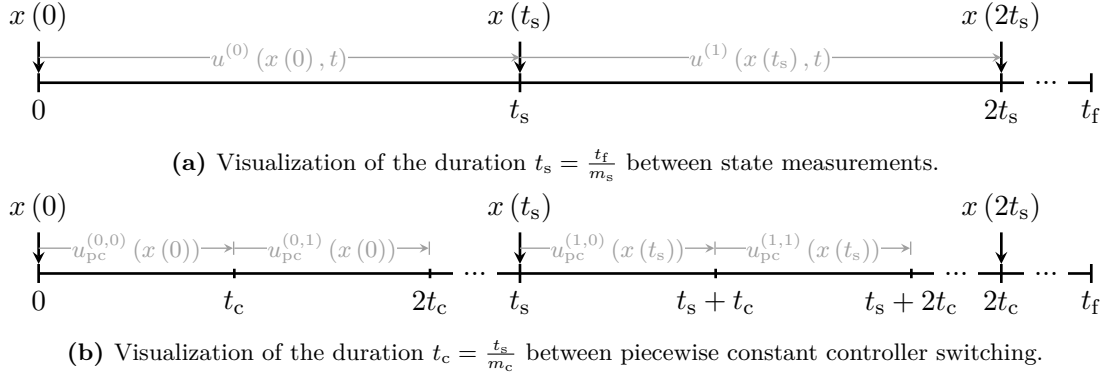
Generally, the synthesis problem in (3.1) is not efficiently solvable for an arbitrary controller  $u_{\text{ctrl}}(x(t), t)$ . Therefore, one has to specify a parameterization for the controller to be synthesized. To that end, we start by subdividing the time horizon  $[0, t_f]$  into  $m_s \in \mathbb{N}_+$  intervals of length  $t_s = \frac{t_f}{m_s}$  (see Fig. 5.1). The duration  $t_s$  denotes the time between state feedback, i.e.,  $u_{\text{ctrl}}(x(t), t) = u^{(i)}(x(it_s), t)$  for  $t \in [i, i+1]t_s$  and  $0 \leq i \leq m_s - 1$ . We further split the time interval  $[i, i+1]t_s$  using  $t_c = \frac{t_s}{m_c}$ , where  $t_c$  is the duration between piecewise constant controllers, i.e.,  $u_{\text{ctrl}}(x(t), t) = u_{\text{pc}}^{(i,j)}(x(it_s))$  for  $t \in it_s + [j, j+1]t_c$  and  $0 \leq j \leq m_c - 1$ .

To enable optimization over these controllers  $u_{\text{pc}}^{(i,j)}(x(it_s))$ , we parameterize them in the controller parameters  $P^{(i,j)} \in \mathbb{R}^{n_u \times a}$  and write  $u(x(it_s), P^{(i,j)})$ , where  $a \in \mathbb{N}_+$  is the number of controller parameters per input dimension (will be discussed in detail in subsequent sections). Let

$$\begin{aligned} P^{(i)} &= \left[ P^{(i,0)T}, P^{(i,1)T}, \dots, P^{(i,m_c-1)T} \right]^T, \\ P &= \left[ P^{(0)T}, P^{(1)T}, \dots, P^{(m_s-1)T} \right]^T. \end{aligned}$$

and denote with  $\mathcal{R}(t, P)$  the closed-loop reachable set at time  $t$  using the controller  $u(x(it_s), P^{(i,j)})$  for  $t \in it_s + [j, j+1]t_c$  with  $0 \leq i \leq m_s - 1$  and  $0 \leq j \leq m_c - 1$ , which is now dependent on  $P$ . Lastly, we assume that the bounded input constraint set  $\mathcal{U} = \langle C_U, d_U \rangle_H$  with  $C_U \in \mathbb{R}^{o_u \times n_u}$  and  $d_U \in \mathbb{R}^{o_u}$ , state constraint set  $\mathcal{X} = \langle C_X, d_X \rangle_H$

## 5 Piecewise Constant Controller Synthesis



**Figure 5.1:** Visualization of the time discretization for discrete-time feedback and piecewise constant control inputs for each sampling interval.

with  $C_{\mathcal{X}} \in \mathbb{R}^{o_{\mathcal{X}} \times n_x}$  and  $d_{\mathcal{X}} \in \mathbb{R}^{o_{\mathcal{X}}}$ , and final state constraint set  $\mathcal{X}_f = \langle C_{\mathcal{X}_f}, d_{\mathcal{X}_f} \rangle_H$  with  $C_{\mathcal{X}_f} \in \mathbb{R}^{o_{\mathcal{X}_f} \times n_x}$  and  $d_{\mathcal{X}_f} \in \mathbb{R}^{o_{\mathcal{X}_f}}$  are given as H-polytopes. We can then rewrite (3.1) as

$$\hat{P} = \arg \min_P \max_{x(\cdot, P) \in \mathcal{R}(\cdot, P)} \left\{ \|x(t_f, P) - x_f\|_1 + \zeta \sum_{i=0}^{m_s-1} \sum_{j=0}^{m_c-1} \left\| u(x(it_s, P), P^{(i,j)}) \right\|_1 \right\}, \quad (5.1a)$$

$$\text{s.t. } \forall i \in \{0, \dots, m_s - 1\} \forall j \in \{0, \dots, m_c - 1\} : u(x(it_s, P), P^{(i,j)}) \subseteq \mathcal{U}, \quad (5.1b)$$

$$\mathcal{R}([0, t_f], P) \subseteq \mathcal{X}, \quad (5.1c)$$

$$\mathcal{R}(t_f, P) \subseteq \mathcal{X}_f. \quad (5.1d)$$

We remark that the inclusion of state constraints means that a feasible solution might not exist.

## 5.2 Generator-Space Control

In this section, we introduce the generator-space control (GSC) approach from [98, 99] to synthesize  $m_s m_c$  piecewise constant controllers while ensuring state and input constraints. Fig. 5.2 visualizes the concept of GSC, while Alg. 1 describes the general procedure which is briefly summarized subsequently:

- (i) First, a reference trajectory – steering the system from the center of the initial set close to the target state  $x_f$  while respecting constraints – is computed, which yields the intermediate target states  $x_f^{(1:m_s)}$  (l. 2).
- (ii) We then iteratively synthesize  $m_c$  controllers for all  $m_s$  steps. We use a moving horizon  $h \in \mathbb{N}_+$  with  $1 \leq h \leq m_s - i$  starting from  $t = it_s$ : This allows us to account for later changes in the nonlinear dynamics  $f$  which might not be visible for  $t \in [0, t_s]$ . The algorithm then iterates as follows:
  - (a) Using the outer approximation of the closed-loop reachable set which is available from the last iteration, we compute the parameterized reachable set, which is an approximation to the closed-loop reachable set for  $t \in it_s + [0, ht_s]$ , parameterized in the controller parameters we wish to determine (l. 5–7).
  - (b) This parameterized reachable set, along with the intermediate target states collected in  $x_f^{(i+(1:h))}$ , is then used to compute the controller parameters subject to input and state constraints by solving a linear program (LP) (l. 8).
  - (c) While we compute the parameters for all controllers until  $t = (i + h)t_s$ , we only use the first  $m_c$  piecewise controllers and compute the closed-loop reachable set using the optimized controllers until  $t = (i + 1)t_s$  (l. 9).

As explained in Sec. 5.2.4 in detail, the efficient inclusion of state constraints in GSC requires the user to provide an adapted state constraint set  $\tilde{\mathcal{X}} \subseteq \mathbb{R}^{n_x}$  and  $\tilde{\mathcal{X}}_f \subseteq \mathbb{R}^{n_x}$ .

---

### Algorithm 1 Generator-space control

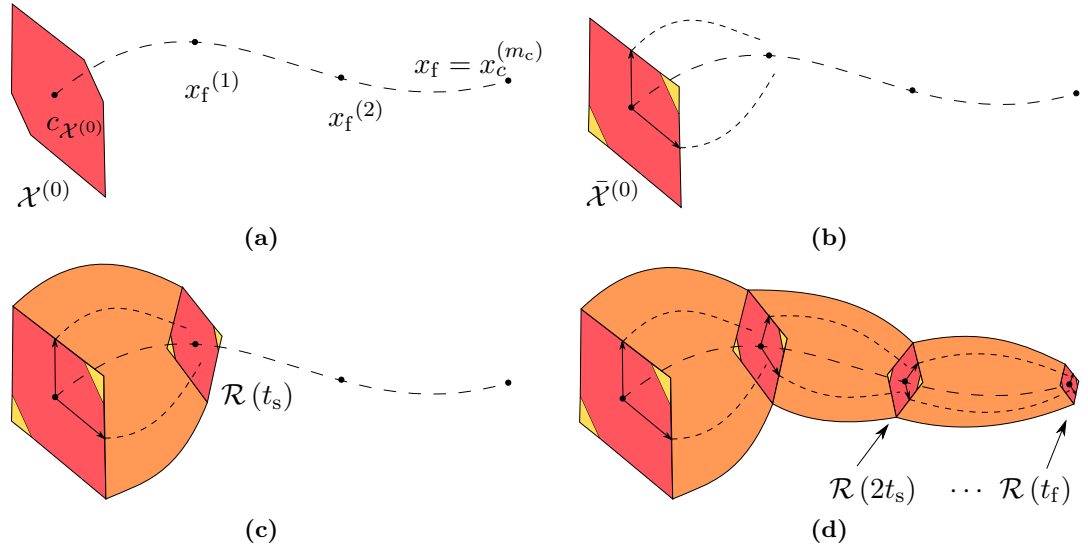
---

```

1: function GENERATORSPACECONTROL( $\mathcal{X}^{(0)}, \mathcal{U}, \mathcal{X}, \tilde{\mathcal{X}}, \mathcal{X}_f, \tilde{\mathcal{X}}_f, \mathcal{W}, m_s, m_c, x_f, h, t_f$ )
2:    $[x_f, u_{\text{ref}}] = \text{REFERENCETRAJECTORY}(\mathcal{X}^{(0)}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, m_s, m_c, x_f, t_f)$   $\triangleright$  Sec. 5.2.1
3:    $\bar{\mathcal{R}} = \mathcal{X}^{(0)}$ 
4:   for  $i = 0; i < m_s; i++$  do
5:      $\bar{\mathcal{R}} = \bar{\mathcal{R}} \downarrow_o$   $\triangleright o \in \mathbb{N}_+$  (reduction order)
6:      $h = \min(h, m_s - i)$ 
7:      $\tilde{\mathcal{R}}^{(i)} = \text{PARAMREACH}(\bar{\mathcal{R}}, h)$   $\triangleright$  Sec. 5.2.3
8:      $\hat{P}^{(i)} = \text{COMPUTECTRL}(\tilde{\mathcal{R}}^{(i)}, \mathcal{U}, \tilde{\mathcal{X}}, \tilde{\mathcal{X}}_f, x_f^{(i+(1:h))}, u_{\text{ref}}^{(im_c+(0:hm_c-1))})$   $\triangleright$  Sec. 5.2.4
9:      $\mathcal{R}([i, i+1]t_s) = \text{REACH}(\bar{\mathcal{R}}, \hat{P}^{(i)}, t_s)$   $\triangleright$  Sec. 5.2.5
10:     $\bar{\mathcal{R}} = \mathcal{R}((i+1)t_s)$ 
11:  end for
12:  return  $\hat{P}, \mathcal{R}([0, t_f])$ 
13: end function

```

---



**Figure 5.2:** Overview of the GSC approach [101, adapted]: First, a reference trajectory is computed (Fig. 5.2a), which is then used to subdivide the overall synthesis problem into  $m_s$  parts and define the intermediate target states and their corresponding control inputs. Using this first intermediate target state  $x_f^{(1)}$ , the synthesis problem (5.1) is solved for  $x_f^{(1)}$  (Fig. 5.2b). The closed-loop reachable set  $\mathcal{R}(t_s)$  for the newly computed controller is then computed (Fig. 5.2c). Applying the same solution process iteratively – where  $\mathcal{R}(t_s)$  is considered the new initial set – we can compute all closed-loop reachable sets, i.e., we obtain  $\mathcal{R}([0, t_f])$ .



We start by introducing the reference trajectory computation in Sec. 5.2.1 and describe the proposed controller template in Sec. 5.2.2. For this template, we then compute its parameterized reachable set in Sec. 5.2.3 and introduce the synthesis of the controller in Sec. 5.2.4, where its closed-loop reachable set computation is then discussed in Sec. 5.2.5. Finally, we briefly discuss both the offline and online complexity of the GSC algorithm in Sec. 5.2.6 and shortly discuss its properties in Sec. 5.2.7.

**Remark:** We will not evaluate the performance of the GSC approach in this section but defer its numerical experiments to the next section, where it will be compared to a novel extension of GSC. We refer the reader to [99, Sec. 4] for a numerical comparison of GSC against other synthesis algorithms.

### 5.2.1 Reference Trajectory

In this section, we want to compute  $m_s m_c$  piecewise constant control inputs such that the center  $c_{\mathcal{X}(0)}$  of the initial set  $\mathcal{X}^{(0)} = \langle c_{\mathcal{X}(0)}, G_{\mathcal{X}(0)} \rangle_Z$  is steered closed to the target state  $x_f$  in time  $t_f$ . These  $m_s m_c$  constant control inputs can be computed by solving

$$u_{\text{ref}}^{(0:m_s m_c - 1)} = \arg \min_{u^{(0:m_s m_c - 1)}} \left( (x(t_f) - x_f)^T Q_{\text{ref}} (x(t_f) - x_f) + \sum_{l=0}^{m_s m_c - 1} u^{(l)T} R_{\text{ref}} u^{(l)} \right), \quad (5.2a)$$

$$\text{s.t. } \forall t \in [k, k+1] t_c : \dot{x}(t) = f(x(t), u^{(k)}, 0) \quad (5.2b)$$

$$x(0) = c_{\mathcal{X}(0)}, \quad (5.2c)$$

$$\forall t \in [0, t_f] : x(t) \in \mathcal{X}, \quad (5.2d)$$

$$u^{(k)} \in \mathcal{U}, \quad (5.2e)$$

$$0 \leq k \leq m_s m_c - 1, \quad (5.2f)$$

where  $Q_{\text{ref}} \in \mathbb{S}_+^{n_x \times n_x}$  and  $R_{\text{ref}} \in \mathbb{S}_+^{n_u \times n_u}$  are user-specified weighting matrices. There exists software, such as the automatic control and dynamic optimization<sup>1</sup> [55] (ACADO) toolbox, which can efficiently solve (5.2). The corresponding state solutions are then given by  $x_{\text{ref}}^{(k+1)} = \xi((k+1)t_c, c_{\mathcal{X}(0)}, u_{\text{ref}}^{(0:k)}, 0)$  for  $0 \leq k \leq m_s m_c - 1$ , where the intermediate target states are given by  $x_f^{(1+i)} = \xi((1+i)t_s, c_{\mathcal{X}(0)}, u_{\text{ref}}^{(0:(1+i)m_c - 1)}, 0)$  for  $0 \leq i \leq m_s - 1$ .

### 5.2.2 Controller Template

Before we can compute a parameterization of the reachable set in the controller parameters, we first need to define the parameterization of the controller. Since we assume the same parameterization for all  $m_s m_c$  piecewise constant controllers, we restrict ourselves to the first piecewise constant controller for  $t \in [0, t_c]$ , i.e., we find the template

<sup>1</sup><https://acado.github.io/>

for  $u(x(0), P^{(0,0)})$ . With slight abuse of notation, we assume for this section that  $P = P^{(0,0)}$  for easier readability.

In Sec. 5.2.1, we already computed (sub-)optimal control inputs for  $\mathcal{X}^{(0)} = \{c_{\mathcal{X}^{(0)}}\}$ , i.e., when the initial set contains only its center. That said, the initial set  $\mathcal{X}^{(0)}$  generally contains infinitely many initial states, and thus computing a control input trajectory  $u_{\text{ref}}^{(0:m_c m_s - 1)}$  for every state in  $\mathcal{X}^{(0)}$  using Sec. 5.2.1 is not feasible. However, using a single input trajectory for every  $x(0) \in \mathcal{X}^{(0)}$  is clearly not ideal, as the optimal input trajectory to reach  $x_f$  generally depends on the chosen initial state.

For the efficient computation of a controller for every initial state, the authors in [99, Sec. 3.2] introduce the controller parameterization

$$\bar{u}(\beta, P) = P_{(:,1)} + P_{(:,2)}\beta = P \begin{bmatrix} 1 \\ \beta \end{bmatrix}, \quad (5.3)$$

where  $P \in \mathbb{R}^{n_u \times (1+l)}$  and  $\beta \in [-1, 1]^l$  are the dependent factors of the generating function  $x^{(0)}(\beta) = c_{\mathcal{X}^{(0)}} + G_{\mathcal{X}^{(0)}}\beta$  of the initial set  $\mathcal{X}^{(0)} = \{x^{(0)}(\beta)\}_{\beta}$ : By making (5.3) dependent on the initial state  $x^{(0)}(\beta)$  through  $\beta$ , we obtain different control inputs for different initial states. We remark here that the optimal control law is assumed to be linear in  $\beta$  due to (5.3); this assumption will be removed in Sec. 5.3. The choice to parameterize the controller template in  $\beta$  instead of the initial state  $x(0)$  directly has numerical advantages since  $\beta \in [-1, 1]^l$ . Additionally, we parameterize the controller in  $P$ , which allows us to later synthesize an optimal controller by finding the optimal value for  $P$ .

For the online application of the controller, however, only the initial state  $x(0)$  is given and thus we require the input parameterized in  $x(0)$  and not  $\beta$ . If  $\mathcal{X}^{(0)}$  is given as a non-degenerate parallelotope, the controller template parameterized in  $x(0)$  then follows by substitution of  $\beta = G_{\mathcal{X}^{(0)}}^{-1}(x(0) - c_{\mathcal{X}^{(0)}})$  into (5.3), yielding (compare to [99, Eq. (15)])

$$u(x(0), P) = \bar{u}(G_{\mathcal{X}^{(0)}}^{-1}(x(0) - c_{\mathcal{X}^{(0)}}), P). \quad (5.4)$$

If  $\mathcal{X}^{(0)}$  is a general zonotope, one can also replace  $\mathcal{X}^{(0)}$  with its parallelotope outer approximation  $\bar{\mathcal{X}}^{(0)} = \mathcal{X}^{(0)} \downarrow_1$  so that the controller template can still be directly parameterized in  $\bar{x}(0) \in \bar{\mathcal{X}}^{(0)}$  ( $\beta$  then describes the dependent factor of the parallelotope outer approximation); however, this comes at the cost of possibly losing some input capacity: We need to ensure that  $\forall x(0) \in \mathcal{X}^{(0)} : u(x(0), P) \subseteq \mathcal{U}$ . If we now parameterize the template in  $\bar{x}(0)$ , this is still enforced, even though we introduce artificial initial states, i.e., there exist  $\bar{x}(0) \notin \mathcal{X}^{(0)}$ . As an alternative to this parallelotope enclosure, the factors  $\beta$  for a corresponding  $x(0)$  can also be efficiently computed online by solving  $x(0) = c_{\mathcal{X}^{(0)}} + G_{\mathcal{X}^{(0)}}\beta$  with  $\|\beta\|_{\infty} \leq 1$  using linear programming, which does not result in a potential loss of input capacity.

**Reference Input for the Center State:** In Sec. 5.2.1, we compute a reference input trajectory for the center  $c_{\mathcal{X}^{(0)}}$  of the initial set. If we wish to use this reference input

directly in the controller template, we have  $u_{\text{ref}}^{(0)} = P \begin{bmatrix} 1 \\ 0 \end{bmatrix} \implies P_{(:,1)} = u_{\text{ref}}^{(0)}$  for  $\beta = 0$  since  $c_{\mathcal{X}^{(0)}} = c_{\mathcal{X}^{(0)}} + G_{\mathcal{X}^{(0)}} 0$ . Thus, the controller template simply returns the input trajectory for the initial state and thus steers the reachable set along the reference trajectory.

### 5.2.3 Parameterized Reachable Set

The solution of (5.1) requires (an approximation of) the closed-loop reachable using the controller template from (5.3), parameterized in the controller parameters. Since the  $m_s$  optimization problems can be solved sequentially (see Alg. 1), we assume without loss of generality that we start from the initial set  $\mathcal{X}^{(0)}$ , generated by  $x^{(0)}(\beta) = c_{\mathcal{X}^{(0)}} + G_{\mathcal{X}^{(0)}}\beta$  for center  $c_{\mathcal{X}^{(0)}} \in \mathbb{R}^{n_x}$  and generator matrix  $G_{\mathcal{X}^{(0)}} \in \mathbb{R}^{n_x \times l}$  with  $\beta \in [-1, 1]^l$ . As discussed in [98], we do not only need this parameterized set for  $t \in [0, t_s]$ , but rather  $t \in [0, ht_s]$  since we optimize over an extended optimization horizon to account for later changes in the nonlinear dynamics  $f$  which might not be visible for  $t \in [0, t_s]$ . That said, we first find the parameterized reachable set for  $t \in [0, t_s]$ , and then extend this to the moving horizon of length  $h$ .

To that end, we linearize and time-discretize the undisturbed nonlinear dynamics  $\dot{x} = f(x, u, 0)$  (see Def. 2.24), yielding (also compare to [99, Sec. 3.1])

$$x^{(j+1)} = c_d^{(j)} + A_d^{(j)} x^{(j)} + B_d^{(j)} u^{(j)}, \quad (5.5)$$

with

$$\begin{aligned} A_d^{(j)} &= e^{A^{(j)}t_c}, \\ B_d^{(j)} &= \int_0^{t_c} e^{A^{(j)}\tau} d\tau B^{(j)}, \\ c_d^{(j)} &= \int_0^{t_c} e^{A^{(j)}\tau} d\tau c^{(j)}, \\ A^{(j)} &= \left. \frac{\partial f(x, u, 0)}{\partial x} \right|_{\substack{x=\bar{x}^{(j)} \\ u=\bar{u}^{(j)}}}, \\ B^{(j)} &= \left. \frac{\partial f(x, u, 0)}{\partial u} \right|_{\substack{x=\bar{x}^{(j)} \\ u=\bar{u}^{(j)}}}, \\ c^{(j)} &= f(\bar{x}^{(j)}, \bar{u}^{(j)}, 0) - A^{(j)}\bar{x}^{(j)} - B^{(j)}\bar{u}^{(j)}, \end{aligned}$$

where  $x^{(0)} \in \mathcal{X}^{(0)}$ , and  $\bar{x}^{(j)} = \frac{1}{2} (x_{\text{ref}}^{(j)} + x_{\text{ref}}^{(j+1)})$  and  $\bar{u}^{(j)} = u_{\text{ref}}^{(j)}$  are the linearization points for  $0 \leq j \leq m_c - 1$  (see Sec. 5.2.1). Evaluation of (5.5) using the corresponding generating functions yields the recursive relation

$$\tilde{r}((j+1)t_c, P, \beta) = c_d^{(j)} + A_d^{(j)} \tilde{r}(jt_c, P, \beta) + B_d^{(j)} \bar{u}(\beta, P^{(0,j)}), \quad (5.6)$$

## 5 Piecewise Constant Controller Synthesis

where  $\tilde{\mathcal{R}}(t, P) = \{\tilde{r}(t, P, \beta)\}_\beta$  and  $\tilde{r}(0, P, \beta) = x^{(0)}(\beta) = c_{\mathcal{X}^{(0)}} + G_{\mathcal{X}^{(0)}}\beta$ . For later convenience, we define

$$\begin{aligned}\tilde{\mathcal{R}}(jt_c, P) &= \left\langle \tilde{c}(jt_c, P), \tilde{G}(jt_c, P) \right\rangle_Z, \\ \left\{ \tilde{u}(\beta, P^{(i,j)}) \right\}_\beta &= \left\langle c_u(P^{(i,j)}), G_u(P^{(i,j)}) \right\rangle_Z.\end{aligned}$$

**Extended Horizon:** We notice from (5.6) that  $\beta$  alone generates  $\tilde{\mathcal{R}}(jt_c, P)$  since both  $\mathcal{X}^{(0)}$  and  $\left\{ \tilde{u}(\beta, P^{(0,j)}) \right\}_\beta$  are generated by  $\beta$ . As a result, the parameterized reachable set at  $t = (k+1)t_s$  for  $0 \leq k \leq h-1$  can be found by simply replacing  $\mathcal{X}^{(0)}$  with  $\tilde{\mathcal{R}}(kt_s, P)$  and executing the aforementioned steps again.

### 5.2.4 Controller Computation

With the parameterized reachable set computed in Sec. 5.2.3, we are now ready to approximate (5.1) as proposed in [99]. We first transform the objective function in (5.1) into a tractable approximation and then briefly discuss how the constraints of (5.1) can be adapted such that the resulting optimization problem can be posed as an LP. Due to the iterative nature of the GSC algorithm, we again assume without loss of generality that  $t = 0$ , i.e., we start at the initial set  $\mathcal{X}^{(0)}$ .

**Objective Function:** We bound (5.1a) over the extended optimization horizon from above as (also compare to [99, Eq. (9), (10)], [98, Lem. 2])

$$\max_{x(\cdot, P) \in \tilde{\mathcal{R}}(\cdot, P)} \sum_{k=1}^h \vartheta_k \left( \left\| x(kt_s, P) - x_f^{(k)} \right\|_1 + \zeta \sum_{j=0}^{m_c-1} \left\| u(x((k-1)t_s, P), P^{(k-1,j)}) \right\|_1 \right) \quad (5.7)$$

$$\leq \sum_{k=1}^h \vartheta_k \left( \left\| \begin{bmatrix} \tilde{c}(kt_c, P) \\ \tilde{G}(kt_c, P)_{(\cdot)} \end{bmatrix} - x_f^{(k)} \right\|_1 + \zeta \sum_{j=0}^{m_c-1} \left\| \begin{bmatrix} c_u(P^{(k-1,j)}) \\ G_u(P^{(k-1,j)})_{(\cdot)} \end{bmatrix} \right\|_1 \right), \quad (5.8)$$

as the cost function, where  $\zeta \in \mathbb{R}_+$  penalizes the overall input size, and  $\vartheta \in \mathbb{R}_{\geq 0}^h$  weighs the size of the reachable set at different points in time over the extended horizon  $h$ . The upper bound from the triangle inequality to arrive at (5.8) corresponds to a maximizing  $x(kt_s, P)$  which in general does not exist, as the following example demonstrates.

**Example 5.1.** Let  $h = 1$ ,  $x_f = 0$ ,  $\vartheta = \vartheta_1 = 1$ ,  $\zeta = 0$ , and  $\tilde{\mathcal{R}}(t_f, P) = \{\tilde{r}(t_f, \beta)\}_\beta$  with  $\tilde{r}(t_f, \beta) = \begin{bmatrix} 2 \\ -2 \end{bmatrix} \beta_1 + \begin{bmatrix} 1 \\ 1 \end{bmatrix} \beta_2$ . Then (5.7) can be equivalently written as

$$\begin{aligned}\max_{\|\beta\|_\infty \leq 1} \|\tilde{r}(t_f, \beta) - x_f\|_1 &= \max_{\|\beta\|_\infty \leq 1} (|2\beta_1 + \beta_2| + |-2\beta_1 + \beta_2|) \\ &= 4.\end{aligned}$$

In contrast, (5.8) yields

$$\begin{aligned} \max_{\|\beta\|_\infty \leq 1} \|\tilde{r}(t_f, \beta) - x_f\|_1 &\leq \left\| \begin{bmatrix} 2 \\ -2 \end{bmatrix} \right\|_1 + \left\| \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\|_1 \\ &= 6. \end{aligned}$$

■

Since the objective function is based on the linearized parameterized reachable set, its approximation quality tends to diminish with an increasing prediction horizon. As a result, one typically chooses  $\vartheta_1 \geq \vartheta_2 \geq \dots \geq \vartheta_h$ , and we have a direct approximation of (5.1a) for  $\vartheta_1 = \vartheta_2 = \dots = \vartheta_{h-1} = 0$ ,  $\vartheta_h = 1$ , and  $h = m_s$ .

**Constraints:** Using the approximated, parameterized reachable set instead of an outer approximation, the constraints in (5.1b) to (5.1d) for the extended horizon  $t \in [0, ht_s]$  are

$$\forall i \in \{0, \dots, h-1\} \forall j \in \{0, \dots, m_c-1\} : \left\{ \bar{u}(\beta, P^{(i,j)}) \right\}_\beta \subseteq \mathcal{U}, \quad (5.9)$$

$$\tilde{\mathcal{R}}([0, ht_s], P) \subseteq \mathcal{X}, \quad (5.10)$$

$$\tilde{\mathcal{R}}(ht_s, P) \subseteq \mathcal{X}_f. \quad (5.11)$$

Using Prop. 4.6, the input constraint  $\left\{ \bar{u}(P, \beta) \right\}_\beta \subseteq \mathcal{U} = \langle C_U, d_U \rangle_H$  can be written as

$$\left\{ \bar{u}(\beta, P^{(i,j)}) \right\}_\beta \subseteq \mathcal{U} \iff \rho_u(C_U, P^{(i,j)}) \leq d_U,$$

which can be encoded as a convex constraint by reformulating the absolute values using auxiliary variables (see Sec. 2.3.6.1) since  $\left\{ \bar{u}(\beta, P^{(i,j)}) \right\}_\beta = \hat{\mathcal{Z}}\left(\left\{ \bar{u}(\beta, P^{(i,j)}) \right\}_\beta\right)$ ; here,  $\rho_u(\cdot, P^{(i,j)})$  denotes the support function of  $\left\{ \bar{u}(\beta, P^{(i,j)}) \right\}_\beta$  and  $\hat{\mathcal{Z}}(\cdot)$  denotes the outer approximation of a polynomial zonotope by a zonotope as described in Sec. 4.1.4.

To check state constraints, we only have the approximate time-point solutions

$$\tilde{x}(it_s + jt_c, P) \in \tilde{\mathcal{R}}(it_s + jt_c, P), \quad 0 \leq i \leq m_s - 1, \quad 0 \leq j \leq m_c,$$

available. While the set of applied inputs  $\left\{ \bar{u}(P, \beta) \right\}_\beta$  is known exactly and thus the input constraints can be checked using a zonotope outer approximation as described above, the set of states is only approximate. Therefore, we cannot directly enforce the state constraints contained in  $\mathcal{X}$ : It might happen that the approximation of the parameterized reachable set is much smaller than an outer approximation of the reachable set, in which case the resulting controller will most likely violate the state constraints. Therefore, let  $\tilde{\mathcal{X}} = \langle C_{\tilde{\mathcal{X}}}, d_{\tilde{\mathcal{X}}} \rangle_H \subseteq \mathbb{R}^{n_x}$  with  $C_{\tilde{\mathcal{X}}} \in \mathbb{R}^{o_x \times n_x}$  and  $d_{\tilde{\mathcal{X}}} \in \mathbb{R}^{o_x}$  denote the user-provided state constraint set which is adapted to the approximated parameterized reachable set (see Sec. 5.2.7 for details).

## 5 Piecewise Constant Controller Synthesis

We now approximate the interval constraint (5.10) with  $m_c h$  time-point constraints and apply Prop. 4.6, yielding

$$\max_{0 \leq k \leq h m_c} \tilde{\rho}_x(C_{\tilde{\mathcal{X}}}, kt_c, P) \leq d_{\tilde{\mathcal{X}}},$$

where  $\tilde{\rho}_x(\cdot, t, P)$  denotes the support function of  $\tilde{\mathcal{R}}(t, P)$ . Finally, the final state constraints can be rewritten using again Prop. 4.6 as

$$\tilde{\rho}_x(C_{\tilde{\mathcal{X}}_f}, ht_s, P) \leq d_{\tilde{\mathcal{X}}_f},$$

where  $\tilde{\mathcal{X}}_f = \langle C_{\tilde{\mathcal{X}}_f}, d_{\tilde{\mathcal{X}}_f} \rangle_H \subseteq \mathbb{R}^{n_x}$  with  $C_{\tilde{\mathcal{X}}_f} \in \mathbb{R}^{o_{\mathcal{X}}_f \times n_x}$  and  $d_{\tilde{\mathcal{X}}_f} \in \mathbb{R}^{o_{\mathcal{X}}_f}$  are user-provided, adapted final state constraints.

**Optimization Problem:** The optimization problem for the extended horizon  $h$  is then given by (also compare to [99, Eq. (12), (13)], [98, Th. 3])

$$\hat{P}^{(0:h-1)} = \arg \min_P \sum_{k=1}^h \vartheta_k \left( \left\| \begin{bmatrix} \tilde{c}(kt_c, P) - x_f^{(k)} \\ \tilde{G}(kt_c, P) \end{bmatrix}_{(\cdot)} \right\|_1 + \zeta \sum_{j=0}^{m_c-1} \left\| \begin{bmatrix} c_u(P^{(k-1,j)}) \\ G_u(P^{(k-1,j)}) \end{bmatrix}_{(\cdot)} \right\|_1 \right), \quad (5.12a)$$

$$\text{s.t.} \quad \max_{\substack{0 \leq i \leq h-1 \\ 0 \leq j \leq m_c-1}} \rho_u(C_U, P^{(i,j)}) \leq d_U, \quad (5.12b)$$

$$\max_{0 \leq k \leq h m_c} \tilde{\rho}_x(C_{\tilde{\mathcal{X}}}, kt_c, P) \leq d_{\tilde{\mathcal{X}}}, \quad (5.12c)$$

$$\tilde{\rho}_x(C_{\tilde{\mathcal{X}}_f}, t_f, P) \leq d_{\tilde{\mathcal{X}}_f}. \quad (5.12d)$$

The absolute values and maximizations can be reformulated using the techniques from Sec. 2.3.6 which then yields an LP. While we obtain optimal controller parameters for all  $h$  steps, we only apply the first  $m_c$  piecewise constant controllers corresponding to  $\hat{P}^{(0)}$  until  $t = t_s$  (see Alg. 1).

### 5.2.5 Closed-Loop Reachable Set

Subsequently, we show how an outer approximation of the closed-loop reachable set can be computed. Without loss of generality, we first focus on the reachable set computation for  $t \in [0, t_s]$ . Given an optimal  $\hat{P}^{(0)}$  as the solution to (5.12) for  $t \in [0, t_s]$ , the controller is (see (5.3))  $\bar{u}(\beta, \hat{P}^{(0,j)})$  for  $0 \leq j \leq m_c - 1$ . The initial set is given by  $\mathcal{X}^{(0)} = \langle c_{\mathcal{X}^{(0)}}, G_{\mathcal{X}^{(0)}} \rangle_Z = \{x^{(0)}(\beta)\}_\beta$ .

We mainly use two different algorithms to compute reachable sets in this thesis: In Sec. 5.2.5.1, we describe the reachable set computation using the conservative linearization approach from [9] described in Sec. 2.7.2, and in Sec. 5.2.5.2 we introduce the reachable set computation for polynomial zonotopes using the conservative polynomialization approach from [5], which we introduced in Sec. 2.7.3.

### 5.2.5.1 Conservative Linearization

In Sec. 2.7, we discussed how the reachable set for given closed-loop system dynamics can be computed (also compare to [98, Sec. 3.3.6]). Thus, we introduce the extended dynamics

$$f_{\text{ext}}(x_{\text{ext}}, w) = \begin{bmatrix} f(x, u, w) \\ 0 \\ 0 \end{bmatrix},$$

with the extended state  $x_{\text{ext}} = [x^T, u^T, \beta^T]^T$  and the generating function of the extended initial set

$$x_{\text{ext}}^{(0)}(\beta) = \begin{bmatrix} x^{(0)}(\beta) \\ \bar{u}(\beta, \hat{P}^{(0,0)}) \\ \beta \end{bmatrix}.$$

Since the computed controller is dependent on  $\beta$ , we include the input as an additional state variable, which allows us to define  $x_{\text{ext}}^{(0)}(\beta)$  such that we always apply the corresponding pair of initial state and applied input. For  $t \in [0, t_c]$  with  $t_c = \frac{t_f}{m_s m_c}$ , we can then compute the extended reachable set  $\mathcal{R}_{\text{ext}}(t) = \{r_{\text{ext}}(t, \tilde{\beta})\}_{\tilde{\beta}}$  using the approach in Sec. 2.7.2. At  $t = t_c$ , we need to halt the reachability computations since the controller and thus the extended dynamics change. Due to reduction operations during reachability analysis and because reachability analysis with zonotopes does not preserve dependencies, the generating function of the extended reachable set is in general now dependent on a new dependent factor  $\tilde{\beta}$  instead of  $\beta$ . Hence, in order to apply the correct control input, we need a mapping from the new dependent factors  $\tilde{\beta}$  to the desired dependent factors  $\beta$ . Fortunately, this is available in the extended generating function of the reachable set at  $t = t_c$ , since by definition of the extended state we have

$$\beta(\tilde{\beta}) = [r_{\text{ext}}(t_c, \tilde{\beta})]_{n_x + n_u + (1:n_x)}. \quad (5.13)$$

Substitution of (5.13) into the controller template then yields  $\bar{u}(\beta(\tilde{\beta}), \hat{P}^{(0,1)})$ , which is now parameterized in the current dependent factor  $\tilde{\beta}$ . Thus, we can define the generating function for the new extended initial set as

$$x_{\text{ext}}^{(1)}(\tilde{\beta}) = \begin{bmatrix} [r_{\text{ext}}(t_c, \tilde{\beta})]_{(1:n_x)} \\ \bar{u}(\beta(\tilde{\beta}), \hat{P}^{(0,1)}) \\ \beta(\tilde{\beta}) \end{bmatrix},$$

from which  $\mathcal{R}_{\text{ext}}(2t_c)$  can be computed. The remaining steps until  $t = t_s$  follow analogously.

### 5.2.5.2 Conservative Polynomialization

Using the approach from Sec. 2.7.3, we can make use of the dependency-preserving operations of polynomial zonotopes (see Sec. 2.6.6). We again define an extended state

## 5 Piecewise Constant Controller Synthesis

$x_{\text{ext}} = \begin{bmatrix} x^T & u^T \end{bmatrix}$  with extended dynamics

$$f_{\text{ext}}(x_{\text{ext}}, w) = \begin{bmatrix} f(x, u, w) \\ 0 \end{bmatrix},$$

and extended initial generating function

$$x_{\text{ext}}^{(0)}(\beta) = \begin{bmatrix} x^{(0)}(\beta) \\ \bar{u}(\beta, \hat{P}^{(0,0)}) \end{bmatrix}.$$

In contrast to Sec. 5.2.5.1, the extended state does not require the inclusion of  $\beta$  due to the dependency-retaining properties of polynomial zonotopes as implemented in CORA at the time of writing. That said, it is in principle possible to also implement dependency-retaining zonotopes.

The reachable set  $\mathcal{R}_{\text{ext}}(t_c) = \{\tilde{r}_{\text{ext}}(t_c, \beta)\}_{\beta} \oplus \langle c_{\text{err}}, G_{\text{err}} \rangle_Z$  is computed by executing the algorithm described in Sec. 2.7.3 until  $t = t_c$  using polynomial zonotopes. Here,  $\tilde{r}_{\text{ext}}(t_c, \beta)$  is the generating function of the reachable set with the  $\beta$  dependency preserved, and  $\{c_{\text{err}} + G_{\text{err}}\delta\}_{\delta}$  represents abstraction errors (see Sec. 2.7.3) and reduction errors (see Sec. 2.8). The new extended initial set is then given by the generating function

$$x_{\text{ext}}^{(1)}(\beta, \delta) = \begin{bmatrix} [\tilde{r}_{\text{ext}}(t_c, \beta) + c_{\text{err}} + G_{\text{err}}\delta]_{(1:n_x)} \\ \bar{u}(\beta, \hat{P}^{(0,1)}) \end{bmatrix},$$

from which  $\mathcal{R}_{\text{ext}}(t)$  for  $t \in [1, 2]t_c$  can be computed. The remaining steps follow analogously.

### 5.2.6 Computational Complexity

We first discuss the offline complexity of the controller synthesis and then briefly describe the complexity of applying the computed controller online. We make the following assumptions.

**Assumption 5.1.** The number of elementary operations  $e \in \mathbb{N}_+$  to evaluate the closed-loop dynamics  $[f_{\text{cl}}]_i : \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \mapsto \mathbb{R}^{n_x}$  for each dimension  $1 \leq i \leq n_x$  as well as any of its derivatives, the number of input constraints  $o_{\mathcal{U}} \in \mathbb{N}_+$ , state constraints  $o_{\mathcal{X}} \in \mathbb{N}_+$ , and final state constraints  $o_{\mathcal{X}_f} \in \mathbb{N}_+$ , the number of generators  $l \in \mathbb{N}_+$  for any reachable set, and the number of dependent factors of the initial set all grow with at most  $O(n)$ , i.e.,

$$\begin{aligned} O(e) &= O(n), \\ O(o_{\mathcal{U}}) &= O(o_{\mathcal{X}}) = O(o_{\mathcal{X}_f}) = O(n), \\ O(l) &= O(n), \end{aligned}$$

where  $n = \max(n_x, n_u)$ . ■



### 5.2.6.1 Offline Complexity

We derive an upper bound on the computational complexity of GSC in  $n$ . Since  $m_s$  is fixed, it suffices to inspect the algorithm for  $m_s = 1$ .

**Reference Trajectory:** Generally, the optimization problem (5.2) to compute the reference trajectory is a smooth non-convex optimization problem, which can be solved (using second-order methods) to a first-order critical point with accuracy  $\epsilon$  using a number of function evaluations polynomially bounded in  $\frac{1}{\epsilon}$  [18], where  $\epsilon > 0$  is the selected accuracy. On the other hand, while tools like ACADO do not provide bounds on the complexity to compute the reference trajectory in Sec. 5.2.1, they solve these optimization problems in practice very efficiently. For this complexity analysis, we simply define the computational complexity as  $O(c_{\text{ref}}(n))$ , where  $c_{\text{ref}}(n)$  is some function of  $n$ .

**Reachable Sets:** Since the parameterized reachable set in Sec. 5.2.3 is given as a simple analytic function, its computational complexity is clearly lower than the complexity of computing the closed-loop reachable set, which is at most  $O(n^5)$  [59, Sec. 4.1.4] when using the conservative polynomialization approach (see Sec. 5.2.5.2).

**Controller Computation:** The complexity of the controller computation is dominated by the solution of the optimization problem – which can be posed as a linear program – whose solution complexity is  $O(c_{\text{LP}}(n)) = O\left(n^\omega \log\left(\frac{n}{\delta}\right) \log^k(n)\right)$ , where  $\delta$  is the relative accuracy and  $k \in \mathbb{N}_+$  is independent of  $n$  [16].

**Overall Complexity:** The computational offline complexity of GSC is thus at most

$$O\left(c_{\text{GSC}}^{(\text{off})}(n)\right) = O\left(c_{\text{ref}}(n) + c_{\text{LP}}(n) + n^5\right). \quad (5.14)$$

### 5.2.6.2 Online Complexity

The output of GSC are the optimal controller parameters  $\hat{P}$  and the closed-loop reachable set  $\mathcal{R}([0, t_f])$ . Due to the sequential nature of GSC and because  $h$ ,  $m_s$ , and  $m_c$  are fixed, it again suffices to restrict the analysis to  $t \in [0, t_c]$ .

If  $\mathcal{X}^{(0)}$  is given as a parallelotope (or enclosed by one), obtaining  $\beta$  for a given  $x(0) \in \mathcal{X}^{(0)}$  amounts to a simple matrix vector multiplication with complexity  $O(n_x^2)$  (see (5.4)) and another matrix vector multiplication with complexity  $O(n_u n_x)$  (see (5.3)), so that the overall online complexity is  $O(n^2)$ . If we do not compute a parallelotope outer approximation of  $\mathcal{R}(it_s)$  and its  $l$  generators, obtaining  $\beta$  requires the solution of an LP with complexity  $O(c_{\text{LP}}(l)) = O(c_{\text{LP}}(n))$  since  $O(l) = O(n)$  by Ass. 5.1. Computing  $\bar{u}(\beta, \hat{P}^{(i,j)})$  is then one matrix-vector multiplication (see (5.3)) with complexity  $O(n_u l) = O(n^2)$ . Thus, the online application of GSC has complexity

$$O\left(c_{\text{GSC}}^{(\text{on})}(n)\right) = \begin{cases} O(c_{\text{LP}}(n)), & \text{no parallelotope outer approximation,} \\ O(n^2), & \text{otherwise} \end{cases}. \quad (5.15)$$

### 5.2.7 Discussion

Because the GSC approach uses a linearized, parameterized reachable set approximation and because the controller template is piecewise linear in the state, the controller synthesis only requires the solution of (5.12) which can be reformulated as an LP and therefore can be solved efficiently. However, this simplicity also means that the parameterized reachable set approximation might not be very accurate for nonlinear systems with large initial sets. Furthermore, a linear controller might not be enough to achieve the desired control performance.

Additionally, we so far simply accepted that the adapted state and final state constraints  $\tilde{\mathcal{X}}$  and  $\tilde{\mathcal{X}}_f$  are given by the user. However, it is still unclear how these sets can be obtained from the original state and final state constraints  $\mathcal{X}$  and  $\mathcal{X}_f$ , respectively. Since we have no information about the accuracy of the parameterized reachable set approximation and because its accuracy is also dependent on the initial set and the input set, finding  $\tilde{\mathcal{X}}$  and  $\tilde{\mathcal{X}}_f$  requires trial and error by the user for each system and each system setup for which a controller needs to be synthesized. As a result, the GSC approach is not trivial to apply for non-experts when state and final state constraints need to be enforced.

### 5.3 Polynomial Generator-Space Control

Due to the linear approximation of the parameterized reachable set and the limitation to linear control laws in GSC, it may be impossible to achieve satisfactory controller performance, especially when the system dynamics are very nonlinear. In this section, we introduce the polynomial generator-space control (PGSC) approach from our work in [37] which is an extension of the GSC approach. Since PGSC subsumes GSC, we again solve (5.1) by synthesizing  $m_s m_c$  piecewise constant controllers (see Fig. 5.1). In contrast to GSC, however, compute the parameterized reachable set using arbitrary system abstractions and allow for an extension to polynomial controller templates.

Alg. 2 shows the general procedure for PGSC. The algorithm iteratively computes the following steps to synthesize  $m_s m_c$  piecewise constant controllers with respect to input, state, and final state constraints:

- (i) Given the center of the current initial set, we compute a reference trajectory from that center to the target state (l. 7).
- (ii) Let  $i \in \{0, \dots, m_s - 1\}$  denote the current iteration. We now compute an approximation to the parameterized reachable set using higher-order abstractions for  $t \in [i, i + h] t_s$  (l. 8).
- (iii) Since we compute the parameterized reachable set using higher-order abstractions, the resulting set is generally no longer a zonotope, making the optimization problem in the next step possibly non-convex. Therefore, we compute an initial guess for the controller parameters using, e.g., the GSC approach (l. 9).
- (iv) Using the updated reference trajectory, the parameterized reachable set, and the computed initial guess for the controller parameters, we solve the synthesis problem to obtain the optimal controller parameters (l. 10).
- (v) Lastly, the outer approximation of the closed-loop reachable set for  $t \in [i, (i + 1)] t_s$  is computed for the newly synthesized optimal controller (l. 11).

Subsequently, we present the described steps of the algorithm in detail: In Sec. 5.3.1, we show how a drift of the closed-loop reachable set away from the reference trajectory for  $m_s > 1$  can be avoided and then introduce the proposed polynomial controller template in Sec. 5.3.2. In Sec. 5.3.3, this template is then used to compute an approximation to the closed-loop reachable set, parameterized in the controller parameters, over which we then optimize in Sec. 5.3.4 to find the optimal controller parameters. The closed-loop reachable set for the resulting controller is then computed in Sec. 5.3.5. We derive the computational offline and online complexity of the approach in Sec. 5.3.6 and prove in Sec. 5.3.7 that PGSC generalizes GSC under mild technical assumptions. We conclude the section with a comparison of the novel PGSC approach to the GSC approach using numerical experiments in Sec. 5.3.8 and a brief discussion the algorithm in Sec. 5.3.9.

---

**Algorithm 2** Polynomial generator-space control
 

---

```

1: function POLYGENSPACECONTROL( $\mathcal{X}^{(0)}, \mathcal{U}, \mathcal{X}, \tilde{\mathcal{X}}, \mathcal{X}_f, \tilde{\mathcal{X}}_f, \mathcal{W}, m_s, m_c, x_f, t_f, h, \kappa$ )
2:    $\bar{\mathcal{R}} = \mathcal{X}^{(0)}$ 
3:   for  $i = 0; i < m_s; i++$  do
4:      $h = \min(h, m_s - i)$ 
5:      $\bar{\mathcal{R}} = \bar{\mathcal{R}} \downarrow_o$  ▷  $o \in \mathbb{N}_+$  (reduction order)
6:      $\bar{\mathcal{X}} = \hat{\mathcal{Z}}(\bar{\mathcal{R}}) \downarrow_1$ 
7:      $[x_f, u_{\text{ref}}] = \text{REFERENCETRAJECTORY}(\bar{\mathcal{X}}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, m_s - i, m_c, x_f, t_f)$  ▷
       Sec. 5.3.1
8:      $\tilde{\mathcal{R}}^{(i)} = \text{PARAMREACH}(\bar{\mathcal{R}}, h, \kappa)$  ▷ Sec. 5.3.3
9:      $\bar{P}^{(i+(1:h))} = \text{INITGUESS}(\bar{\mathcal{X}}, \mathcal{U}, \mathcal{X}, \tilde{\mathcal{X}}, \mathcal{X}_f, \tilde{\mathcal{X}}_f, \mathcal{W}, m_s - i, m_c, x_f^{(\cdot)}, h, t_f)$  ▷ e.g.
       GSC (Sec. 5.2)
10:     $\hat{P}^{(i)} = \text{COMPUTECTRL}(\bar{P}^{(i+(1:h))}, \tilde{\mathcal{R}}^{(i)}, \mathcal{U}, \tilde{\mathcal{X}}, \tilde{\mathcal{X}}_f, x_f, u_{\text{ref}}, h, \kappa)$  ▷ Sec. 5.3.4
11:     $\mathcal{R}([i, i+1]t_s) = \text{REACH}(\bar{\mathcal{R}}, \hat{P}^{(i)})$  ▷ Sec. 5.3.5
12:     $\bar{\mathcal{R}} = \mathcal{R}((i+1)t_s)$ 
13:  end for
14:  return  $\hat{P} = [\hat{P}^{(0)T}, \dots, \hat{P}^{(m_s-1)T}]^T, \mathcal{R}([0, t_f])$ 
15: end function
    
```

---

### 5.3.1 Reference Trajectory

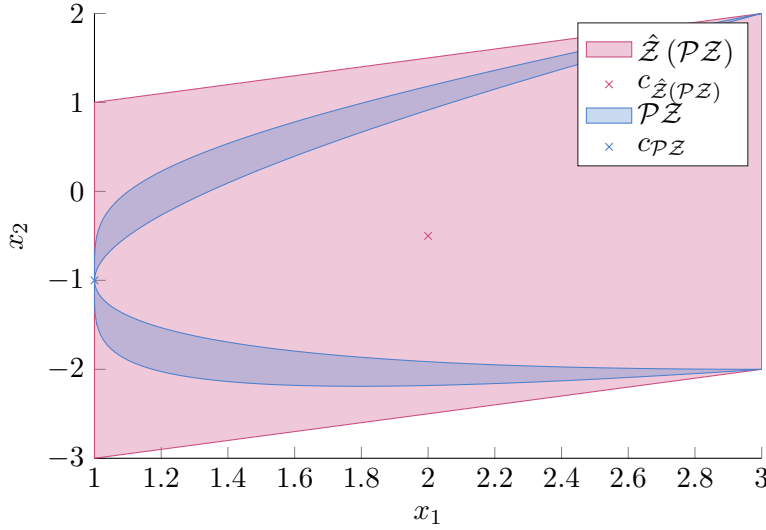
For the general computation procedure of the reference trajectory, we refer to Sec. 5.2.1. In this section, we discuss the necessary adaptations for  $\mathcal{X}^{(0)}$  being a polynomial zonotope, and then argue why an iterative update of the reference trajectory can be beneficial.

**Starting Point:** Let again be  $t = 0$  without loss of generality. In GSC,  $\mathcal{X}^{(0)}$  is given as a zonotope and thus the reference trajectory starts at the center of the initial set. While it need not be the case for nonlinear systems in general, in practice this often means that the reference trajectory runs centrally through the reachable set and ends centrally in the final, closed-loop reachable set at  $t = ht_s$ . As a result, it makes sense to use the reference trajectory at  $t = kt_c$  for  $0 \leq k \leq hm_c$  as linearization points for the computation of the abstraction (also see Sec. 5.3.3), based on which the parameterized reachable set is formed, to reduce abstraction errors.

In PGSC,  $\mathcal{X}^{(0)}$  is generally given as a polynomial zonotope. If now at least one exponent vector in the exponent matrix of  $\mathcal{X}^{(0)}$  has all-even exponents, its starting point no longer represents the center of the initial set due its even exponent vector, as the following example demonstrates.

**Example 5.2** (Center Shift of a Polynomial Zonotope). Given is the polynomial zonotope

$$\mathcal{PZ} = \langle c_{\mathcal{PZ}}, \cdot, \cdot \rangle_{\mathcal{PZ}} = \left\langle \left[ \begin{array}{c} 1 \\ -1 \end{array} \right], \left[ \begin{array}{cc} 0 & 2 \\ 2 & 1 \end{array} \right], \left[ \begin{array}{c} 0 \\ 0 \end{array} \right], \left[ \begin{array}{cc} 1 & 2 \\ 1 & 4 \end{array} \right] \right\rangle_{\mathcal{PZ}},$$



**Figure 5.3:** Comparison of the starting point of a polynomial zonotope to the center of its corresponding zonotope outer approximation using Prop. 4.4 as an approximation to the geometric center.

with its zonotope outer approximation according to Prop. 4.4 given by

$$\hat{\mathcal{Z}}(\mathcal{PZ}) = \langle c_{\mathcal{Z}}, \cdot \rangle_Z = \left\langle \begin{bmatrix} 2 \\ -\frac{1}{2} \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ 2 & \frac{1}{2} \end{bmatrix} \right\rangle_Z.$$

As is clearly visible in Fig. 5.3, the geometric center is better approximated by the center  $c_{\mathcal{Z}}$  of the zonotope outer approximation than by the starting point  $c_{\mathcal{PZ}}$  of the polynomial zonotope. ■

Thus, to obtain a reference trajectory that is roughly centered in the convex hull of the reachable set – which is beneficial to reduce the abstraction error during reachability analysis (see Sec. 2.7) – we compute the parallelotope outer approximation

$$\bar{\mathcal{X}}^{(0)} = \langle c_{\bar{\mathcal{X}}^{(0)}}, G_{\bar{\mathcal{X}}^{(0)}} \rangle_Z = \hat{\mathcal{Z}}(\mathcal{X}^{(0)}) \downarrow_1,$$

and use its center  $c_{\bar{\mathcal{X}}^{(0)}}$  as an approximation to the geometric center of  $\mathcal{X}^{(0)}$  for the computation of the reference trajectory.

**Update:** Finding a reference trajectory from a given initial state to the target state using optimal control as described in Sec. 5.2.1 works very well in practice: Not only is the solution process efficient, but we also obtain an input trajectory that steers us very close to the target state if such a trajectory exists under input and state constraints. Thus, it often makes sense to use the reference input directly in the controller template as described in Sec. 5.2.2: This not only ensures that the controller steers all initial

states close to the target state, but also reduces the number of optimization variables in the synthesis problem.

To that end, we assume that  $t = 0$  and let the reference trajectory from the geometric center  $c^{(0)}$  of the initial set to the target state  $x_f$  be given by  $u_{\text{ref}}^{(0:m_s m_c - 1)}$  with intermediate target states  $x_f^{(1+i)}$  for  $0 \leq i \leq m_s - 1$ . Now denote with  $u(c^{(0)}, P^{(0,j)})$  the controller template parameterized in  $P^{(0,j)}$  (see Sec. 5.3.2 for details) and let the first  $m_c$  controllers be chosen such that  $u(c^{(0)}, P^{(0,j)}) = u_{\text{ref}}^{(j)}$  for  $0 \leq j \leq m_c - 1$ . Thus, the exact state trajectory brings  $c^{(0)}$  to the first intermediate target state  $x_f^{(1)}$ . However, reachability analysis does not integrate the dynamics exactly, such that the geometric center  $c^{(1)}$  of the closed-loop reachable set at  $t = t_s$  is generally not equal to  $x_f^{(1)}$ . If we now synthesize the controllers for the next step again such that  $u(c^{(1)}, P^{(1,j)}) = u_{\text{ref}}^{(m_c+j)}$ , even the exact state trajectory from  $c^{(1)}$  under  $u_{\text{ref}}^{(m_c+(0:m_c-1))}$  does not reach  $x_f^{(2)}$  exactly since  $u_{\text{ref}}^{(m_c+(0:m_c-1))}$  steers  $x_f^{(1)}$  to  $x_f^{(2)}$  and  $c^{(1)} \neq x_f^{(1)}$ . To avoid this shift, we update the reference trajectory in each iteration, starting from the geometric center (see Alg. 2, l. 7).

### 5.3.2 Controller Template

Without loss of generality, we again start at  $t = 0$ , and define with slight abuse of notation  $P = P^{(0,0)}$  to avoid unnecessary indices.

Since  $\mathcal{X}^{(0)} = \{x^{(0)}(\beta)\}_{\beta}$  is in general represented as a polynomial zonotope, finding  $\beta$  for a given  $x(0) \in \mathcal{X}^{(0)}$  for the online application of the controller by solving

$$x(0) = x^{(0)}(\beta), \quad \|\beta\|_{\infty} \leq 1, \quad (5.16)$$

is in general hard since  $\beta$  now is the solution of a polynomial equation under the constraint  $\|\beta\|_{\infty} \leq 1$ . Thus, we compute a parallelotope outer approximation  $\tilde{\mathcal{X}}^{(0)} = \langle c_{\tilde{\mathcal{X}}^{(0)}}, G_{\tilde{\mathcal{X}}^{(0)}} \rangle_Z = \hat{\mathcal{Z}}(\mathcal{X}^{(0)}) \downarrow_1 \supseteq \mathcal{X}^{(0)}$  with corresponding dependent factor  $\tilde{\beta} \in [-1, 1]^{n_x}$  to obtain the analytic dependence  $\tilde{\beta} = G_{\tilde{\mathcal{X}}^{(0)}}^{-1}(\bar{x}(0) - c_{\tilde{\mathcal{X}}^{(0)}})$  for  $\bar{x}(0) \in \tilde{\mathcal{X}}^{(0)}$ . Because  $\bar{x}(0) \in \tilde{\mathcal{X}}^{(0)} \implies x(0) \in \mathcal{X}^{(0)}$ , we use

$$\begin{aligned} c_{\tilde{\mathcal{X}}^{(0)}} + G_{\tilde{\mathcal{X}}^{(0)}} \tilde{\beta} &= x^{(0)}(\beta) \\ \iff \tilde{\beta} &= G_{\tilde{\mathcal{X}}^{(0)}}^{-1}(x^{(0)}(\beta) - c_{\tilde{\mathcal{X}}^{(0)}}), \end{aligned} \quad (5.17)$$

to obtain an analytic dependence between  $\tilde{\beta}$  and  $\beta$ . Next, we derive the controller template in  $\tilde{\beta}$ ; by substituting  $\tilde{\beta}$  from (5.17) into said template, we also obtain a template in  $\beta$ .

While the linear controller template of the GSC approach in (5.3) is simple, it is numerically ill-conditioned: If, e.g.,  $\mathcal{U}$  is near-degenerate, different row dimensions of  $P$  may have values over multiple orders of magnitude. For example, the thrust force and the ruder angle as inputs to an aircraft model may have wildly different magnitudes. As

a result, a numerically well-conditioned linear controller template of the GSC approach can be written as

$$\tilde{u}_{\text{lin}}(\tilde{\beta}, P) = c_{\bar{U}} + G_{\bar{U}}P \begin{bmatrix} 1 \\ \tilde{\beta} \end{bmatrix}, \quad (5.18)$$

where we used a parallelotope outer approximation  $\bar{U} = \langle c_{\bar{U}}, G_{\bar{U}} \rangle_Z \supseteq \mathcal{U}$ . Instead of a direct parameterization of the input in  $P$ , we use  $\bar{U} = \{\bar{u}(\alpha)\}_\alpha$  with  $\bar{u}(\alpha) = c_{\bar{U}} + G_{\bar{U}}\alpha$  to parameterize the dependent factor  $\alpha$  in  $P$ , which achieves a parameterization with evenly scaled  $P$ : We have  $\mathcal{U} \subseteq \bar{U} \subseteq \{\tilde{u}_{\text{lin}}(\tilde{\beta}, P)\}_{\tilde{\beta}, P}$  since  $[-1, 1]^{n_u} \subseteq \left\{P \begin{bmatrix} 1 \\ \tilde{\beta} \end{bmatrix}\right\}_{\tilde{\beta}, P}$ . Thus, it

suffices to choose  $P \in [-1, 1]^{n_u \times (1+n_x)}$  to generate all admissible inputs. Note that  $\bar{U}$  is used for scaling purposes only; we still ensure the original input constraints using  $\mathcal{U}$ . We compute a parallelotope outer approximation of  $\mathcal{U}$  since a zonotope outer approximation of  $\mathcal{U}$  would increase the dimension of  $\alpha$  in  $\bar{u}(\alpha)$  beyond  $n_u$ . Thus, the row dimension of  $P$  would also grow, which in turn increases the number of optimization variables. If  $\mathcal{U}$  is also available as a zonotope,  $\bar{U}$  can be obtained using reduction techniques (see Sec. 2.8). Otherwise,  $\bar{U}$  can be obtained using, e.g., Prop. 4.3. A natural extension of (5.18) to realize polynomial controllers then is

$$\tilde{u}(\tilde{\beta}, P) = c_{\bar{U}} + G_{\bar{U}} \left( \sum_{k=1}^a P_{(:,k)} \tilde{\beta}^{E_{(:,k)}} \right), \quad (5.19)$$

where  $E \in \mathbb{N}^{n_x \times a}$  is a matrix of exponents with the assumption of  $E_{(:,1)} = 0$  to model constant offsets,  $a \in \mathbb{N}_+$  is the number of monomials in  $E$ , and  $P \in [-1, 1]^{n_u \times a}$ . Similarly to the linear template in (5.18), the polynomial template in (5.19) also has uniform scaling in  $P$ , i.e.,  $P \in [-1, 1]^{n_u \times a}$  since  $[-1, 1]^{n_u} \subseteq \left\{ \sum_{k=1}^a P_{(:,k)} \tilde{\beta}^{E_{(:,k)}} \right\}_{\tilde{\beta}, P}$ . As a more convenient alternative to defining the exponent matrix directly, we further define the controller order  $\kappa \in \mathbb{N}_+$ : For a given controller order, we use the exponent matrix  $E(\kappa) \in \mathbb{N}^{n_x \times o(\kappa)}$ , which collects all  $o(\kappa) = \binom{\kappa+n_x-1}{n_x-1}$  unique exponent vectors up to the total order  $\kappa$ , i.e.,  $\forall k \in \{1, \dots, o(\kappa)\} : 1_{n_x}^T E(\kappa)_{(:,k)} \leq \kappa$ , where  $E(\kappa)_{(:,1)} = 0$ .

Substitution of (5.17) into (5.19) then yields the controller template in  $x(0)$  as

$$u(x(0), P) = \tilde{u} \left( G_{\bar{\mathcal{X}}(0)}^{-1} (x(0) - c_{\bar{\mathcal{X}}(0)}), P \right), \quad (5.20)$$

and the controller template in  $\beta$  is

$$\bar{u}(\beta, P) = u \left( x^{(0)}(\beta), P \right). \quad (5.21)$$

**Reference Input for Center State:** If we wish to use the reference input in the controller directly, i.e., set the respective parameters such that  $u(c_{\bar{\mathcal{X}}(0)}, P) = u_{\text{ref}}^{(0)}$  (see Sec. 5.3.1 for why  $c_{\bar{\mathcal{X}}(0)}$ ), it follows that  $\tilde{u}(0, P) = u_{\text{ref}}^{(0)}$  and thus

$$\begin{aligned} c_{\bar{U}} + G_{\bar{U}}P_{(:,1)} &= u_{\text{ref}}^{(0)} \\ \iff P_{(:,1)} &= G_{\bar{U}}^{-1} \left( u_{\text{ref}}^{(0)} - c_{\bar{U}} \right), \end{aligned}$$

since  $E_{(:,1)} = 0$  by assumption.

**Remark:** As stated initially, the derivations in this section assumed that we start from the initial set  $\mathcal{X}^{(0)}$ . If the controller is to be synthesized from  $t = it_s$  for  $1 \leq i \leq m_s - 1$ , one can simply replace the initial set with the final, closed-loop reachable set at  $t = it_s$  that is obtained from the closed-loop reachable set computation using the previous controller (see Sec. 5.3.5).

### 5.3.3 Parameterized Reachable Set

In this section, we use a higher-order polynomial abstraction of our system to compute a parameterized reachable set (see also Sec. 2.7.3). This is in contrast to the GSC approach, where only a linear, time-discretized approximation is used to compute the parameterized reachable set. Similarly to GSC, we compute the parameterized reachable set based on the non-disturbed, nominal system dynamics.

In order to retain the dependency of  $\mathcal{X}^{(0)}$  and the controller template on their dependent factor  $\beta$ , we extend the state by the input, which yields the generating function

$$x_{\text{ext}}^{(0)}(P, \beta) = \begin{bmatrix} x^{(0)}(\beta) \\ \bar{u}(\beta, P^{(0,0)}) \end{bmatrix}, \quad (5.22)$$

of the extended initial set for  $x_{\text{ext}} = [x^T, u^T]^T$  with the extended flow given by

$$f_{\text{ext}}(x_{\text{ext}}) = \begin{bmatrix} f(x, u, 0) \\ 0 \end{bmatrix}.$$

Applying reachability analysis as described in Sec. 2.7.3 while interpreting both  $\beta$  and  $P$  in (5.22) as dependent factors for  $t \in [0, t_c]$  yields

$$\mathcal{R}(t, P) = \mathcal{D}(t, P) \oplus \mathcal{Z}_{\text{err}} = \left\{ [r_{\text{ext}}(t, P, \delta)]_{(1:n_x)} \right\}_{\delta}, \quad (5.23)$$

i.e., we generate  $\mathcal{R}(t, P)$  over  $\delta$ , which collects all dependent and independent factors except for  $P$ : While we treat  $P$  as a dependent factor during the computation of the parameterized reachable set to retain the dependency of the generating function  $r_{\text{ext}}(t, P, \delta)$  on  $P$ , we then generate  $r_{\text{ext}}(t, P, \delta)$  only over  $\delta$  to obtain  $\mathcal{R}(t, P)$  as an expression with  $P$  as an argument. We choose the linearization point  $\bar{x}_{\text{ext}}^{(k)} = [\bar{x}^{(k)T}, \bar{u}^{(k)T}]^T$  for the computation of  $\hat{\mathcal{R}}(t, P)$  with  $0 \leq k \leq m_r - 1$ ,  $m_r \in \mathbb{N}_+$  denoting the number of reachability steps in  $t \in [0, t_c]$ , and reachability step length  $r = \frac{t_c}{m_r}$ , where

$$\bar{x}^{(k)} = \frac{1}{2} \left( \xi(kr, c_{\bar{\mathcal{X}}^{(0)}}, u_{\text{ref}}^{(0)}, 0) + \xi((k+1)r, c_{\bar{\mathcal{X}}^{(0)}}, u_{\text{ref}}^{(0)}, 0) \right), \quad (5.24)$$

$$\bar{u}^{(k)} = u_{\text{ref}}^{(0)}. \quad (5.25)$$

Hence we ensure that the linearization during reachability analysis also occurs along the reference trajectory as in Sec. 5.2.3; specifically, we have  $\xi(t_c, c_{\bar{\mathcal{X}}^{(0)}}, u_{\text{ref}}^{(0)}, 0) = x_{\text{ref}}^{(1)}$ , where  $\bar{\mathcal{X}}^{(0)} = \hat{\mathcal{Z}}(\mathcal{X}^{(0)}) \downarrow_1 = \langle c_{\bar{\mathcal{X}}^{(0)}}, G_{\bar{\mathcal{X}}^{(0)}} \rangle_Z$ . The reachable set approximation  $\hat{\mathcal{R}}([0, t_c], P)$  is



then directly available since reachability analysis returns time-point and time-interval solutions.

However, by treating  $P$  as a dependent factor in the computation of  $r_{\text{ext}}(t, P, \delta)$ , we essentially compute the reachable set for all closed-loop controllers parameterized by  $P \in [-1, 1]^{n_u \times a}$ . Since the result of reachability analysis generally worsens for large initial sets (since the abstraction used in reachability analysis is only locally accurate around the linearization points), increasing the size of the bounded input  $\mathcal{U}$  causes an increasingly large outer approximation in reachability analysis which is detrimental to the approximation quality of  $\hat{\mathcal{R}}(t, P)$ , as demonstrated in the following example with  $m_s = 1$  for ease of presentation.

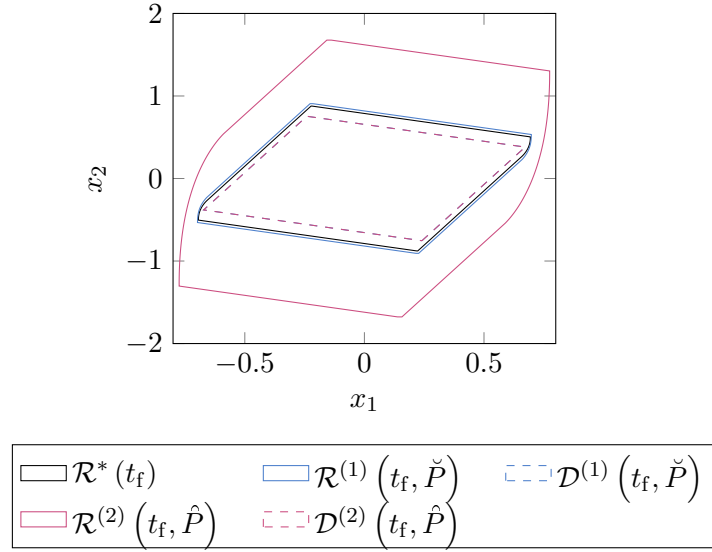
**Example 5.3.** Let us consider a controlled Van-der-Pol oscillator with the dynamics  $f(x, u) = [x_2, (1 - x_1^2)x_2 - x_1 + u]^T$ ,  $\mathcal{X}^{(0)} = \left\{ \left[ \begin{array}{c} 0.5\beta_1 \\ 0.3\beta_2 \end{array} \right] \right\}_\beta$ ,  $t_f = 0.5$ ,  $\bar{u}(\beta, P^{(0,j)}) = G_{\bar{u}}(p_1^{(j)}\beta_1 + p_2^{(j)}\beta_2)$  for  $P^{(0,j)} = [p_1^{(j)}, p_2^{(j)}]$  with  $0 \leq j \leq m_c - 1 = 4$ ,  $E = I_2$ , and the parallelotope outer approximation  $\bar{\mathcal{U}} = \langle 0, G_{\bar{u}} \rangle_Z$ .

Let a tight outer approximation  $\mathcal{R}^*(t_f)$  of the closed-loop reachable set for the controller  $\bar{u}^{*(j)}(\beta) = \frac{1}{2}(\beta_1 + \beta_2)$  be given. We now compare this outer approximation to the parameterized reachable sets

$$\begin{aligned} \mathcal{R}^{(1)}(t, \check{P}) &= \mathcal{D}^{(1)}(t, \check{P}) \oplus \mathcal{Z}_{\text{err}}^{(1)}, \\ \mathcal{R}^{(2)}(t, \hat{P}) &= \mathcal{D}^{(2)}(t, \hat{P}) \oplus \mathcal{Z}_{\text{err}}^{(2)}, \end{aligned}$$

evaluated at two fixed controller parameters  $\check{P} \in [-1, 1]^{m_c n_u \times a}$  and  $\hat{P} \in [-1, 1]^{m_c n_u \times a}$ , with increasing input capacities  $\bar{\mathcal{U}}^{(1)} = \langle 0, \frac{1}{5} \rangle_Z$  and  $\bar{\mathcal{U}}^{(2)} = \langle 0, 2 \rangle_Z$ : The controller templates are hereby given by  $\bar{u}^{(1)}(\beta, P^{(0,j)}) = \frac{1}{5}(p_1^{(j)}\beta_1 + p_2^{(j)}\beta_2)$  and  $\bar{u}^{(2)}(\beta, P^{(0,j)}) = 2(p_1^{(j)}\beta_1 + p_2^{(j)}\beta_2)$ , from which  $\check{P} = 1_{m_c n_u \times a}$  and  $\hat{P} = 0.1 \cdot 1_{m_c n_u \times a}$  to replicate  $\bar{u}^{*(j)}(\beta)$  follow. Fig. 5.4 visualizes these sets. Evidently, increasing the input capacity clearly increases the bloating of  $\mathcal{R}^{(2)}(t, \check{P})$  compared to  $\mathcal{R}^{(1)}(t, \hat{P})$ . In contrast,  $\mathcal{D}^{(1)}(t, \check{P})$  and  $\mathcal{D}^{(2)}(t, \hat{P})$  are practically identical. ■

Thus, while using  $\mathcal{R}(t, P)$  directly as an approximation for the parameterized reachable set can become quite inaccurate for increasing control capacity, the set  $\mathcal{D}(t, P)$  remains even for larger input capacities a good approximation as the example demonstrated. Therefore, we propose to use  $\tilde{\mathcal{R}}(t, P) = \mathcal{D}(t, P)$  as an approximation to the parameterized reachable set in order to avoid bloating effects for large input sets. The approximated reachable set for  $t \in [1, 2]t_c$  follows by choosing  $\tilde{\mathcal{R}}(t_c, P)$  as the initial set for the next step with corresponding input  $\left\{ \bar{u}(\beta, P^{(0,1)}) \right\}_\beta$  and executing the above steps again, where the remaining steps until  $t = ht_s$  follow analogously.



**Figure 5.4:** Approximation of the tight outer approximation  $\mathcal{R}^*(t_f)$  of the reachable set by  $\mathcal{R}^{(1)}(t, \check{P}) = \mathcal{D}^{(1)}(t, \check{P}) \oplus \mathcal{Z}_{\text{err}}^{(1)}$  for  $\check{P} = 1_{m_c n_u \times a}$  and  $\mathcal{R}^{(2)}(t, \hat{P}) = \mathcal{D}^{(2)}(t, \hat{P}) \oplus \mathcal{Z}_{\text{err}}^{(2)}$  for  $\hat{P} = 0.1 \cdot 1_{m_c n_u \times a}$ . Note that  $\mathcal{D}^{(1)}(t, \check{P})$  is not visible as it is visually identical to  $\mathcal{D}^{(2)}(t, \hat{P})$ .

### 5.3.4 Controller Computation

We are now ready to pose the synthesis problem using the parameterized reachable set computed previously. Similarly to Sec. 5.2, we also consider an extended optimization horizon  $h$ . Without loss of generality, we further assume that we start from  $t = 0$ .

Let the parameterized reachable set  $\tilde{\mathcal{R}}(t, P) = \{\tilde{r}(t, P, \beta)\}_{\beta}$  for  $t \in [0, ht_s]$  from Sec. 5.3.3 be given. Since  $\tilde{\mathcal{R}}(t, P)$  is represented as a polynomial zonotope, there are dependencies between different generators. Since range bounding of polynomials is NP-hard in general [33], we compute a zonotope outer approximation  $\tilde{\mathcal{Z}}(t, P) = \hat{\mathcal{Z}}(\tilde{\mathcal{R}}(t, P)) = \langle \tilde{c}(t, P), \tilde{G}(t, P) \rangle_{\mathcal{Z}}$  of the approximated parameterized reachable set to be able to efficiently bound the maximization in (5.1a). We note that this zonotope outer approximation only operates on  $\beta$ , i.e.,  $\tilde{\mathcal{Z}}(t, P)$  still captures the nonlinearity of its center and generators on  $P$ .

**Objective Function:** It is clear that the parameterized zonotope  $\tilde{\mathcal{Z}}(t, P)$  is in structure identical to the parameterized reachable set of the GSC approach – with the exception that it is now nonlinearly parameterized in  $P$  – which yields an upper bound on the objective function analogously to the derivation in (5.8).

**Constraints:** Using the approximated, parameterized reachable set instead of an outer approximation, the constraints in (5.1b) to (5.1d) for the extended horizon  $t \in [0, ht_s]$

can be written as

$$\forall i \in \{0, \dots, m_s - 1\} \forall j \in \{0, \dots, m_c - 1\} : \left\{ \bar{u} \left( \beta, P^{(i,j)} \right) \right\} \subseteq \mathcal{U}, \quad (5.26)$$

$$\tilde{\mathcal{R}}([0, ht_s], P) \subseteq \mathcal{X}, \quad (5.27)$$

$$\tilde{\mathcal{R}}(t_s, P) \subseteq \mathcal{X}_f. \quad (5.28)$$

To efficiently bound the range of the controller template in  $\beta$ , we use the support function  $\rho_u^{(\mathcal{Z})}(\cdot, P^{(i,j)})$  of the zonotope outer approximation  $\hat{\mathcal{Z}}\left(\left\{u\left(P^{(i,j)}, \beta\right)\right\}_\beta\right) = \left\langle c_u\left(P^{(i,j)}\right), G_u\left(P^{(i,j)}\right) \right\rangle_{\mathcal{Z}}$  and replace the input constraint in (5.26) using Prop. 4.6 by

$$\max_{\substack{0 \leq i \leq h-1 \\ 0 \leq j \leq m_c-1}} \rho_u^{(\mathcal{Z})}\left(C_{\mathcal{U}}, P^{(i,j)}\right) \leq d_{\mathcal{U}}.$$

For the formulation of state constraints, time-point and time-interval solutions  $\tilde{\mathcal{Z}}(t, P)$  for  $t \in [0, ht_s]$  are available since we use reachability analysis for the computation of  $\tilde{\mathcal{R}}(t, P)$ . Let  $m_r \in \mathbb{N}_+$  be the number of steps for reachability analysis for the computation of  $\tilde{\mathcal{R}}(t, P)$  with  $t \in it_s + [j, j+1]t_c$  for  $0 \leq i \leq m_s - 1$  and  $0 \leq j \leq m_c - 1$ . Further, denote with  $\tilde{\rho}_x^{(\mathcal{Z})}(\cdot, t, P)$  the corresponding support function of the zonotope outer approximation  $\hat{\mathcal{Z}}(\tilde{\mathcal{R}}(t, P))$ : By using a zonotope outer approximation of the reachable set, we can conservatively approximate the polynomial zonotope-in-polytope constraints in (5.27) and (5.28) using Prop. 4.6, which yields

$$\begin{aligned} \max_{0 \leq k \leq m_r h m_c - 1} \tilde{\rho}_x^{(\mathcal{Z})}\left(C_{\tilde{\mathcal{X}}}, [k, k+1]r, P\right) &\leq d_{\tilde{\mathcal{X}}}, \\ \tilde{\rho}_x^{(\mathcal{Z})}\left(C_{\tilde{\mathcal{X}}_f}, ht_s, P\right) &\leq d_{\tilde{\mathcal{X}}_f}, \end{aligned}$$

where  $\tilde{\mathcal{X}} = \langle C_{\tilde{\mathcal{X}}}, d_{\tilde{\mathcal{X}}} \rangle_H \subseteq \mathbb{R}^{n_x}$  with  $C_{\tilde{\mathcal{X}}} \in \mathbb{R}^{o_{\mathcal{X}} \times n_x}$  and  $d_{\tilde{\mathcal{X}}} \in \mathbb{R}^{o_{\mathcal{X}}}$ ,  $\tilde{\mathcal{X}}_f = \langle C_{\tilde{\mathcal{X}}_f}, d_{\tilde{\mathcal{X}}_f} \rangle_H \subseteq \mathbb{R}^{n_x}$  with  $C_{\tilde{\mathcal{X}}_f} \in \mathbb{R}^{o_{\mathcal{X}_f} \times n_x}$  and  $d_{\tilde{\mathcal{X}}_f} \in \mathbb{R}^{o_{\mathcal{X}_f}}$  denote the adapted state and final state constraints required due to the approximate nature of  $\tilde{\mathcal{R}}(t, P)$ , and  $r = \frac{t_f}{m_r m_s m_c}$ . Similarly to Sec. 5.2.4, these adapted constraints are necessary since we do not know how accurate the approximated parameterized reachable set is.

**Optimization Problem:** Collecting the objective function and the constraints finally yields the synthesis problem

$$\hat{P}^{(0:h-1)} = \arg \min_P \tilde{J}(P), \quad (5.29a)$$

$$\text{s.t. } \tilde{g}(P) \leq 0, \quad (5.29b)$$

with

$$\tilde{J}(P) = \sum_{k=1}^h \vartheta_k \left( \left\| \begin{bmatrix} \tilde{c}(kt_s, P) - x_f^{(k)} \\ \tilde{G}(kt_s, P) \end{bmatrix}_{(\cdot)} \right\|_1 + \zeta \sum_{j=0}^{m_c-1} \left\| \begin{bmatrix} c_u(P^{(k-1,j)}) \\ G_u(P^{(k-1,j)}) \end{bmatrix}_{(\cdot)} \right\|_1 \right), \quad (5.30)$$

$$\tilde{g}(P) = \begin{bmatrix} \max_{\substack{0 \leq j \leq m_c-1 \\ 0 \leq i \leq h-1}} \rho_u^{(\mathcal{Z})}(C_{\mathcal{U}}, P^{(i,j)}) - d_{\mathcal{U}} \\ \max_{0 \leq k \leq m_r h m_c - 1} \tilde{\rho}_x^{(\mathcal{Z})}(C_{\tilde{\mathcal{X}}}, [k, k+1]r, P) - d_{\tilde{\mathcal{X}}} \\ \tilde{\rho}_x^{(\mathcal{Z})}(C_{\tilde{\mathcal{X}}}, ht_s, P) - d_{\tilde{\mathcal{X}}} \end{bmatrix}, \quad (5.31)$$

where  $\vartheta \in \mathbb{R}_{\geq 0}^h$  weighs each term of the objective function and where absolute values and maximization expressions can be reformulated (see Sec. 2.3.6) to yield a differentiable optimization problem. Depending on the chosen abstraction order and controller order, (5.29) is no longer expressible as an LP or even a convex optimization problem (see Sec. 5.3.7 for details).

Strictly speaking, (5.29) cannot be written as the abstracted problem in (3.2) directly. However, let  $s \in \mathbb{R}^{o_{\mathcal{U}}+o_{\mathcal{X}}+o_{\mathcal{X}_f}}$  be a slack variable  $s \geq 0$  to relax (5.29) to

$$\hat{P}^{(0:h-1)} = \arg \min_P \tilde{J}(P) + \sigma \|s\|_1, \quad (5.32a)$$

$$\text{s.t. } \tilde{g}(P) \leq s, \quad (5.32b)$$

where  $\sigma \in \mathbb{R}_+$  is a large enough factor such that minimizing the constraint violation, i.e., avoiding  $s > 0$ , is always prioritized over minimizing  $\tilde{J}(P)$  (also see [50]). Assuming (5.29) has a feasible global optimizer  $\hat{P}$ , i.e.,  $\tilde{g}(\hat{P}) \leq 0$ , the feasible point  $(\hat{P}, 0)$  of (5.32) is the global optimum of (5.32) since  $\tilde{J}(\hat{P}) < \tilde{J}(\hat{P}) + \sigma \|s\|$  for any  $s > 0$  due to  $\sigma$  being chosen large enough, and since (5.29) and (5.32) are identical for  $s = 0$ .

### 5.3.5 Closed-Loop Reachable Set

Since we use nonlinear abstractions, which involve quadratic (or even higher-order) maps, for reachability analysis, it makes sense to use polynomial zonotopes for reachability analysis since they are closed under these maps. To compute the closed-loop reachable set, the procedure in Sec. 5.2.5.2 can be used.

### 5.3.6 Computational Complexity

We again differentiate between the complexity of the offline controller synthesis and the application of said controller online. We only discuss the complexity of the first iteration starting at  $t = 0$  since the total number of iterations  $m_s$  is fixed. For simplicity, let  $\kappa \geq 1$  for the order of the controller template introduced in Sec. 5.3.2 and assume that the number of monomials  $a$  of the controller template grows at least with  $O(n)$ , i.e.,  $O(n+a) = O(a)$ . Additionally, we again make use of Ass. 5.1.

### 5.3.6.1 Offline Complexity

Subsequently, we derive an upper bound on the computational complexity of the proposed PGSC approach in  $n$ . Due to the similarities between GSC and PGSC (also see Sec. 5.3.7), we only present the complexity analysis for parts of the algorithm that are different from GSC.

**Controller Template:** Without loss of generality, we inspect the complexity for the controller template parameterized in  $P^{(0,0)}$  and set –with slight abuse of notation –  $P = P^{(0,0)}$  to avoid unnecessary indices. The computation of

$$\bar{u}(\beta, P) = \tilde{u}(\tilde{\beta}(\beta), P) = c_{\bar{u}} + G_{\bar{u}} \left( \sum_{k=1}^a P_{(:,k)} \tilde{\beta}(\beta)^{E_{(:,k)}} \right),$$

from Sec. 5.3.2 requires the polynomial substitution of  $\tilde{\beta}(\beta) = G_{\bar{\mathcal{X}}^{(0)}}^{-1} (x^{(0)}(\beta) - c_{\bar{\mathcal{X}}^{(0)}})$ . With the controller template  $\{\tilde{u}(\tilde{\beta}, P)\}_{\tilde{\beta}}$  given by the user with controller order  $\kappa$ , we compute its composition with  $\{\tilde{\beta}(\beta)\}_{\beta}$  using Prop. A.2.

The computation of  $\{\tilde{\beta}(\beta)\}_{\beta}$  has complexity  $O(n_x^\omega + n_u^2 l) = O(n^3)$  due the computation of the inverse  $G_{\bar{\mathcal{X}}^{(0)}}^{-1}$  and the matrix multiplication [61, Prop. 8], where  $\omega \in \mathbb{R}_+$  is the current complexity exponent for matrix multiplication. The complexity of the composition  $\{\bar{u}(\beta, P)\}_{\beta}$  from  $\{\tilde{u}(\tilde{\beta}, P)\}_{\tilde{\beta}}$  and  $\{\tilde{\beta}(\beta)\}_{\beta}$  is then  $O((\kappa(n_x + an) + an_u) l^\kappa) = O(an^{1+\kappa})$  according to Prop. A.2.

**Reachable Sets:** We compute reachable sets using the conservative polynomialization approach from [5] as described in Sec. 2.7.3. Thus, the computation of the parameterized reachable set is  $O(n^5 + n^2 a^2 + n^3 a \log n)$  according to Prop. A.1 and the complexity of the closed-loop reachable set computation can be bounded from above by  $O(n^5)$  [59, Sec. 4.1.4].

**Controller Computation:** The number of objective and constraint evaluations can be polynomially bounded in  $\frac{1}{\epsilon}$  for second-order optimization methods [18], where  $\epsilon > 0$  is the specified accuracy of a first-order critical point of (5.29); thus, the number of iterations required is independent of  $n$  for a given  $\epsilon$ . As a result, it is sufficient to derive the complexity of one objective and constraint function evaluation. Since we check a fixed number of different constraint types, i.e., state, input, and final state constraints, and the number of constraints for each grows at most with  $O(n)$  by assumption, we limit the analysis here to final state constraints. Thus, let  $C_{\bar{\mathcal{X}}_f} \tilde{\mathcal{R}}(t_f, P) = \left\{ c + \sum_{k=1}^l M_{(:,k)} \prod_{o=1}^{an_u} p_o^{E_{(o,k)}} \prod_{o=1}^m \beta_o^{B_{(o,k)}} \right\}_{\beta}$  be the final constraint set obtained from the parameterized reachability analysis, where  $c \in \mathbb{R}^{o_{\mathcal{X}_f}}$  is the starting point,  $M \in \mathbb{R}^{o_{\mathcal{X}_f} \times l}$  is the generator matrix,  $\begin{bmatrix} E \\ B \end{bmatrix} \in \mathbb{N}_{\geq 0}^{(an_u+m) \times l}$  is the exponent matrix,  $p = P_{(:,\cdot)}$ ,  $\beta \in [-1, 1]^m$ ,

## 5 Piecewise Constant Controller Synthesis

and we have  $O(o_{\mathcal{X}_f}) = O(l) = O(m) = O(n)$  due to Ass. 5.1. First, we compute the collected polynomial zonotope  $\langle \cdot, N(P), 0, B \rangle_{PZ}$  with

$$N_{(:,j)}(P) = \sum_{k \in \mathcal{K}^{(j)}} M_{(:,k)} \prod_{o=1}^{an_u} p_o^{E_{(o,k)}},$$

$$\mathcal{K}^{(j)} = \left\{ k \in \{1, \dots, l\} \mid \forall o \in \{1, \dots, an_u\} : B_{(o,j)} = B_{(o,k)} \right\},$$

for  $1 \leq j \leq \tilde{l}$ , where  $\tilde{l}$  is the number of unique columns in  $B$  with  $\tilde{l} < l$  and  $\mathcal{K}^{(j)}$  contains all generator/exponent vector column indices that have the same exponents in  $\beta$ , and where it further holds that  $\sum_{j=1}^{\tilde{l}} |\mathcal{K}^{(j)}| = l$  by construction. The complexity of computing this collected polynomial zonotope is identical to the complexity of computing all  $\mathcal{K}^{(j)}$ , which is at most  $O(an_u \tilde{l} \log \tilde{l}) = O(an^2 \log n)$  [59, Prop 2.5.2]. We then compute its zonotope outer approximation  $\hat{\mathcal{Z}}(\langle \cdot, N(P), 0, B \rangle_{PZ}) = \langle \cdot, G(P) \rangle_Z$  with complexity  $O(n^2)$  [59, Prop. 3.1.4], where  $G(P) \in \mathbb{R}^{o_{\mathcal{X}_f} \times \tilde{l}}$  denotes its generator matrix and where we ignore the center as it can be included as a generator with an all-zero exponent vector. From Sec. 2.3.6.1, it is clear that the complexity of computing and evaluating the objective and constraint functions of the smooth reformulation of (5.29) for final state constraints only will be dominated by evaluating and computing the  $O(o_{\mathcal{X}_f} \tilde{l}) = O(n_x^2)$  auxiliary constraints since  $G(P) \in \mathbb{R}^{o_{\mathcal{X}_f} \times \tilde{l}}$ . Further, the complexity of computing and then evaluating these auxiliary constraints will be dominated by the computation and evaluation of all Hessian matrices of  $G_{ij}(P)$ .

To evaluate  $\nabla_{P_{(i)}}^2 G_{ij}(P)$ , we form  $\mathcal{G}_{(:,j)} = \{G_{(:,j)}(P)\}_P = \langle \cdot, \bar{G}^{(j)}, 0, \bar{E} \rangle_{PZ}$  for  $1 \leq j \leq \tilde{l}$ , i.e., a polynomial zonotope that represents  $G_{(:,j)}(P)$  with  $\bar{G}^{(j)} \in \mathbb{R}^{o_{\mathcal{X}_f} \times |\mathcal{K}^{(j)}|}$  and  $\bar{E}^{(j)} \in \mathbb{N}^{an_u \times |\mathcal{K}^{(j)}|}$ : Any polynomial zonotope  $G_{(:,j)}(P)$  has at most  $|\mathcal{K}^{(j)}|$  generators and any of its  $|\mathcal{K}^{(j)}|$  exponent vectors is one of the at most  $l$  unique exponent vectors in  $E \in \mathbb{N}_{\geq 0}^{an_u \times l}$ . Thus, we first compute the Hessian matrices of all  $l$  possible monomials; the Hessian matrix of any particular  $\nabla_{P_{(i)}}^2 G_{ij}(P)$  then follows as a product of the  $|\mathcal{K}^{(j)}|$  corresponding monomial Hessian matrices with the generator values in  $\bar{G}_{(i,:)}^{(j)}$ .

To evaluate the Hessian matrix of a monomial  $m^{(j)}(P) = \prod_{k=1}^{an_u} p_k^{E_{(k,j)}}$  for a fixed  $p = P_{(i)}$ , we first evaluate all  $l$  monomials  $m^{(j)}(P)$ , which has complexity  $O(ln_u a) = O(an^2)$  since its evaluation requires  $n_u a$  exponentiations and  $n_u a - 1$  multiplications. Given  $m^{(j)}(P)$ , any second-order derivative can be found as

$$\frac{\partial^2 m^{(j)}(P)}{\partial p_k \partial p_o} = \begin{cases} \frac{E_{(k,j)} E_{(o,j)}}{p_k p_o} m^{(j)}(P), & E_{(k,j)} \geq 1 \wedge E_{(o,j)} \geq 1 \wedge p_k p_o \neq 0, \\ 0, & \text{otherwise} \end{cases}.$$

Since computing any second-order derivative from a given  $m^{(j)}(P)$  has complexity  $O(1)$  by inspection of the above equation, computing all Hessian matrices of  $m^{(j)}(P)$  for a

given  $P$  has complexity  $O\left(l(n_u a)^2\right) \stackrel{\text{Ass. 5.1}}{=} O(a^2 n^3)$  since each matrix has  $O\left((n_u a)^2\right)$  elements.

Computing  $\nabla_{P(i)}^2 G_{ij}(P)$  for a given  $i$  and  $j$  now requires  $|\mathcal{K}^{(j)}|$  scalar multiplications of the generator values in  $\tilde{G}_{(i,:)}^{(j)}$  with the pre-computed Hessian matrices of the monomials and a subsequent summation of these  $|\mathcal{K}^j|$  matrices, which both have complexity  $O\left(|\mathcal{K}^{(j)}|(n_u a)^2\right)$ . Repeating this for all  $o_{\mathcal{X}_f}$  rows and all  $\tilde{l}$  columns of  $G(P)$ , the overall complexity is bounded from above by  $O\left(o_{\mathcal{X}_f} \tilde{l} (n_u a)^2\right) \stackrel{\text{Ass. 5.1}}{=} O(a^2 n^4)$  since  $\sum_{j=1}^{\tilde{l}} |\mathcal{K}^{(j)}| = \tilde{l}$ .

**Overall Complexity:** The overall offline complexity of PGSC is thus at most

$$\begin{aligned} O\left(c_{\text{PGSC}}^{(\text{off})}(n)\right) &= O\left(c_{\text{ref}}(n) + an^{1+\kappa} + n^5 + a^2 n^2 + an^3 \log n + a^2 n^4\right) \\ &= O\left(c_{\text{ref}}(n) + an^{1+\kappa} + a^2 n^4\right). \end{aligned} \quad (5.33)$$

### 5.3.6.2 Online Complexity

For the evaluation of the controller in (5.20), we first compute  $\tilde{\beta}$  using one matrix vector multiplication with complexity  $O(n_x^2)$  and then evaluate the control law in (5.19) over  $\tilde{\beta}$ , which has complexity  $O(an_x)$  [62, Sec. 3.3], followed by one last matrix-vector multiplication with complexity  $O(n_u^2)$ .

The online complexity of deploying PGSC is thus

$$O\left(c_{\text{PGSC}}^{(\text{on})}(n)\right) = O(an). \quad (5.34)$$

### 5.3.7 Generalization of Generator-Space Control

In this section, we show that the PGSC approach generalizes the GSC approach under mild technical assumptions.

**Proposition 5.1** (PGSC as Generalization of GSC). *Let  $\kappa = 1$  (linear controller),  $\pi = 1$  (linear abstraction),  $r = 1$  (number of reachability steps), omit any reduction operations, let  $\mathcal{X}^{(0)}$  be given as a non-degenerate parallelotope, and consider only input and final state constraints. Then the PGSC and GSC approach are identical.*

*Proof.* Comparing (5.3) and (5.21), it is obvious that both controllers are identical for  $\kappa = 1$  since  $\mathcal{X}^{(0)} = \tilde{\mathcal{X}}^{(0)}$  by assumption, and therefore  $\tilde{\beta} = \beta$  and  $\tilde{u}_{\text{lin}}(\tilde{\beta}, P^{(i,j)}) = \tilde{u}(\tilde{\beta}, P^{(i,j)}) = \bar{u}(\beta, P^{(i,j)})$  (see (5.18), (5.19), (5.21)) for all  $0 \leq i \leq m_s - 1$  and  $0 \leq j \leq m_c - 1$  if their parameterized reachable sets are identical. Because (5.12) and (5.29) are structurally identical, it suffices to show that both parameterized reachable sets are equal.

## 5 Piecewise Constant Controller Synthesis

Let  $\tilde{\mathcal{R}}(t, P) = \{\tilde{r}(t, P, \beta)\}_\beta$  denote the parameterized reachable set where  $\tilde{\mathcal{R}}(0, P) = \mathcal{X}^{(0)}$ . The extended generating function at step  $t = it_s$  following Sec. 5.3.3 is

$$x_{\text{ext}}^{(i)}(P, \beta) = \begin{bmatrix} \tilde{r}(it_s, P, \beta) \\ \bar{u}(\beta, P^{(i,0)}) \end{bmatrix},$$

and the linear abstraction of the extended flow  $\dot{x}_{\text{ext}} = f_{\text{ext}}(x_{\text{ext}})$  with  $x_{\text{ext}} = [x^T, u^T]^T$  while ignoring abstraction errors – projected onto the first  $n_x$  dimensions – yields

$$\tilde{r}(it_s + t_c, P, \beta) = \left( e^{A_{\text{ext}}^{(0)} t_c} x_{\text{ext}}^{(i)}(P, \beta) + \int_0^{t_c} e^{A_{\text{ext}}^{(0)} \tau} d\tau \begin{bmatrix} c^{(0)} \\ 0 \end{bmatrix} \right)_{(1:n_x)}, \quad (5.35)$$

where

$$A_{\text{ext}}^{(0)} = \left. \frac{df_{\text{ext}}(x_{\text{ext}})}{dx_{\text{ext}}} \right|_{x_{\text{ext}} = \bar{x}_{\text{ext}}^{(0)}} = \begin{bmatrix} A^{(0)} & B^{(0)} \\ 0 & 0 \end{bmatrix},$$

with  $\bar{x}_{\text{ext}}^{(0)} = [\bar{x}^{(0)T}, \bar{u}^{(0)T}]^T$ , where  $\bar{x}^{(0)}$  and  $\bar{u}^{(0)}$  are given as in (5.24) and (5.25), and  $A^{(0)}$ ,  $B^{(0)}$ , and  $c^{(0)}$  are defined as in Sec. 5.2.3. Further

$$\begin{aligned} e^{A_{\text{ext}}^{(0)} t} &= \sum_{k=0}^{\infty} \frac{(A_{\text{ext}}^{(0)} t)^k}{k!} = I_{n_x+n_u} + \sum_{k=1}^{\infty} \frac{\begin{bmatrix} (A^{(0)} t)^k & A^{(0)k-1} t^k B^{(0)} \\ 0 & 0 \end{bmatrix}}{k!} \\ &= \begin{bmatrix} \sum_{k=0}^{\infty} \frac{(A^{(0)} t)^k}{k!} & \sum_{k=1}^{\infty} \frac{A^{(0)k} t^k B^{(0)}}{k!} \\ 0 & I_{n_u} \end{bmatrix} = \begin{bmatrix} e^{A^{(0)} t} & \int_0^t e^{A^{(0)} \tau} d\tau \\ 0 & I_{n_u} \end{bmatrix}, \end{aligned}$$

which, substituted into (5.35), yields

$$\tilde{r}(it_s + t_c, P, \beta) = c_d^{(0)} + A_d^{(0)} \tilde{r}(it_s, P, \beta) + B_d^{(0)} \bar{u}(\beta, P^{(i,0)}), \quad (5.36)$$

where  $A_d^{(0)}$ ,  $B_d^{(0)}$ , and  $c_d^{(0)}$  are defined as in Sec. 5.2.3. Thus, (5.36) is identical to the parameterized reachable set computation in (5.6), and applying the remaining steps iteratively concludes the proof.  $\square$

Prop. 5.1 excludes state constraints since GSC only approximately enforces state constraint satisfaction at discrete points in time while PGSC approximately enforces state constraint satisfaction over time intervals.

### 5.3.8 Experiments

In this section, we compare the PGSC approach to the GSC approach introduced in Sec. 5.2 from [99]. Since manually guessing the adapted state and final state constraints is in general non-trivial, we choose to omit both state and final state constraints here, but show the final state constraints  $\mathcal{X}_f = \{-c_{\mathcal{X}^{(0)}} + x_f\} \oplus \mathcal{X}^{(0)}$  with  $\mathcal{X}^{(0)} = \langle c_{\mathcal{X}^{(0)}}, G_{\mathcal{X}^{(0)}} \rangle_Z$



for reference (helpful when result is used in a maneuver automaton; see Sec. 3.1); we will enforce state constraints on the benchmarks shown here in the next section. Lastly, we do not penalize input energy used here, i.e., we set  $\zeta = 0$ , for all benchmarks in this section.

To measure the control performance of each approach, we introduce

$$\text{csize}(\mathcal{R}(t_f)) = \left\| \left[ \begin{array}{c} c - x_f \\ G^{(\cdot)} \end{array} \right] \right\|_1, \quad (5.37)$$

for a given outer approximation  $\mathcal{R}(t_f)$  of the closed-loop reachable set with the zonotope outer approximation  $\mathcal{Z} = \hat{\mathcal{Z}}(\mathcal{R}(t_f)) = \langle c, G \rangle_{\mathcal{Z}}$ : We measure how much  $\mathcal{Z}$  differs from the target state  $x_f \in \mathbb{R}^{n_x}$  by measuring both the deviation of its center from  $x_f$  and the size of  $\mathcal{Z}$  by measuring the absolute sum of its generators.

### 5.3.8.1 Bicycle (Single Track)

While kinematic models are simpler, they do not consider tire slip and thus cannot model effects such as oversteering and understeering. We introduce a variant of the single-track model<sup>2</sup> – illustrated in Fig. 5.5 – given by

$$\dot{x}_1 = \frac{\mu}{x_4 l} \left( \left( \frac{c_r(u_2) l_r - c_f(u_2) l_f}{x_4} - 1 \right) x_3 + (c_r(u_2) + c_f(u_2)) x_1 + c_f(u_2) u_1 \right), \quad (5.38a)$$

$$\dot{x}_2 = x_3, \quad (5.38b)$$

$$\dot{x}_3 = -\frac{\mu m}{I_z (l_r + l_f)} \left( \frac{l_f^2 c_f(u_2) + l_r^2 c_r(u_2)}{x_4} x_3 + (l_r c_r(u_2) - l_f c_f(u_2)) x_1 + l_f c_f(u_2) u_1 \right), \quad (5.38c)$$

$$\dot{x}_4 = u_2, \quad (5.38d)$$

$$\dot{x}_5 = x_4 \cos(x_1 + x_2), \quad (5.38e)$$

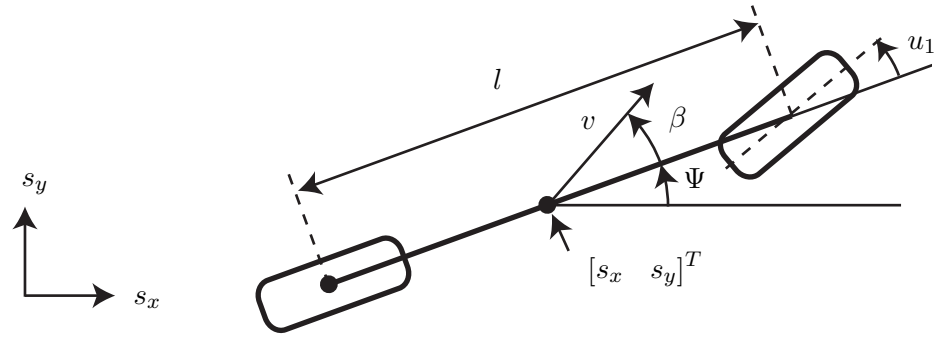
$$\dot{x}_6 = x_4 \sin(x_1 + x_2), \quad (5.38f)$$

with

$$\begin{aligned} c_f(u_2) &= C_{S,f}(gl_r - u_2 h_{cg}), \\ c_r(u_2) &= C_{S,r}(gl_f + u_2 h_{cg}), \\ l &= l_r + l_f, \end{aligned}$$

and where  $x = [\beta, \Psi, \dot{\Psi}, v, s_x, s_y]^T$  is the state,  $\beta$  is the slip angle at the vehicle center,  $\Psi$  and  $\dot{\Psi}$  are the yaw angle and yaw rate,  $v$  is the longitudinal velocity of the vehicle,  $[s_x, s_y]^T$  are the two spatial coordinates,  $u_1$  is the steering angle of the front

<sup>2</sup>Sec. 7 ([https://gitlab.lrz.de/tum-cps/commonroad-vehicle-models/-/blob/master/vehicleModels\\_commonRoad.pdf](https://gitlab.lrz.de/tum-cps/commonroad-vehicle-models/-/blob/master/vehicleModels_commonRoad.pdf))



**Figure 5.5:** States and inputs of the single-track model<sup>4</sup>. Shown are the slip angle  $\beta$ , the yaw angle  $\Psi$ , longitudinal velocity  $v$ , spatial coordinates  $[s_x \ s_y]^T$ , the steering angle  $u_1$ , and the wheelbase  $l$ .

wheels, and  $u_2$  is the longitudinal acceleration of the vehicle. The benchmark parameters are listed in Tab. 5.1 and have been chosen appropriately<sup>3</sup>.

We compute a feedforward controller for the undisturbed system, i.e., we set  $m_s = 1$  and  $\mathcal{W} = \{0\}$ . Fig. 5.6 shows a comparison between the performance of the GSC approach and the PGSC approach, where computations times were 27.5s and 36.9s, respectively, and  $\text{csize}(\mathcal{R}_{\text{GSC}}(t_f)) = 0.92$  and  $\text{csize}(\mathcal{R}_{\text{PGSC}}(t_f)) = 0.65$ . Clearly, the PGSC approach outperforms the GSC approach significantly as it reduces the overall cost significantly compared to GSC: While both approaches keep the yaw rate and velocity small (see Fig. 5.6,  $[x_3, \ x_4]$ ) and fail to keep the slip angle small (see Fig. 5.6,  $x_1$ ), the PGSC approach steers the set of initial positions more closely to the target state in the two spatial dimensions ( $[x_5, \ x_6]$ ).

### 5.3.8.2 Controlled Van-der-Pol Oscillator

We present a controlled version of the Van-der-Pol oscillator benchmark [9, Sec. VII]

$$\dot{x}_1 = x_2, \quad (5.39a)$$

$$\dot{x}_2 = (1 - x_1^2)x_2 - x_1 + u + w. \quad (5.39b)$$

We compute a piecewise constant, discrete-time feedback controller with  $m_s = 2$  and  $m_c = 4$ , where all parameters for the benchmark are listed in Tab. 5.2.

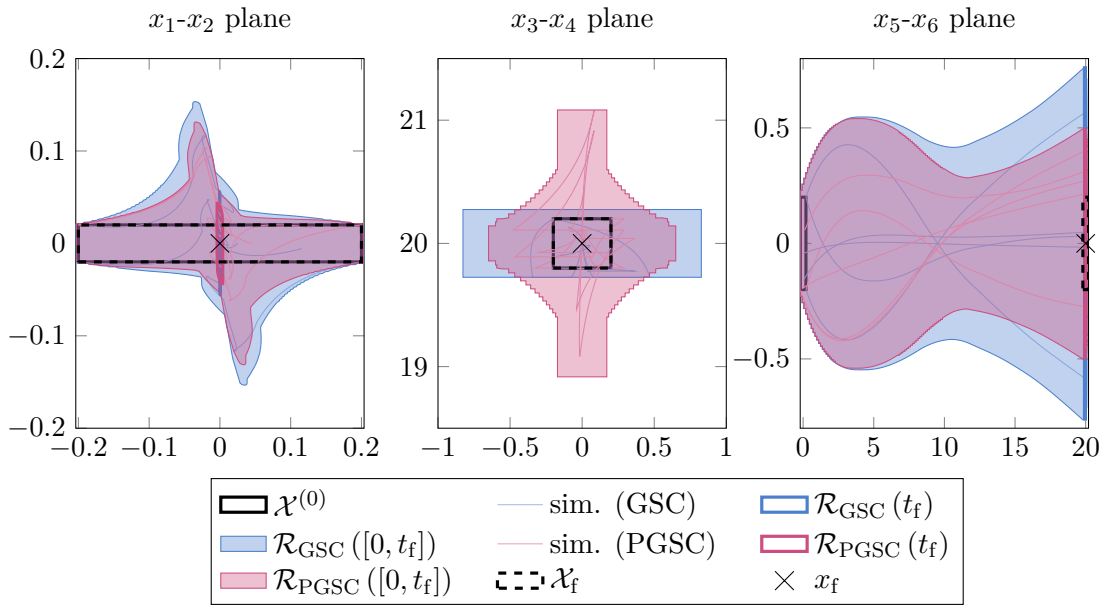
Fig. 5.7 shows the reachable set using the GSC approach and the PGSC approach, where the computation took 15.50s and 20.44s, respectively. Clearly, PGSC achieves a better control result as the final reachable set is contained within the final reachable

<sup>3</sup>Tab. 4, vehicle identifier 1 ([https://gitlab.lrz.de/tum-cps/commonroad-vehicle-models/-/blob/master/vehicleModels\\_commonRoad.pdf](https://gitlab.lrz.de/tum-cps/commonroad-vehicle-models/-/blob/master/vehicleModels_commonRoad.pdf))

<sup>4</sup>Fig. 4, vehicle identifier 1 ([https://gitlab.lrz.de/tum-cps/commonroad-vehicle-models/-/blob/master/vehicleModels\\_commonRoad.pdf](https://gitlab.lrz.de/tum-cps/commonroad-vehicle-models/-/blob/master/vehicleModels_commonRoad.pdf))

**Table 5.1:** Benchmark parameters for the bicycle model.

Parameter	Value
$m_s$	1
$m_c$	6
$\mathcal{X}^{(0)}$	$\langle c_{\mathcal{X}^{(0)}}, \text{diag}([0.2, 0.02, 0.2, 0.2, 0.2, 0.2]) \rangle_Z$
$c_{\mathcal{X}^{(0)}}$	$[0, 0, 0, 20, 0, 0]^T$
$\mathcal{U}$	$\langle 0, \text{diag}([0.4, 9.81]) \rangle_Z$
$\mathcal{W}$	$\{0\}$
$t_f$	1 s
$x_f$	$[0, 0, 0, 20, 20, 0]^T$
$\kappa$	1


**Figure 5.6:** Comparison of the GSC approach and the PGSC approach for the single-track (bicycle) model. For a fair comparison, both approaches use a linear controller template.

**Table 5.2:** Benchmark parameters for the Van-der-Pol model.

Parameter	Value
$m_s$	2
$m_c$	4
$\mathcal{X}^{(0)}$	$\left\langle \left[ 1.4, 2.4 \right]^T, \text{diag} \left( \left[ 0.6, 0.6 \right] \right) \right\rangle_Z$
$\mathcal{U}$	$\langle 0, 2 \rangle_Z$
$t_f$	1 s
$x_f$	$\left[ 1.93, -0.47 \right]^T$
$\kappa$	1

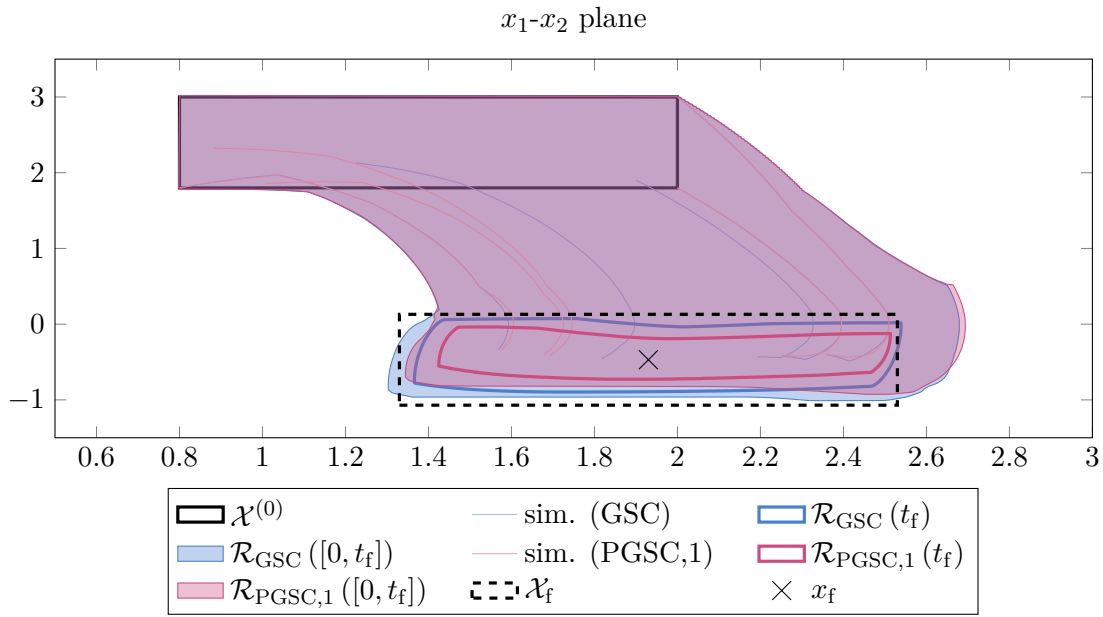
set computed by GSC. The combined deviations from the target state and size of final reachable set are  $\text{csize} \left( \hat{\mathcal{Z}} \left( \mathcal{R}_{\text{GSC}} \left( t_f \right) \right) \right) = 1.22$  and  $\text{csize} \left( \mathcal{R}_{\text{PGSC},1} \left( t_f \right) \right) = 1.09$  (see (5.37)).

Further, Fig. 5.8 shows the improved performance when using a quadratic controller ( $\kappa = 2$ ) over the same linear controller ( $\kappa = 1$ ) for this benchmark using PGSC, where the computation of the quadratic controller took 18.52 s. Here, we have  $\text{csize} \left( \mathcal{R}_{\text{PGSC},2} \left( t_f \right) \right) = 0.79$ .

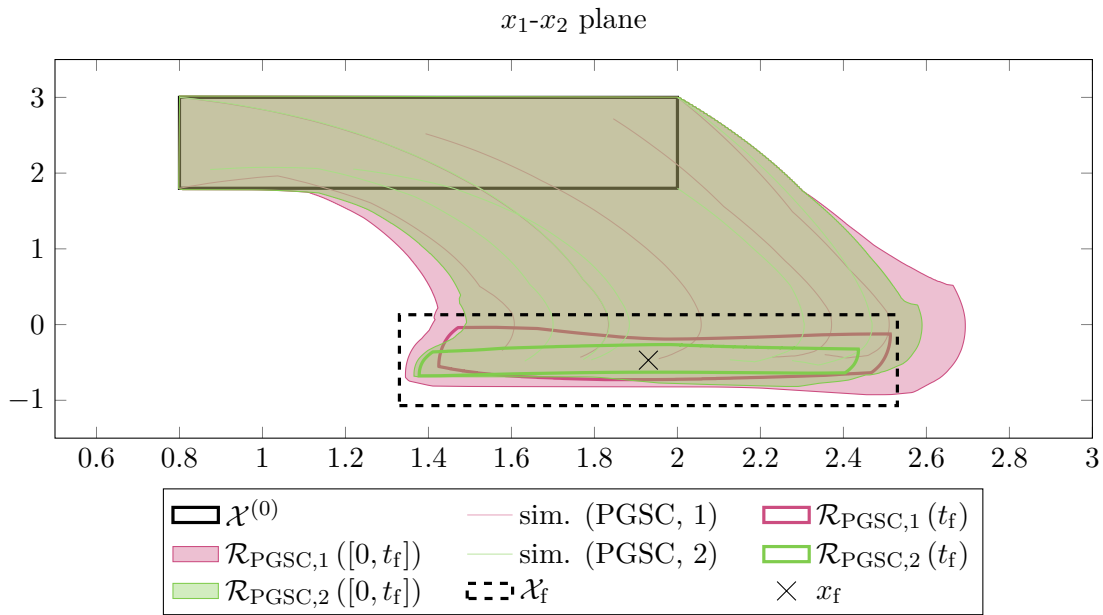
### 5.3.9 Discussion

The numerical experiments in Sec. 5.3.8 demonstrate improved performance of PGSC over GSC, both by using a more accurate parameterized reachable set and thus achieving a better control performance for the same linear controller template as well as by using, e.g., a quadratic controller over a linear one.

While both approaches work reasonably well for input constraints, the fact that they both only approximate the parameterized reachable set during the optimization of the controller – without any regard for the accuracy of this approximation – means that state and final state constraints can seldom be enforced as is; rather, they require the user to specify adapted state and final state constraints. Since this can become quite involved for non-experts, it is a major downside of both approaches. In the next section, we thus propose a novel synthesis algorithm that does not require the manual specification of these adapted state constraints.



**Figure 5.7:** Comparison of the GSC approach and the PGSC approach for the Van-der-Pol benchmark. For a fair comparison, both approaches use a linear controller template.



**Figure 5.8:** Comparison of a linear and quadratic controller template using the PGSC approach for the Van-der-Pol benchmark.

## 5.4 Iterative Polynomial Generator-Space Control

While the GSC approach used only linear abstractions for the computation of the parameterized reachable set and was limited to linear control laws, the PGSC approach extended this concept to polynomial abstraction orders for reachability analysis and polynomial control laws. While the reachable set can in principle be arbitrarily closely approximated using the conservative polynomialization approach by choosing the abstraction order large enough, in practice, we are often limited to at most quadratic abstractions due to extended computation times or the fact that higher-order abstractions generate even more generators during reachability analysis, which in turn increases the number of generators needing to be reduced; this effect can mostly or even completely negate the benefit of a higher abstraction order to begin with.

However, we still find it necessary to better approximate the parameterized reachable set for strongly nonlinear system dynamics. Therefore, we introduce the iterative polynomial generator-space control (iPGSC) approach, an iterative trust-region approach for the synthesis of formally verified controllers under input and state constraints according to (5.1), which is part of our work in [35]. Since (5.1) is generally not efficiently solvable, we approximate it using a non-differentiable cost function in the controller parameters. By locally approximating this non-differentiable cost function with a differentiable optimization problem, we can iteratively update our controller parameters based on this local approximation, where we control its accuracy by restricting the admissible controller parameters for the optimization problem to a bounded trust-region. In contrast to GSC and PGSC, iPGSC does not require manual tuning of the state constraints. Alg. 3 describes the general solution procedure, where the most important steps are briefly summarized subsequently:

- (i) We compute the reference trajectory at the start of each iteration, which yields the intermediate target states necessary for the optimization of the controller parameters (Alg. 3, l. 7).
- (ii) We compute an initial guess for the controller parameters using GSC (Alg. 3, l. 8).
- (iii) We compute the closed-loop reachable set of the controller, parameterized by the initial guess  $\bar{P}$  (Alg. 4, l. 3) and evaluate its initial cost  $J(\bar{P})$  (Alg. 4, l. 4). The algorithm then iteratively executes the following steps:
  - (a) Since the proposed cost is non-differentiable and its evaluation requires the computation of the closed-loop reachable set, we first compute the parameterized reachable set locally for a region of controller parameters  $\mathcal{P}_\gamma(\bar{P})$  around the current controller parameters  $\bar{P}$ , whose radius is bounded by the trust-region radius  $\gamma \in (0, 1]$  (Alg. 4, l. 6). To form this approximation, we use the available outer approximation of the closed-loop reachable set of the current controller parameterized by  $\bar{P}$ .
  - (b) Using this local parameterized reachable set, we approximate the cost function with a differentiable optimization problem, where the accuracy of this

trust-region subproblem is controlled by the trust-region radius  $\gamma$  (Alg. 4, l. 7). Solving the trust-region subproblem, we obtain the optimal controller parameters  $\hat{P}$ . Because we use the local parameterized reachable set for the trust-region subproblem, we avoid the recomputation of the closed-loop reachable set in each optimization iteration.

- (c) We then evaluate the cost of the new controller parameterized in  $\hat{P}$  using an outer approximation of its closed-loop reachable set (Alg. 4, l. 8–9).
  - (d) By comparing the actual decrease in cost against the decrease in cost predicted by the trust-region subproblem, we adapt the trust-region radius to ensure an accurate approximation of the parameterized reachable set in the next iteration (Alg. 4, l. 10): If the ratio is close enough to one, the approximated trust-region subproblem is accurate and we enlarge  $\gamma$  to allow for faster progress. Otherwise, we shrink the trust-region radius to ensure an accurate approximation in the next iteration.
  - (e) We then compare the cost of the current controller to the cost of the newly computed controller: If the cost decreased, we accept the step and save these new controller parameters (Alg. 4, l. 11–19). If the algorithm converged (Alg. 4, l. 11), we return the current controller parameters.
- (iv) Once the trust-region algorithm converges, we have the optimal controller parameters as well as the corresponding closed-loop reachable set available. If we have not yet reached the final iteration, we again start at (i). Otherwise, the algorithm terminates.

---

**Algorithm 3** Iterative polynomial generator-space control
 

---

```

1: function ITERPOLYGENSPACECONTROL( $\mathcal{X}^{(0)}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, \mathcal{W}, m_s, m_c, x_f, h, \kappa$ )
2:    $\bar{\mathcal{R}} = \mathcal{X}^{(0)}$ 
3:   for  $i = 0; i < m_s; i++$  do
4:      $h = \min(h, m_s - i)$ 
5:      $\bar{\mathcal{R}} = \bar{\mathcal{R}} \downarrow_o$  ▷  $o \in \mathbb{N}_+$  some sensible order
6:      $\bar{\mathcal{X}} = \hat{\mathcal{Z}}(\bar{\mathcal{R}}) \downarrow_1$ 
7:      $x_f^{(\cdot)} = \text{REFERENCETRAJECTORY}(\bar{\mathcal{X}}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, m_s - i, m_c, x_f, t_f)$  ▷ Sec. 5.3.1
8:      $\bar{P}^{(\cdot)} = \text{INITGUESS}(\bar{\mathcal{R}}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, \mathcal{W}, m_s - i, m_c, x_f^{(\cdot)}, h, t_f, \kappa)$  ▷ e.g. Sec. 5.2
9:      $[\hat{P}^{(i)}, \mathcal{R}([i, i+1]t_s)] = \text{COMPUTECTRL}(\bar{P}^{(\cdot)}, \bar{\mathcal{R}}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, x_f, u_{\text{ref}}, \kappa)$  ▷ Alg. 4
10:     $\bar{\mathcal{R}} = \mathcal{R}((i+1)t_s)$ 
11:  end for
12:  return  $\hat{P} = [\hat{P}^{(0)T}, \dots, \hat{P}^{(m_s-1)T}]^T, \mathcal{R}([0, t_f])$ 
13: end function
    
```

---

The following example illustrates the solution procedure, which is visualized in Fig. 5.9, where we set  $h = m_s = m_c = 1$  and only consider input constraints for illustrative purposes.

---

**Algorithm 4** Trust-region controller computation

---

```

1: function COMPUTECTRL( $\bar{P}^{(\cdot)}, \bar{\mathcal{X}}, h, x_f^{(\cdot)}$ )
2:   Given:  $\mu, \sigma, \gamma = 1$ 
3:    $\bar{\mathcal{R}} = \text{REACH}(\bar{\mathcal{X}}, \bar{P}^{(\cdot)}, ht_s)$  ▷ Sec. 5.2.5.2
4:    $\bar{J} = \text{COST}(\bar{\mathcal{R}})$  ▷ Sec. 5.4.1
5:   for  $k = 1; k \leq l_{\max}; k++$  do
6:      $[\tilde{\mathcal{R}}, \gamma] = \text{PARAMREACH}(\bar{P}^{(\cdot)}, \gamma, \bar{\mathcal{R}})$  ▷ Sec. 5.4.2
7:      $\hat{P}^{(\cdot)} = \text{SOLVETRSubPROBLEM}(\tilde{\mathcal{R}}, \bar{P}^{(\cdot)}, \gamma, \bar{\mathcal{R}}, x_f^{(\cdot)})$  ▷ Sec. 5.4.3
8:      $\hat{\mathcal{R}} = \text{REACH}(\bar{\mathcal{X}}, \hat{P}^{(\cdot)}, ht_s)$  ▷ Sec. 5.3.5
9:      $\hat{J} = \text{COST}(\hat{\mathcal{R}})$  ▷ Sec. 5.4.1
10:     $\gamma = \text{ADAPTRADIUS}(\hat{J}, \tilde{\mathcal{R}}, \hat{P})$  ▷ Sec. 5.4.4
11:    if  $\hat{J} < \bar{J}$  then
12:       $\nu = \text{tol}(\hat{J}, \bar{J})$  ▷ Def. 2.1
13:       $\bar{J} = \hat{J}$ 
14:       $\bar{\mathcal{R}} = \hat{\mathcal{R}}$ 
15:       $\bar{P}^{(\cdot)} = \hat{P}^{(\cdot)}$ 
16:      if  $\nu \leq \mu$  then
17:        break
18:      end if
19:    end if
20:  end for
21:  return  $[\bar{P}^{(0)}, \bar{\mathcal{R}}]$ 
22: end function

```

---



**Example 5.4** (Solve progress for a 1D system [35, Ex. 1, adapted]). Let the final, exact parameterized reachable set at  $t = t_f$  for a 1-dimensional system be given by

$$\mathcal{R}(t_f, P) = \left\{ -1 + 2p_1^2 + 4p_0^3 + (3 + 10p_0p_1^3)\beta + (2 - 2p_1)\beta^2 \right\}_\beta,$$

where  $\beta \in [-1, 1]$  is the dependent factor of the corresponding initial set,  $\bar{u}(\beta, P) = p_0 + p_1\beta$  with  $P = P^{(0)} = P^{(0,0)} = [p_0, p_1]$  is the controller template as given in Sec. 5.3.2,  $x_f = 0$ , and  $\mathcal{U} = [-1, 1] = \left\langle [1, -1]^T, [1, 1]^T \right\rangle_H = \langle C\mathcal{U}, d\mathcal{U} \rangle_H$ . Using Prop. 4.4, the zonotope outer approximation of the reachable set required to form (5.42) is given by

$$\begin{aligned} \hat{\mathcal{Z}}(\mathcal{R}(t_f, P)) &= \left\{ -1 + 2p_1^2 + 4p_0^3 + (3 + 10p_0p_1^3)\nu_1 + (2 - 2p_1)\left(\frac{1}{2} + \frac{1}{2}\nu_2\right) \right\}_\nu \\ &= \left\{ -p_1 + 2p_1^2 + 4p_0^3 + (3 + 10p_0p_1^3)\nu_1 + (1 - p_1)\nu_2 \right\}_\nu \\ &= \left\langle -p_1 + 2p_1^2 + 4p_0^3, [3 + 10p_0p_1^3, 1 - p_1] \right\rangle_Z \\ &= \langle c(t_f, P), G(t_f, P) \rangle_Z. \end{aligned}$$

Since we only consider input constraints and because the controller template  $\bar{u}(P, \beta)$  is linear in  $\beta$ , we have  $\hat{\mathcal{Z}}(\{\bar{u}(\beta, P)\}_\beta) = \langle p_0, p_1 \rangle_Z$  and thus

$$\rho_u^{(\mathcal{Z})}(C\mathcal{U}, [0, t_f], P) - d\mathcal{U} = \rho_u^{(\mathcal{Z})}\left(\left[\begin{array}{c} 1 \\ -1 \end{array}\right], [0, t_f], P\right) - \left[\begin{array}{c} 1 \\ 1 \end{array}\right] = \left[\begin{array}{c} p_0 + |p_1 - 1| \\ -p_0 + |p_1 - 1| \end{array}\right],$$

so that  $g(P)$  in (5.41) reduces to

$$g(P) = |p_0| + |p_1| - 1.$$

Thus, the cost function is given by

$$J(P) = \left| -p_1 + 2p_1^2 + 4p_0^3 \right| + \left| 3 + 10p_0p_1^3 \right| + |1 - p_1| + 100 \max(0, |p_0| + |p_1| - 1),$$

for  $\sigma = 100$ . Fig. 5.9 visualizes the contour lines of  $J(P)$  and its approximation  $\tilde{J}(P)$  over the iterations.

Here, we assume that an approximation of the exact reachable set is available as its quadratic Taylor expansion with the initial trust-region radius  $\gamma = \frac{1}{4}$ ; in practice, the exact reachable set is not available and we thus instead compute an approximation to it (see Sec. 5.4.2). For the initial controller parameter guess  $\bar{P} = 0$ , the generating function of the parameterized, time-point reachable set for  $P \in \mathcal{P}_{\frac{1}{4}}(0) = \left[-\frac{1}{2}, \frac{1}{2}\right]^{1 \times 2}$  (see (5.45))

then is

$$\begin{aligned}\tilde{\mathcal{R}}(t_f, P) &= \left\{ \begin{aligned} & \left( -1 + 2\bar{p}_1^2 + 4\bar{p}_0^3 + (3 + 10\bar{p}_0\bar{p}_1^3)\beta + (2 - 2\bar{p}_1)\beta^2 \right) \Big|_{\bar{P}=0} \\ & + \left( [12\bar{p}_0^2 + 10\bar{p}_1^3\beta, \quad 4\bar{p}_1 + 30\bar{p}_0\bar{p}_1^2\beta - 2\beta^2] (P - \bar{P}) \right) \Big|_{\bar{P}=0} \\ & + \frac{1}{2} \left( (P - \bar{P})^T \begin{bmatrix} 24\bar{p}_0, & 30\bar{p}_1^2\beta \\ 30\bar{p}_1^2\beta, & 4 + 60\bar{p}_0\bar{p}_1\beta \end{bmatrix} (P - \bar{P}) \right) \Big|_{\bar{P}=0} \end{aligned} \right\}_\beta \\ &= \left\{ -1 + 3\beta + 2\beta^2 - 2p_1\beta^2 + 2p_1^2 \right\}_\beta,\end{aligned}$$

and thus

$$\hat{\mathcal{Z}}(\tilde{\mathcal{R}}(t_f, P)) = \left\langle -p_1 + 2p_1^2, [3, \quad 1 - p_1] \right\rangle_Z.$$

The approximated cost function for  $P \in \left[-\frac{1}{2}, \frac{1}{2}\right]^{1 \times 2}$  for the first step thus is

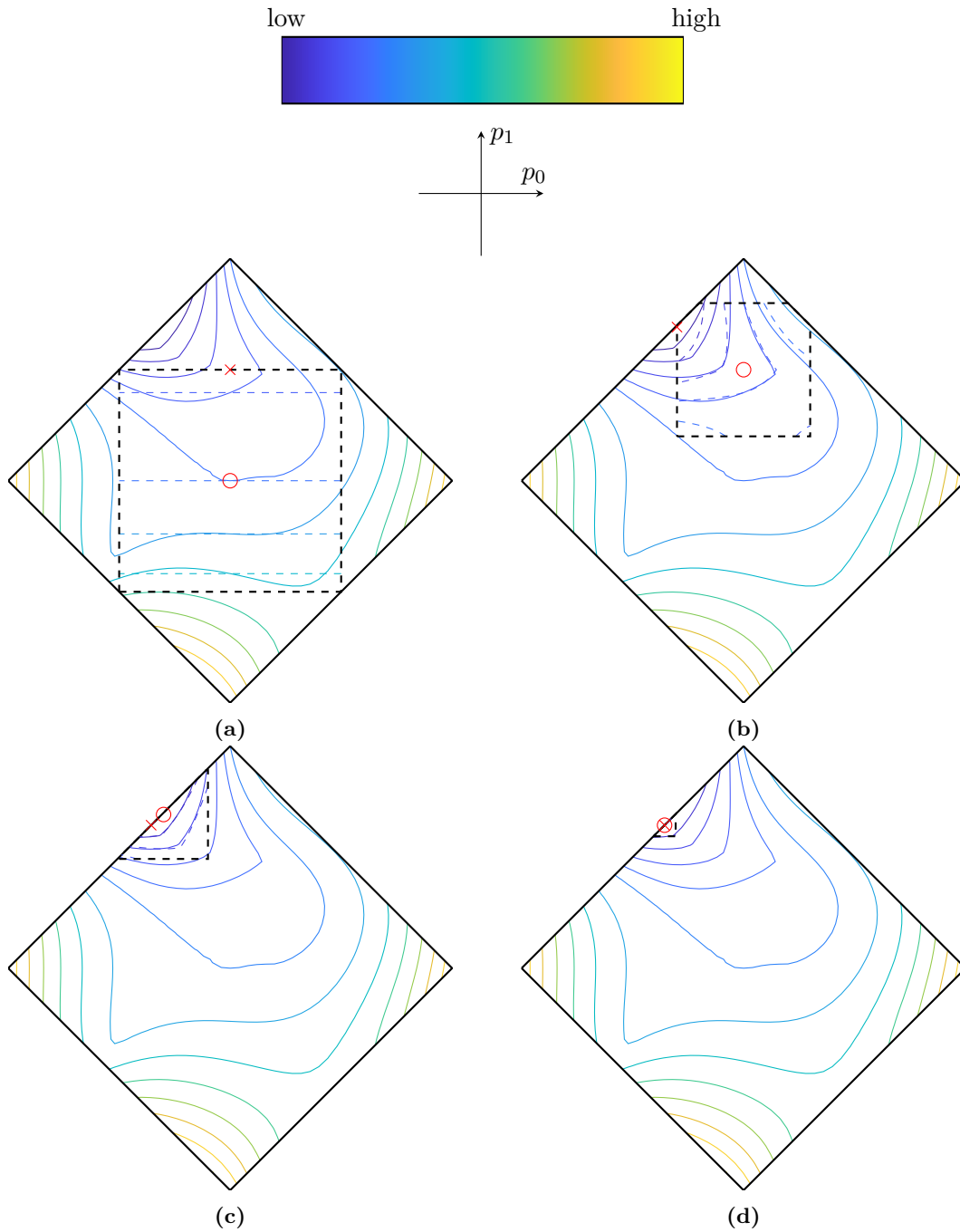
$$\tilde{J}(P) = \left| -p_1 + 2p_1^2 \right| + |1 - p_1| + 3 + 100 \max(0, |p_0| + |p_1| - 1).$$

For the first iteration,  $\tilde{J}(P)$  clearly does not approximate the exact cost function  $J(P)$  particularly well, even though  $J(\hat{P}) < J(\bar{P})$ , where  $\bar{P}$  and  $\hat{P}$  are the initial and optimized controller parameters, respectively (see Fig. 5.9a). Thus, we accept the step but shrink the trust-region radius  $\gamma$  to ensure more accurate approximations in the next steps (see Figs. 5.9b and 5.9c) in order to eventually arrive at the minimum (see Fig. 5.9d). ■

The remainder of this section is structured as follows: We first describe how we approximate (5.1) using a cost function in Sec. 5.4.1 and then derive the computation of the parameterized reachable set in Sec. 5.4.2. To avoid repeated reachable set computations for the cost of the controller and its non-differentiability, we derive the trust-region subproblem – a differentiable optimization problem – as an approximation to this cost in Sec. 5.4.3, whose solution then yields a new controller candidate. To ensure that the trust-region subproblem accurately approximates the cost of the newly computed controller candidate, we describe the tuning of the trust-region radius in Sec. 5.4.4 such that the trust-region subproblem approximates the controller cost arbitrarily closely. We discuss the offline and online complexity of the iPGSC approach in Sec. 5.4.5, demonstrate its applicability using numerical experiments in Sec. 5.4.6, and finally discuss the algorithm in Sec. 5.4.7.

### 5.4.1 Cost of the Controller

Since we generally cannot solve (5.1) efficiently, we instead propose a cost function as a tractable approximation of (5.1). Inspection of (5.1) yields the following observations: We want to minimize the size of the reachable set, centered at the target state, but also minimize the control effort. Further, input, state, and final state constraints must



**Figure 5.9:** Visualization of the iterative controller computation for Ex. 5.4 from [35]. Shown is the set of feasible parameters  $\mathcal{P}$  (outer black diamond), the contour lines of the exact cost  $J(P)$  (solid) and the current approximation of the cost  $\tilde{J}(P)$  (dashed) based on the approximated, parameterized reachable set for  $P \in \mathcal{P}_\gamma(\bar{P})$  (dashed black box), the current initial guess of the controller parameters  $\bar{P}$  (red circle), and the optimizer  $\hat{P} \in \mathcal{P}_\gamma(\bar{P})$  for the current approximated optimization problem (red x).

## 5 Piecewise Constant Controller Synthesis

not be violated. In this section, we thus first transform the objective function and the constraints of (5.1) separately to tractable expressions and then use both to define the cost as an approximation to (5.1). To that end, we first require the following definitions.

Let

$$\begin{aligned}\hat{\mathcal{Z}}(\mathcal{R}(t, P)) &= \langle c(t, P), G(t, P) \rangle_Z, \\ \hat{\mathcal{Z}}\left(\left\{\bar{u}(\beta, P^{(i,j)})\right\}_\beta\right) &= \left\langle c_u(P^{(i,j)}), G_u(P^{(i,j)}) \right\rangle_Z,\end{aligned}$$

where  $\mathcal{R}(t, P)$  is an outer approximation of the closed-loop reachable set for the controller  $\bar{u}(\beta, P^{(i,j)})$  with  $0 \leq i \leq m_s - 1$  and  $0 \leq j \leq m_c - 1$ . We define the number of reachability steps  $m_r \in \mathbb{N}_+$  for the computation of  $\mathcal{R}(t, P)$  in each time interval  $[k, k+1]\delta$  for  $0 \leq k \leq m_s m_c m_r - 1$ . Lastly, we define the support functions over any arbitrary time interval  $[\underline{t}, \bar{t}]$  as

$$\begin{aligned}\rho_x^{(\mathcal{Z})}(\cdot, [\underline{t}, \bar{t}], P) &= \max_{k \in \mathcal{I}_{[\underline{t}, \bar{t}]}^{(m_s m_c m_r)}} \rho_x^{(\mathcal{Z})}\left(\cdot, [k, k+1] \frac{t_c}{m_r}, P\right), \\ \rho_u^{(\mathcal{Z})}(\cdot, [\underline{t}, \bar{t}], P) &= \max_{k \in \mathcal{I}_{[\underline{t}, \bar{t}]}^{(m_s m_c)}} \rho_u^{(\mathcal{Z})}(\cdot, [k, k+1] t_c, P),\end{aligned}$$

with

$$\mathcal{I}_{[\underline{t}, \bar{t}]}^{(n)} = \left\{ k \in \{0, \dots, n-1\} \mid [k, k+1] \frac{t_f}{n} \cap [\underline{t}, \bar{t}] \neq \emptyset \right\},$$

where  $\rho_x^{(\mathcal{Z})}(\cdot, [k, k+1] \frac{t_c}{m_r}, P)$  denotes the support function of  $\hat{\mathcal{Z}}(\mathcal{R}([k, k+1] t_c, P))$ ,  $\rho_u^{(\mathcal{Z})}(\cdot, [k, k+1] t_c, P)$  denotes the support function of  $\hat{\mathcal{Z}}\left(\left\{\bar{u}(\beta, P^{(i,j)})\right\}_\beta\right)$  for  $k = im_c + j$  with  $0 \leq i \leq m_s - 1$  and  $0 \leq j \leq m_c - 1$ , and we have  $0 \leq \underline{t} \leq \bar{t} \leq t_f$  for  $\underline{t} \in \mathbb{R}_{\geq 0}$  and  $\bar{t} \in \mathbb{R}_{\geq 0}$ .

**Objective:** We bound the objective in (5.1a) from above – while considering the extended optimization horizon – by

$$\begin{aligned}& \max_{x(t, P) \in \mathcal{R}(t, P)} \left\{ \sum_{k=1}^h \left( \left\| x(kt_s, P) - x_f^{(k)} \right\|_1 + \zeta \sum_{j=0}^{m_c-1} \left\| u(x(kt_s, P), P^{(k-1, j)}) \right\|_1 \right) \right\} \\ & \leq \sum_{k=1}^h \vartheta_k \left( \left\| \begin{bmatrix} c(kt_s, P) - x_f^{(k)} \\ [G(kt_s, P)]_{(\cdot)} \end{bmatrix} \right\|_1 + \zeta \sum_{j=0}^{m_c-1} \left\| \begin{bmatrix} c_u(P^{(k-1, j)}) \\ [G_u(P^{(k-1, j)})]_{(\cdot)} \end{bmatrix} \right\|_1 \right),\end{aligned}\quad (5.40)$$

where  $\vartheta \in \mathbb{R}_{\geq 0}^h$  weighs the contribution of the extended horizon. The inequality follows from a zonotope outer approximation of the reachable set and the input set as well as multiple applications of the triangle inequality.

**Constraints:** We can conservatively transform the constraints (5.1b) to (5.1d) – again considering the extended optimization horizon – to

$$g(P) = \begin{bmatrix} \rho_u^{(\mathcal{Z})}(C_{\mathcal{U}}, [0, ht_s], P) - d_{\mathcal{U}} \\ \rho_x^{(\mathcal{Z})}(C_{\mathcal{X}}, [0, ht_s], P) - d_{\mathcal{X}} \\ \rho_x^{(\mathcal{Z})}(C_{\mathcal{X}_f}, ht_s, P) - d_{\mathcal{X}_f} \end{bmatrix} \leq 0, \quad (5.41)$$

since  $\mathcal{U} = \langle C_{\mathcal{U}}, d_{\mathcal{U}} \rangle_H$ ,  $\mathcal{X} = \langle C_{\mathcal{X}}, d_{\mathcal{X}} \rangle_H$ , and  $\mathcal{X}_f = \langle C_{\mathcal{X}_f}, d_{\mathcal{X}_f} \rangle_H$  are  $H$ -polytopes and we can hence apply Prop. 4.6.

**Cost:** Using (5.40) and (5.41), we propose the cost function

$$J(P) = \sum_{k=1}^h \vartheta_k \left( \left\| \begin{bmatrix} c(k t_s, P) - x_f^{(k)} \\ [G(k t_s, P)]_{(\cdot)} \end{bmatrix} \right\|_1 + \zeta \sum_{j=0}^{m_c-1} \left\| \begin{bmatrix} c_u(P^{(k-1, j)}) \\ [G_u(P^{(k-15, j)})]_{(\cdot)} \end{bmatrix} \right\|_1 \right) + \sigma \|\max(0, g(P))\|_1, \quad (5.42)$$

where  $\sigma \in \mathbb{R}_+$ : The first term in (5.42) penalizes the deviation of the center from the corresponding target state and the size of the reachable set by penalizing each generator of the reachable set, whereas the second term penalizes the size of the applied input. The third term then penalizes all constraint violations, where  $\sigma$  is an exact penalty multiplier and is chosen such that feasibility is always prioritized over minimizing the first and second term (see [50]).

### 5.4.2 Parameterized Reachable Set

In Sec. 2.7.3, we remarked that it is in principle possible to represent the exact reachable set arbitrarily closely if the abstraction order and the number of steps during reachability analysis are both large enough. For PGSC, we used this fact as a motivation to compute the parameterized reachable set via the abstraction described in Sec. 2.7.3, where we ignored the abstraction error term since it vanishes for a large enough abstraction order and number of steps for reachability analysis. In practice, however, we are often limited to a quadratic abstraction for the computation of the parameterized reachable set due to the increasing computational effort for a large abstraction order and a large number of time steps during reachability analysis. That said, we would still like to compute the parameterized reachable set accurately. Since computing the reachable set exactly is hard in general [68], it is unreasonable to expect that the parameterized reachable set can become accurate with respect to the exact reachable set. In this section, we thus derive an approximation to the parameterized reachable set which is accurate with respect to a tight enough outer approximation of the exact reachable set. For simplicity, we assume without loss of generality that we start from  $t = 0$ , and denote  $P = P^{(0:h-1)}$  with slight abuse of notation.

We want to find an approximation  $\tilde{\mathcal{R}}(t, P)$  that does not need to be recomputed for each  $P$  and is  $\epsilon$ -accurate with respect to a tight outer approximation  $\mathcal{R}(t, P)$  of the

reachable set for  $\epsilon > 0$ , i.e.

$$\forall P \in [-1, 1]^{hm_c n_u \times a} \forall t \in [0, ht_s] : d_H \left( \mathcal{R}(t, P), \tilde{\mathcal{R}}(t, P) \right) \leq \epsilon, \quad (5.43)$$

where  $d_H(\cdot, \cdot)$  denotes the Hausdorff distance between two sets (see Def. 2.16). Because both  $\mathcal{R}(t, P)$  and  $\tilde{\mathcal{R}}(t, P)$  are reachable sets represented by polynomial zonotopes, they are always closed and thus  $d_H \left( \mathcal{R}(t, P), \tilde{\mathcal{R}}(t, P) \right) = 0 \iff \mathcal{R}(t, P) = \mathcal{D}(t, P)$  [103, Chap. 2]. Thus, for a given accuracy  $\epsilon$ , we want to find an approximation  $\tilde{\mathcal{R}}(t, P)$  whose Hausdorff distance to the tight outer approximation  $\mathcal{R}(t, P)$  is at most  $\epsilon$ .

In Ex. 5.3, we saw that the dependency-preserving part  $\mathcal{D}(t, P)$  (see (5.23)) used for PGSC often is a reasonable approximation of the parameterized reachable set up to a scaling factor. When only minimizing the size of the final reachable set, a rough approximation of the parameterized reachable set is often sufficient. However, once we also consider state constraints in the synthesis problem, the absolute size of the approximated parameterized reachable set becomes highly relevant. In PGSC and GSC, this was solved by introducing the adapted state constraints  $\tilde{\mathcal{X}}$ , which have to be provided by the user.

We avoid these adapted constraints by noticing that for the current controller parameters  $\bar{P}$  (see Alg. 4), a tight outer approximation  $\mathcal{R}(t, \bar{P})$  of the reachable set at  $\bar{P}$  for all  $t \in [0, ht_s]$  is available. Thus, we construct an approximation to the parameterized reachable set as

$$\tilde{\mathcal{R}}(t, P) = \mathcal{R}(t, \bar{P}) \oplus_e \mathcal{D}(t, P) \oplus_e \left( -\mathcal{D}(t, \bar{P}) \right), \quad (5.44)$$

where  $\oplus_e$  denotes the exact addition of polynomial zonotopes (see Sec. 2.6.6): Since  $\mathcal{D}(t, \bar{P}) \oplus_e \left( -\mathcal{D}(t, \bar{P}) \right) = 0$  by definition, we recover  $\mathcal{R}(t, \bar{P})$  by evaluating (5.44) for the current controller parameters  $\bar{P}$ . For any  $P \neq \bar{P}$ , we then use  $\mathcal{D}(t, P) \oplus_e \left( -\mathcal{D}(t, \bar{P}) \right)$  to approximate the unknown difference  $\mathcal{R}(t, P) \oplus_e \left( -\mathcal{R}(t, \bar{P}) \right)$ . In order to control the approximation quality of (5.44), we limit the controller parameters to the trust region  $P \in \mathcal{P}_\gamma(\bar{P})$ , where

$$\mathcal{P}_\gamma(\bar{P}) = \left\{ \bar{P} + 2\gamma M \mid M \in [-1, 1]^{hm_c n_u \times a} \right\}. \quad (5.45)$$

Here,  $\gamma \in (0, 1]$  is the trust-region radius that restricts the range of possible  $P$  around the controller parameters  $\bar{P}$ . To keep  $\gamma$  within  $(0, 1]$  for consistency with later chapters, we include the factor 2 in (5.45) so that  $\forall \bar{P} \in [-1, 1]^{hm_c n_u \times a} : [-1, 1]^{hm_c n_u \times a} \subseteq \mathcal{P}_1(\bar{P})$ , i.e.,  $\gamma = 1$  means that

$$\begin{aligned} \forall i \in \{0, \dots, h-1\} \forall j \in \{0, \dots, m_c-1\} \forall u \in \mathcal{U} \exists P \in \mathcal{P}_1(\bar{P}) \exists \beta \in [-1, 1]^l : \\ \bar{u}(\beta, P^{(i,j)}) = u, \end{aligned}$$

where  $\beta$  is the dependent factor of the generating function  $x^{(0)}(\beta)$  of the initial set  $\mathcal{X}^{(0)}$ . For  $P \in \mathcal{P}_\gamma(\bar{P})$ , it thus holds that

$$\lim_{\gamma \rightarrow 0} \tilde{\mathcal{R}}(t, P) = \mathcal{R}(t, \bar{P}),$$

and therefore we can ensure (5.43) for  $P \in \mathcal{P}_\gamma(\bar{P})$ , i.e.

$$\exists \gamma \in (0, 1] \forall P \in \mathcal{P}_\gamma(\bar{P}) \forall t \in [0, ht_s] : d_H(\mathcal{R}(t, P), \tilde{\mathcal{R}}(t, P)) \leq \epsilon, \quad (5.46)$$

for any  $\epsilon > 0$  and independent of the choice of  $\bar{P}$ . We defer the tuning of  $\gamma$  to achieve (5.46) to Sec. 5.4.4.

### 5.4.3 Trust-Region Subproblem

We now propose the trust-region subproblem as an approximation to (5.42). With a slight abuse of notation, we again set  $P = P^{(0:h-1)}$  and propose the trust-region subproblem

$$\hat{P}, \hat{s} = \arg \min_{P, s} \tilde{J}_{\text{TR}}(P, s), \quad (5.47a)$$

$$\text{s.t. } \tilde{g}(P) \leq s, \quad (5.47b)$$

$$P \in \mathcal{P}_\gamma(\bar{P}) \cap [-1, 1]^{hm_c n_u \times a}, \quad (5.47c)$$

$$s \geq 0, \quad (5.47d)$$

with

$$\tilde{J}_{\text{TR}}(P, s) = \sum_{k=1}^h \vartheta_k \left( \left\| \begin{bmatrix} \tilde{c}(kt_s, P) - x_f^{(k)} \\ \tilde{G}(kt_s, P) \end{bmatrix}_{(\cdot)} \right\|_1 + \zeta \sum_{j=0}^{m_c-1} \left\| \begin{bmatrix} c_u(P^{(k-1,j)}) \\ G_u(P^{(k-1,j)}) \end{bmatrix}_{(\cdot)} \right\|_1 \right) + \sigma \|s\|_1, \quad (5.48)$$

as an approximation to the cost in (5.42), where  $\tilde{g}(P)$  is defined as in (5.41) but with  $\rho_x^{(\mathcal{Z})}(\cdot, t, P)$  replaced by the support function  $\tilde{\rho}_x^{(\mathcal{Z})}(\cdot, t, P)$  of the zonotope outer approximation  $\tilde{\mathcal{Z}}(t, P) = \hat{\mathcal{Z}}(\tilde{\mathcal{R}}(t, P))$ :

- (5.47a) and (5.48): We avoid the need for an outer approximation of the reachable set in (5.42) by replacing it with the approximated parameterized reachable set  $\tilde{\mathcal{Z}}(t, P)$  (see Sec. 5.4.2). Further, we replace the constraint violation in (5.42) with the slack variable  $s$  (see next bullet point for more details).
- (5.47b): We replace the support function  $\rho_x^{(\mathcal{Z})}$  of the zonotope outer approximation with the support function  $\tilde{\rho}_x^{(\mathcal{Z})}$  of  $\hat{\mathcal{Z}}(\tilde{\mathcal{R}}(t, P))$ , i.e., the zonotope outer approximation of the parameterized reachable set. Further, we include input, state, and final state constraints in (5.47b) but relax them with  $s \in \mathbb{R}_{\geq 0}^{o_u + o_x + o_{x_f}}$  so that (5.47)

always has a feasible solution even if  $\nexists P \in \mathcal{P}_\gamma(\bar{P}) \forall k \in \{1, \dots, o_U + o_X + o_{X_f}\} : \tilde{g}_k(P) \leq 0$  holds true, i.e., no feasible solution of (5.47) for  $s = 0$  exists. We force a reduction in constraint violation for infeasible  $P$  by penalizing  $s$  in the objective function, where the exact penalty multiplier  $\sigma \gg 1$  is chosen large enough (see, e.g., [50] for a discussion on penalty multipliers) so that the minimization of constraint violation is always prioritized over minimizing the objective function for  $s = 0$ .

- (5.47c): To control the accuracy of the approximated, parameterized reachable set, we restrict  $P$  to its trust region, i.e.  $P \in \mathcal{P}_\gamma(\bar{P}) \cap [-1, 1]^{hm_c n_u \times a}$  (see Sec. 5.4.2 for details).

The trust-region subproblem in (5.47) can be reformulated as a smooth optimization problem (see Sec. 3.2.2) and thus allows us to use off-the-shelf optimization codes, such as the interior point optimizer<sup>5</sup> [111] (IPOPT). Further, regularity of (5.47) follows from Prop. 3.2. As it turns out, we can express the objective value (5.48) at a critical point in  $\hat{P}$  only, which allows the definition of the approximated cost  $\tilde{J}(P)$ .

**Proposition 5.2** (Stationary Objective Value). *At a critical point of the smooth reformulation of (5.47) with optimizer  $\hat{P}$ , the optimal objective value – only dependent on  $\hat{P}$  – is given by*

$$\tilde{J}(\hat{P}) = \tilde{J}_{\text{TR}}(\hat{P}, \max(0, \tilde{g}(\hat{P}))). \quad (5.49)$$

*Proof.* Follows from Prop. 3.1. □

**Remark:** Because  $\tilde{\mathcal{R}}(t, \bar{P}) = \mathcal{R}(t, \bar{P})$  (see (5.44)), the cost  $J(P)$  from (5.42) and the approximated cost  $\tilde{J}(P)$  as defined in (5.49) are equal at  $P = \bar{P}$ , i.e.,  $\tilde{J}(\bar{P}) = J(\bar{P})$ . Thus, choosing  $\bar{s} = \max(0, \tilde{g}(\bar{P}))$  as in Prop. 5.2, the initial guess  $(\bar{P}, \bar{s})$  for solving the smooth reformulation of (5.47) ensures that we start optimizing from the actual controller cost.

#### 5.4.4 Tuning of the Trust-Region Radius

Previously, we derived an approximation to the controller cost  $J(\hat{P})$  in (5.42), where  $\hat{P}$  is a critical point of (5.47). In this section, we show how the value of the trust-region radius can be adapted after each iteration of Alg. 4 so that the objective value  $\tilde{J}(\hat{P})$  of (5.47) at a critical point  $\hat{P}$  approximates  $J(\hat{P})$  arbitrarily closely.

After the trust-region subproblem in (5.47) is solved to a critical point  $\hat{P}$ , we have both the approximated cost  $\tilde{J}(\hat{P})$  (see Prop. 5.2) and the cost  $J(\hat{P})$ , calculated from the outer approximation of the closed-loop reachable set, available (see Alg. 4). This inspires the following tuning rule.

<sup>5</sup><https://github.com/coin-or/Ipopt>



**Theorem 5.1** (Tuning of Trust-Region Radius). *Let  $J(\hat{P})$  and  $\tilde{J}(\hat{P})$  be defined as in (5.42) and (5.49), respectively, where  $\hat{P}$  is a critical point of the smooth reformulation of the trust-region subproblem (5.47). If we adapt the trust-region radius  $\gamma$  according to*

$$\gamma \leftarrow \min \left( 1, v \left( \left| \tilde{J}(\hat{P}) - J(\hat{P}) \right| \right) \gamma \right), \quad (5.50)$$

where  $v : [0; \infty) \mapsto \mathbb{R}_+$  is an arbitrary, monotonically decreasing function with  $v(0) = \bar{c}$ ,  $v(\psi) = 1$ , and  $\forall r \geq \psi : v(r) = \underline{c}$ , where  $\frac{1}{\underline{c}} > \bar{c} > 1$  and  $0 \leq \psi < \epsilon$ , then

$$\left| \tilde{J}(\hat{P}) - J(\hat{P}) \right| \leq \epsilon,$$

is achieved after a finite number of iterations for  $\epsilon > 0$ .

*Proof.* Comparing (5.42) with (5.49), it is clear that  $\lim_{\gamma \rightarrow 0} \left| \tilde{J}(\hat{P}) - J(\hat{P}) \right| = 0$  since  $\lim_{\gamma \rightarrow 0} d_H \left( \hat{\mathcal{Z}} \left( \mathcal{R}(t, \hat{P}) \right), \hat{\mathcal{Z}} \left( \tilde{\mathcal{R}}(t, \hat{P}) \right) \right) = 0$  due to  $\lim_{\gamma \rightarrow 0} \hat{P} = \bar{P}$  and  $\tilde{\mathcal{R}}(t, \bar{P}) = \mathcal{R}(t, \bar{P})$ . Therefore, following (5.50) for a finite number of steps yields a trust-region radius  $\gamma > 0$  such that  $\left| \tilde{J}(\hat{P}) - J(\hat{P}) \right| \leq \epsilon$ , which concludes the proof.  $\square$

## 5.4.5 Computational Complexity

We again differentiate between the complexity of the offline controller synthesis and the complexity of applying said controller online. We again make use of Ass. 5.1. For simplicity, let  $\kappa \geq 1$  for the order of the controller template introduced in Sec. 5.3.2 and assume that the number of monomials  $a$  of the controller template grows at least with  $O(n)$ , i.e.,  $O(n + a) = O(a)$ .

### 5.4.5.1 Offline Complexity

We first inspect the complexity of the different components of the iPGSC approach. Since  $m_s$  is fixed, it suffices to inspect the algorithm for  $m_s = 1$ .

**Controller Template:** The complexity of the polynomial composition dominates with  $O(an^{1+\kappa})$  (see Sec. 5.3.6.1).

**Reachable Set Computations:** The combined complexity for the closed-loop reachable set computation and the parameterized reachable set computation is bounded from above by  $O(n^5 + n^2a^2 + n^3a \log n)$  (see Sec. 5.3.6.1).

**Controller Computation:** Since the maximum number of steps  $l_{\max} \in \mathbb{N}_+$  in Alg. 4 is fixed, it suffices to check the complexity of one trust-region iteration. As we already considered the complexity for the reachable set computations above and because the trust-region tuning has lower (or equal) computational complexity by visual inspection, it suffices to derive the solution complexity of the trust-region subproblem in (5.47).

Computing the optimal controller parameters  $\hat{P}$  requires the solution of the smooth reformulation of the non-convex optimization problem in (5.47). Following the same argument as for the complexity of the controller computation in Sec. 5.3.6.1, the computational complexity of solving the trust-region subproblem is at most  $O(a^2n^4)$ .

**Overall Complexity:** The offline complexity of the iPGSC approach is identical to the PGSC approach, i.e.

$$O\left(c_{\text{iPGSC}}^{(\text{off})}(n)\right) = O\left(c_{\text{PGSC}}^{(\text{off})}(n)\right) \stackrel{(5.33)}{=} O\left(c_{\text{ref}}(n) + an^{1+\kappa} + a^2n^4\right). \quad (5.51)$$

#### 5.4.5.2 Online Complexity

The online complexity of iPGSC is identical to that of PGSC, i.e.

$$O\left(c_{\text{iPGSC}}^{(\text{on})}(n)\right) = O\left(c_{\text{PGSC}}^{(\text{on})}(n)\right) \stackrel{(5.34)}{=} O(an). \quad (5.52)$$

### 5.4.6 Experiments

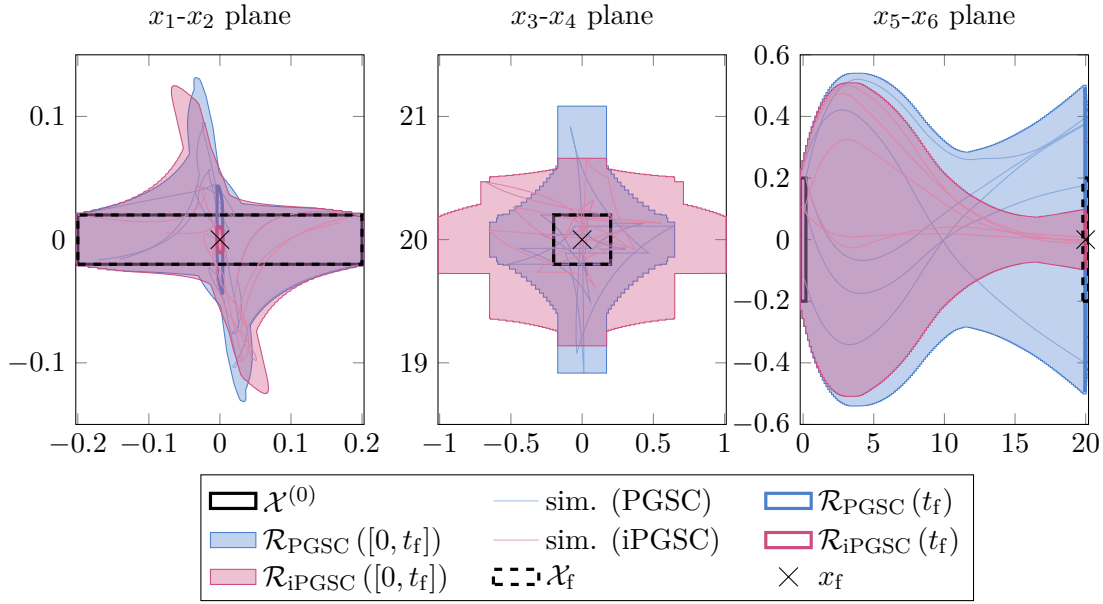
In this section, we demonstrate the applicability of the proposed iPGSC approach by comparing it against the PGSC approach introduced in Sec. 5.3. For that comparison, we revisit both the controlled Van-der-Pol oscillator and the bicycle benchmark described in Sec. 5.3.8. We set  $\zeta = 0$  for all benchmarks in this section.

#### 5.4.6.1 Bicycle (Single Track)

We consider the bicycle model with the dynamics given in (5.38) and parameters given as in Tab. 5.1. Additionally, we explicitly enforce the final state constraint  $\mathcal{X}_f = \{-c_{\mathcal{X}^{(0)}} + x_f\} \oplus \mathcal{X}^{(0)}$  for the iPGSC approach now (also see Sec. 5.3.9 for why we only enforce it here).

Fig. 5.10 compares the performance of PGSC against iPGSC, where runtimes were 34.7s and 196s, respectively. The deviation and size of the final reachable sets is  $\text{csize}(\mathcal{R}_{\text{PGSC}}(t_f)) = 0.65$  and  $\text{csize}(\mathcal{R}_{\text{iPGSC}}(t_f)) = 0.17$ . Clearly, the PGSC approach is outperformed by iPGSC even with  $\mathcal{X}_f$  enforced, albeit at the cost of significantly increased computation times, which are mainly due to the more frequent computation of reachable set outer approximations.

In Fig. 5.11, the progress of the iPGSC trust-region algorithm is visualized over the iterations. First, as remarked at the end of Sec. 5.4.3, the approximated cost equals the formally correct controller cost at the current controller guess  $\bar{P}^{(i-1)}$ . Second, we see how the trust-region radius  $\gamma$  is decreased as the iterations at  $i \in \{1, 3\}$  are rejected to recover a “good enough” approximation of the parameterized reachable set for the next iteration.



**Figure 5.10:** Comparison of the PGSC approach and the iPGSC approach for the single-track (bicycle) model. A linear controller template is used.

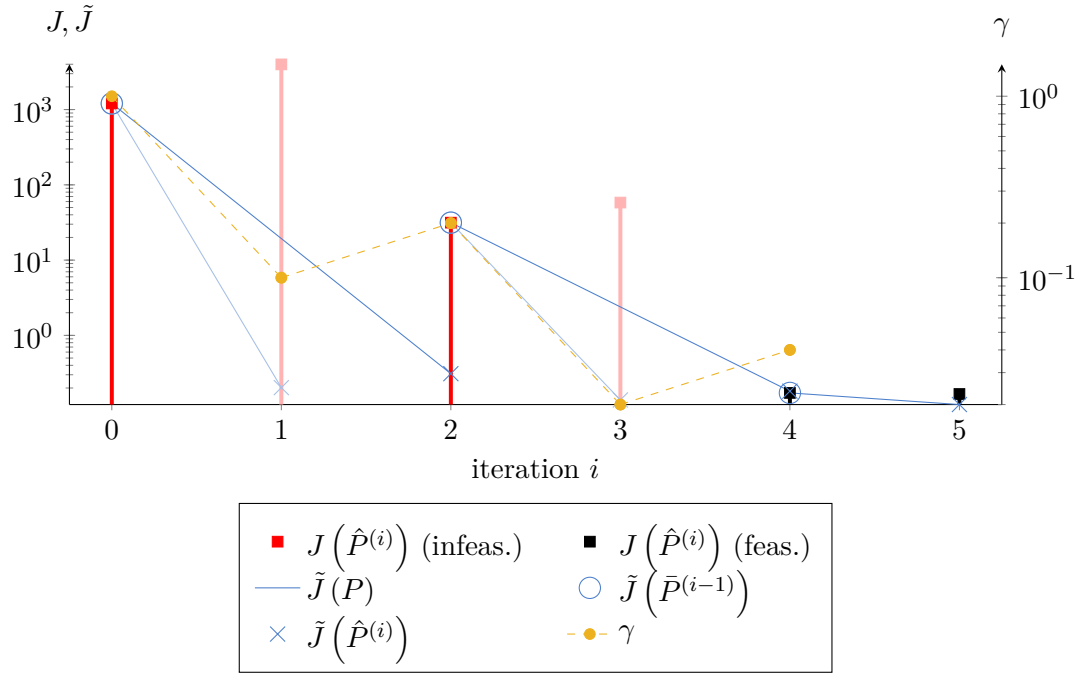
#### 5.4.6.2 Controlled Van-der-Pol Oscillator

We consider the Van-der-Pol oscillator model described in (5.39) with benchmark parameters given as in Tab. 5.2 but with a quadratic controller, i.e.,  $\kappa = 2$ .

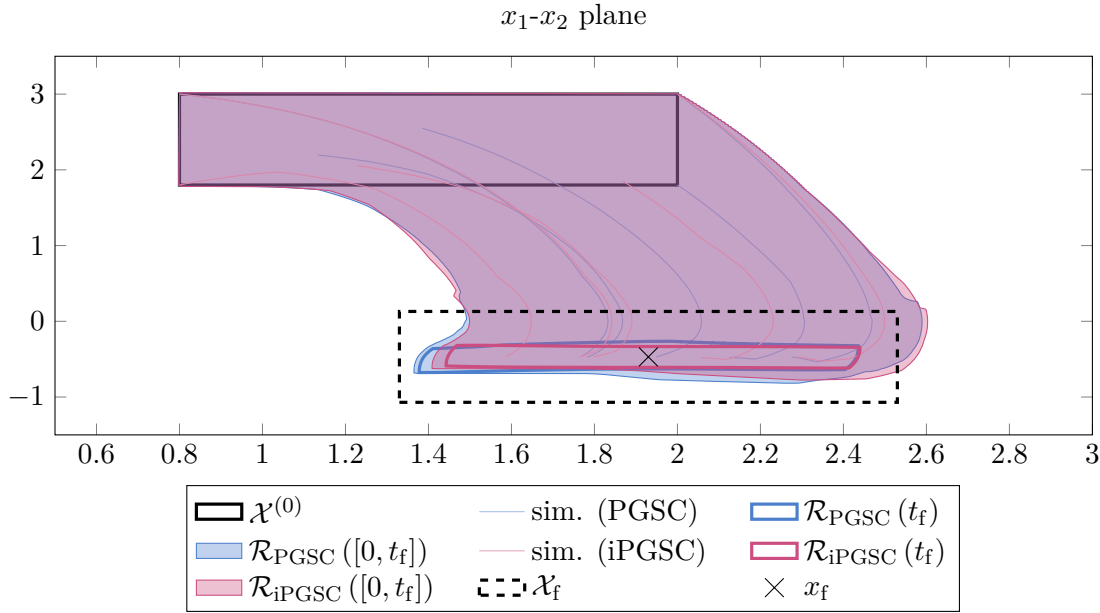
Fig. 5.12 compares the control performance of PGSC against the proposed iPGSC approach, where PGSC and iPGSC took 19.68s and 146.66s, respectively, and further  $\text{csize}(\mathcal{R}_{\text{PGSC}}(t_f)) = 0.79$  and  $\text{csize}(\mathcal{R}_{\text{iPGSC}}(t_f)) = 0.67$ . Clearly, iPGSC achieves a better control result, even though we additionally enforce final state constraints  $\mathcal{X}_f = \{-c_{\mathcal{X}^{(0)}} + x_f\} \oplus \mathcal{X}^{(0)}$  for the iPGSC approach: By approximating the parameterized reachable set only locally around the current controller parameters, the iPGSC uses a more accurate approximation of the closed-loop reachable set during the optimization of the controller parameters. In contrast, the PGSC approach computes its parameterized reachable set only as a global approximation of the closed-loop reachable set.

#### 5.4.7 Discussion

As the numerical experiments indicate, iPGSC synthesizes controllers with better performance compared to the PGSC approach (and thus the GSC approach by extension). Further, iPGSC uses state constraints directly as provided by the user without the need for manual constraint tightening or bloating since we control the accuracy of the parameterized reachable set automatically using the trust-region radius. Furthermore, its design naturally allows to improve existing controllers by providing them as an initial guess to the algorithm. Because we only accept iterations that improve the control per-



**Figure 5.11:** Visualization of the solver progress for the PGSC approach (bicycle benchmark). Values between  $\tilde{J}(\bar{P}^{(i-1)})$  and  $\tilde{J}(\hat{P}^{(i)})$  for  $1 \leq i \leq 5$  – where  $\bar{P}^{(i-1)}$  and  $\hat{P}^{(i)}$  are the initial and critical points of the trust-region subproblem – are linearly interpolated ( $\hat{P}^{(0)} = \bar{P}^{(0)}$ ). The opaque red vertical bars at  $i \in \{1, 3\}$  visualize a rejected iteration. As one might expect, the trust-region radius  $\gamma$  is shrunk at each rejected iteration to ensure a more accurate approximation for the next step.



**Figure 5.12:** Comparison of the PGSC approach and the iPGSC approach for the Van-der-Pol model. Both approaches use a quadratic controller template.

formance, the iPGSC algorithm can guarantee that the control performance of the newly computed controller is at least as good as the one provided as an initial guess.

That said, all these advantages come at the cost of increased computational effort, which is mainly caused by the need to compute outer approximations of the reachable sets. With that in mind, the increased computational effort for iPGSC only occurs offline and still has polynomial complexity (see Sec. 5.4.5.1) so that it should always be the first piecewise constant synthesis algorithm to try on any system if very fast offline computation times are not required.

## 5.5 Summary

Starting from the problem statement in Sec. 5.1, we reviewed existing approaches for the piecewise constant controller synthesis for disturbed nonlinear systems in Sec. 5.2 and extended existing work to include polynomial controllers and allow parameterized reachable sets of arbitrary abstraction order in Sec. 5.3 to improve the performance of the synthesized controllers. Furthermore, we introduced a novel trust-region algorithm for piecewise constant controller synthesis in Sec. 5.4 that can directly consider state constraints without the need for manual tuning.

Generator-space control (GSC) introduced in [99] synthesizes piecewise constant feed-forward controllers that are linear in the initial state under input and state constraints. With the initial set given as a zonotope, a zonotope is used to represent the controller, where the center and generator matrix are parameters that need to be computed. An

approximation to the closed-loop reachable set using that controller is computed by linearizing the undisturbed nonlinear dynamics along a previously computed reference trajectory. Because of this linearization, we obtained an analytic expression for the parameterized reachable set in the controller parameters. Using this parameterized reachable set, a linear program (LP) is formulated which minimizes the distance of this approximated reachable set from the target state and further minimizes the size of the final, approximated reachable set as well as the control energy. Due to the parameterization of the controller, input constraints can be exactly represented as linear constraints. Since the parameterized reachable set is only approximately computed, state constraints can only be considered approximately and require additional user input. Sub-optimal choices of this input may require additional runs of the algorithm.

To better approximate the reachable set, we described the polynomial generator-space control (PGSC) approach from our work in [37]: Here, we computed the parameterized reachable set using a slightly modified version of reachability analysis from [5] with a polynomial abstraction order to achieve better results. In contrast to the GSC approach [99] where only linear control laws are synthesized, we proposed a controller which is polynomial in the initial state and showed that PGSC generalizes the GSC approach. Numerical experiments demonstrated the improved control performance due to the increased abstraction order for the parameterized reachable set and the polynomial controller compared to the GSC approach from previous work [99].

To avoid the need for the manual tuning of state constraints by the user as required for GSC and PGSC, we introduced the iterative polynomial generator-space control (iPGSC) approach as part of our work in [35] to solve the synthesis problem. To that end, we first defined a cost function in the controller parameters as an approximation to the general synthesis problem, penalizing the deviation from the target state, the size of the reachable set, the input energy, and the constraint violations. Since this cost function is non-differentiable and because its evaluation requires the computation of the closed-loop reachable set, we optimized over this cost iteratively: We approximate the non-differentiable cost in the neighborhood of the current controller parameters by the differentiable trust-region subproblem, whose accuracy is ensured by restricting its domain to a bounded trust region around the current controller parameters. The size of the trust region is then updated by comparing the change in cost between the newly synthesized controller and the current controller to the change predicted by the trust-region subproblem. If the new controller improves the controller cost, we accept the step and update the current controller parameters accordingly. The advantages of this novel synthesis algorithm were demonstrated in numerical experiments.

## 6 Piecewise Constant Controller Synthesis with Continuous State Feedback

Previously, we considered the synthesis of controllers which are piecewise constant in time. Due to their simple structure, they are easy to implement and can be applied efficiently online. Additionally, they also allow the incorporation of state feedback at discrete points in time. In practice, however, it is often necessary to shrink the time between input updates, i.e., increase the number of piecewise constant controllers that need to be computed, e.g., when a system is dominated by disturbances. Thus, larger disturbance sets may require a larger number of piecewise constant controllers, increasing the computational effort required, both offline and online. In this chapter, we combine piecewise constant control with continuous state feedback in order to continuously counteract disturbances without the need for a larger number of piecewise constant controllers.

To that end, we first introduce the synthesis problem for this combined controller in Sec. 6.1 and present a solution approach from previous work [101],[98, Sec. 3.6] in Sec. 6.2. We then propose a novel synthesis approach for the combined controller in Sec. 6.3 which – for the first time – synthesizes the feedforward and feedback controller simultaneously.

### 6.1 Problem Statement

In this chapter, we synthesize controllers for  $t \in \tau^{(j)} = [j, j + 1] \frac{t_f}{m_c}$  with  $0 \leq j \leq m_c - 1$  of the form (as proposed in [98, Eq. (3.43)])

$$u(t, x(t), P, K) = u_{\text{ff}}(x(0), P^{(0,j)}) + K^{(j)}(x(t, P, K) - x_{\text{ff}}(t, P)), \quad (6.1)$$

where  $x(0) \in \mathcal{X}^{(0)}$ ,  $K(t) = K^{(j)} \in \mathbb{R}^{n_u \times n_x}$  for  $t \in \tau^{(j)}$ , the controller parameters of the piecewise constant feedforward controller  $u_{\text{ff}}(x(0), P^{(0,j)})$  (see Chap. 5) are given by  $P = [P^{(0,0)T}, P^{(0,1)T}, \dots, P^{(0,m_c-1)T}] \in [-1, 1]^{m_c n_u \times a}$ , and

$$x_{\text{ff}}(t, P) = \xi(t, x(0), u_{\text{ff}}(x(0), P^{(0,\cdot)}), 0), \quad (6.2)$$

is the feedforward trajectory that one obtains by starting from the initial state  $x(0)$  and applying the  $m_c$  piecewise constant feedforward controllers: Piecewise constant controllers are able to solve a given reach-avoid problem efficiently when little or no disturbance is involved (see Chap. 5). Thus, to steer the undisturbed system close to the target state, the feedforward controller is applied. Since the actual system is affected by

## 6 Piecewise Constant Controller Synthesis with Continuous State Feedback

disturbances, the continuous feedback controller is applied on the difference between the nominal feedforward state and the actual state – caused by disturbances – to minimize this difference. We sometimes omit time indices and dependencies for readability, e.g., we write  $K$  instead of  $K(t)$ . The feedforward part is a piecewise constant controller as described in Chap. 5 with  $m_s = 1$  and therefore  $P = P^{(0)}$ . To simplify notation, we further set  $P^{(j)} = P^{(0,j)}$  for the remainder of this chapter. We again assume that the bounded input constraints  $\mathcal{U} = \langle C_U, d_U \rangle_H$  with  $C_U \in \mathbb{R}^{o_U \times n_U}$  and  $d_U \in \mathbb{R}^{o_U}$ , the state constraints  $\mathcal{X} = \langle C_X, d_X \rangle_H$  with  $C_X \in \mathbb{R}^{o_X \times n_X}$  and  $d_X \in \mathbb{R}^{o_X}$ , and the final state constraints  $\mathcal{X}_f = \langle C_{X_f}, d_{X_f} \rangle_H$  with  $C_{X_f} \in \mathbb{R}^{o_{X_f} \times n_X}$  and  $d_{X_f} \in \mathbb{R}^{o_{X_f}}$  are given as H-polytopes.

Let  $\mathcal{R}(t, P, K)$  denote an outer approximation of the closed-loop reachable set of the combined controller  $u(t, x(t), P, K)$  from (6.1). Substitution of (6.1) into the general synthesis problem in (3.1) then yields an optimization problem for the combined synthesis of a piecewise constant feedforward controller with continuous state feedback, i.e. (also compare to )

$$\hat{P}, \hat{K} = \arg \min_{P, K} \max_{x(t, P, K) \in \mathcal{R}(t, P, K)} \left\{ \|x(t_f, P, K) - x_f\|_1 + \zeta \int_0^{t_f} \|u(\tau, x(\tau), P, K)\|_1 d\tau \right\} \quad (6.3a)$$

$$\text{s.t. } \forall t \in [0, t_f] : \mathcal{S}_u(t, P, K) \subseteq \mathcal{U}, \quad (6.3b)$$

$$\mathcal{R}([0, t_f], P, K) \subseteq \mathcal{X}, \quad (6.3c)$$

$$\mathcal{R}(t_f, P, K) \subseteq \mathcal{X}_f, \quad (6.3d)$$

where

$$\mathcal{S}_u(t, P, K) = \mathcal{S}_{u_{\text{ff}}}(t, P) \oplus_e \mathcal{S}_{\Delta u}(t, P, K), \quad (6.4)$$

follows from (6.1) with

$$\mathcal{S}_{u_{\text{ff}}}(t, P) = \left\{ u_{\text{ff}}(x(0), P^{(0,j)}) \mid x(0) \in \mathcal{X}^{(0)} \right\}, \quad (6.5)$$

$$\mathcal{S}_{\Delta u}(t, P, K) = K^{(j)}(\mathcal{R}(t, P, K) \oplus_e (-\mathcal{R}_{x_{\text{ff}}}(t, P))), \quad (6.6)$$

for  $t \in \tau^{(j)}$  with  $0 \leq j \leq m_c - 1$ . Here,  $\mathcal{R}_{x_{\text{ff}}}(t, P)$  denotes the reachable set of the feedforward state from (6.2).



## 6.2 Reachset Optimal Control

We already established that continuous state feedback is beneficial to quickly counteract disturbances for systems where the disturbances are large or the dynamics of the system are changing quickly. That said, the control goal is not only rejecting disturbances as efficiently as possible but also steering the system close to the target state. Therefore, the authors in [101] combine the GSC feedforward synthesis approach with a subsequent optimization over the feedback matrices in order to minimize the effect of disturbances, which we review in this section. The general structure of the reachset optimal control (ROC) approach is detailed in Alg. 5 and visualized in Fig. 6.1. The most important steps can be summarized as follows:

- (i) We start by computing a reference trajectory steering the center of the initial set close to the target state (l. 2).
- (ii) Since the feedforward and the feedback synthesis problems are solved separately, we compute the feedforward controller using GSC with adapted state and input constraints (l. 3).
- (iii) Next, we compute an initial guess of the feedback matrices for the feedback synthesis (l. 4).
- (iv) Given the feedforward controller from (ii), we then optimize the feedback matrices to obtain a feedback control law that minimizes the deviation of the disturbed, closed-loop state from the undisturbed feedforward state. However, we do not directly optimize over the feedback matrices; rather, we parameterize all feedback matrices using LQR control, i.e., the feedback matrices become functions of the LQR weighting matrices  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $R \in \mathbb{S}_{++}^{n_u \times n_u}$  (also see Sec. 6.2.1); here,  $Q$  and  $R$  are kept constant over all time steps. This reduces the feedback synthesis to an optimization over the weighting matrices  $Q$  and  $R$ . In each optimization iteration, we compute a formally verified outer approximation to the closed-loop reachable set for the combined controller (see Sec. 6.2.2) and then use this reachable set to evaluate the objective function and check constraint satisfaction (l. 5).

---

### Algorithm 5 Reachset Optimal Control

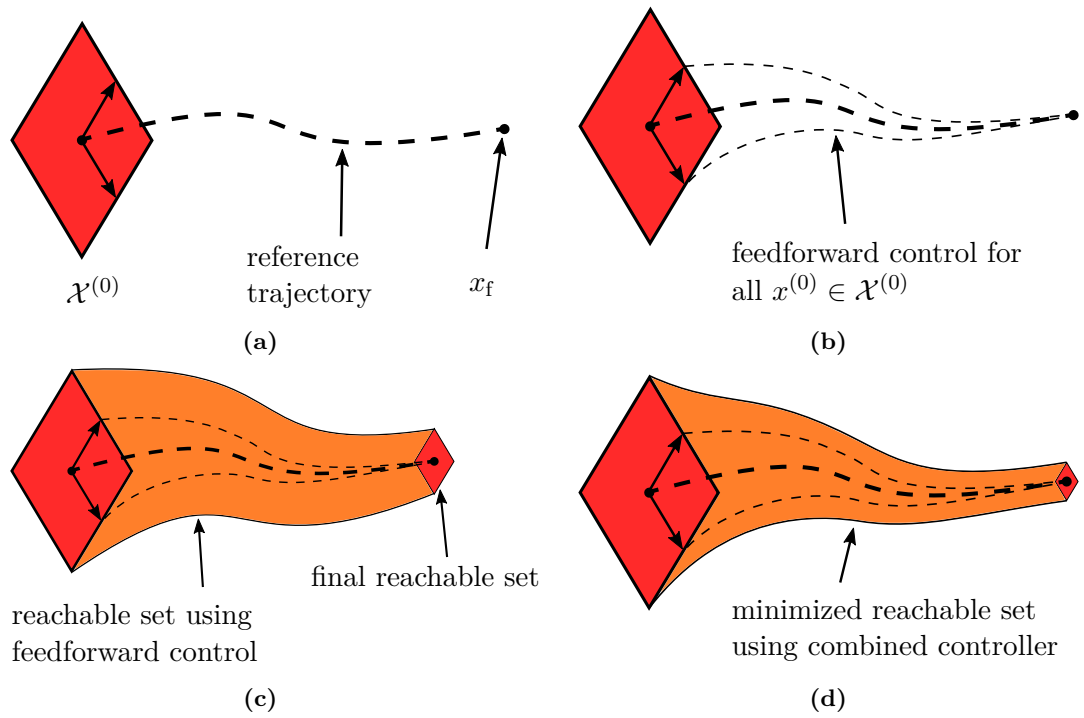
---

```

1: function REACHSETOPTIMALCONTROL( $\mathcal{X}^{(0)}, \mathcal{U}, \tilde{\mathcal{U}}, \mathcal{X}, \tilde{\mathcal{X}}, \mathcal{X}_f, \tilde{\mathcal{X}}_f, \mathcal{W}, m_c, x_f$ )
2:    $[x_{\text{ref}}, u_{\text{ref}}] = \text{REFERENCETRAJECTORY}(\mathcal{X}^{(0)}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, 1, m_c, x_f, t_f) \quad \triangleright \text{Sec. 5.2.1}$ 
3:    $\hat{P} = \text{GENERATORSPACECONTROL}(\mathcal{X}^{(0)}, \tilde{\mathcal{U}}, \tilde{\mathcal{X}}, \tilde{\mathcal{X}}_f, \{0\}, 1, m_c, x_f, 1) \quad \triangleright \text{Sec. 6.2.3}$ 
4:    $[\bar{Q}, \bar{R}] = \text{INITGUESS}(\dots) \quad \triangleright \text{Sec. 6.2.5}$ 
5:    $[\hat{K}, \mathcal{R}([0, t_f])] = \text{OPTPROB}(\hat{P}, \bar{Q}, \bar{R}, \mathcal{X}^{(0)}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, \mathcal{W}, m_c, t_f, x_{\text{ref}}, u_{\text{ref}}) \quad \triangleright (6.16)$ 
6:   return  $[\hat{P}, \hat{K}, \mathcal{R}([0, t_f])]$ 
7: end function

```

---



**Figure 6.1:** Overview of the ROC approach [101, adapted]. First, a reference trajectory is computed (Fig. 6.1a). For the adapted input and state constraints, the feedforward synthesis problem for the target state  $x_f$  is solved using the GSC approach as described in Sec. 5.2 (Fig. 6.1b). Since the feedforward synthesis problem does not consider disturbances, the resulting closed-loop reachable set under disturbances might be too large (Fig. 6.1c). Thus, we optimize over the feedback controller to minimize the effect of the disturbance (Fig. 6.1d).

We structure this section as follows: In Sec. 6.2.1, we briefly introduce the parameterization of the feedback matrices using LQR control, and then discuss the closed-loop reachable set computation for a given combined controller in Sec. 6.2.2. In Sec. 6.2.3, we discuss how the GSC approach can be used as a feedforward controller and then describe the feedback controller synthesis in Sec. 6.2.4. Next, we briefly comment on the role of initial guesses in Sec. 6.2.5. In Sec. 6.2.6, we then analyze the computational complexity of ROC and finally discuss important properties of the algorithm in Sec. 6.2.7. Note that we defer numerical experiments of ROC to Sec. 6.3, where we compare ROC against our novel combined synthesis approach. We refer the reader to [101, Sec. VI] for a comparison of ROC to other state-of-the-art approaches.

### 6.2.1 Feedback Matrix Parameterization

Directly optimizing over  $K^{(j)}$  for  $0 \leq j \leq m_c - 1$  requires  $m_c n_x n_u$  variables since  $K^{(j)} \in \mathbb{R}^{n_u \times n_x}$ . In this section, we therefore introduce a parameterization for the feedback matrices using LQR control (see Sec. 2.4) as proposed in [101, Sec. IV.B.], [98, Sec. 3.6.3] to reduce the number of optimization variables.

Since LQR computes a control law for an LTI system, we first linearize the nonlinear dynamics  $f$  along the reference trajectory to obtain  $m_c$  LTI systems. For  $t \in \tau^{(j)} = [j, j+1] \frac{t_f}{m_c}$ , we obtain the LTI system  $(A^{(j)}, B^{(j)})$  as

$$A^{(j)} = \left. \frac{\partial f(x, u, w)}{\partial x} \right|_{\substack{x=x_{\text{lin}}^{(j)} \\ u=u_{\text{lin}}^{(j)} \\ w=0}},$$

$$B^{(j)} = \left. \frac{\partial f(x, u, w)}{\partial u} \right|_{\substack{x=x_{\text{lin}}^{(j)} \\ u=u_{\text{lin}}^{(j)} \\ w=0}},$$

where

$$x_{\text{lin}}^{(j)} = \frac{1}{2} (x_{\text{ref}}^{(j)} + x_{\text{ref}}^{(j+1)}),$$

$$u_{\text{lin}}^{(j)} = u_{\text{ref}}^{(j)},$$

and the reference inputs and reference trajectory are computed as described in Sec. 5.2.1 with  $m_s = 1$ . To the best of our knowledge, the reason why the LTI systems  $(A^{(j)}, B^{(j)})$  are an appropriate choice for the computation of the LQR feedback matrix is not given by the original authors; we provide more details in Sec. 6.3.1 as we use this parameterization in a later algorithm. Assume that  $(A^{(j)}, B^{(j)})$  is controllable per Corr. 2.1; the corresponding LQR matrix  $K^{(j)}$  of the LTI system  $(A^{(j)}, B^{(j)})$  with given  $Q$  and  $R$  can then be computed as described in Sec. 2.4. Intuitively, the assumption of controllability of the linearized dynamics may seem quite restrictive; however, trying to control the nonlinear dynamics when even the linearized dynamics (which are only valid approximations in a small neighbor around  $x_{\text{lin}}^{(j)}$  and  $u_{\text{lin}}^{(j)}$ ) are not controllable seems rather

ambitious. While keeping  $Q$  and  $R$  constant over the entire horizon reduces the number of optimization variables, it may also reduce the overall performance of the controller as we limit the range of possible feedback matrices. If this extra performance is required, one may introduce separate weighting matrices  $Q$  and  $R$  for each of the  $m_c$  LTI systems. Alternatively, one can also directly optimize over  $K^{(j)}$ , in which case controllability is not required; however, this means that one may also generate non-stabilizing feedback matrices.

## 6.2.2 Closed-Loop Reachable Set

Since ROC optimizes directly over outer approximations of the closed-loop reachable set, we discuss the closed-loop reachable set computation for a combined controller of the form (6.1) for the conservative linearization approach (see Sec. 2.7.2) for a given  $K^{(j)}$  (also compare to [98, Sec. 3.6.4]). The reachable set computation using conservative polynomialization approach follows analogously and – since CORA implements dependency preservation for polynomial zonotopes – does not require the dependent factors  $\beta$  in the extended state definition.

To retain the dependency between each initial state and its corresponding control input, we define the extended state  $x_{\text{ext}} = [x^T, x_{\text{ref}}^T, \Delta\tilde{x}_{\text{ff}}^T, \beta^T]^T$  so that the extended dynamics for  $t \in \tau^{(j)}$  with  $0 \leq j \leq m_c - 1$  are given by

$$f_{\text{ext}}(x, u, w) = \begin{bmatrix} f(x, \bar{u}(\beta, \hat{P}^{(j)}) + K^{(j)}(Q, R)(x - (x_{\text{ref}} + \Delta\tilde{x}_{\text{ff}})), w) \\ f(x_{\text{ref}}, u_{\text{ref}}, 0) \\ A\Delta\tilde{x}_{\text{ff}} + B(\bar{u}(\beta, \hat{P}^{(j)}) - u_{\text{lin}}) \\ 0 \end{bmatrix}, \quad (6.7)$$

where the feedforward controller  $\bar{u}(\beta, \hat{P}^{(j)})$ , parameterized in  $\beta$ , is defined as in (5.3). Here,  $\hat{P}$  denotes the optimal controller parameters that define the feedforward controller as synthesized by GSC. The generating function of the extended initial set then is

$$x_{\text{ext}}^{(0)}(\beta) = \begin{bmatrix} x^{(0)}(\beta) \\ c_{\mathcal{X}^{(0)}} \\ G_{\mathcal{X}^{(0)}}\beta \\ \beta \end{bmatrix},$$

where  $\mathcal{X}^{(0)} = \{x^{(0)}(\beta)\}_{\beta} = \langle c_{\mathcal{X}^{(0)}}, G_{\mathcal{X}^{(0)}} \rangle_Z$ . Note that we do not use the feedforward state

$$x_{\text{ff}}(t) = \xi(t, x(0), u_{\text{ff}}(x(0), \hat{P}^{(j)}), 0),$$

directly, but rather the same linearized approximation  $\tilde{x}_{\text{ff}} = x_{\text{ref}} + \Delta\tilde{x}_{\text{ff}}$  that is used to compute the parameterized reachable set for the GSC approach. By doing so, we try to bring the closed-loop reachable set of the combined controller as close as possible to the parameterized reachable set, evaluated at the optimal controller parameter  $\hat{P}$ . Since we

will replace  $x_{\text{ff}}$  by  $\tilde{x}_{\text{ff}}$  during the feedback controller synthesis (see Sec. 6.2.4) and during the online deployment of GSC (see Sec. 6.2.6.2), all formal guarantees still hold.

Applying the reachability algorithm from Sec. 2.7.2 to the extended system dynamics in (6.7) for  $t \in [0, t_c]$  yields an outer approximation of the closed-loop reachable set, given by

$$\mathcal{R}_{\text{ext}}(t, Q, R) = \left\{ \begin{bmatrix} r_x(t, Q, R, \tilde{\beta}) \\ r_{x_{\text{ref}}}(t, \tilde{\beta}) \\ r_{\Delta\tilde{x}_{\text{ff}}}(t, \tilde{\beta}) \\ \beta(\tilde{\beta}) \end{bmatrix} \right\}_{\tilde{\beta}},$$

where its generating function is no longer dependent on  $\beta$  (e.g. due to reduction operations during reachability analysis) but  $\tilde{\beta}$ . Since this dependency is required for the evaluation of the dynamics in (6.7), we extended the state with  $\beta$ , making  $\beta(\tilde{\beta})$  available. The next extended initial set at  $t = t_c$  is then given by

$$x_{\text{ext}}^{(1)}(\tilde{\beta}, Q, R) = \begin{bmatrix} r_x\left(\frac{t_f}{m_c}, Q, R, \tilde{\beta}\right) \\ r_{x_{\text{ref}}}\left(\frac{t_f}{m_c}, Q, R, \tilde{\beta}\right) \\ r_{\Delta\tilde{x}_{\text{ff}}}\left(\frac{t_f}{m_c}, \tilde{\beta}\right) \\ \beta(\tilde{\beta}) \end{bmatrix},$$

and hence  $\mathcal{R}(t, Q, R) = \left\{ r_x(t, Q, R, \tilde{\beta}) \right\}_{\tilde{\beta}}$  for  $t \in [0, t_f]$  follows by the iterative application of the above steps.

### 6.2.3 Feedforward Controller Computation

Because the feedforward and feedback synthesis are executed sequentially, one can simply use the GSC approach described in Sec. 5.2 to compute the feedforward controller parameters. However, we cannot directly use the input constraint set  $\mathcal{U}$  as is: If we compute the feedforward controller for the input constraint set  $\mathcal{U}$ , the feedforward controller may use up all of the input capacity, leaving no available input capacity for the continuous feedback controller. As a result, the user needs to specify a tightened input constraint set  $\tilde{\mathcal{U}} \subseteq \mathcal{U}$  in addition to the adapted state and final state constraint sets  $\tilde{\mathcal{X}}$  and  $\tilde{\mathcal{X}}_f$ , which are necessary for non-conservative constraint checking.

### 6.2.4 Feedback Controller Computation

With the feedforward controller fixed, we only describe the synthesis of the feedback matrices here. To pose an approximation of (6.3), we subsequently introduce an approximation of its objective function and an approximation of its constraints as proposed in [98]. To that end, we first introduce all required generating functions.

## 6 Piecewise Constant Controller Synthesis with Continuous State Feedback

From Sec. 6.2.2, the generating function of the extended reachable set is

$$\mathcal{R}_{\text{ext}}(t, Q, R) = \left\{ \left[ \begin{array}{c} r_x(t, Q, R, \tilde{\beta}) \\ r_{x_{\text{ref}}}(t, \tilde{\beta}) \\ r_{\Delta \tilde{x}_{\text{ff}}}(t, \tilde{\beta}) \\ \beta(\tilde{\beta}) \end{array} \right] \right\}_{\tilde{\beta}}. \quad (6.8)$$

Let  $m_r \in \mathbb{N}_+$  denote the number of steps during reachability analysis per interval  $\tau^{(j)}$  for  $0 \leq j \leq m_c - 1$ , and let  $\tau^{(k,j)} = j \frac{t_f}{m_c} + [k, k+1] \Delta t$  for  $0 \leq k \leq m_r - 1$  with  $\Delta t = \frac{t_f}{m_c m_r}$ . Using (6.4) and (6.8), the generating function for the set of applied inputs  $\hat{\mathcal{S}}_u(t, Q, R) = \left\{ \hat{s}_u(t, Q, R, \tilde{\beta}) \right\}_{\tilde{\beta}} = \mathcal{S}_u(t, \hat{P}, Q, R)$  for a fixed  $\hat{P}$  and  $t \in \tau^{(k,j)}$  (or  $t = \tau^{(k,j)}$  since reachability analysis returns both time-point and time-interval solutions) can be computed via

$$\hat{s}_u(t, Q, R, \tilde{\beta}) = \bar{u}(\beta(\tilde{\beta}), \hat{P}^{(j)}) + K^{(j)}(Q, R) \left( r_x(t, Q, R, \tilde{\beta}) - r_{\tilde{x}_{\text{ff}}}(t, Q, R, \tilde{\beta}) \right),$$

where

$$r_{\tilde{x}_{\text{ff}}}(t, Q, R, \tilde{\beta}) = r_{x_{\text{ref}}}(t, \tilde{\beta}) + r_{\Delta \tilde{x}_{\text{ff}}}(t, Q, R, \tilde{\beta}).$$

We define

$$\begin{aligned} \hat{\mathcal{Z}}(\hat{\mathcal{S}}_u(t, Q, R)) &= \langle c_u(t, Q, R), G_u(t, Q, R) \rangle_Z, \\ \hat{\mathcal{Z}}\left(\left\{ r_x(t, Q, R, \tilde{\beta}) \right\}_{\tilde{\beta}}\right) &= \langle c(t, Q, R), G(t, Q, R) \rangle_Z. \end{aligned}$$

**Objective Function:** Let  $\hat{u}(t, x, Q, R) = u(t, x, \hat{P}, Q, R)$ . The objective function in (6.3a) can be approximated as (also see [98, Sec. 3.5.3])

$$\begin{aligned} & \max_{x(t, Q, R) \in \mathcal{R}(t, Q, R)} \left\{ \|x(t_f, Q, R) - x_f\|_1 + \zeta \int_0^{t_f} \|\hat{u}(\xi, x, Q, R)\|_1 d\xi \right\} \\ & \leq \|c(t_f, Q, R) - x_f\|_1 + \left\| [G(t_f, Q, R)]_{(\cdot)} \right\|_1 + \max_{x(t, Q, R) \in \mathcal{R}(t, Q, R)} \zeta \int_0^{t_f} \|\hat{u}(\xi, x, Q, R)\|_1 d\xi \\ & \leq \|c(t_f, Q, R) - x_f\|_1 + \left\| [G(t_f, Q, R)]_{(\cdot)} \right\|_1 + \max_{\|\tilde{\beta}\|_\infty \leq 1} \sum_{k=0}^{m_r-1} \sum_{j=0}^{m_c-1} \Delta t \zeta \left\| \hat{s}_u(\tau^{(k,j)}, Q, R, \tilde{\beta}) \right\|_1 \\ & \leq \left\| \begin{bmatrix} c(t_f, Q, R) - x_f \\ [G(t_f, Q, R)]_{(\cdot)} \end{bmatrix} \right\|_1 + \Delta t \zeta \sum_{k=0}^{m_r-1} \sum_{j=0}^{m_c-1} \left\| \begin{bmatrix} c_u(\tau^{(k,j)}, Q, R) \\ [G_u(\tau^{(k,j)}, Q, R)]_{(\cdot)} \end{bmatrix} \right\|_1, \end{aligned} \quad (6.9)$$

where the first inequality follows from the triangle inequality, the second inequality follows from applying the worst-case value during the entire interval of length  $\Delta t$ , and the last inequality again follows from the triangle inequality.

**Constraints:** Denote with  $\rho_u^{(\mathcal{Z})}(\cdot, \tau^{(k,j)}, Q, R)$  the support function of  $\hat{\mathcal{Z}}(\hat{\mathcal{S}}_u(t, Q, R))$ . For the input constraint from (6.3b) with  $\mathcal{U} = \langle C_U, d_U \rangle_H$ , it holds that

$$\forall t \in [0, t_f] : \hat{\mathcal{S}}_u(t, Q, R) \subseteq \mathcal{U} \quad (6.10)$$

$$\stackrel{Prop. 4.4}{\iff} \forall t \in [0, t_f] : \hat{\mathcal{Z}}(\hat{\mathcal{S}}_u(t, Q, R)) \subseteq \mathcal{U} \quad (6.11)$$

$$\iff \bigcup_{\substack{0 \leq k \leq m_r \\ 0 \leq j \leq m_c - 1}} \hat{\mathcal{Z}}(\hat{\mathcal{S}}_u(\tau^{(k,j)}, Q, R)) \subseteq \mathcal{U} \quad (6.12)$$

$$\stackrel{Prop. 4.6}{\iff} \rho_u^{(\mathcal{Z})}(C_U, [0, t_f], Q, R) \leq d_U, \quad (6.13)$$

where

$$\rho_u^{(\mathcal{Z})}(C_U, [0, t_f], Q, R) = \max_{\substack{0 \leq k \leq m_r \\ 0 \leq j \leq m_c - 1}} \rho_u^{(\mathcal{Z})}(C_U, \tau^{(k,j)}, Q, R).$$

Further, denote with  $\rho_x^{(\mathcal{Z})}(\cdot, t, Q, R)$  the support function of  $\hat{\mathcal{Z}}(\{r_x(t, Q, R, \tilde{\beta})\}_{\tilde{\beta}})$ . Then

$$\begin{aligned} & \mathcal{R}([0, t_f], \hat{P}, Q, R) \subseteq \mathcal{X} \\ & \stackrel{Prop. 4.4}{\iff} \bigcup_{\substack{0 \leq k \leq m_r \\ 0 \leq j \leq m_c - 1}} \hat{\mathcal{Z}}(\{r_x(\tau^{(k,j)}, Q, R, \tilde{\beta})\}_{\tilde{\beta}}) \subseteq \mathcal{X} \\ & \stackrel{Prop. 4.6}{\iff} \rho_x^{(\mathcal{Z})}(C_{\mathcal{X}}, [0, t_f], Q, R) \leq d_{\mathcal{X}}, \end{aligned} \quad (6.14)$$

where

$$\rho_x^{(\mathcal{Z})}(C_{\mathcal{X}}, [0, t_f], Q, R) = \max_{\substack{0 \leq k \leq m_r - 1 \\ 0 \leq j \leq m_c - 1}} \rho_x^{(\mathcal{Z})}(C_{\mathcal{X}}, \tau^{(k,j)}, Q, R),$$

and

$$\begin{aligned} & \mathcal{R}(t_f, Q, R) \subseteq \mathcal{X}_f \\ & \stackrel{Prop. 4.4}{\iff} \hat{\mathcal{Z}}(\{r_x(t_f, Q, R, \tilde{\beta})\}_{\tilde{\beta}}) \subseteq \mathcal{X}_f \\ & \stackrel{Prop. 4.6}{\iff} \rho_x^{(\mathcal{Z})}(C_{\mathcal{X}_f}, t_f, Q, R) \leq d_{\mathcal{X}_f}. \end{aligned} \quad (6.15)$$

**Optimization Problem:** Collecting (6.9) and (6.13) to (6.15) yields the optimization problem

$$\min_{Q,R} \left\| \begin{bmatrix} c(t_f, Q, R) - x_f \\ [G(t_f, Q, R)]_{(:,)} \end{bmatrix} \right\|_1 + \Delta t \zeta \sum_{k=0}^{m_r-1} \sum_{j=0}^{m_c-1} \left\| \begin{bmatrix} c_u(\tau^{(k,j)}, Q, R) \\ [G_u(\tau^{(k,j)}, Q, R)]_{(:,)} \end{bmatrix} \right\|_1, \quad (6.16a)$$

$$\text{s.t. } \rho_u^{(\mathcal{Z})}(C_{\mathcal{U}}, [0, t_f], Q, R) \leq d_{\mathcal{U}}, \quad (6.16b)$$

$$\rho_x^{(\mathcal{Z})}(C_{\mathcal{X}}, [0, t_f], Q, R) \leq d_{\mathcal{X}}, \quad (6.16c)$$

$$\rho_x^{(\mathcal{Z})}(C_{\mathcal{X}_f}, t_f, Q, R) \leq d_{\mathcal{X}_f}, \quad (6.16d)$$

$$Q \in \left\{ \text{diag}(q) \mid q \in [q, \bar{q}] \right\}, \quad (6.16e)$$

$$R \in \left\{ \text{diag}(r) \mid r \in [r, \bar{r}] \right\}, \quad (6.16f)$$

where we start from  $Q = \text{diag}(\bar{q})$  and  $R = \text{diag}(\bar{r})$  with  $[q, \bar{q}] \subseteq \mathbb{R}_+^{n_x}$  and  $[r, \bar{r}] \subseteq \mathbb{R}_+^{n_u}$ . To avoid semi-definite constraints for  $Q$  and  $R$ , the authors in [101] restrict the weighting matrices to positive diagonal matrices. Since the computation of  $K^{(j)}(Q, R)$  for  $0 \leq m_c - 1$  in dependence of  $Q$  and  $R$  is invariant when multiplying  $Q$  and  $R$  with some positive scalar, i.e.,  $K^{(j)}(Q, R) = K^{(j)}(sQ, sR)$  for  $s > 0$ , the authors in [101] further propose to set  $q_1 = \bar{q}_1 = 1$ .

## 6.2.5 Initial Guess For Feedback Synthesis

Since the feedforward controller is fixed and  $K^{(j)}(Q, R)$  for  $0 \leq j \leq m_c - 1$  is a function of  $Q$  and  $R$ , it suffices to find an initial guess for  $Q$  and  $R$ . Because solving (6.16) requires the evaluation of the closed-loop reachable set for each optimization iteration (see Sec. 6.2.4 for details), providing good initial guesses  $\bar{Q} \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $\bar{R} \in \mathbb{S}_{++}^{n_u \times n_u}$  can save a lot of computation time. Such an initial guess computation is outside of the scope of this thesis; that said, a possible initial guess – an idea from the code of the authors of [101] – can be obtained by solving (6.16), but replacing the outer approximations of the closed-loop reachable set for the nonlinear dynamics with approximations using reachability analysis of a linearized sequence of LTI systems. This has the advantage that linear reachability analysis is much more efficient than reachability analysis for nonlinear systems.

## 6.2.6 Computational Complexity

We again use Ass. 5.1. We distinguish between the complexity of the offline controller synthesis and the complexity of applying the controller online, both in  $n$ .

### 6.2.6.1 Offline Complexity

Since we already discussed the offline complexity of GSC (see Sec. 5.2.6.1), we focus on the complexity of computing the continuous feedback controller.



**Closed-loop Reachability Analysis:** Computing the reachable set for a given combined controller has a complexity of at most  $O(n^5)$  [59, Sec. 4.1.4].

**Feedback Optimization Problem:** The computation of the closed-loop reachable set in each optimization iteration when solving (6.16) naturally dominates the computational complexity with  $O(n^5)$  (if the conservative polynomialization algorithm is used; see Sec. 2.7.3). Since (6.16) cannot be reformulated as a smooth optimization problem (see Sec. 6.2.7 for details), there is no straightforward way to bound the number of optimization iterations.

**Overall Complexity:** Since the optimization problem to solve is not differentiable, we cannot make a statement about the overall complexity of ROC.

### 6.2.6.2 Online Complexity

Evaluation of the combined controller requires the evaluation of the GSC feedforward controller, one matrix vector multiplication between the feedback matrix and  $x - \tilde{x}_{\text{ff}}$ , and the forward simulation of  $\tilde{x}_{\text{ff}}$ . If a fixed-step ODE solver is used, the online computational complexity for the forward simulation of  $\tilde{x}_{\text{ff}}$  is at most  $O\left(b\left(n_x e + c_{\text{GSC}}^{(\text{on})}(n)\right)\right) \stackrel{\text{Ass. 5.1}}{=} O\left(n^2 + c_{\text{GSC}}^{(\text{on})}(n)\right)$ , where  $b \in \mathbb{N}_+$  collects both the number of fixed steps as well as the number of fixed function evaluations per step: We have  $n_x$  dimensions, the evaluation of  $f$  in each dimension requires  $e$  elementary operations, and we need to evaluate the feedforward controller  $b$  times for the computation of the feedforward input.

With the complexity of the matrix-vector multiplication of at most  $O(n_u n_x) = O(n^2)$ , the online complexity of the ROC approach is

$$O\left(c_{\text{ROC}}^{(\text{on})}(n)\right) = O\left(c_{\text{GSC}}^{(\text{on})}(n) + n^2\right) \stackrel{(5.15)}{=} \begin{cases} O(n^2), & \mathcal{X}^{(0)} = \mathcal{X}^{(0)} \downarrow_1, \\ O(c_{\text{LP}}(n) + n^2), & \text{otherwise} \end{cases}. \quad (6.17)$$

### 6.2.7 Discussion

Using the GSC approach for the feedforward synthesis, ROC is conceptually an intuitive extension for the inclusion of a continuous feedback term: By directly optimizing over outer approximations of the reachable set, constraints can be easily checked using Prop. 4.6. With an appropriate choice of both the tightened feedforward input constraint set as well as the adapted state and final state constraint sets, ROC significantly improves the control result compared to GSC when the disturbance sets get larger [101, Sec. VI].

However, this direct optimization over the reachable set in (6.16) requires the recomputation of an outer approximation of the extended reachable set in each optimization iteration. Further, the absolute values in (6.16) cannot be resolved as done previously: Such a reformulation requires that the same generator for different values of  $Q$  and  $R$  is

used for each resulting inequality. Since we have no way of identifying the same generator for different values of  $Q$  and  $R$ , this approach cannot be used. Since a reformulation of absolute values is not possible, differentiability of (6.16) is strictly speaking not given. Thus, applying most gradient-based solvers to (6.16) will yield no theoretical convergence guarantee. Furthermore, gradients for (6.16) cannot be computed analytically and thus gradient information has to be approximated using finite differences (see [53, Chap. 5] for details on finite differences) if a gradient-based solver is used. Additionally, ROC requires the user to choose the adapted constraint sets – such as adapted state constraints and tightened feedforward input constraints – appropriately, which makes the approach non-trivial to apply for non-experts. Lastly, ROC first computes the feedforward controller and only then synthesizes the continuous feedback controller. Since the choice of the feedforward controller directly influences the feedback controller, this separation into two synthesis steps might result in a sub-optimal controller (also see Sec. 6.3.10).

### 6.3 Iterative Polynomial Reachset Optimal Control

The ROC approach first synthesizes a feedforward controller and then computes the continuous feedback control law while keeping the feedforward controller fixed. Not only does this generally result in a sub-optimal solution, but it also introduces additional algorithm parameters, such as the adapted input and state constraints, which have to be set appropriately to arrive at a performant controller; this can limit the applicability of ROC by non-experts. To avoid these additional parameters, we present the iterative polynomial reachset optimal control (iPROC) approach from our work in [35], which solves the synthesis problem for the combined controller simultaneously, i.e., without separation of the feedforward and feedback controller synthesis: Since (6.3) is generally not efficiently solvable, we approximate (6.3) – similarly to iPGSC – using a non-differentiable cost function in the combined controller parameters. By locally approximating this non-differentiable cost function with a differentiable optimization problem, we can iteratively update our controller parameters based on this local approximation, where we control its accuracy by restricting the admissible controller parameters for the optimization problem to a bounded trust region.

As it turns out, parameterizing  $K^{(j)} = K^{(j)}(z)$  for  $0 \leq j \leq m_c - 1$  in the controller parameters  $z = \left[ P_{(\cdot)}^T, Q_{(\cdot)}^T, R_{(\cdot)}^T \right]^T$  using LQR control – as used in Sec. 6.2 and originally proposed in [101, Sec. IV.B] – is also beneficial here, where  $P \in [-1, 1]^{m_c n_u \times a}$  are the feedforward controller parameters of the controller template from (5.20), and  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $R \in \mathbb{S}_{++}^{n_u \times n_u}$  are the feedback controller parameters. For now, we simply take this parameterization as a given and defer its motivation and introduction to Sec. 6.3.1. Fig. 6.2 visualizes the general idea of iPROC, where each step of Alg. 6 is briefly described subsequently.

- (i) Starting from an initial guess for both the feedforward controller parameters  $\bar{P}$  and feedback controller parameters  $\bar{Q}$  and  $\bar{R}$  (l. 3), we first compute an outer approximation to the undisturbed feedforward reachable set for the current feedforward controller (l. 5).
- (ii) Using this feedforward outer approximation, we then construct a locally accurate approximation of the undisturbed feedforward reachable set around the current controller parameters (l. 6).
- (iii) We evaluate the cost of the initial controller, parameterized by the initial controller parameters, using an outer approximation of the closed-loop reachable set of this initial controller (l. 7–8). We then iteratively update the current controller parameters by executing the following steps:
  - (a) Starting from the current controller parameters, we solve the trust-region subproblem – which is a locally accurate, differentiable approximation of the controller cost – to obtain a new controller candidate (l. 10): To ensure accuracy, we restrict the domain of controller parameters  $z$  to a bounded trust region, whose size is determined by the feedforward and feedback trust-region radii  $\gamma$  and  $\eta \in (0, 1]$ .

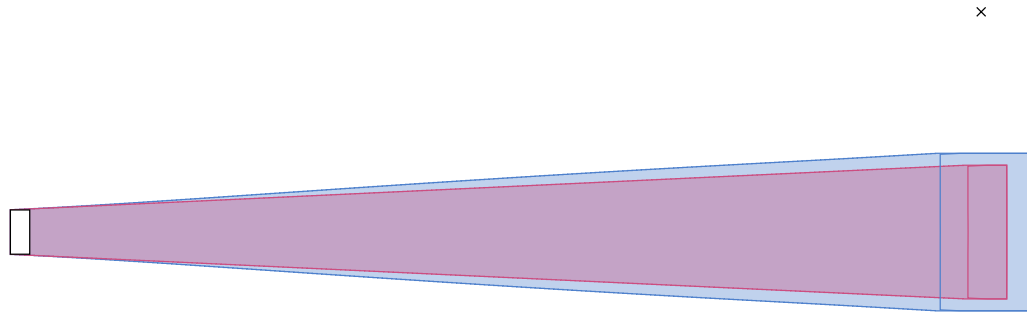
- (b) We try to compute an outer approximation to the closed-loop reachable set of the new controller candidate (l. 13–18): For very large trust-region radii, the inaccurate approximation of the reachable set may lead to a controller for which the closed-loop reachable set computation encounters numerical errors. In that case, we shrink the trust-region radii and restart the iteration.
- (c) We then compute the cost of the newly synthesized controller (l. 19). (l. 20).
- (d) To ensure the accuracy of the trust-region subproblem, we compare the difference in the controller cost to the difference predicted by the trust-region subproblem and tune the trust-region radii accordingly (l. 20).
- (e) If the newly synthesized controller has a lower cost than the current controller candidate, the step is accepted and we update all necessary values (l. 23–24). We terminate the algorithm if the change from the previous best parameters to the current best parameters is small enough (l. 26). Otherwise, we recompute an updated parameterized feedforward reachable set (l. 29) and continue with (iii-a).

The remainder of this section explains iPROC in more detail: We start by introducing the feedback matrix parameterization using LQR control in Sec. 6.3.1 and describe the closed-loop reachable set computation of the combined controller in Sec. 6.3.2. We then derive the controller cost as an approximation to (6.3) in Sec. 6.3.3. Because the evaluation of the cost is computationally expensive, we compute a locally accurate approximation to the parameterized feedforward reachable set in Sec. 6.3.4, which is then used in Sec. 6.3.5 to construct the trust-region subproblem in each iteration to locally approximate the controller cost. In Sec. 6.3.6, we then derive tuning rules for the trust-region radii, which limit the domain of the controller parameters to a local neighborhood around the current controller parameters, to guarantee an accurate approximation of the controller cost. Since the trust-region subproblem is a differentiable optimization problem, we demonstrate the computation of necessary derivatives using differentials in Sec. 6.3.7. Since the proposed algorithm strongly benefits from good initial guesses for the controller parameters, we briefly discuss possible initial guesses in Sec. 6.3.8 and then derive the computational complexity, both online and offline, for the iPROC approach in Sec. 6.3.9. Finally, we demonstrate the applicability of iPROC with a comparison to the ROC approach in numerical experiments in Sec. 6.3.10, and discuss advantages of the algorithm in Sec. 6.3.11.

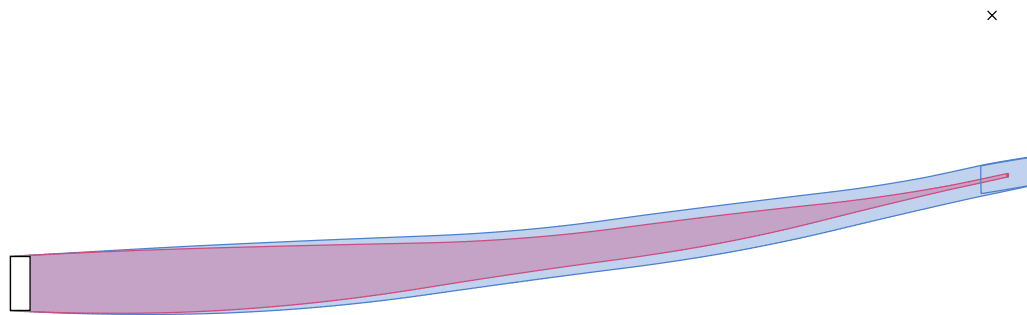
### 6.3.1 Feedback Matrix Parameterization

Directly optimizing over  $K^{(j)}$  for  $0 \leq j \leq m_c - 1$  requires  $m_c n_x n_u$  variables since  $K^{(j)} \in \mathbb{R}^{n_u \times n_x}$ . In Sec. 6.2.1, this same parameterization has already been introduced; in this section, we focus on the necessary, additional requirements to use this parameterization for iPROC. First, we briefly motivate the choice of LTI systems for the LQR control approach. Then, we further show that the LQR feedback matrix is continuously differentiable in the controller parameters  $z$ , which is a necessary requirement to later optimize over them (see Sec. 6.3.5).

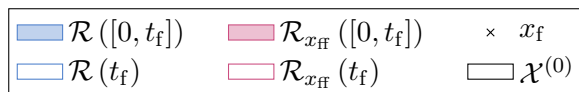
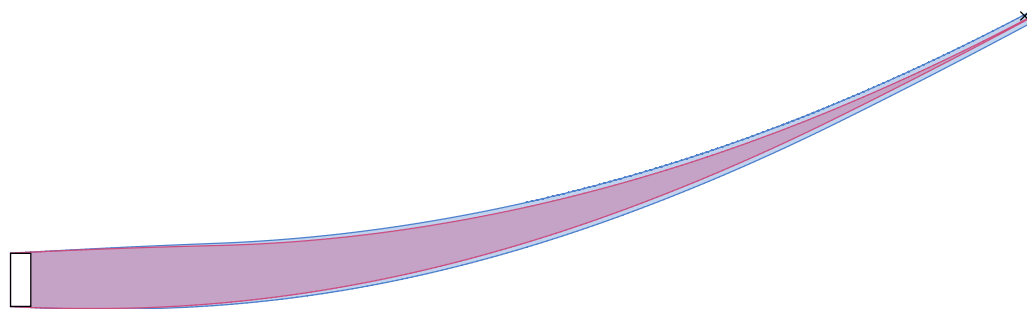
Initial reachable set



Reachable set during optimization



Final reachable set



**Figure 6.2:** Visualization of the iPROC approach. Starting from an initial guess for both the feedforward controller parameters as well as state and input weighting matrices of the LQR formulation, the feedforward controller steers the reachable set closer to the target state while the feedback controller simultaneously minimizes the effect of the disturbance.

---

**Algorithm 6** Iterative Polynomial Reachset Optimal Control
 

---

```

1: function ITERPOLYREACHSETOPTIMALCONTROL( $\mathcal{X}^{(0)}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f, \mathcal{W}, m_c, x_f, \kappa$ )
2:    $\gamma = \eta = 1$ 
3:    $[\bar{P}, \bar{Q}, \bar{R}] = \text{INITGUESS}(\dots)$  ▷ Sec. 6.3.8
4:    $\bar{z} = [\bar{P}_{(\cdot)}^T, \bar{Q}_{(\cdot)}^T, \bar{R}_{(\cdot)}^T]^T$ 
5:    $\bar{\mathcal{R}}_{x_{\text{ff}}} = \text{REACHFF}(\mathcal{X}^{(0)}, \bar{P}, t_f, \kappa)$  ▷ see Sec. 5.2.5.2 with  $m_s = h = 1$ 
6:    $[\tilde{\mathcal{R}}_{x_{\text{ff}}}, \gamma] = \text{REACHPARAMFF}(\mathcal{X}^{(0)}, \bar{P}, \gamma, \bar{\mathcal{R}}_{x_{\text{ff}}}, t_f, \kappa)$  ▷ Sec. 5.4.2
7:    $[\bar{\mathcal{R}}_{\text{ext}}, \bar{K}] = \text{REACH}(\mathcal{X}^{(0)}, \bar{z}, \bar{\mathcal{R}}_{x_{\text{ff}}}, t_f, \kappa)$  ▷ Sec. 6.3.2
8:    $\bar{J} = \text{COST}(\bar{\mathcal{R}}_{\text{ext}})$  ▷ Sec. 6.3.3
9:   for  $k = 1; k \leq l_{\text{max}}; k++$  do
10:      $\hat{z} = \text{TRSUBPROBLEM}(\bar{z}, \gamma, \eta, \tilde{\mathcal{R}}_{x_{\text{ff}}}, \bar{\mathcal{R}}_{\text{ext}}, \mathcal{U}, \mathcal{X}, \mathcal{X}_f)$  ▷ Sec. 6.3.5
11:     try
12:        $[\hat{\mathcal{R}}_{\text{ext}}, \hat{K}] = \text{REACH}(\mathcal{X}^{(0)}, \hat{z}, \tilde{\mathcal{R}}_{x_{\text{ff}}}, t_f, \kappa)$  ▷ Sec. 6.3.2
13:     catch
14:        $\gamma = \frac{1}{2}\gamma$ 
15:        $\eta = \frac{1}{2}\eta$ 
16:        $[\tilde{\mathcal{R}}_{x_{\text{ff}}}(t, P), \gamma] = \text{REACHPARAMFF}(\mathcal{X}^{(0)}, \bar{P}, \gamma, \bar{\mathcal{R}}_{x_{\text{ff}}}, t_f, \kappa)$  ▷ Sec. 5.4.2
17:     continue
18:     end try
19:      $\hat{J} = \text{COST}(\hat{\mathcal{R}}_{\text{ext}})$  ▷ Sec. 6.3.3
20:      $[\gamma, \eta] = \text{TUNERADII}(\gamma, \eta, \tilde{\mathcal{R}}_{x_{\text{ff}}}, \bar{\mathcal{R}}_{\text{ext}}, \hat{\mathcal{R}}_{\text{ext}}, \bar{J}, \hat{J})$  ▷ Sec. 6.3.6
21:     if  $\hat{J} < \bar{J}$  then
22:        $\nu = \text{tol}(\hat{J}, \bar{J})$ 
23:        $\bar{z} = \hat{z}, \bar{J} = \hat{J}, \bar{\mathcal{R}}_{\text{ext}} = \hat{\mathcal{R}}_{\text{ext}}, \bar{K} = \hat{K}$ 
24:        $\bar{\mathcal{R}}_{x_{\text{ff}}} = [\bar{\mathcal{R}}_{\text{ext}}]_{(n_x + (1:n_x))}$ 
25:       if  $\nu < \mu$  then
26:         break
27:       end if
28:     end if
29:      $[\tilde{\mathcal{R}}_{x_{\text{ff}}}(t, P), \gamma] = \text{REACHPARAMFF}(\mathcal{X}^{(0)}, \bar{P}, \gamma, \bar{\mathcal{R}}_{x_{\text{ff}}}, t_f, \kappa)$  ▷ Sec. 5.4.2
30:   end for
31:   return  $\bar{P}, \bar{K}, [\bar{\mathcal{R}}_{\text{ext}}]_{(1:n_x)}$ 
32: end function

```

---

### 6.3 Iterative Polynomial Reachset Optimal Control

**LTI systems:** We now derive the LTI systems required for the computation of the feedback matrices using LQR control. We split the state

$$x(t, z) = x_{\text{ff}}(t, P) + \Delta x(t, z), \quad (6.18)$$

where  $P$  are the feedforward parameters, and  $z = [P_{(\cdot)}^T, Q_{(\cdot)}^T, R_{(\cdot)}^T]^T$  are the combined controller parameters with LQR weighting matrices  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $R \in \mathbb{S}_{++}^{n_u \times n_u}$ . Since the feedback controller tries to follow the undisturbed feedforward state  $x_{\text{ff}}$  and  $x = x_{\text{ff}} + \Delta x$ , we ideally try to steer  $\Delta x$  to the origin. For the  $m_c$  feedback matrices  $K(t, z) = K^{(j)}(z)$  with  $t \in \tau^{(j)}$  and  $0 \leq j \leq m_c - 1$  that need to be computed, we can thus compute the  $m_c$  LTI systems required for LQR control by performing a first-order Taylor expansion of  $\Delta \dot{x} = f(x, u, w) - f(x_{\text{ff}}, u_{\text{ff}}, 0)$  around  $x = x_{\text{ff}}$ ,  $u = u_{\text{ff}}$  and  $w = 0$ , which yields (arguments omitted where clear from the context for readability)

$$\begin{aligned} \Delta \dot{x} &= \dot{x} - \dot{x}_{\text{ff}} = f(x, u, w) - f(x_{\text{ff}}, u_{\text{ff}}, 0) \\ &\approx f(x_{\text{ff}}, u_{\text{ff}}, 0) + \frac{\partial f(x, u, w)}{\partial x} \Big|_{\substack{x=x_{\text{ff}} \\ u=u_{\text{ff}} \\ w=0}} (x - x_{\text{ff}}) \\ &\quad + \frac{\partial f(x, u, w)}{\partial u} \Big|_{\substack{x=x_{\text{ff}} \\ u=u_{\text{ff}} \\ w=0}} (u - u_{\text{ff}}) + \frac{\partial f(x, u, w)}{\partial w} \Big|_{\substack{x=x_{\text{ff}} \\ u=u_{\text{ff}} \\ w=0}} w \\ &\quad - f(x_{\text{ff}}, u_{\text{ff}}, 0) \\ &\approx A_{\text{cl}}^{(j)}(z) \Delta x + E^{(j)}(P) w, \end{aligned} \quad (6.19)$$

since  $u(t, x(0), z) = u_{\text{ff}}(x(0), P^{(j)}) + K^{(j)}(z) \Delta x(t, z)$  for  $t \in \tau^{(j)}$  with  $0 \leq j \leq m_c - 1$  (see (6.1)) and where

$$A^{(j)}(P) = \frac{\partial f(x, u, w)}{\partial x} \Big|_{\substack{x=x_{\text{lin}}^{(j)}(P) \\ u=u_{\text{lin}}^{(j)}(P) \\ w=0}}, \quad (6.20)$$

$$B^{(j)}(P) = \frac{\partial f(x, u, w)}{\partial u} \Big|_{\substack{x=x_{\text{lin}}^{(j)}(P) \\ u=u_{\text{lin}}^{(j)}(P) \\ w=0}}, \quad (6.21)$$

$$E^{(j)}(P) = \frac{\partial f(x, u, w)}{\partial w} \Big|_{\substack{x=x_{\text{lin}}^{(j)}(P) \\ u=u_{\text{lin}}^{(j)}(P) \\ w=0}}, \quad (6.22)$$

$$A_{\text{cl}}^{(j)}(z) = A^{(j)}(P) + B^{(j)}(P) K^{(j)}(z), \quad (6.23)$$

with

$$\begin{aligned} x_{\text{lin}}^{(j)}(P) &= \frac{1}{2} (\tilde{c}_{x_{\text{ff}}}(j t_c, P) + \tilde{c}_{x_{\text{ff}}}((j+1) t_c, P)), \\ u_{\text{lin}}^{(j)}(P) &= c_{u_{\text{ff}}}(P^{(j)}). \end{aligned}$$

Here,

$$\begin{aligned}\hat{\mathcal{Z}}\left(\left\{\bar{u}\left(\beta, P^{(j)}\right)\right\}_{\beta}\right) &= \left\langle c_{u_{\text{ff}}}\left(P^{(j)}\right), G_{u_{\text{ff}}}\left(P^{(j)}\right)\right\rangle_Z, \quad 0 \leq j \leq m_c - 1, \\ \hat{\mathcal{Z}}\left(\tilde{\mathcal{R}}_{x_{\text{ff}}}(t, P)\right) &= \left\langle \tilde{c}_{x_{\text{ff}}}(t, P), \tilde{G}_{x_{\text{ff}}}(t, P)\right\rangle_Z,\end{aligned}$$

are the zonotope outer approximations of  $\left\{\bar{u}\left(\beta, P^{(j)}\right)\right\}_{\beta}$  and  $\tilde{\mathcal{R}}_{x_{\text{ff}}}(t, P)$  from (5.21) and (5.44) for  $m_s = 1$ , whose centers  $\tilde{c}_{x_{\text{ff}}}(t, P)$  and  $c_{u_{\text{ff}}}\left(P^{(j)}\right)$  are used to define the linearization points for the linearized matrices in (6.20) to (6.22). As a result, we use the feedforward input and the parameterized feedforward reachable set to extract their centers as approximations for a reference trajectory parameterized in  $P$ ; the zonotope outer approximation is a preferred approximation of the geometric center (see Ex. 5.2). Note that in (6.19), we replaced the implicit dependence of the linearized dynamics on the initial state  $x_{\text{ff}}(0) = x(0)$  with  $\tilde{c}_{x_{\text{ff}}}(0, P)$ : This approximation is motivated by the fact that the feedforward controller generally minimizes the size of the undisturbed reachable set and the feedback controller minimizes the effect of the disturbance.

If controllability for all  $P$  can be assumed, the LQR parameterization  $K(t, z) = K^{(j)}(z)$  using the weighting matrices  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $R \in \mathbb{S}_{++}^{n_u \times n_u}$  – which are constant for the entire time horizon – and the system matrices  $A^{(j)}(P)$  and input matrices  $B^{(j)}(P)$  is then well-defined and motivated as follows: By design, LQR then yields the stable closed-loop system matrices  $A_{\text{cl}}^{(j)}(z)$  (see (6.23)) and thus tries to bring  $\Delta x$ , which is approximately described by the  $m_c$  systems  $\left(A_{\text{cl}}^{(j)}(z), E^{(j)}(P)\right)$ , to the origin. While keeping  $Q$  and  $R$  constant over the entire horizon reduces the number of optimization variables, it may also reduce the overall performance of the controller as we limit the range of possible feedback matrices. If this extra performance is required, one may introduce separate weighting matrices for each of the  $m_c$  LTI systems. Alternatively, one can also directly optimize over  $K^{(j)}$ , in which case controllability is not required; however, this means that one may also generate non-stabilizing feedback matrices.

**Differentiability:** In Prop. 3.1, we showed that the abstracted synthesis problem can be replaced by a smooth optimization problem which shares its minimum. Since Prop. 3.1 requires all involved functions to be sufficiently smooth, we must show that the LQR feedback matrix  $K^{(j)}(z)$  is continuously differentiable with respect to  $z$ .

**Lemma 6.1** (Continuous Differentiability of LQR Matrix). *Let  $(A(P), B(P))$  for  $P \in [-1, 1]^{m_c n_u \times a}$  denote an  $n_x$ -dimensional controllable LTI system, where  $A(P)$  and  $B(P)$  are  $k$ -times continuously differentiable in each matrix element. Then the LQR gain matrix  $K(z)$  for the system  $(A(P), B(P))$  and weight matrices  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $R \in \mathbb{S}_{++}^{n_u \times n_u}$  is  $k$ -times continuously differentiable with respect to the collected variables  $z = \left[P_{(\cdot)}^T, Q_{(\cdot)}^T, R_{(\cdot)}^T\right]^T$ .*

*Proof.* We first prove the differentiability of  $X$  with respect to  $z$ . The algebraic Riccati equation (see (2.25)) is

$$F(z) = A(P)^T X(z) + X(z) A(P) - X(z) B(z) R^{-1} B(P)^T X(z) + Q = 0.$$



### 6.3 Iterative Polynomial Reachset Optimal Control

From the implicit function theorem [27, Th. 1B.1, Prop. 1B.5], we know that  $X_{(\cdot)}(z)$  exists, is unique and  $k$ -times differentiable in a neighborhood around  $z$  if  $\det\left(\frac{dF_{(\cdot)}}{dX_{(\cdot)}}\right) \neq 0$ . Since  $K = -R^{-1}B^T X$  (see (2.24)), we have

$$\begin{aligned} dF &= A^T dX + dXA - dXBR^{-1}B^T X - XBR^{-1}B^T dX \\ &= (A + BK)^T dX + dX(A + BK) \\ \stackrel{(2.1)}{\iff} dF_{(\cdot)} &= \left( I_{n_x} \boxtimes (A + BK)^T + (A + BK)^T \boxtimes I_{n_x} \right) dX_{(\cdot)}, \end{aligned}$$

and thus

$$\frac{dF_{(\cdot)}}{dX_{(\cdot)}} = A_{\text{cl}}^T \boxplus A_{\text{cl}}^T,$$

where we used the definition of the Kronecker sum and  $A_{\text{cl}} = A + BK$ . Since  $A_{\text{cl}}$  is stable by design, i.e.,  $\text{real}(\lambda(A_{\text{cl}})) < 0$ , it follows that  $\text{real}\left(\lambda\left(A_{\text{cl}}^T \boxplus A_{\text{cl}}^T\right)\right) < 0$  [70, Th. 13.16] and thus  $\det\left(\frac{dF_{(\cdot)}}{dX_{(\cdot)}}\right) = \det\left(A_{\text{cl}}^T \boxplus A_{\text{cl}}^T\right) \neq 0$ . As a result,  $X(z)$  exists, is unique, and  $k$ -times continuously differentiable in a neighborhood around  $z$ , and since we made no assumptions about  $z$ , the result follows for all  $P \in [-1, 1]^{m_c m_s n_u \times a}$ ,  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$ , and  $R \in \mathbb{S}_{++}^{n_u \times n_u}$ .

The differentiability of  $K$  then directly follows from (2.24) and the fact that  $R \in \mathbb{S}_{++}^{n_u \times n_u}$ .  $\square$

**Remark:** As it turns out, one can use Lem. 6.1 to derive an analytic approximation of the LQR matrix with arbitrary precision, which we now briefly sketch. The second-order Taylor expansion of  $K_{ij}(z)$  for  $1 \leq i \leq n_u$  and  $1 \leq j \leq n_x$  around some expansion point  $\bar{z} \in \mathbb{R}^{n_z}$  is given by

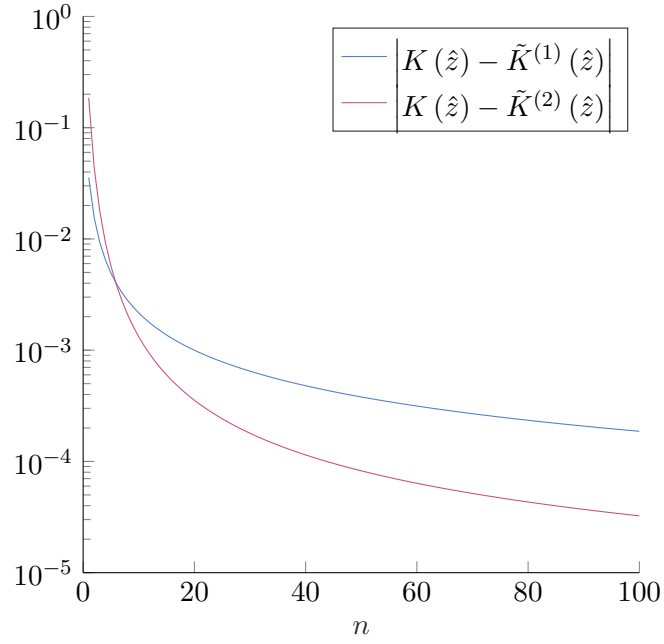
$$K_{ij}(z) = K_{ij}(\bar{z}) + J_{K_{ij}}(\bar{z})(z - \bar{z}) + \frac{1}{2}(z - \bar{z})^T H_{K_{ij}}(\bar{z})(z - \bar{z}) + R_{ij}^{(3)}(z - \bar{z}), \quad (6.24)$$

where  $R_{ij}^{(3)}(z - \bar{z})$  is the remainder of the second-order Taylor expansion in (6.24) and [29, Corr. 7.2]

$$\lim_{z - \bar{z} \rightarrow 0} \frac{R_{ij}^{(3)}(z - \bar{z})}{\|z - \bar{z}\|_2^3} = 0. \quad (6.25)$$

Thus, if we want to compute  $K(\hat{z})$  for some  $\hat{z} \in \mathbb{R}^{n_z}$  starting from  $\bar{z}$ , we apply (6.24) recursively  $n \in \mathbb{N}_+$  times, yielding

$$\begin{aligned} K_{ij}(z^{(k+1)}) &= K_{ij}(z^{(k)}) + J_{K_{ij}}(z^{(k)})(z^{(k+1)} - z^{(k)}) \\ &\quad + \frac{1}{2}(z^{(k+1)} - z^{(k)})^T H_{K_{ij}}(z^{(k)})(z^{(k+1)} - z^{(k)}) + R_{ij}^{(3)}(z^{(k+1)} - z^{(k)}) \\ &= K_{ij}(z^{(k)}) + J_{K_{ij}}(z^{(k)})\Delta z + \frac{1}{2}\Delta z^T H_{K_{ij}}(z^{(k)})\Delta z + R_{ij}^{(3)}(\Delta z), \end{aligned}$$



**Figure 6.3:** Visualization of the decrease in error for the recursive approximation of the LQR feedback matrix as defined in (6.26) and (6.27).

where  $\forall k \in \{1, \dots, n-1\} : \Delta z = z^{(k+1)} - z^{(k)} = \frac{\hat{z} - \bar{z}}{n}$ ,  $\hat{z} = z^{(n)}$  and  $J_{K_{ij}}(z^{(k)})$  and  $H_{K_{ij}}(z^{(k)})$  can be computed using Lem. 6.1. Since  $R_{ij}^{(3)}(\Delta z)$  tends to 0 at least as fast as  $\|\Delta z\|_2^3$  for  $\Delta z \rightarrow 0$  (see (6.25)) and increasing  $n$  decreases  $\|\Delta z\|_2$ , the total error  $nR_{ij}^{(3)}(\Delta z)$  shrinks for increasing  $n$ . We define

$$K_{ij}^{(1)}(z^{(k+1)}) = K_{ij}(z^{(k)}) + J_{K_{ij}}(z^{(k)}) \Delta z, \quad (6.26)$$

$$K_{ij}^{(2)}(z^{(k+1)}) = K_{ij}(z^{(k)}) + J_{K_{ij}}(z^{(k)}) \Delta z + \frac{1}{2} \Delta z^T H_{K_{ij}}(z^{(k)}) \Delta z. \quad (6.27)$$

The following example demonstrates the expected decrease in error with increasing  $n$ .

**Example 6.1.** Let an LTI system  $(A, B)$  with  $A = 1$  and  $B = 1$ , where  $n_x = 1$  and  $n_u = 1$ , be given. Further, we set  $\bar{Q} = \bar{R} = 1$  so that  $\bar{z} = [1, 1]^T$  since we assume that  $A$  and  $B$  are constant and thus omit the feedforward parameters in  $z$ . Fig. 6.3 visualizes the error of the first-order and second-order approximations in (6.26) and (6.27) with  $\hat{z} = [10, 5]^T$  for increasing  $n$ . As expected, the error decreases for increasing  $n$  and the second-order approximation achieves a smaller error for the same  $n$  and large enough  $n$ . ■

### 6.3.2 Closed-Loop Reachable Set

During each iteration of Alg. 6, an outer approximation to the extended closed-loop reachable set needs to be computed. The combined controller for  $t \in \tau^{(j)}$  with  $0 \leq j \leq m_c - 1$  is given by (see (6.1))

$$u(t, x, z) = u_{\text{ff}}(x(0), P^{(j)}) + K^{(j)}(z)(x(t) - x_{\text{ff}}(t)),$$

where  $u_{\text{ff}}(x(0), P^{(j)})$  is defined as in (5.20).

To retain dependencies of the combined state, feedforward state, and feedforward input on the initial state during the closed-loop reachable set computation, we define the extended state  $x_{\text{ext}} = [x^T, x_{\text{ff}}^T, u_{\text{ff}}^T]^T$  with the extended flow given by

$$f_{\text{ext}}(x_{\text{ext}}, w) = \begin{bmatrix} f(x, u_{\text{ff}} + K(t, z)(x - x_{\text{ff}}), w) \\ f(x_{\text{ff}}, u_{\text{ff}}, 0) \\ 0 \end{bmatrix},$$

and the extended initial set

$$\mathcal{X}_{\text{ext}}^{(0)}(P) = \left\{ \begin{bmatrix} x^{(0)}(\beta) \\ x^{(0)}(\beta) \\ \bar{u}(\beta, P^{(0)}) \end{bmatrix} \right\}_{\beta}.$$

Reachability analysis then yields

$$\mathcal{R}_{\text{ext}}(t, z) = \left\{ \begin{bmatrix} r_x(t, \beta, \xi, z) \\ r_{x_{\text{ff}}}(t, \beta, \xi, P) \\ \bar{u}(\beta, P^{(0)}) \end{bmatrix} \right\}_{\beta, \xi},$$

for  $t \in [0, 1]t_c$ , where  $r_x$  and  $r_{x_{\text{ff}}}$  denote the respective generating functions and the dependent factor  $\xi$  is caused by the abstraction and reduction errors as well as the convex combination required for the computation of the time-interval reachable sets. The extended initial set for  $t \in [1, 2]t_c$  then is

$$\mathcal{X}_{\text{ext}}^{(1)}(z) = \left\{ \begin{bmatrix} r_x(t_c, \beta, \xi, z) \\ r_{x_{\text{ff}}}(t_c, \beta, \xi, P) \\ \bar{u}(\beta, P^{(1)}) \end{bmatrix} \right\}_{\beta, \xi},$$

from which the remaining  $m_c - 1$  steps until  $t = t_f$  follow analogously.

### 6.3.3 Cost of the Controller

In this section, we approximate (6.3) with a cost function. To that end, let the extended closed-loop reachable set  $\mathcal{R}_{\text{ext}}(t, z)$  at  $z$  (see Sec. 6.3.2), containing an outer approximation of the closed-loop reachable set  $\mathcal{R}(t, z)$  and an outer approximation of

the closed-loop feedforward reachable set  $\mathcal{R}_{x_{\text{ff}}}(t, P)$ , be available for  $m_r \in \mathbb{N}_+$  reachability steps, i.e., interval sets for  $t \in [k, k+1] \frac{t_c}{m_r}$  and time point sets for  $t = k \frac{t_c}{m_r}$  with  $0 \leq k \leq m_c m_r - 1$  are given.

The remainder of this section is structured as follows: First, we review the split of the state and the input, which simplifies all subsequent computations. Then, we approximate the objective function and constraints of (6.3) and form the cost as an approximation to (6.3).

**State and input separation:** Let the split of the state be defined as in (6.18), where  $P$  are the feedforward parameters, and  $z = \begin{bmatrix} P_{(\cdot)}^T & Q_{(\cdot)}^T & R_{(\cdot)}^T \end{bmatrix}^T$  are the combined controller parameters (for details see Sec. 6.3.1). Further, let

$$\Delta x(t, z) \in \mathcal{S}_{\Delta x}(t, z) = \mathcal{R}(t, z) \oplus_e (-\mathcal{R}_{x_{\text{ff}}}(t, P)), \quad (6.28)$$

collect all  $\Delta x$  as defined in (6.18) (we use  $\mathcal{S}_{\Delta x}$  instead of  $\mathcal{R}_{\Delta x}$  since a reachable set requires a flow equation according to Def. 2.25). By definition of the exact sum, we have

$$x(t, z) \in \mathcal{R}(t, z) = \mathcal{R}_{x_{\text{ff}}}(t, P) \oplus_e \mathcal{S}_{\Delta x}(t, z) \subseteq \mathcal{R}_{x_{\text{ff}}}(t, P) \oplus \mathcal{S}_{\Delta x}(t, z), \quad (6.29)$$

where the last inclusion follows from the fact that ignoring dependencies always results in an outer approximation. By removing these dependencies,  $\mathcal{R}_{x_{\text{ff}}}(t, P)$  can be computed independently of  $\mathcal{S}_{\Delta x}(t, z)$ . Naturally, the removal of these dependencies introduces conservatism but we argue that this does not cause a large outer approximation: Without any disturbances,  $\mathcal{S}_{\Delta x}(t, z) = \emptyset$  and thus  $\mathcal{R}(t, z) = \mathcal{R}_{x_{\text{ff}}}(t, P)$ . As the disturbance set grows, the size of  $\mathcal{S}_{\Delta x}(t, z)$  is mainly affected by the size of the disturbance set (see Sec. 6.3.4 for the derivation of the approximated flow of  $\Delta x$ ).

Similarly, one can define an input separation following (6.4) as

$$\mathcal{S}_u(t, z) = \mathcal{S}_{u_{\text{ff}}}(t, P) \oplus_e \mathcal{S}_{\Delta u}(t, z) \subseteq \mathcal{S}_{u_{\text{ff}}}(t, P) \oplus \mathcal{S}_{\Delta u}(t, z), \quad (6.30)$$

where the set of feedforward inputs is given by  $\mathcal{S}_{u_{\text{ff}}}(t, z) = \left\{ \bar{u} \left( \beta, P^{(j)} \right) \right\}$  for  $t \in \tau^{(j)}$  with  $0 \leq j \leq m_c - 1$  and  $\mathcal{S}_{\Delta u}(t, z)$  contains all input deviations  $\Delta u(t, z) = K(t, z) \Delta x(t, z)$  as defined in (6.5) and (6.6).

**Objective:** Based on the separation of the reachable set into the feedforward reachable set  $\mathcal{R}_{x_{\text{ff}}}(t, P)$  and the set of deviation states  $\mathcal{S}_{\Delta x}(t, z)$  from (6.29) as well as the corresponding separation of the combined input into the feedforward input set  $\mathcal{S}_{u_{\text{ff}}}(t, P)$  and the set of deviation inputs  $\mathcal{S}_{\Delta u}(t, z)$  from (6.30), we bound the objective function

in (6.3a) from above by

$$\begin{aligned}
 & \max_{x(t,z) \in \mathcal{R}(t,z)} \left\{ \|x(t_f, z) - x_f\|_1 + \zeta \int_0^{t_f} \|u(\tau, x, z)\|_1 d\tau \right\} \\
 = & \max_{\substack{x_{\text{ff}}(t,P) \in \mathcal{R}_{x_{\text{ff}}}(t,P) \\ \Delta x(t,z) \in \mathcal{S}_{\Delta x}(t,z) \\ \Delta u(t,z) \in \mathcal{S}_{\Delta u}(t,z) \\ \|\beta\|_\infty \leq 1}} \left\{ \|x_{\text{ff}}(t_f, P) + \Delta x(t, z) - x_f\|_1 \right. \\
 & \left. + \zeta \sum_{j=0}^{m_c-1} \sum_{k=0}^{m_r-1} \int_{jt_c+kr}^{jt_c+(k+1)\delta} \left\| \bar{u}(\beta, P^{(j)}) + \Delta u(\tau, z) \right\|_1 d\tau \right\} \\
 \leq & \max_{\|\xi\|_\infty \leq 1} \|c_{x_{\text{ff}}}(t_f, P) - x_f + G_{x_{\text{ff}}}(t_f, P) \xi\|_1 + \max_{\Delta x(t_f, z) \in \mathcal{Z}_{\Delta x}(t_f, z)} \|\Delta x(t_f, z)\|_1 \\
 & + \zeta \sum_{j=0}^{m_c-1} \sum_{k=0}^{m_r-1} \int_{jt_c+kr}^{jt_c+(k+1)\delta} \left( \max_{\|\xi\|_\infty \leq 1} \|c_{u_{\text{ff}}}(P^{(j)}) + G_{u_{\text{ff}}}(P^{(j)}) \xi\|_1 \right. \\
 & \left. + \max_{\Delta u(\tau, z) \in \mathcal{Z}_{\Delta u}(\tau, z)} \|\Delta u(\tau, z)\|_1 \right) d\tau \\
 \leq & \left\| \begin{bmatrix} c_{x_{\text{ff}}}(t_f, P) - x_f \\ [G_{x_{\text{ff}}}(t_f, P)]_{(\cdot)} \end{bmatrix} \right\|_1 + 1_{n_x}^T \max \left( \rho_{\Delta x}^{(\mathcal{Z})}(-I_{n_x}, t_f, z), \rho_{\Delta x}^{(\mathcal{Z})}(I_{n_x}, t_f, z) \right) \\
 & + \zeta \delta \sum_{j=0}^{m_c-1} \sum_{k=0}^{m_r-1} \left( \left\| \begin{bmatrix} c_{u_{\text{ff}}}(P^{(j)}) \\ [G_{u_{\text{ff}}}(P^{(j)})]_{(\cdot)} \end{bmatrix} \right\|_1 \right. \\
 & \left. + 1_{n_u}^T \max \left( \rho_{\Delta u}^{(\mathcal{Z})}(-I_{n_u}, \tau^{(j,k)}, z), \rho_{\Delta u}^{(\mathcal{Z})}(I_{n_u}, \tau^{(j,k)}, z) \right) \right),
 \end{aligned}$$

where

$$\mathcal{Z}_{x_{\text{ff}}}(t, P) = \hat{\mathcal{Z}}(\mathcal{R}_{x_{\text{ff}}}(t, P)) = \langle c_{x_{\text{ff}}}(t, P), G_{x_{\text{ff}}}(t, P) \rangle_Z, \quad (6.31a)$$

$$\mathcal{Z}_{\Delta x}(t, z) = \hat{\mathcal{Z}}(\mathcal{R}(t, z) \oplus_e (-\mathcal{R}_{x_{\text{ff}}}(t, P))), \quad (6.31b)$$

$$\mathcal{Z}_{u_{\text{ff}}}(t, P) = \hat{\mathcal{Z}}\left(\left\{ \bar{u}(\beta, P^{(j)}) \right\}_\beta\right) = \langle c_{u_{\text{ff}}}(P^{(j)}), G_{u_{\text{ff}}}(P^{(j)}) \rangle_Z, \quad (6.31c)$$

$$\mathcal{Z}_{\Delta u}(t, z) = K(t, z) \mathcal{Z}_{\Delta x}(t, z), \quad (6.31d)$$

are the zonotope outer approximations of  $\mathcal{R}_{x_{\text{ff}}}(t, P)$ ,  $\mathcal{S}_{u_{\text{ff}}}(t, P)$ ,  $\mathcal{S}_{\Delta x}(t, z)$ , and  $\mathcal{S}_{\Delta u}(t, z)$  with  $\rho_{x_{\text{ff}}}^{(\mathcal{Z})}(\cdot, t, z)$ ,  $\rho_{\Delta x}^{(\mathcal{Z})}(\cdot, t, z)$ ,  $\rho_{u_{\text{ff}}}^{(\mathcal{Z})}(\cdot, t, z)$ , and  $\rho_{\Delta u}^{(\mathcal{Z})}(\cdot, t, z)$  denoting the corresponding support functions: The first inequality follows from the separation of  $x_{\text{ff}}$  and  $\Delta x$  as well as  $u_{\text{ff}}$  and  $\Delta u$  using the triangle inequality, where we replaced all sets with their respective zonotope outer approximations as defined in (6.31). The second inequality then follows from applying the triangle inequality again to separate the center and the generators for states and inputs (also see Sec. 5.2.4). Furthermore, it holds that

$$\max_{\Delta x(t,z) \in \mathcal{Z}_{\Delta x}(t,z)} \|\Delta x(t_f, z)\|_1 = 1_{n_x}^T \max \left( \rho_{\Delta x}^{(\mathcal{Z})}(-I_{n_x}, t_f, z), \rho_{\Delta x}^{(\mathcal{Z})}(I_{n_x}, t_f, z) \right),$$

which follows from  $0 \in \mathcal{Z}_{\Delta x}(t, z)$  and the central symmetry of  $\mathcal{Z}_{\Delta x}(t, z)$  since it is a zonotope. Similarly, it holds that

$$\max_{\Delta u(\tau, z) \in \mathcal{Z}_{\Delta u}(\tau, z)} \|\Delta u(\tau, z)\|_1 = 1_{n_u}^T \max \left( \rho_{\Delta u}^{(\mathcal{Z})}(-I_{n_u}, \tau, z), \rho_{\Delta u}^{(\mathcal{Z})}(I_{n_u}, \tau, z) \right).$$

**Constraints:** We conservatively encode the constraints of (6.3) as

$$\forall t \in [0, t_f] : \mathcal{Z}_{u_{\text{ff}}}(t, P) \oplus \mathcal{Z}_{\Delta u}(t, z) \subseteq \mathcal{U}, \quad (6.32)$$

$$\mathcal{Z}_{x_{\text{ff}}}([0, t_f], z) \oplus \mathcal{Z}_{\Delta x}([0, t_f], z) \subseteq \mathcal{X}, \quad (6.33)$$

$$\mathcal{Z}_{x_{\text{ff}}}(t_f, z) \oplus \mathcal{Z}_{\Delta x}(t_f, z) \subseteq \mathcal{X}_f, \quad (6.34)$$

where we used (6.29) and (6.30) and replaced all sets with their respective zonotope outer approximations from (6.31). Let

$$\rho_x^{(\mathcal{Z})}(\cdot, \tau, z) = \max_{\bar{\tau} \in \mathcal{T}_{\tau}(m_c m_r)} \left( \rho_{x_{\text{ff}}}^{(\mathcal{Z})}(\cdot, \bar{\tau}, z) + \rho_{\Delta x}^{(\mathcal{Z})}(\cdot, \bar{\tau}, z) \right),$$

$$\rho_u^{(\mathcal{Z})}(\cdot, \tau, z) = \max_{\bar{\tau} \in \mathcal{T}_{\tau}(m_c m_r)} \left( \rho_{u_{\text{ff}}}^{(\mathcal{Z})}(\cdot, \bar{\tau}, z) + \rho_{\Delta u}^{(\mathcal{Z})}(\cdot, \bar{\tau}, z) \right),$$

with

$$\mathcal{T}_{\tau}(n) = \left\{ \bar{\tau} = [o, o+1] \frac{t_f}{n} \mid 0 \leq o \leq n-1, \tau \cap \bar{\tau} \neq \emptyset \right\}, \quad (6.35)$$

for some interval  $\tau \in [0, t_f]$ . Following Prop. 4.6, the set containment constraints in (6.32) to (6.34) can thus be encoded by

$$g(z) = \begin{bmatrix} \rho_u^{(\mathcal{Z})}(C_{\mathcal{U}}, [0, t_f], z) - d_{\mathcal{U}} \\ \rho_x^{(\mathcal{Z})}(C_{\mathcal{X}}, [0, t_f], z) - d_{\mathcal{X}} \\ \rho_x^{(\mathcal{Z})}(C_{\mathcal{X}_f}, t_f, z) - d_{\mathcal{X}_f} \end{bmatrix} \leq 0, \quad (6.36)$$

since the bounded input constraints  $\mathcal{U} = \langle C_{\mathcal{U}}, d_{\mathcal{U}} \rangle_H$ , state constraints  $\mathcal{X} = \langle C_{\mathcal{X}}, d_{\mathcal{X}} \rangle_H$ , and final state constraints  $\mathcal{X}_f = \langle C_{\mathcal{X}_f}, d_{\mathcal{X}_f} \rangle_H$  are all given as H-polytopes.

**Cost:** Thus, we propose the cost function (also see (5.42))

$$\begin{aligned} J(z) &= \|c_{x_{\text{ff}}}(t_f, P) - x_f\|_1 + \left\| [G_{x_{\text{ff}}}(t_f, P)]_{(\cdot)} \right\|_1 \\ &\quad + 1_{n_x}^T \max \left( \rho_{\Delta x}^{(\mathcal{Z})}(-I_{n_x}, t_f, z), \rho_{\Delta x}^{(\mathcal{Z})}(I_{n_x}, t_f, z) \right) \\ &\quad + \zeta \delta \sum_{j=0}^{m_c-1} \sum_{k=0}^{m_r-1} \left( \left\| \begin{bmatrix} c_{u_{\text{ff}}}(P^{(j)}) \\ G_{u_{\text{ff}}}(P^{(j)})_{(\cdot)} \end{bmatrix} \right\|_1 \right. \\ &\quad \left. + 1_{n_u}^T \max \left( \rho_{\Delta u}^{(\mathcal{Z})}(-I_{n_u}, \tau^{(j,k)}, z), \rho_{\Delta u}^{(\mathcal{Z})}(I_{n_u}, \tau^{(j,k)}, z) \right) \right) \\ &\quad + \sigma \|\max(0, g(z))\|_1, \end{aligned} \quad (6.37)$$

where we penalize the constraint violation by choosing the exact penalty multiplier  $\sigma \in \mathbb{R}_+$  large enough such that constraint violations always dominate all other terms in (6.37) (see Sec. 5.4.1).

### 6.3.4 Parameterized Reachable Set

In Sec. 5.4.2, we described how the feedforward trust-region radius  $\gamma$  can be used to construct an approximation of the feedforward reachable set for the iPGSC approach that can be brought arbitrarily close to the tight outer approximation of the closed-loop feedforward reachable set. In this section, we derive how this approximation can be used to approximate the closed-loop reachable set for the combined controller – which is required to solve (6.3) – without recomputing outer approximations of the closed loop reachable set during each optimization iteration.

Using the state separation from (6.29) allows us to compute the parameterized feedforward reachable set of  $x_{\text{ff}}$  independently of  $\Delta x$ . Since the parameterized feedforward reachable set can be computed as described in Sec. 5.4.2, we first derive an approximated flow for  $\Delta x$ , and then discuss how this flow can be used to find an efficient approximation for the set of deviation vectors  $\mathcal{S}_{\Delta x}(t, z)$ . That said, we only need the support function of this approximation in all directions which are normal to the half-spaces of the polytopic input, state, and final state constraints (see (6.37)); since  $\mathcal{R}_{\text{ext}}(t, \bar{z})$  is available for the current controller parameters  $\bar{z}$  (see Alg. 6) and thus  $\rho_{\Delta x}^{(\mathcal{Z})}(\cdot, t, \bar{z})$  can be efficiently computed, we next construct an approximation for the support function of  $\Delta x$  that only approximates the unknown difference from  $\rho_{\Delta x}^{(\mathcal{Z})}(\cdot, t, \bar{z})$  to  $\rho_{\Delta x}^{(\mathcal{Z})}(\cdot, t, z)$ .

**Approximated reachable set:** Let the sequence of LTI systems for  $\Delta x$  as derived in Sec. 6.3.1 with  $A_{\text{cl}}(t, z) = A_{\text{cl}}^{(k)}(z)$  and  $E(t, P) = E^{(k)}(P)$  for  $0 \leq k \leq m_c m_d$  be given, where instead of  $m_c$  LTI systems, we linearize  $m_c m_d$  times: One may introduce  $m_d \in \mathbb{N}_+$  to obtain a better approximation of the original nonlinear dynamics. Since these LTI systems only approximate the flow for  $\Delta x$ , we can only compute an approximation of  $\mathcal{S}_{\Delta x}(t, z)$ . To that end, we compute the reachable set of the linearized flow in (6.19) using reachability analysis for linear systems (see Sec. 2.7.1)

$$\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{Z})}(t + \delta, z) = e^{A_{\text{cl}}(t, z)\delta} \tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{Z})}(t, z) \oplus \int_0^\delta e^{A_{\text{cl}}(t, z)(\delta - \tau)} E(t, P) \mathcal{W} d\tau, \quad (6.38)$$

for  $t + \delta \leq t_f$  with  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{Z})}(0, z) = \{0\}$ . In contrast to  $\mathcal{S}_{\Delta x}(t, z)$ , its approximation  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{Z})}(t, z)$  has a flow and thus is a reachable set. Since  $\mathcal{W}$  is given as a zonotope and the dynamics consist of  $m_c m_d$  LTI systems, one can compute an arbitrarily tight outer approximation of the reachable set in (6.38) [114, Sec. IV.A]. However, since  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{Z})}(t, z)$  will be formed in each optimization iteration of the trust-region subproblem problem later, its growing representation size due to the repeated Minkowski sums (see Sec. 2.6.5) when increasing  $m_c$  or  $m_d$  means additional computational effort, especially since we want to compute both Jacobian and Hessian matrices for the optimization of the trust-region subproblem. To maintain a constant representation size, we thus propose to evaluate (6.38) using ellipsoids. With  $\hat{\mathcal{W}} = \hat{\mathcal{E}}(\mathcal{W})$  being an ellipsoid outer approximation of  $\mathcal{W}$  according to Th. 4.1, we compute the ellipsoidal reachable set of

the linearized flow of the  $m_c m_d$  LTI systems with [65, Lem. 2.2]

$$\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t + \delta, z) = e^{A_{cl}(t,z)\delta} \tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t, z) \oplus \int_0^\delta e^{A_{cl}(t,z)\phi} E(t, P) \mathcal{E}(0, Q_W) d\phi. \quad (6.39)$$

Since the representation size of ellipsoids remains constant, there is no need to manually approximate the integral in (6.39) and we can directly use numerical integration methods (see Sec. 6.3.7 for more details on computing (6.39)). However, while ellipsoids are closed under linear maps, the Minkowski sum of two ellipsoids is not an ellipsoid in general (see Sec. 2.6.3). While efficient methods for the outer approximation of the Minkowski sum exist [47], finding the minimum-volume ellipsoid is not required for our purposes since we only need a good enough approximation. To save computation time, we instead propose to approximate this Minkowski sum by simply adding their respective shape matrices. While this does not result in an outer approximation, we argue that it is a good approximation of the Minkowski sum up to a scaling factor: First, for two ellipsoids  $\langle 0, U \rangle_E$  with  $U \in \mathbb{S}_{++}^{n \times n}$  and  $\langle 0, V \rangle_E$  with  $V \in \mathbb{S}_{++}^{n \times n}$ , it holds that  $\langle 0, U \rangle_E \oplus \langle 0, V \rangle_E \subseteq \langle 0, \frac{1}{\alpha_1} U + \frac{1}{\alpha_2} V \rangle_E$  for  $\alpha_1 + \alpha_2 = 1$  and  $\alpha_1, \alpha_2 > 0$  [28, Th. 4.2]. Choosing  $\alpha_1 = \alpha_2 = \frac{1}{2}$  thus means that an outer approximation of the Minkowski sum is given by the sum of their shape matrices, scaled by a scalar factor (here 2). Second, we can interpret the two ellipsoids  $\langle 0, U \rangle_E$  and  $\langle 0, V \rangle_E$  as confidence ellipsoids of two zero-mean, independently distributed Gaussian random variables with covariance matrices  $Q$  and  $U$ , respectively [81]. Then  $\langle 0, U \rangle_E$  and  $\langle 0, V \rangle_E$  describe an area of the support of the two underlying random variables such that any realization of these random variables lies within  $\langle 0, U \rangle_E$  and  $\langle 0, V \rangle_E$  with some probability. It follows that the covariance matrix of the sum of these two random variables is given by the sum of the covariance matrices of each random variable.

**Approximated support function:** So denote with  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t, z) = \mathcal{E}(0, D_{\Delta x}(t, z))$  the propagation of  $\bar{W}$  through (6.39), where we replace the Minkowski sum of two ellipsoids with the sum of their respective shape matrices. We now derive the approximated support function  $\tilde{\rho}_{\Delta x}(\cdot, t, z)$  of  $\rho_{\Delta x}^{(\mathcal{Z})}(\cdot, t, z)$ .

We start with the support function of  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t, z)$ , which is given by

$$\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z) = \sqrt{l^T D_{\Delta x}(t, z) l} + \epsilon, \quad (6.40)$$

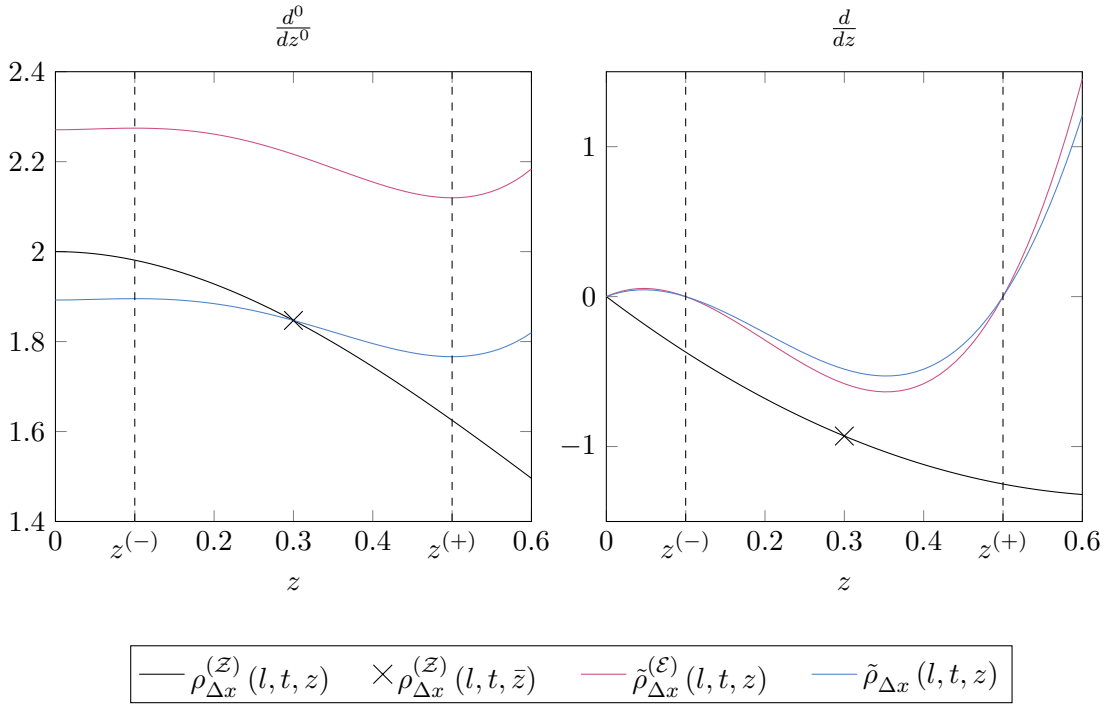
where  $l \in \mathbb{R}^{n_x}$  and  $\epsilon \ll 1$  to achieve differentiability for all  $z$  since

$$\forall l \in \mathbb{R}^{n_x} : l^T D_{\Delta x}(t, z) l \geq 0,$$

due to  $D_{\Delta x}(t, z) \succeq 0$ . Using the support function  $\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)$  of  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t, z)$  and the known support function value  $\rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z})$  at the current controller parameters  $\bar{z}$ , we are now ready to define the approximation of the support function  $\rho_{\Delta x}^{(\mathcal{Z})}(l, t, z)$  using the support function  $\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)$  of the ellipsoidal reachable set  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t, z)$  as

$$\tilde{\rho}_{\Delta x}(l, t, z) = \rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z}) + \chi_{\Delta x}(l, t) \left( \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z) - \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, \bar{z}) \right), \quad (6.41)$$





**Figure 6.4:** Visualization (1D) of the approximation of the support function  $\rho_{\Delta x}^{(\mathcal{Z})}(l, t, z)$ : Generally, we assume that  $\frac{d}{dz}\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)$  approximates the  $\frac{d}{dz}\rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z})$  reasonably well for  $z \in [z^{(-)}, z^{(+)}] = [\bar{z} - r, \bar{z} + r]$  with a given radius  $r \in \mathbb{R}_+$ ; for illustrative purposes, we here thus assume  $\text{sign}(\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)) = \text{sign}(\rho_{\Delta x}^{(\mathcal{Z})}(l, t, z))$  for  $z \in [z^{(-)}, z^{(+)}]$  (see right plot). With appropriate scaling, we obtain the reasonable approximation  $\rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z})$  with  $\tilde{\rho}_{\Delta x}(l, t, \bar{z}) = \rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z})$  (see left plot).

where  $\chi_{\Delta x}(l, t) = \frac{\rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z})}{\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, \bar{z})}$ : Since  $\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)$  is calculated based on various approximations we made previously, we do not assume that its magnitude is accurate. Instead, we define the scaling function  $\chi_{\Delta x}(l, t)$  which is chosen such that the scaled approximation is identical to the support function of the outer approximation, i.e.

$$\rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z}) = \chi_{\Delta x}(l, t) \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, \bar{z}),$$

holds for  $z = \bar{z}$ . Using the scaled approximation  $\chi_{\Delta x}(l, t) \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)$ , we then only approximate the unknown difference

$$\rho_{\Delta x}^{(\mathcal{Z})}(l, t, z) - \rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z}) \approx \chi_{\Delta x}(l, t) \left( \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z) - \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, \bar{z}) \right),$$

so that we recover  $\tilde{\rho}_{\Delta x}(l, t, \bar{z}) = \rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z})$  by substitution of  $z = \bar{z}$  into (6.41). Fig. 6.4 visualizes this approximation using a 1D example.

An approximation of the support function for the set of input deviations  $S_{\Delta u}(t, z)$  can be derived analogously. Since  $\mathcal{S}_{u_{\text{ff}}}(t, z)$  is known, it only remains to compute an

approximation of the support function of  $\mathcal{S}_{\Delta u}(t, z)$  to approximate the input constraint in (6.3b). It follows from (6.6) that  $S_{\Delta u}(t, z) = K(t, z) S_{\Delta x}(t, z)$ , and thus

$$\tilde{\mathcal{S}}_{\Delta u}^{(\mathcal{E})}(t, z) = K(t, z) \tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t, z) = \left\langle 0, K(t, z) D_{\Delta x}(t, z) K(t, z)^T \right\rangle_E.$$

With

$$\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, t, z) = \sqrt{K(t, z) D_{\Delta x}(t, z) K(t, z)^T + \epsilon}, \quad (6.42)$$

given analogously to (6.40), we can – equivalently to (6.41) – define

$$\tilde{\rho}_{\Delta u}(l, t, z) = \rho_{\Delta u}^{(\mathcal{Z})}(l, t, \bar{z}) + \chi_{\Delta u}(l, t) \left( \tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, t, z) - \tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, t, \bar{z}) \right), \quad (6.43)$$

as an approximation to the support function  $\rho_{\Delta u}^{(\mathcal{Z})}(\cdot, t, z)$  of  $\mathcal{Z}_{\Delta u}(t, z)$  as defined in (6.31d) with  $\chi_{\Delta u}(l, t) = \frac{\rho_{\Delta u}^{(\mathcal{Z})}(l, t, \bar{z})}{\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, t, \bar{z})}$ .

Combining (6.29) with (6.41) and (6.30) with (6.43) finally yields

$$\tilde{\rho}_x(l, t, z) = \tilde{\rho}_{x_{\text{ff}}}^{(\mathcal{Z})}(l, t, P) + \tilde{\rho}_{\Delta x}(l, t, z), \quad (6.44)$$

$$\tilde{\rho}_u(l, t, z) = \rho_{u_{\text{ff}}}^{(\mathcal{Z})}(l, t, P) + \tilde{\rho}_{\Delta u}(l, t, z), \quad (6.45)$$

where we used (2.26),  $\tilde{\rho}_{x_{\text{ff}}}^{(\mathcal{Z})}(l, t, P)$  denotes the support function of the feedforward reachable set approximation  $\hat{\mathcal{Z}}(\tilde{\mathcal{R}}_{x_{\text{ff}}}(t, P))$  according to Sec. 5.4.2, and  $\rho_{u_{\text{ff}}}^{(\mathcal{Z})}(l, t, P)$  denotes the support function of  $\mathcal{Z}_{u_{\text{ff}}}(t, z)$  as defined in (6.31c). Choosing  $m_d$  large enough further motivates the use of (6.44) and (6.45) to obtain the approximated support functions for time intervals (instead of time points in (6.44) and (6.45))

$$\begin{aligned} \tilde{\rho}_{\Delta x}(l, [k, k+1] \delta, z) &\approx \tilde{\rho}_{\Delta x}(l, (k+1) \delta, z), \\ \tilde{\rho}_{\Delta u}(l, [k, k+1] \delta, z) &\approx \tilde{\rho}_{\Delta u}(l, (k+1) \delta, z). \end{aligned}$$

For later convenience, we introduce the definition of the support function approximations over any time interval  $\tau \subseteq [0, t_f]$  as

$$\tilde{\rho}_x(l, \tau, z) = \max_{\bar{\tau} \in \mathcal{T}_{\tau}(m_c m_d)} \left( \max_{\tilde{\tau} \in \mathcal{T}_{\bar{\tau}}(m_c m_r)} \tilde{\rho}_{x_{\text{ff}}}^{(\mathcal{Z})}(l, \tilde{\tau}, P) + \tilde{\rho}_{\Delta x}(l, \bar{\tau}, z) \right), \quad (6.46)$$

$$\tilde{\rho}_u(l, \tau, z) = \max_{\bar{\tau} \in \mathcal{T}_{\tau}(m_c)} \left( \rho_{u_{\text{ff}}}^{(\mathcal{Z})}(l, \tilde{\tau}, P) + \max_{\tilde{\tau} \in \mathcal{T}_{\bar{\tau}}(m_c m_d)} \tilde{\rho}_{\Delta u}(l, \bar{\tau}, z) \right), \quad (6.47)$$

where  $\mathcal{T}_{\tau}(\cdot)$  is defined in (6.35), we assumed that  $m_r \geq m_d$  with  $m_r$  being the number of reachability steps per duration  $t_c$  and  $m_d$  being the number of additional linearization points for (6.20) to (6.22) per duration  $t_c$ .

### 6.3.5 Controller Computation

Ideally, we want to optimize over the controller cost in (6.37). However, its evaluation requires the computation of an outer approximation of the closed-loop extended reachable set for each controller parameter value  $z$ , so that direct optimization over (6.37) is computationally expensive. For the efficient, simultaneous synthesis of the combined controller, we thus need to find a smooth approximation to (6.37). In Sec. 6.3.4, we already derived expressions for the approximated support function for both the combined state and input set. In this section, we use these expressions to approximate (6.37).

**Trust Region:** As demonstrated in Sec. 5.4, the trust-region radius  $\gamma$  is required to restrict  $P$  to a local, compact neighborhood around the current controller guess  $\bar{P}$ . Similarly, we now introduce the trust-region radius  $\eta \in (0, 1]$  to restrict  $Q$  and  $R$  to local neighborhoods of  $\bar{Q} \in \mathcal{M}^{(Q)}$  and  $\bar{R} \in \mathcal{M}^{(R)}$ , respectively, i.e.

$$\begin{aligned} Q &\in \mathcal{M}_\eta^{(Q)}(\bar{Q}) \cap \mathcal{M}^{(Q)}, \\ R &\in \mathcal{M}_\eta^{(R)}(\bar{R}) \cap \mathcal{M}^{(R)}, \end{aligned}$$

where  $\mathcal{M}^{(Q)} \subset \mathbb{S}_{++}^{n_x \times n_x}$  and  $\mathcal{M}^{(R)} \subset \mathbb{S}_{++}^{n_u \times n_u}$  are bounded, given sets from which  $Q$  and  $R$  are chosen, and where

$$\mathcal{M}_\eta^{(Q)}(\bar{Q}) = \left\{ \bar{Q} + h_Q(\eta M) \mid \|M\|_\infty \leq 1 \right\}, \quad (6.48)$$

$$\mathcal{M}_\eta^{(R)}(\bar{R}) = \left\{ \bar{R} + h_R(\eta N) \mid \|N\|_\infty \leq 1 \right\}, \quad (6.49)$$

are bounded sets defined by the trust region  $\eta$ . Here,  $h_Q : [-1, 1]^{p \times q} \mapsto \mathbb{S}^{n_x \times n_x}$  and  $h_R : [-1, 1]^{r \times s} \mapsto \mathbb{S}^{n_u \times n_u}$  can be any sufficiently smooth functions which generate symmetric matrices. That said, this creates semi-definite constraints in the final optimization problem which can significantly impact solver times. Therefore, we restrict  $Q$  and  $R$  to diagonal matrices, i.e., we set

$$\begin{aligned} h_Q(\omega) &= \text{diag}(\bar{q} - \underline{q}) \omega, \\ h_R(\xi) &= \text{diag}(\bar{r} - \underline{r}) \xi, \\ \mathcal{M}^{(Q)} &= \left\{ \text{diag}(\tilde{q}) \mid \tilde{q} \in [\underline{q}, \bar{q}] \right\}, \\ \mathcal{M}^{(R)} &= \left\{ \text{diag}(\tilde{r}) \mid \tilde{r} \in [\underline{r}, \bar{r}] \right\}, \end{aligned}$$

with  $\bar{q} \geq \underline{q} \in \mathbb{R}_+^{n_x}$ ,  $\bar{r} \geq \underline{r} \in \mathbb{R}_+^{n_u}$ , and dependent factors  $\omega \in [-1, 1]^{n_x}$  and  $\xi \in [-1, 1]^{n_u}$ . Additionally, we set  $\underline{q}_1 = \bar{q}_1 = 1$  to eliminate the invariance of the LQR computation under equal scaling of  $Q$  and  $R$ . We remark that more general choices are possible while still avoiding positive definite constraints, e.g., one can generate symmetric, strictly diagonally dominant matrices with real positive entries, which is sufficient to ensure positive definiteness (follows from the Gershgorin Circle Theorem [11, Th. 0]).

**Optimization Problem:** Starting from the current controller parameters collected in  $\bar{z} = [\bar{P}_{(\cdot)}^T, \bar{Q}_{(\cdot)}^T, \bar{R}_{(\cdot)}^T]^T$ , an approximation to the controller cost in (6.37) is given by

$$\min_{z,s} \tilde{J}_{\text{TR}}(z, s), \quad (6.50a)$$

$$\text{s.t. } \tilde{g}(z) \leq s, \quad (6.50b)$$

$$P \in \mathcal{P}_\gamma(\bar{P}) \cap [-1, 1]^{m_c n_u \times a}, \quad (6.50c)$$

$$Q \in \mathcal{M}_\eta^{(Q)}(\bar{Q}) \cap \mathcal{M}^{(Q)}, \quad (6.50d)$$

$$R \in \mathcal{M}_\eta^{(R)}(\bar{R}) \cap \mathcal{M}^{(R)}, \quad (6.50e)$$

$$s = [s_U^T, s_{\mathcal{X}}^T, s_{\mathcal{X}_f}^T]^T \geq 0, \quad (6.50f)$$

with

$$\begin{aligned} \tilde{J}_{\text{TR}}(z, s) = & \|\tilde{c}_{x_{\text{ff}}}(t_f, P) - x_f\|_1 + \left\| \left[ \tilde{G}_{x_{\text{ff}}}(t_f, P) \right]_{(\cdot)} \right\|_1 \\ & + 1_{n_x}^T \max(\tilde{\rho}_{\Delta x}(-I_{n_x}, t_f, z), \tilde{\rho}_{\Delta x}(I_{n_x}, t_f, z)) \\ & + \zeta \tilde{\delta} \sum_{j=0}^{m_c-1} \sum_{k=0}^{m_d-1} \left( \left\| \begin{bmatrix} c_{u_{\text{ff}}}(P^{(j)}) \\ G_{u_{\text{ff}}}(P^{(j)})_{(\cdot)} \end{bmatrix} \right\|_1 \right. \\ & \left. + 1_{n_u}^T \max(\tilde{\rho}_{\Delta u}(-I_{n_u}, \tau^{(j,k)}, z), \tilde{\rho}_{\Delta u}(I_{n_u}, \tau^{(j,k)}, z)) \right) \\ & + \sigma \|s\|_1, \end{aligned} \quad (6.51)$$

and

$$\tilde{g}(z) = \begin{bmatrix} \tilde{\rho}_u(C_U, [0, t_f], z) - d_U \\ \tilde{\rho}_x(C_{\mathcal{X}}, [0, t_f], z) - d_{\mathcal{X}} \\ \tilde{\rho}_x(C_{\mathcal{X}_f}, t_f, z) - d_{\mathcal{X}_f} \end{bmatrix}, \quad (6.52)$$

where  $\sigma \in \mathbb{R}_+$  and  $\tilde{\delta} = \frac{t_f}{m_c m_d}$ :

- (6.50a) and (6.51): We replace the outer approximations used in (6.37) with their respective approximations derived previously. Additionally, we have  $m_c m_d$  time intervals instead of  $m_c m_r$  as in Sec. 5.4.1 and thus adapt the sum accordingly. Lastly, we replace the constraint violation with a slack variable (see next bullet point for more details).
- (6.50b): Instead of incorporating the constraints in the objective function, we enforce input, state, and final state constraints with  $\tilde{g}(z) \leq s$  directly; however, we relax these constraints with  $s \geq 0$  so that (6.50) always has a feasible solution, even if there is no feasible solution for  $s = 0$ . To still ensure feasibility, we penalize  $s$  in the objective function with  $\sigma \in \mathbb{R}_+$ .

- (6.50c) to (6.50e): Because we replaced tight outer approximations of the closed-loop reachable sets and input sets with their respective approximations, (6.50) is only accurate in a neighborhood of  $\bar{z}$ . We therefore restrict  $z$  to a bounded set around  $\bar{z}$ .

Since (6.50) can be written as the abstracted synthesis problem, it can be reformulated as a smooth non-convex optimization problem according to Prop. 3.1 and the KKT conditions are always necessary for this smooth reformulation due to Prop. 3.2. We show in Sec. 6.3.6 that (6.50) can approximate (6.37) arbitrarily closely. For this proof, we require a cost value of the trust-region solution only dependent on  $\hat{z}$ , which we introduce subsequently.

**Proposition 6.1** (Stationary Objective Value). *At a first-order critical point of the smooth reformulation of (6.50) with optimizer  $\hat{z}$ , the optimal objective value is given by*

$$\tilde{J}(\hat{z}) = \tilde{J}_{\text{TR}}(\hat{z}, \max(0, \tilde{g}(\hat{z}))), \quad (6.53)$$

where  $\tilde{g}(\hat{z})$  is given by (6.52).

*Proof.* Follows from Prop. 3.1. □

### 6.3.6 Tuning of Trust-Region Radii

Since the trust-region subproblem approximates the cost to find a new controller candidate, it is paramount to analyze how close this approximation is to the real cost; otherwise, optimizing the trust-region problem might not decrease the actual cost of the controller. However, if the trust-region subproblem and thus its approximated objective value (see Prop. 6.1) is close enough to the real cost, optimizing (6.50) – and thus reducing the approximated cost in (6.53) – also reduces the real cost in (6.37). In this section, we therefore show that the trust-region radii  $\gamma$  and  $\eta$  can be tuned so that the approximated cost described in Prop. 5.2 stays arbitrarily close to the actual controller cost from (6.37).

To that end, we first define  $\tilde{J}_{\Delta x}(z)$  equivalently to  $\tilde{J}(z)$  but replace all approximated feedforward sets and corresponding feedforward support functions with their respective formally correct versions using the closed-loop feedforward reachable set available at  $\bar{z}$ . Intuitively,  $\tilde{J}_{\Delta x}(z)$  then contains no inaccuracies (compared to  $J(z)$ ) from approximating the feedforward reachable set and hence is only inaccurate due to inaccurate support function approximations for  $\Delta x$  and  $\Delta u$ . We are now ready to state the main result of this section.

**Theorem 6.1** (Accurate Trust-Region Subproblem). *Let  $\hat{z}$  denote a first-order critical point of the reformulation of the trust-region subproblem in (6.50). Then*

$$\left| J(\hat{z}) - \tilde{J}(\hat{z}) \right| \leq \epsilon, \quad (6.54)$$

## 6 Piecewise Constant Controller Synthesis with Continuous State Feedback

is achieved after a finite number of trust-region iterations with given tolerance  $\epsilon > 0$  if we adapt the trust-region radii according to

$$\gamma \leftarrow \min(1, v(\max(e_{\text{ff},\gamma}, e_{\Delta,\gamma}))\gamma), \quad (6.55)$$

$$\eta \leftarrow \min(1, v(e_{\Delta,\eta})\eta), \quad (6.56)$$

with

$$e_{\text{ff},\gamma} = \left| \tilde{J}(\hat{z}) - \tilde{J}_{\Delta x}(\hat{z}) \right|, \quad (6.57)$$

$$e_{\Delta,\gamma} = \left| \tilde{J}_{\Delta x}(\hat{P}, \bar{Q}, \bar{R}) - J(\hat{P}, \bar{Q}, \bar{R}) \right|, \quad (6.58)$$

$$e_{\Delta,\eta} = \left| \tilde{J}_{\Delta x}(\hat{z}) - \tilde{J}_{\Delta x}(\hat{P}, \bar{Q}, \bar{R}) - (J(\hat{z}) - J(\hat{P}, \bar{Q}, \bar{R})) \right|, \quad (6.59)$$

where  $v : (0; \infty) \mapsto \mathbb{R}_+$ , is an arbitrary, monotonically decreasing function with  $v(0) = \bar{c}$ ,  $v(\psi) = 1$ ,  $\lim_{r \rightarrow \infty} v(r) = \underline{c}$ ,  $\frac{1}{\underline{c}} > \bar{c} > 1$ , and  $0 \leq \psi < \epsilon$ .

*Proof.* The proof is structured as follows: We bound (6.54) from above by using (6.57) to (6.59), and then show that all latter expressions can be made arbitrarily small by following the tuning rules in (6.55) and (6.56).

It follows from the triangle inequality that

$$\begin{aligned} & \left| \tilde{J}(\hat{z}) - J(\hat{z}) \right| \\ &= \left| \tilde{J}(\hat{z}) - \tilde{J}_{\Delta x}(\hat{z}) + \tilde{J}_{\Delta x}(\hat{P}, \bar{Q}, \bar{R}) - J(\hat{P}, \bar{Q}, \bar{R}) \right. \\ & \quad \left. + \tilde{J}_{\Delta x}(\hat{z}) - \tilde{J}_{\Delta x}(\hat{P}, \bar{Q}, \bar{R}) - (J(\hat{z}) - J(\hat{P}, \bar{Q}, \bar{R})) \right| \\ &\leq e_{\text{ff},\gamma} + e_{\Delta,\gamma} + e_{\Delta,\eta}. \end{aligned} \quad (6.60)$$

Since  $\tilde{J}_{\Delta x}(\hat{z})$  uses the outer approximation of the feedforward reachable set, it follows that  $\lim_{\gamma \rightarrow 0} e_{\text{ff},\gamma} = 0$  since  $\lim_{\gamma \rightarrow 0} \tilde{\mathcal{R}}_{x_{\text{ff}}}(t, \hat{P}) = \tilde{\mathcal{R}}_{x_{\text{ff}}}(t, \bar{P}) = \mathcal{R}_{x_{\text{ff}}}(t, \bar{P})$  (see (5.44)). Since  $\lim_{\gamma \rightarrow 0} \tilde{J}_{\Delta x}(\hat{P}, \bar{Q}, \bar{R}) = \tilde{J}_{\Delta x}(\bar{z})$ ,  $\tilde{J}_{\Delta x}(\bar{z}) = J(\bar{z})$  by construction, and  $\lim_{\gamma \rightarrow 0} J(\hat{P}, \bar{Q}, \bar{R}) = J(\bar{z})$ , we further have  $\lim_{\gamma \rightarrow 0} e_{\Delta,\gamma} = 0$ . Lastly, it holds that  $\lim_{\eta \rightarrow 0} \tilde{J}_{\Delta x}(\hat{z}) = \tilde{J}_{\Delta x}(\hat{P}, \bar{Q}, \bar{R})$  and  $\lim_{\eta \rightarrow 0} J(\hat{z}) = J(\hat{P}, \bar{Q}, \bar{R})$  because  $\lim_{\eta \rightarrow 0} \hat{z} = \left[ \hat{P}_{(\cdot)}^T, \bar{Q}_{(\cdot)}^T, \bar{R}_{(\cdot)}^T \right]^T$ , and therefore it holds that  $\lim_{\eta \rightarrow 0} e_{\Delta,\eta} = 0$ .

Hence we can always achieve  $e_{\text{ff},\gamma} \leq \delta$ ,  $e_{\Delta,\gamma} \leq \delta$ , and  $e_{\Delta,\eta} \leq \delta$  for  $\delta > 0$  by following (6.55) and (6.56) for a finite number of steps, yielding  $\left| \tilde{J}(\hat{z}) - J(\hat{z}) \right| \leq 3\delta$ , which concludes the proof since  $\delta$  is arbitrary.  $\square$

Evaluating (6.57) to (6.59) requires – in addition to the extended closed-loop reachable set already available at  $z = \bar{z}$  and  $z = \hat{z}$  – the extended closed-loop reachable set at  $z = \left[ \hat{P}_{(\cdot)}^T, \bar{Q}_{(\cdot)}^T, \bar{R}_{(\cdot)}^T \right]^T$ : While the feedforward trust-region radius  $\gamma$  ensures that the feedforward approximation remains accurate, it also indirectly contributes to

inaccuracies in the approximated support functions of  $\Delta x$  and  $\Delta u$ . Thus, even if  $\gamma$  is small enough so that the approximation of the feedforward reachable set is accurate, the approximated support function for  $\Delta x$  and  $\Delta u$  might still be too inaccurate in  $P$ . As a result, we consider both the accuracy of the feedforward reachable set as well as the accuracy of the feedforward parameters in the disturbance support function approximation (see (6.55)) for the tuning of  $\gamma$ . Put differently, we bound

$$e_{\Delta} = \left| \tilde{J}_{\Delta x}(\hat{z}) - J(\hat{z}) \right| \leq e_{\Delta, \gamma} + e_{\Delta, \eta},$$

from above, which allows us to assess whether the error  $e_{\Delta, \gamma}$  due to  $\gamma$  is too large or the error  $e_{\Delta, \eta}$  due to  $\eta$  is too large. Alternatively, one can skip this separation by setting  $e_{\text{ff}, \gamma} = e_{\Delta, \eta} = e_{\Delta}$  and then apply (6.55) and (6.56), albeit at the cost of possibly unnecessarily shrinking either  $\gamma$  or  $\eta$ . In this thesis, we choose  $v$  from Th. 6.1 as a linear function interpolating between  $\underline{c}$  and  $\bar{c}$ ; that said, a better choice for  $v$  may be possible.

### 6.3.7 Feedback Derivative Information

For the efficient implementation of the trust-region subproblem in Sec. 6.3.5, it is beneficial to provide derivative information for the objective function and the constraints. Since the parameterized feedforward input sets are generated by the polynomial controller template and the approximations to the parameterized feedforward reachable sets are computed according to Sec. 5.4.2, both can be expressed as polynomial zonotopes and thus their derivatives are straightforward to compute. Thus, we focus on the computation of the first-order and second-order derivatives for the approximated support function  $\tilde{\rho}_{\Delta u}(l, t, z)$  and  $\tilde{\rho}_{\Delta x}(l, t, z)$  with respect to  $z \in \mathbb{R}^{n_z}$ .

**Deviation state:** For convenience, we restate (6.40) and (6.41) here:

$$\begin{aligned} \tilde{\rho}_{\Delta x}(l, t, z) &= \rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z}) + \chi_{\Delta x}(l, t) \left( \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z) - \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, \bar{z}) \right), \\ \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z) &= \sqrt{l^T D_{\Delta x}(t, z) l} + \epsilon. \end{aligned}$$

Thus, only the derivative of  $\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)$  is required since  $\rho_{\Delta x}^{(\mathcal{Z})}(l, t, \bar{z})$  and  $\chi_{\Delta x}(l, t)$  are not dependent on  $z$ . For the derivative of  $\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)$ , we require the derivative of each matrix element of the positive-definite shape matrix  $D_{\Delta x}(t, z)$ . The corresponding ellipsoid is given by  $\tilde{\mathcal{R}}_{\Delta}^{(\mathcal{E})}(t, z) = \langle 0, D_{\Delta x}(t, z) \rangle_E$  and we repeat the propagation equation for  $\bar{\mathcal{W}}$  from (6.38) here for convenience:

$$\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t + \delta, z) = e^{A_{\text{cl}}(t, z)\delta} \mathcal{R}_{\Delta x}^{(\mathcal{E})}(t, z) \oplus \int_0^{\delta} e^{A_{\text{cl}}(t, z)(\delta - \tau)} E(t, P) \bar{\mathcal{W}} d\tau, \quad (6.61)$$

where  $A_{\text{cl}}(t, z) = A(t, P) + B(t, P)K(t, z)$ .

In order to compute the Hessian matrix  $H_{\tilde{\rho}_{\Delta x}^{(\mathcal{E})}}(l, t, z)$  of  $\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, t, z)$ , we use differentials as introduced in Sec. 2.2. For the sake of simplicity and because  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t + \delta, z)$  is defined recursively, we assume that  $\tilde{\mathcal{R}}_{\Delta x}^{(\mathcal{E})}(t, z)$  with first-order differential tensor  $dD_{\Delta x}(t, z)$

and second-order differential tensor  $d^2D_{\Delta x}(t, z)$  are given and we want to compute derivative information for  $D_{\Delta x}(\hat{t}, z)$  with  $\hat{t} = t + \delta$ . From (6.20) to (6.22), it is clear that  $A(t, P) = A^{(k)}(P)$ ,  $B(t, P) = B^{(k)}(P)$ , and  $E(t, P) = E^{(k)}(P)$  for  $0 \leq k \leq m_c m_d - 1$  can be precomputed, and therefore we assume that (arguments omitted for readability) the first-order differential tensors  $dA \in \mathbb{R}^{1 \times n_z \times n_x \times n_x}$ ,  $dB \in \mathbb{R}^{1 \times n_z \times n_x \times n_u}$ , and  $dE \in \mathbb{R}^{1 \times n_z \times n_x \times n_w}$  as well as second-order differential tensors  $d^2A \in \mathbb{R}^{n_z \times n_z \times n_x \times n_x}$ ,  $d^2B \in \mathbb{R}^{n_z \times n_z \times n_x \times n_u}$ , and  $d^2E \in \mathbb{R}^{n_z \times n_z \times n_x \times n_w}$  are available (see Sec. 2.2). These can be computed by computing the Hessian matrix of each element of  $A$ ,  $B$ , and  $E$ . For this example, we directly express all results with differential tensors instead of differentials (see Sec. 2.2 for details).

We start by applying the differential operator to  $\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z)$ , yielding

$$\begin{aligned} \tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) &= \sqrt{l^T D_{\Delta x}(\hat{t}, z) l + \epsilon}, \\ d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) &= \frac{l^T dD_{\Delta x}(\hat{t}, z) l}{2\sqrt{l^T D_{\Delta x}(\hat{t}, z) l + \epsilon}} \\ &= \frac{l^T dD_{\Delta x}(\hat{t}, z) l}{2\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z)} \end{aligned} \quad (6.62)$$

$$\begin{aligned} d^2\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) &= \frac{l^T d^2D_{\Delta x}(\hat{t}, z) l}{2\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z)} - \frac{l^T dD_{\Delta x}(\hat{t}, z) l}{2\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z)^2} d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z), \\ &= \frac{l^T d^2D_{\Delta x}(\hat{t}, z) l}{2\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z)} - \frac{1}{\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z)} d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) \\ &= \frac{l^T d^2D_{\Delta x}(\hat{t}, z) l - 2d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z)}{2\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z)}, \end{aligned} \quad (6.63)$$

where  $d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) \in \mathbb{R}^{1 \times n_z \times 1 \times 1}$  and  $d^2\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) \in \mathbb{R}^{n_z \times n_z \times 1 \times 1}$ . Ignoring trailing singleton dimensions, we thus have

$$J_{\tilde{\rho}_{\Delta x}^{(\mathcal{E})}}(l, \hat{t}, z) = d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z), \quad (6.64)$$

$$H_{\tilde{\rho}_{\Delta x}^{(\mathcal{E})}}(l, \hat{t}, z) = \frac{1}{2} \left( d^2\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) + \left( d\tilde{\rho}_{\Delta x}^{(\mathcal{E})}(l, \hat{t}, z) \right)^T \right), \quad (6.65)$$

due to Th. 2.4 and Th. 2.5. However, the differential tensors  $dD_{\Delta x}(\hat{t}, z)$  and  $d^2D_{\Delta x}(\hat{t}, z)$  still need to be computed.

With the shape matrix  $D_{\Delta x}(t, z)$  as defined in (6.61), its differential tensors – when using the sum of the shape matrices as an approximation of the Minkowski sum as



### 6.3 Iterative Polynomial Reachset Optimal Control

discussed in Sec. 6.3.4 and the definition of the linear map of ellipsoids from Sec. 2.6.3 – are then given by

$$\begin{aligned}
D_{\Delta x}(\hat{t}, z) &= e^{A_{\text{cl}}(t,z)\delta} D_{\Delta x}(t, z) \left( e^{A_{\text{cl}}(t,z)} \right)^T + V(t, z), \\
dD_{\Delta x}(\hat{t}, z) &= de^{A_{\text{cl}}(t,z)\delta} D_{\Delta x}(t, z) \left( e^{A_{\text{cl}}(t,z)} \right)^T + e^{A_{\text{cl}}(t,z)\delta} dD_{\Delta x}(t, z) \left( e^{A_{\text{cl}}(t,z)} \right)^T \\
&\quad + e^{A_{\text{cl}}(t,z)\delta} D_{\Delta x}(t, z) \left( de^{A_{\text{cl}}(t,z)} \right)^T + dV(t, z), \\
d^2 D_{\Delta x}(\hat{t}, z) &= d^2 e^{A_{\text{cl}}(t,z)\delta} D_{\Delta x}(t, z) \left( e^{A_{\text{cl}}(t,z)} \right)^T + 2de^{A_{\text{cl}}(t,z)\delta} D_{\Delta x}(t, z) \left( de^{A_{\text{cl}}(t,z)} \right)^T \\
&\quad + e^{A_{\text{cl}}(t,z)\delta} d^2 D_{\Delta x}(t, z) \left( de^{A_{\text{cl}}(t,z)} \right)^T \\
&\quad + e^{A_{\text{cl}}(t,z)\delta} D_{\Delta x}(t, z) \left( d^2 e^{A_{\text{cl}}(t,z)} \right)^T + d^2 V(t, z),
\end{aligned}$$

where

$$\begin{aligned}
V(t, z) &= \int_0^\delta U(\tau, t, z) d\tau, \tag{6.66} \\
U(\tau, t, z) &= M(\tau, t, z) \bar{W} M(\tau, t, z)^T, \\
M(\tau, t, z) &= e^{A_{\text{cl}}(t,z)(t-\tau)} E(t, P),
\end{aligned}$$

such that  $\int_0^\delta e^{A_{\text{cl}}(t,z)(\delta-\tau)} E(t, P) \bar{W} d\tau = \langle 0, V(t, z) \rangle_E$  with  $\bar{W} = \langle 0, \bar{W} \rangle_E$  and  $\bar{W} \in \mathbb{S}_{++}^{n_w \times n_w}$ . That said, we still require the differential tensors of  $V(t, z)$ .

Using Prop. 2.1, the differential tensors of  $V(t, z)$  are simply

$$dV(t, z) = d \int_0^\delta U(\tau, t, z) d\tau = \int_0^\delta dU(\tau, t, z) d\tau, \tag{6.67}$$

$$d^2 V(t, z) = \int_0^\delta d^2 U(\tau, t, z) d\tau, \tag{6.68}$$

$$\begin{aligned}
dU(\tau, t, z) &= dM(\tau, t, z) \bar{W} M(\tau, t, z)^T + M(\tau, t, z) \bar{W} dM(\tau, t, z)^T, \\
d^2 U(\tau, t, z) &= d^2 M(\tau, t, z) \bar{W} M(\tau, t, z)^T + 2dM(\tau, t, z) \bar{W} dM(\tau, t, z)^T \\
&\quad + M(\tau, t, z) \bar{W} d^2 M(\tau, t, z)^T, \\
dM(\tau, t, z) &= de^{A_{\text{cl}}(t,z)(t-\tau)} E(t, P) + e^{A_{\text{cl}}(t,z)(t-\tau)} dE(t, P), \\
d^2 M(\tau, t, z) &= d^2 e^{A_{\text{cl}}(t,z)(t-\tau)} E(t, P) + 2de^{A_{\text{cl}}(t,z)(t-\tau)} dE(t, P) \\
&\quad + e^{A_{\text{cl}}(t,z)(t-\tau)} d^2 E(t, P),
\end{aligned}$$

with  $d\bar{W} = 0$  and  $d^2 \bar{W} = 0$ . We can now evaluate  $V(t, z)$ ,  $dV(t, z)$ , and  $d^2 V(t, z)$  in (6.66) to (6.68) by numerical integration for a given  $t$  and  $z$ . It now only remains to compute the differential tensors of  $e^{A_{\text{cl}}(t,z)t}$ .

Since the differential tensors of  $e^{A_{\text{cl}}(t,z)t}$  also require the the differential tensors of  $A_{\text{cl}}(t, z)$ , we first compute  $dA_{\text{cl}} = dA + dBK + BdK$  and  $d^2 A_{\text{cl}} = d^2 A + d^2 BK + dBdK + Bd^2 K$  by applying the chain rule from Prop. 2.1, where  $dK \in \mathbb{R}^{1 \times n_z \times n_u \times n_x}$  and

## 6 Piecewise Constant Controller Synthesis with Continuous State Feedback

$d^2K \in \mathbb{R}^{n_z \times n_z \times n_u \times n_x}$  are computed according to Prop. A.3 and Prop. A.4, respectively. The differential tensor of the matrix exponential with  $dt = 0$  can then be computed by using its definition, i.e.

$$e^{A_{\text{cl}}t} = \sum_{k=0}^{\infty} \frac{(A_{\text{cl}})^k t^k}{k!},$$

so that

$$\begin{aligned} de^{A_{\text{cl}}t} &= \sum_{k=1}^{\infty} \frac{t^k \left( dA_{\text{cl}}A_{\text{cl}}^{k-1} + A_{\text{cl}}d\left(A_{\text{cl}}^{k-1}\right) \right)}{k!}, \\ d^2e^{A_{\text{cl}}t} &= \sum_{k=1}^{\infty} \frac{t^k \left( d^2A_{\text{cl}}A_{\text{cl}}^{k-1} + 2dA_{\text{cl}}d\left(A_{\text{cl}}^{k-1}\right) + A_{\text{cl}}d^2\left(A_{\text{cl}}^{k-1}\right) \right)}{k!}, \end{aligned}$$

where we used

$$\begin{aligned} d\left(A_{\text{cl}}^{i+1}\right) &= dA_{\text{cl}}A_{\text{cl}}^i + A_{\text{cl}}d\left(A_{\text{cl}}^i\right), \\ d^2\left(A_{\text{cl}}^{i+1}\right) &= d^2A_{\text{cl}}A_{\text{cl}}^i + 2dA_{\text{cl}}d\left(A_{\text{cl}}^i\right) + A_{\text{cl}}d^2\left(A_{\text{cl}}^i\right), \end{aligned}$$

for  $i \geq 1$  and where the infinite series of the matrix exponential can be truncated after a finite number of steps since its convergent for any  $A_{\text{cl}}t$  [49, Prop. 2.1].

**Deviation input:** Analogously to the deviation state, we only require the derivatives of  $\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, t, z)$  as defined in (6.42), here restated for convenience:

$$\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, t, z) = \sqrt{l^T D_{\Delta u}(t, z)^T l + \epsilon},$$

where

$$\begin{aligned} D_{\Delta u}(t, z) &= K(t, z) D_{\Delta x}(t, z) K(t, z), \\ dD_{\Delta u}(t, z) &= K(t, z) dD_{\Delta x}(t, z) K(t, z) + 2dK(t, z) D_{\Delta x}(t, z) K(t, z), \\ d^2D_{\Delta u}(t, z) &= K(t, z) d^2D_{\Delta x}(t, z) K(t, z)^T + 2d^2K(t, z) D_{\Delta x}(t, z) K(t, z)^T \\ &\quad + 4dK(t, z) dD_{\Delta x}(t, z) K(t, z)^T + 2dK(t, z) D_{\Delta x}(t, z) dK(t, z)^T. \end{aligned}$$

Its derivatives are then – analogously to (6.62) and (6.63) – given by

$$\begin{aligned} d\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, \hat{t}, z) &= \frac{l^T dD_{\Delta u}(\hat{t}, z) l}{2\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, \hat{t}, z)}, \\ d^2\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, \hat{t}, z) &= \frac{l^T d^2D_{\Delta u}(\hat{t}, z) l - 2d\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, \hat{t}, z) d\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, \hat{t}, z)}{2\tilde{\rho}_{\Delta u}^{(\mathcal{E})}(l, \hat{t}, z)}, \end{aligned}$$

so that its Jacobian and Hessian matrix follow identically to (6.64) and (6.65) using Th. 2.4 and Th. 2.5.

### 6.3.8 Initial Guess

Due to the structure of the iPROC approach (also see Alg. 6), good initial guesses may drastically reduce the number of overall iterations. As a result, we shortly describe possible initial guess computations for both the feedforward controller parameters and the feedback controller parameters.

We initialize the feedforward controller parameters by executing the GSC approach described in Sec. 5.2. If polynomial feedforward controllers are to be used, PGSC can also be employed to compute an initial guess. While iPGSC may also be used, it is not as suitable for the initial guess computation since there is no point in fine tuning  $P$  for the given input capacity without considering the input capacity needed for the feedback controller.

For the feedback controller parameters  $Q$  and  $R$ , one may simply set  $Q = I_{n_x}$  and  $R = I_{n_u}$ . Since the feedback matrices are parameterized as LQR gain matrices, any choice for  $Q$  and  $R$  results in (locally) asymptotically stable systems for  $\Delta x = x - x_{\text{ff}}$ , which describes the deviation of the closed-loop, disturbed state from the undisturbed feedforward state.

### 6.3.9 Computational Complexity

We again use Ass. 5.1 and distinguish between the offline and the online computational complexity in  $n$ . For simplicity, let  $\kappa \geq 1$  hold for the order of the controller template introduced in Sec. 5.3.2 and let

$$n_z = an_u + n_x + n_u, \quad (6.69)$$

$$O(a) \geq O(n), \quad (6.70)$$

so that  $O(n_z) = O(an)$ .

We provide an upper bound on the polynomial complexity of the offline computation in Sec. 6.3.9.1 and then briefly describe the online application of the controller in Sec. 6.3.9.2.

#### 6.3.9.1 Offline Complexity

The maximum number of iterations in Alg. 6 is fixed. As a result, it suffices to derive the complexity for a single trust-region iteration. As the numerical experiments in Sec. 6.3.10 will demonstrate, this does not really affect the performance of iPROC as the number of iterations in practice is often small and independent of the state dimension.

**Controller Template:** The dominating complexity is given by  $O(an^{1+\kappa})$  according to Sec. 5.3.6.1.

**Reachable Sets:** The combined complexity of computing the closed-loop reachable set and the parameterized feedforward reachable set can be bounded from above by  $O(n^5 + n^2a^2 + n^3a \log n)$  (see Sec. 5.3.6.1).

**Trust-Region Subproblem:** Solving (6.50) involves the precomputation of the Jacobian and Hessian matrices with complexity  $O(a^2n^4)$  (see Sec. 5.3.6.1) of the feedforward reachable set approximation and the computation of the objective and constraint functions during optimization. Since the number of function evaluations  $v(\epsilon)$  can be polynomially bounded in  $\frac{1}{\epsilon}$  for the accuracy  $\epsilon > 0$  of a first-order critical point if second-order methods are used [18], we focus on the complexity of one evaluation of the objective and constraint function subsequently.

We now analyze the complexity of evaluating the approximations to the state deviation and input deviation support functions (and their corresponding derivatives) from (6.41) and (6.43). For a worst-case complexity analysis, we only consider the computation of the Hessian matrix of these approximations. For their computation, we require the feedback matrix using LQR control with complexity  $O(n_x^3)$  (see Sec. 2.4), its second-order tensor with complexity  $O(n_x^{2\omega} + n_x^4 n_z^2) = O(a^2 n^6)$  (see Prop. A.4) using (6.69), (6.70), and  $2\omega \leq 6$ , and the second-order tensor of the approximated shape matrix  $D_{\Delta x}$  which is computed via (6.39). For the computation of  $D_{\Delta x}$ , we evaluate (6.39)  $m_c m_d$  times, where one evaluation is dominated by the complexity of computing the second-order tensor of the matrix integral in (6.39). Here, computing the second-order tensor of  $e^{A_{\text{cl}}(t,z)\delta}$  dominates with complexity  $O(\vartheta n_x^\omega n_z^2) \stackrel{(6.69)}{=} O(a^2 n^{\omega+2})$ , where  $\vartheta \in \mathbb{N}_+$  denotes the finite number of terms from the infinite series of  $e^{A_{\text{cl}}(t,z)\delta}$  to compute the matrix exponential accurately enough; the complexity for each of the  $\vartheta$  terms is  $O(n_x^\omega n_z^2)$ , which follows from the complexity of multiplying two matrices but where each element multiplication is the outer product of two  $n_z$ -dimensional vectors (also see Sec. 6.3.7). To derive the computational complexity of solving the integral in (6.39), we solve the ODE corresponding to this integral with complexity  $O(b\vartheta n_x^\omega n_z^2) \stackrel{(6.69)}{=} O(a^2 n^{\omega+2})$ , where  $b \in \mathbb{N}_+$  collects the fixed number of function evaluations per ODE step and the total number of ODE steps: For instance, the classical Runge-Kutta integration scheme requires four function evaluations per step and the number of ODE steps can be assumed fixed if a solver with a fixed step size is used. Evaluation of (6.39)  $m_c m_d$  times over  $v(\epsilon)$  optimization iterations and computing the second-order differential tensor for the LQR gain matrix  $m_c$  times thus yields the overall complexity

$$O\left(v(\epsilon) m_c \left(n_x^{2\omega} + n_x^4 n_z^2 + m_d b \vartheta n_x^\omega n_z^2\right)\right) = O\left(a^2 n^6\right),$$

for one objective and constraint function evaluation of (6.50).

**Overall Complexity:** The offline complexity of iPROC can thus be bounded from above by

$$\begin{aligned} O\left(c_{\text{iPROC}}^{(\text{off})}(n)\right) &= O\left(c_{\text{GSC}}^{(\text{off})}(n) + an^{1+\kappa} + n^5 + n^2 a^2 + n^3 a \log n + a^2 n^6\right) \\ &\stackrel{(5.14)}{=} O\left(c_{\text{ref}}(n) + an^{1+\kappa} + a^2 n^6\right), \end{aligned} \quad (6.71)$$

if the GSC approach is used to compute an initial guess for the feedforward parameters (see Sec. 6.3.8) and where  $O(c_{\text{ref}}(n))$  denotes the complexity of the reference trajectory computation as detailed in Sec. 5.2.1.

### 6.3.9.2 Online Complexity

The feedforward controller evaluation has the same complexity as the PGSC approach, i.e.  $O\left(c_{\text{PGSC}}^{(\text{on})}(n)\right) \stackrel{(5.34)}{=} O(an)$  since we use the same controller template. The feedback controller evaluation is dominated by the simulation of the feedforward state with complexity  $O\left(b\left(n_x e + c_{\text{PGSC}}^{(\text{on})}(n)\right)\right) = O(an)$  using (5.34) and Ass. 5.1, where  $b \in \mathbb{N}_+$  collects the total number of function evaluations for a fixed step-size ODE solver: Evaluating the dynamics  $f$  requires  $O(e) = O(n)$  elementary operations  $e \in \mathbb{N}_+$  by Ass. 5.1 for each dimension and we need to evaluate the feedforward controller for all  $b$  function evaluations. Thus, the overall online complexity can be bounded from above by

$$O\left(c_{\text{iPROC}}^{(\text{on})}(n)\right) = O\left(c_{\text{PGSC}}^{(\text{on})}(n) + an\right) \stackrel{(5.34)}{=} O(an).$$

### 6.3.10 Experiments

In this section, we compare the performance of ROC to our novel iPROC synthesis algorithm. Since derivatives for the optimization of the synthesis problem of ROC cannot be computed analytically (see Sec. 6.3.7), ROC uses parallel computing to numerically approximate the derivative information using finite differences (see [53, Chap. 5] for details on finite differences): For the processor used in this thesis, this means that ROC uses eight processing cores compared to a single core used by iPROC. Put differently, ROC uses around eight times more computational power than iPROC. That said, the experiments will demonstrate that iPROC is still competitive, even with much less computational power. We set  $\zeta = 0$  for all benchmarks in this section and use an optimality tolerance  $\mu = 10^{-2}$  and controller order  $\kappa = 1$  for the iPROC algorithm.

#### 6.3.10.1 Space Rendezvous

We present the space rendezvous benchmark proposed in [38, p. 3.6], [20] for the *approaching* mode. The dynamics of the controlled spacecraft, given in a coordinate system relative to the target as shown in Fig. 6.5, are given by

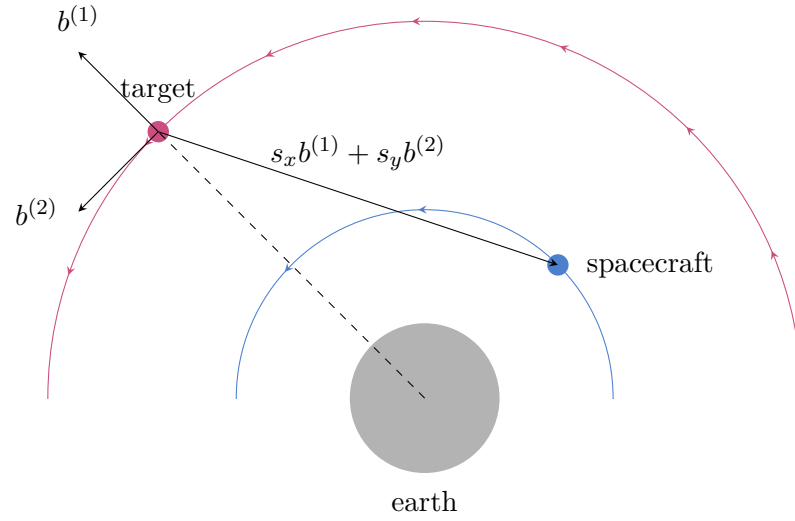
$$\dot{x}_1 = x_3, \tag{6.72a}$$

$$\dot{x}_2 = x_4, \tag{6.72b}$$

$$\dot{x}_3 = n^2 x_1 + 2n x_4 + \frac{c}{r^2} - \frac{c}{r_c^2} (r + x_1) + \frac{u_1}{m_c}, \tag{6.72c}$$

$$\dot{x}_4 = n^2 x_2 - 2n x_3 - \frac{c}{r_c^3} x_2 + \frac{u_2}{m_c}, \tag{6.72d}$$

where  $c = 1.4350 \times 10^{18} \text{ m}^3 \text{ min}^{-2}$ ,  $r = 42\,164 \times 10^3 \text{ m}$ ,  $m_c = 500 \text{ kg}$ ,  $n = \sqrt{\frac{c}{r^3}}$ , and  $r_c = \sqrt{(r + s_x)^2 + s_y^2}$ . Let  $x = [s_y, s_x, v_x, v_y]^T$  denote the state with spatial coordinates  $[s_x, s_y]$  and velocity  $[v_x, v_y]$ , and denote with  $u = [F_x, F_y]^T$  the input to the system, where  $F_x$  and  $F_y$  are thrust forces. The parameters of the benchmark are shown



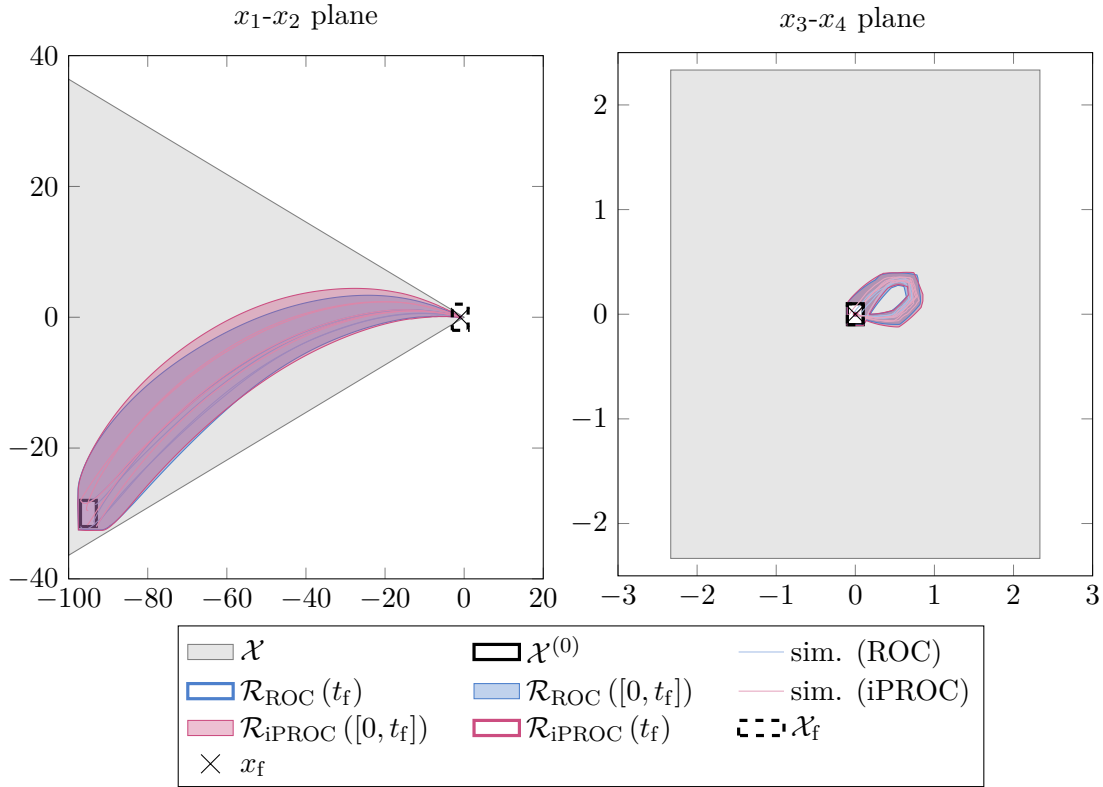
**Figure 6.5:** Depiction of the relative coordinate system for the spacecraft model. Here,  $b^{(1)} \in \mathbb{R}^2$  and  $b^{(2)} \in \mathbb{R}^2$  denote the basis vectors of the relative coordinate system in some global coordinate system.

in Tab. 6.1. For ROC, we additionally set  $\tilde{\mathcal{U}} = 0.75\mathcal{U}$  and  $\tilde{\mathcal{X}} = \mathcal{X}$ . We note that the original authors of the benchmark in [20] proposed  $\mathcal{V}_{\text{orig}} = \{v \in \mathbb{R}^2 \mid \sqrt{v_x^2 + v_y^2} \leq 3.3\}$  as state velocity constraints. Generally, an arbitrarily tight inner approximation of  $\mathcal{V}_{\text{orig}}$  can be computed by first inner approximating  $\mathcal{V}_{\text{orig}}$  using Th. 4.4 and then using [8, Th. 7] for the exact conversion from a zonotope to an H-polytope; however, the resulting H-polytope has potentially many halfspaces. For this benchmark, a crude inner approximation of  $\mathcal{V}_{\text{orig}}$  suffices: Looking at Fig. 6.6, the velocity constraints in  $[x_3, x_4]$  are easily satisfied. Therefore, we compute a parallelotope inner approximation of the velocity sphere constraint as  $\mathcal{V} = \check{\mathcal{Z}}(\langle 0, 3.3^2 I_2 \rangle_E)_1$  using Th. 4.4.

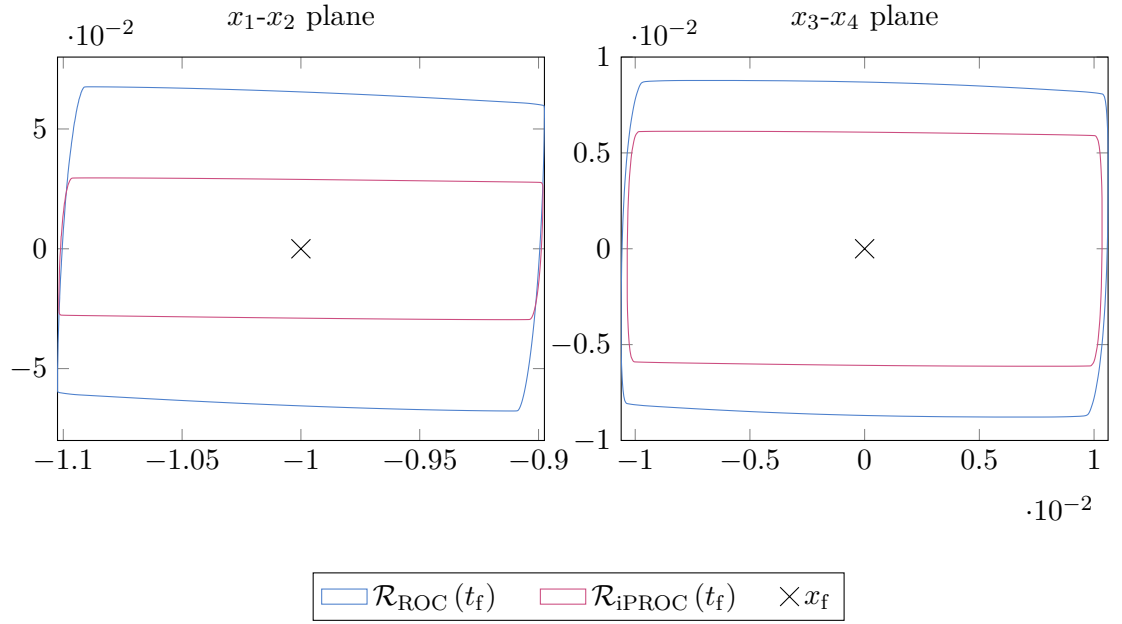
Both ROC and iPROC find a feasible solution in 553 s and 732 s, respectively, where  $\text{csize}(\mathcal{R}_{\text{ROC}}(t_f)) = 0.19$  and  $\text{csize}(\mathcal{R}_{\text{iPROC}}(t_f)) = 0.15$ . Fig. 6.6 shows the solution for both approaches, with a more detailed comparison of their final reachable sets in Fig. 6.7. Additionally, Fig. 6.8 shows the progress of iPROC over the iterations, which validates the result of Th. 6.1 that the approximated trust-region problem can approximate the formally verified controller cost arbitrarily closely.

**Table 6.1:** Parameters for the space rendezvous benchmark.

Parameter	Value
$m_c$	5
$\mathcal{X}^{(0)}$	$\left\langle \left[ -95, -30, 0, 0 \right]^T, \text{diag} \left( \left[ 2, 2, 0.1, 0.1 \right] \right) \right\rangle_Z$
$\mathcal{U}$	$\langle 0, 2 \rangle_Z$
$\mathcal{W}$	$\langle 0, 0.1I_2 \rangle_Z$
$\mathcal{X}$	$\mathcal{S} \times \mathcal{V}$
$\mathcal{S}$	$\left\langle \left[ \begin{array}{cc} -1, & 0 \\ \tan 20^\circ, & -1 \\ \tan 20^\circ, & 1 \end{array} \right], \left[ \begin{array}{c} 100 \\ 0 \end{array} \right] \right\rangle_H$
$\mathcal{V}$	$\tilde{\mathcal{Z}} \left( \langle 0, 3.3^2 I_2 \rangle_E \right)_1 = \langle 0, G \rangle_Z = \left\langle \left[ \begin{array}{c} G^{-1} \\ -G^{-1} \end{array} \right], 1_4 \right\rangle_H$
$t_f$	200 s
$x_f$	$\left[ -1, 0, 0, 0 \right]^T$
$\kappa$	1



**Figure 6.6:** Comparison of the ROC approach and the iPROC approach for the space rendezvous benchmark. Fig. 6.7 shows the final, closed-loop reachable sets of both approaches in more detail.



**Figure 6.7:** Comparison of the final closed-loop reachable set of ROC and iPROC for the space rendezvous benchmark.

### 6.3.10.2 Kinematic Single-Track Model (Car)

We consider a variant of the kinematic single-track model<sup>1</sup>

$$\dot{x}_1 = u_1 + w_1, \quad (6.73a)$$

$$\dot{x}_2 = u_2 + w_2, \quad (6.73b)$$

$$\dot{x}_3 = x_1 \cos(x_2), \quad (6.73c)$$

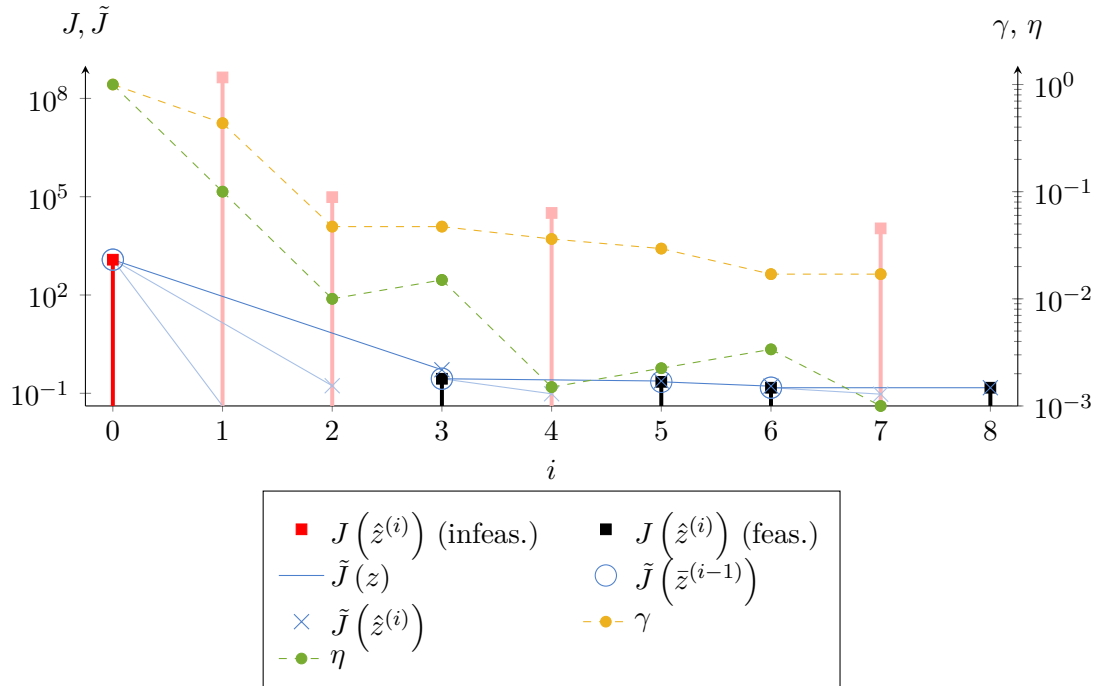
$$\dot{x}_4 = x_1 \sin(x_2), \quad (6.73d)$$

where  $x = [v, \Psi, s_x, s_y]^T$  is the state of the system,  $v$  is the longitudinal velocity of the vehicle,  $\Psi \in [-\pi, \pi]$  is the vehicle heading,  $[s_x, s_y]^T$  denote the spatial coordinates, and  $u = [u_1, u_2]^T = [a, \phi]^T$  is the input to the system, where  $a$  is the longitudinal acceleration and  $\phi$  is the steering rate. All benchmark parameters are shown in Tab. 6.2. To arrive at a feasible solution for ROC, we additionally set  $\tilde{\mathcal{U}} = \text{diag}([0.3, 0.7])\mathcal{U}$ , weigh the final reachable set in the objective function of the combined controller in (6.9) with  $Q_{\text{obj}} = \text{diag}([0.5, 10, 1, 1])$ , and weigh the undisturbed, final reachable set in the objective function of the feedforward computation with  $Q_{\text{obj}}^{(\text{ff})} = \text{diag}([1, 1000, 1, 1])^2$ .

<sup>1</sup>Sec. 5, [https://gitlab.lrz.de/tum-cps/commonroad-vehicle-models/-/blob/master/vehicleModels\\_commonRoad.pdf](https://gitlab.lrz.de/tum-cps/commonroad-vehicle-models/-/blob/master/vehicleModels_commonRoad.pdf)

<sup>2</sup>These values are mostly obtained by trial and error.





**Figure 6.8:** Visualization of the solver progress for the iPGSC approach (space rendezvous benchmark). Values between  $\tilde{J}(\bar{z}^{(i-1)})$  and  $\tilde{J}(\hat{z}^{(i)})$  for  $1 \leq i \leq 8$  – where  $\bar{z}^{(i-1)}$  and  $\hat{z}^{(i)}$  are the initial and critical points of the trust-region subproblem – are linearly interpolated ( $\hat{z}^{(0)} = \bar{z}^{(0)}$ ). Opaque values at  $i \in \{1, 2, 4, 7\}$  visualize rejected iterations.

**Table 6.2:** Benchmark parameters for the car model.

Parameter	Value
$m_c$	5
$c_{\mathcal{X}^{(0)}}$	$[20, 0, 0, 0]^T$
$\mathcal{X}^{(0)}$	$\langle c_{\mathcal{X}^{(0)}}, \text{diag}([0.2, 0.02, 0.2, 0.2]) \rangle_Z$
$\mathcal{U}$	$\langle 0, \text{diag}([9.81, 0.4]) \rangle_Z$
$\mathcal{W}$	$\langle 0, \text{diag}([2, 0.08]) \rangle_Z$
$\mathcal{X}_f$	$\{-c_{\mathcal{X}^{(0)}} + x_f\} \oplus \mathcal{X}^{(0)}$
$t_f$	1 s
$x_f$	$[20, 0.2, 19.87, 1.99]^T$
$\kappa$	1

Fig. 6.9 depicts the reachable sets of the ROC and the iPROC approach for the car benchmark, where Fig. 6.10 shows the final reachable sets of both approaches in more detail. The sizes of the final reachable sets using (5.37) are  $\text{csize}(\mathcal{R}_{\text{ROC}}(t_f)) = 0.39$  and  $\text{csize}(\mathcal{R}_{\text{iPROC}}(t_f)) = 0.13$ , where computation times for ROC and iPROC were 487 s and 497 s, respectively. While both ROC and iPROC find feasible solutions, ROC requires manual tuning (see, e.g.,  $\tilde{\mathcal{U}}$  above) while performing worse, i.e., finding a larger final reachable set.

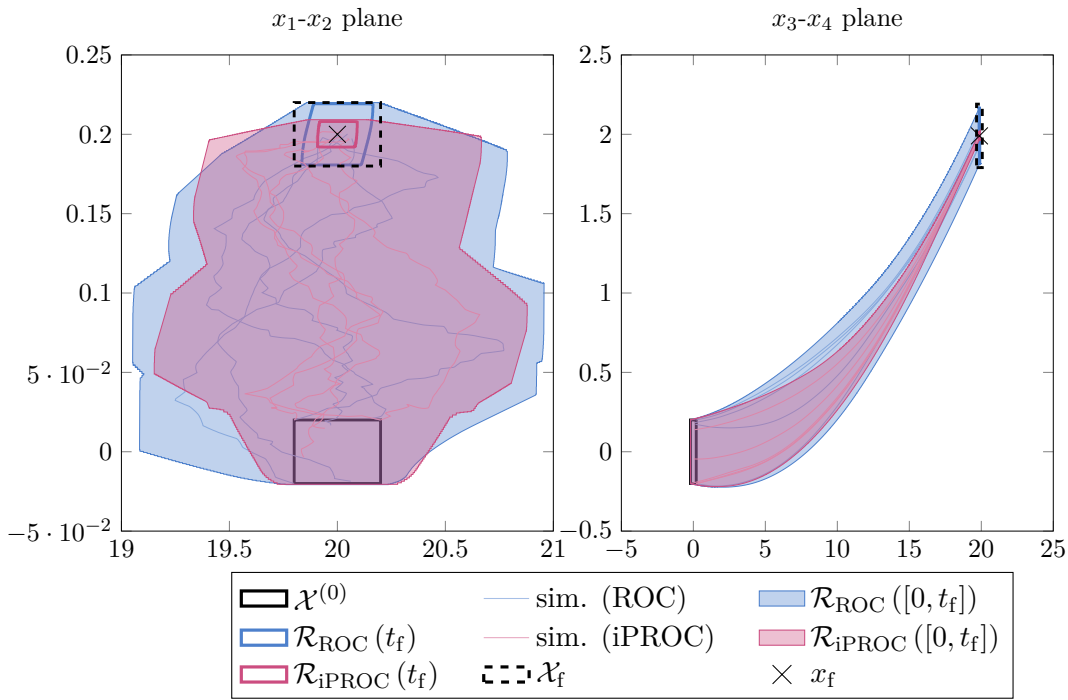
### 6.3.10.3 Water Tanks

We now compare the scalability of ROC and iPROC for a model of  $n \in \mathbb{N}_+$  water tanks [9] as shown in Fig. 6.11. The dynamics are given by

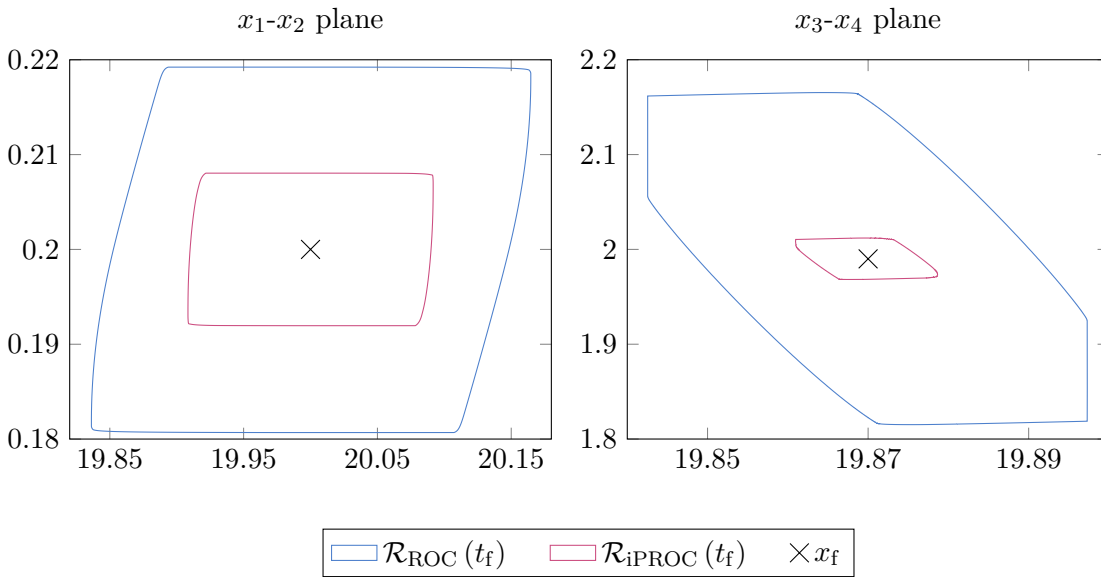
$$\begin{aligned}\dot{x}_1 &= u + w - k\sqrt{2gx_1}, \\ \dot{x}_i &= k\left(\sqrt{2gx_{i-1}} - \sqrt{2gx_i}\right),\end{aligned}$$

for  $2 \leq i \leq n$ , where  $x_k$  is the water level of the  $k$ -th tank for  $1 \leq k \leq n$ ,  $u$  is the inflow into the first tank,  $w$  is the uncontrollable inflow into the first tank, and we set  $k = 0.015$  and  $g = 9.81$ . The parameters for this benchmark are listed in Tab. 6.3 and Fig. 6.12 shows the reachable sets for  $n = 4$  tanks.

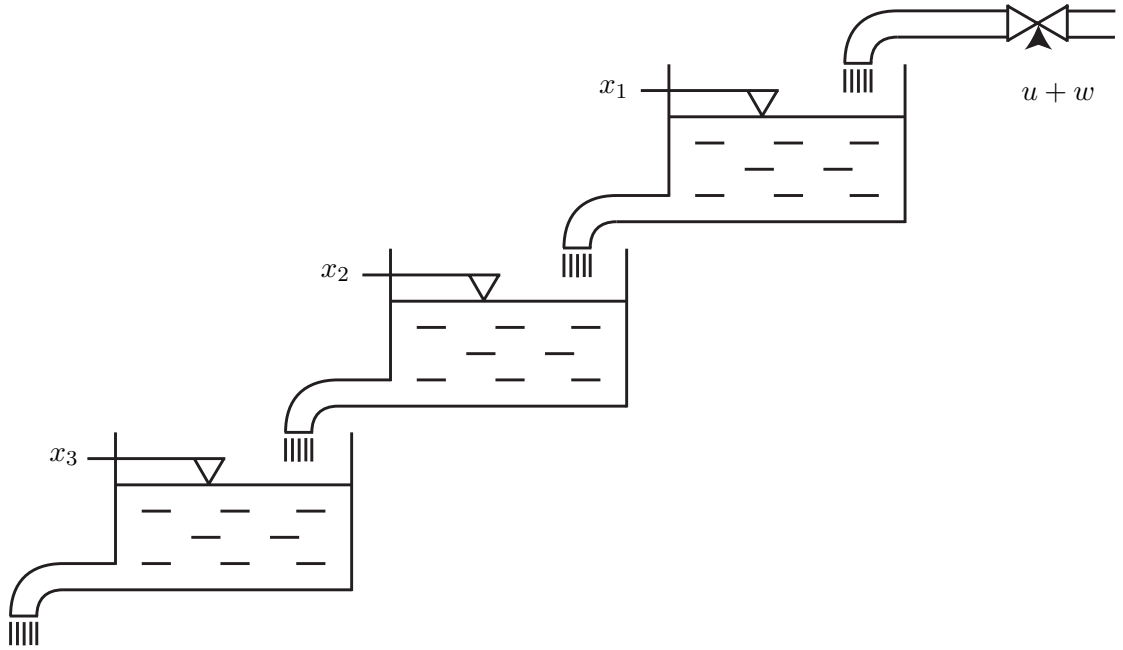
Tab. 6.4 shows the control performance using the metric given in (5.37) and the computation times for both ROC and iPROC for up to eight tanks, where ROC uses parallel computing and thus more computational resources. That said, the iPROC approach seems to scale better with the system dimension in this benchmark than ROC while also performing better overall.



**Figure 6.9:** Comparison of the ROC approach and the iPROC approach for the kinematic single-track model (car).



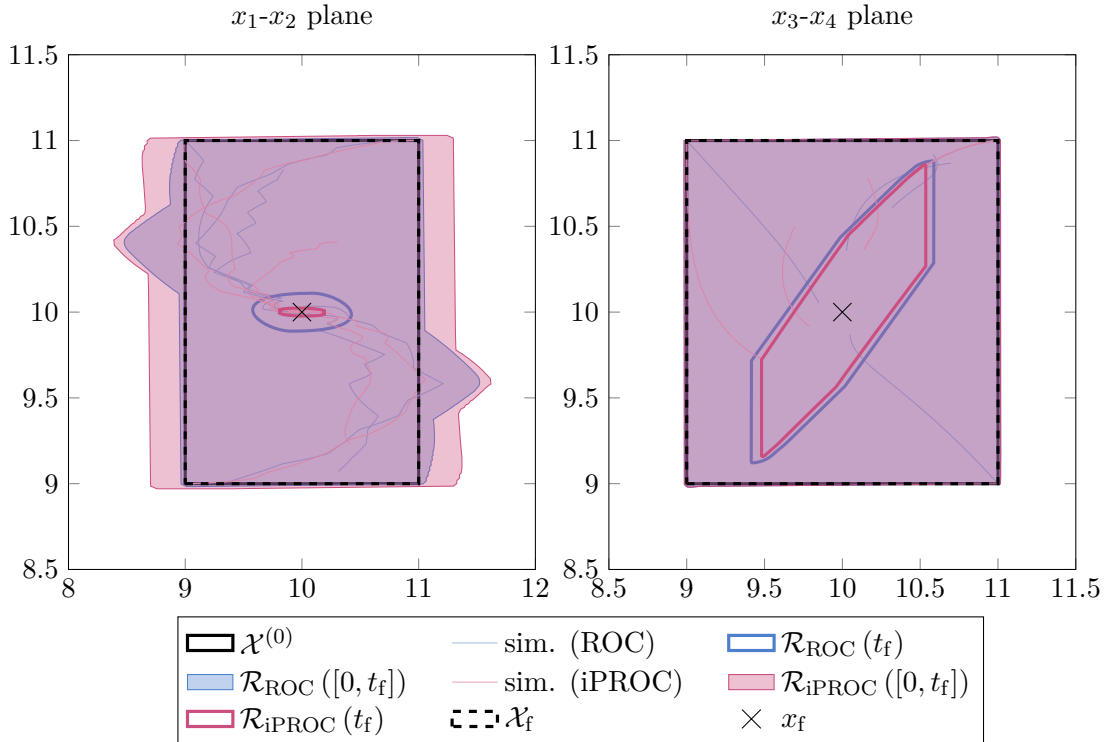
**Figure 6.10:** Comparison of the final closed-loop reachable set of the ROC approach and the iPROC approach for the car benchmark.



**Figure 6.11:** Model of the water tank, where  $x_i$  denotes the water level of tank  $i$  for  $1 \leq i \leq 3$  and  $u + w$  denotes the inflow rate caused by the input and the disturbance [9, Fig. 5, adapted].

**Table 6.3:** Benchmark parameters for the  $n$ -dimensional water tank model.

Parameter	Value
$m_c$	2
$c_{\mathcal{X}^{(0)}}$	$10 \cdot \mathbf{1}_n$
$\mathcal{X}^{(0)}$	$\langle c_{\mathcal{X}^{(0)}}, I_n \rangle_Z$
$\mathcal{U}$	$\langle 1, 1 \rangle_Z$
$\mathcal{W}$	$\langle 0, 0.02 \rangle_Z$
$\mathcal{X}_f$	$\{-c_{\mathcal{X}^{(0)}} + x_f\} \oplus \mathcal{X}^{(0)}$
$t_f$	120 s
$x_f$	$10 \cdot \mathbf{1}_n$
$\kappa$	1



**Figure 6.12:** Comparison of the ROC approach and the iPROC approach for the  $n = 4$  tank benchmark. The final state constraint  $\mathcal{X}_f$  is only shown for reference but not included in the controller synthesis.

**Table 6.4:** Scalability comparison of iPROC and ROC for the tank benchmark using 2 to 8 tanks with  $m_c = 2$  and no final state constraints.

nr. of tanks	2	4	6	8
time [s] (iPROC)	116	218	473	1056
time [s] (ROC)	30	102	383	1209
csize ( $\mathcal{R}(t_f)$ ) (iPROC)	0.17	1.62	3.66	5.77
csize ( $\mathcal{R}(t_f)$ ) (ROC)	0.48	2.00	4.00	6.06
nr. of iter. (iPROC)	10	8	8	8

### 6.3.11 Discussion

Compared to ROC, iPROC synthesizes the combined controller jointly and thus does not require the user to provide tightened feedforward input and adapted state constraints or weighting matrices for the feedforward or feedback synthesis. Additionally, the numerical experiments in Sec. 6.3.10 demonstrate that iPROC can synthesize more performant controllers than ROC, albeit at the cost of higher computation times, especially for lower-dimensional systems. Nevertheless, Tab. 6.4 indicates that iPROC scales better than ROC for an increasing number of system states in practice. Overall, iPROC is an almost one-click synthesis approach for the synthesis of a combined controller which requires minimal user input.

That said, the tuning of the trust-region radii can still be improved: While Th. 6.1 provides a general framework for how the radii can be adjusted, choosing the monotonically decreasing function and its parameters adequately so that iPROC works well over a range of systems is not trivial. Additionally, while not discussed here, it might further be possible to choose the  $m_c$  controller pairs of feedforward and feedback controller and the feedforward controller order  $\kappa$  automatically. Furthermore, it may be possible to obtain an initial guess for  $P$ ,  $Q$ , and  $R$  by solving the synthesis problem for the combined controller for a finite set of initial states using optimal control.

## 6.4 Summary

In this chapter, we described the reachset optimal control (ROC) approach from [101] and introduced our novel iterative polynomial reachset optimal control (iPROC) algorithm from [35], which – for the first time – synthesizes the feedforward and feedback controller simultaneously and thus avoids the introduction of additional algorithm parameters.

ROC synthesizes a combined controller, consisting of a piecewise constant feedforward controller, which is linear in the initial state, and a continuous state feedback term. Using the GSC approach, the feedforward controller is first synthesized separately. The feedforward controller then is kept fixed and the feedback controller is optimized: To evaluate the objective function and the constraints, an outer approximation of the closed-loop reachable set has to be computed in each iteration. To reduce the number of optimization variables for the feedback controller, LQR control with diagonal weighting matrices for LQR are used. In summary, ROC allows for a straightforward extension of the feedforward synthesis to the combined synthesis of a feedforward controller with continuous state feedback by combining reachability analysis with optimization.

However, even though the feedforward controller has direct influence on the feedback controller, ROC keeps the feedforward controller constant during the feedback synthesis, yielding suboptimal results. Furthermore, this sequential optimization introduces additional algorithm parameters that need to be tuned by the user. Additionally, the feedback synthesis recomputes an outer approximation to the closed-loop reachable set in each optimization iteration, and thus is computationally expensive. As a result, we proposed the novel iPROC approach, which – for the first time – jointly synthesizes the feedforward and feedback controller iteratively: Since the original synthesis prob-

lem is in general hard to solve, we approximate it with a cost function. Because the evaluation of this controller cost is computationally expensive as it requires reachability analysis, we avoid optimizing over it directly and rather propose a smooth optimization problem in the controller parameters, which approximates this cost and can be solved with any off-the-shelf optimization solver. That said, this optimization problem approximates the cost accurately only in a local neighborhood of these parameters. Thus we restrict the domain of the approximated optimization problem to a trust region around the current controller parameters and iteratively update them whenever the cost decreases. We adapt the size of the trust region in each iteration to ensure the accuracy of the approximated parameterized reachable set. Hence, iPROC enables the simultaneous optimization over the combined controller parameters and avoids the introduction of additional algorithm parameters that need to be tuned by the user. The effectiveness of iPROC was demonstrated using numerical benchmarks.





# 7 Conclusion and Future Directions

## 7.1 Conclusion

In this thesis, we developed different control synthesis algorithms which combine reachability analysis and optimization theory to obtain non-conservative and formally verified results. To formulate efficient algorithms, we further developed new set conversions between zonotopes and ellipsoids. The presented algorithms all operate on disturbed nonlinear systems, respect input and state constraints, and have polynomial complexity in the state dimension, making them applicable to a wide range of tasks, such as autonomous driving and human-robot interaction tasks.

We first introduced the general synthesis problem and defined the abstracted synthesis problem, which is constructed such that all relevant optimization problems from subsequent chapters can be generalized to this problem. We showed that the abstracted synthesis problem can be reformulated as a smooth non-convex optimization problem. Because many solvers find solutions to an optimization problem by solving the associated Karush-Kuhn-Tucker (KKT) conditions, we furthermore proved that any critical point of said optimization problem is regular. As a result, the first-order KKT conditions are always necessary conditions for all novel optimization problems in subsequent chapters, and thus many off-the-shelf optimization solvers can be used.

For an efficient and numerically stable way to pose and solve these synthesis problems, we further derived novel conversions between sets, such as ellipsoid-zonotope conversions from our work in [36]. These conversions are necessary since not all set operations are closed under all set operations. As a result, the same operation for two different set representations may have different tightness and computational efficiency. Since we synthesize controllers under input and state constraints, we furthermore reviewed existing methods to check the containment of two sets, which are used in subsequent chapters to, e.g., check the state constraint satisfaction of a given reachable set.

Next, we first focused on the synthesis of piecewise constant controllers. Here, we introduced the generator-space control (GSC) approach from previous work in [99] for the synthesis of piecewise constant controllers for disturbed nonlinear systems under constraints. While GSC is computationally efficient since it only requires the solution of a single linear program (LP) for the synthesis of its controller, it only optimizes over the reachable set of the linearized dynamics. Thus, we proposed polynomial generator-space control (PGSC) from our work in [37] as an extension of GSC, which synthesizes piecewise constant controllers that are polynomially dependent on the initial state based on a polynomial abstraction of the nonlinear dynamics. However, since both GSC and PGSC only approximate the reachable set without any formal guarantees, the inclusion of state constraints always requires additional input from the user. To solve this problem, we

introduced the iterative polynomial generator-space control (iPGSC) approach, allowing us to synthesize polynomial, piecewise constant controllers under input and state constraints without the need for additional user input.

However, piecewise constant controllers often cannot counteract disturbances efficiently when disturbances are large or act quickly, resulting in larger reachable sets. Thus, we first described the reachset optimal control (ROC) algorithm from [101], which synthesizes a piecewise constant feedforward controller with continuous state feedback for nonlinear disturbed systems under input and state constraints. However, because this synthesis is executed sequentially, i.e., the feedforward controller is computed first and then fixed in the subsequent optimization of the feedback control term, the resulting controller may be suboptimal. Furthermore, this sequential computation approach produces additional algorithm parameters that require manual tuning by the user. Thus, we proposed the iterative polynomial reachset optimal control (iPROC) approach from our work in [35] which – for the first time – synthesizes a feedforward and a feedback controller simultaneously by solving a single optimization problem. To avoid extensive reachability computations, we approximated the synthesis problem and used trust-region radii to restrict the domain of the controller parameters for this approximated optimization problem to ensure accuracy.

### 7.2 Future Directions

We presented approaches for the formally verified controller synthesis of disturbed nonlinear systems under input and state constraints. Subsequently, we propose possible directions for future research to extend this work.

In PGSC and iPGSC, we use a piecewise constant controller which is polynomially dependent on the initial state. While this can achieve better control performance than linear controllers as demonstrated in Sec. 5.3.8, it requires the user to choose the appropriate controller order. Formulating the synthesis problem for a single initial state, optimal control may be able to quickly solve the synthesis problem for a set of possible controller orders to find a suitable choice. Additionally, one might be able to directly identify exponent vectors of the controller template which contribute most to good control performance. Such a sparse controller template could avoid an unnecessarily large number of controller parameters.

While our proposed approaches, specifically iPGSC and iPROC, already reduce the number of algorithm parameters that need to be set by the user, the number of piecewise controllers for iPGSC and the number of feedforward and feedback controller pairs in iPROC still needs to be chosen by the user. Here, optimal control for a single initial state could be used to determine the effect of a varying number of controllers on the control performance: By formulating the synthesis problem for a single initial state, optimal control methods can be used to quickly solve the synthesis problem for a given number of controllers. Since these methods are very efficient, evaluating the optimal control problem for a larger number of controllers might be feasible.

Furthermore, we assumed that the time horizon is split uniformly, i.e., each controller is always active for the same amount of time. However, there might be instances where a large number of controllers is only required for a short time interval – e.g., due to quickly changing system dynamics – so that subdividing the time horizon uniformly to achieve the required number of controllers for this short time interval is quite inefficient. Here, ideas similar to numerically integrating differential equations might be able to identify points in time where a higher number of controllers is required: Adaptive step-size solvers for ODEs are able to automatically decrease the time step size when the dynamics require a finer time resolution. A similar approach might be able to identify a non-uniform time discretization for the number of controllers.

Lastly, we approximated the synthesis problem in iPROC to avoid the computation of the reachable set in each optimization iteration. While this approximation is less expensive to compute than the reachable set in practice, it still evaluates the Riccati equation and its derivatives in each optimization iteration to compute the LQR feedback matrix for the feedback controller. To further reduce the computational effort and speed up the optimization, one may approximate the trust-region subproblem using yet another, inner trust-region approximation: We initially compute the Jacobian and Hessian matrices of the outer trust-region subproblem and then form an inner approximation by assuming that these computed Jacobian and Hessian matrices are constant. By again restricting the domain for the controller parameters of this inner approximation to a trust region, we can control the accuracy of this inner approximation. As a result, we avoid the recomputation of the derivative information in each optimization iteration of the trust-region subproblem.



# A Appendix

## A.1 Abstracted Synthesis Problem

### A.1.1 Reformulation Equivalence

*Proof of Prop. 3.1.* We first prove the second claim by deriving the optimizers  $\hat{s}$ ,  $\hat{w}$ ,  $\hat{y}$ ,  $\hat{V}$ ,  $\hat{T}$ , and  $\hat{p}$  as functions of  $\hat{z}$ , i.e., in dependence of  $\hat{z}$  only. We note that the KKT conditions used in this proof are always necessary for (3.5) due to Prop. 3.2.

**Value of  $\hat{w}$ :** We start by analyzing the complementary slackness conditions of (3.5b), given by

$$\hat{\lambda}_i^{(+)} (+h_i(\hat{z}) - \hat{w}_i) = 0, \quad (\text{A.1a})$$

$$\hat{\lambda}_i^{(-)} (-h_i(\hat{z}) - \hat{w}_i) = 0, \quad (\text{A.1b})$$

where  $\hat{\lambda}^{(+)} \in \mathbb{R}_{\geq 0}^{n_w}$  and  $\hat{\lambda}^{(-)} \in \mathbb{R}_{\geq 0}^{n_w}$  are the corresponding optimal Lagrange multipliers with  $1 \leq i \leq n_w$ . Further, the stationary condition with respect to  $w$  is given by

$$l - \hat{\lambda}^{(+)} - \hat{\lambda}^{(-)} + M^T \sum_{k=1}^{n_k} \sum_{m=1}^{n_s} F_{(m)}^{(k)T} \hat{\chi}_{(m)}^{(k)} = 0, \quad (\text{A.2})$$

where  $\hat{\chi}_{(m)}^{(k)} \in \mathbb{R}_{\geq 0}^q$  is the optimal Lagrange multiplier of the  $(km)$ -th constraint in (3.5d).

Now assume that  $\hat{w}_i > |h_i(\hat{z})|$ . It then holds that  $\hat{\lambda}_i^{(+)} = \hat{\lambda}_i^{(-)} = 0$  due to (A.1). Since (A.2) must hold for general  $l \geq 0$  and  $M \in \mathbb{R}_{\geq 0}^{n_r \times n_w}$ , it must also hold for  $l = 1_{n_s}$  and  $M = 0$ , from which  $\hat{\lambda}_i^{(+)} + \hat{\lambda}_i^{(-)} = 1$  follows, and thus by contradiction, it must hold that

$$\hat{w} = |h(\hat{z})|, \quad (\text{A.3})$$

since  $\hat{w} \geq |h(\hat{z})|$  (see (3.5b)).

**Value of  $\hat{y}$ :** The complementary slackness conditions of (3.5c) and the stationarity condition in  $y$  are given by

$$\hat{\zeta}_i \left( \bar{R}_{(i,:)} \rho(\hat{z}) - y_i \right) = 0, \quad (\text{A.4})$$

$$\hat{\zeta}_i \left( \bar{R}_{(i,:)} \rho(\hat{z}) - y_i \right) = 0, \quad (\text{A.5})$$

$$1 - \hat{\zeta}_i - \hat{\zeta}_i = 0, \quad (\text{A.6})$$

## A Appendix

for  $1 \leq i \leq n_y$ , where  $\hat{\zeta} \in \mathbb{R}_{\geq 0}^{n_y}$  and  $\hat{\bar{\zeta}} \in \mathbb{R}_{\geq 0}^{n_y}$  are the optimal Lagrange multipliers of the respective constraints in (3.5c). It follows from (A.6) that  $\hat{\zeta}_i + \hat{\bar{\zeta}}_i = 1$ . Assuming  $\hat{y}_i > \max\left(\underline{R}_{(i,:)}\rho(\hat{z}), \bar{R}_{(i,:)}\rho(\hat{z})\right)$  implies  $\hat{\zeta}_i = \hat{\bar{\zeta}}_i = 0$  due to (A.4) and (A.5), and thus

$$\hat{y}_i = \max\left(\underline{R}_{(i,:)}\rho(\hat{z}), \bar{R}_{(i,:)}\rho(\hat{z})\right) \quad (\text{A.7})$$

follows by contradiction and due to (3.5c).

**Value of  $\hat{s}$ :** The relevant KKT conditions in  $s$  are

$$\sigma - \hat{\delta}_m - \hat{\vartheta}_m = 0, \quad (\text{A.8})$$

$$\hat{\delta}_m(\hat{g}_m - \hat{s}_m) = 0, \quad (\text{A.9})$$

$$-\hat{\vartheta}_m\hat{s}_m = 0, \quad (\text{A.10})$$

for  $1 \leq m \leq n_s$ , where  $\hat{\delta} \in \mathbb{R}_{\geq 0}^{n_s}$  and  $\hat{\vartheta} \in \mathbb{R}_{\geq 0}^{n_s}$  are the optimal Lagrange multipliers of the constraints in (3.5g) and (3.5i) and we define  $\hat{g} = \hat{p} - b$ . Summing (A.9) and (A.10) yields  $\hat{\delta}_m\hat{g}_m - (\hat{\delta}_m + \hat{\vartheta}_m)\hat{s}_m = 0$ , where a subsequent substitution of  $\hat{\delta}_m + \hat{\vartheta}_m = \sigma$  from (A.8) and solving for  $\hat{s}_m$  results in

$$\hat{s}_m = \frac{\hat{\delta}_m\hat{g}_m}{\sigma}. \quad (\text{A.11})$$

Now let  $\mathcal{I} = \{i \in \{1, \dots, n_s\} \mid \hat{g}_i \leq 0\}$ . Then  $\forall i \in \mathcal{I} : \hat{\delta}_i\hat{g}_i = 0$  since necessarily  $\hat{\delta}_i = 0$  for  $\hat{g}_i < 0$  due to  $s_i \geq 0$  and (3.5g) (statement trivially holds for  $\hat{g}_i = 0$ ). For  $\mathcal{J} = \{1, \dots, n_s\} \setminus \mathcal{I}$ , we further have  $\forall j \in \mathcal{J} : \hat{g}_j > 0 \xrightarrow{(3.5g)} \hat{s}_j > 0 \xrightarrow{(A.10)} \hat{\vartheta}_j = 0 \xrightarrow{(A.8)} \hat{\delta}_j = \sigma$ . Thus, (A.11) can be written as

$$\hat{s}_m = \begin{cases} \hat{g}_m, & \hat{g}_m > 0, \\ 0, & \text{otherwise} \end{cases} \quad (\text{A.12})$$

for  $1 \leq m \leq n_s$ . To express  $\hat{s}$  in  $\hat{z}$  only, it remains to derive  $\hat{g} = \hat{p} - b$  in  $\hat{z}$  only. Thus, we assume  $\hat{g}_m > 0$  for the remainder of this proof.

**Value of  $\hat{p}$ :** The stationarity condition in  $p$  and the complementary slackness condition of (3.5f) are given by

$$-1_{n_k}^T \hat{\gamma}_{(m)} + \hat{\delta}_m = 0, \quad (\text{A.13})$$

$$\left[\hat{\gamma}_{(m)}\right]_k \left(\hat{V}_{km} + \hat{T}_{km} - \hat{p}_m\right) = 0, \quad (\text{A.14})$$

for  $1 \leq m \leq n_s$  and  $1 \leq k \leq n_k$ , where  $\hat{\gamma}_{(m)} \in \mathbb{R}_{\geq 0}^{n_k}$  is the optimal Lagrange multiplier of the  $m$ -th constraint in (3.5f). Now assume that  $\hat{p}_m > \max\left(\hat{V}_{(:,m)} + \hat{T}_{(:,m)}\right)$  for  $\hat{g}_m > 0$ .

## A.1 Abstracted Synthesis Problem

It then holds that  $\hat{\gamma}_{(m)} = 0$  due to (A.14), and thus  $\hat{\delta}_m = 1_{n_k}^T \hat{\gamma}_{(m)} = 0$  due to (A.13); however, we previously proved that  $\hat{\delta}_m = \sigma$  for  $\hat{g}_m > 0$ , and therefore

$$\hat{p}_m = \max \left( \hat{V}_{(:,m)} + \hat{T}_{(:,m)} \right) \text{ for } \hat{g}_m > 0, \quad (\text{A.15})$$

follows by contradiction.

**Value of  $\hat{V}$ :** Since  $\hat{p}_m = \max \left( \hat{V}_{(:,m)} + \hat{T}_{(:,m)} \right)$  for  $\hat{g}_m > 0$ , there must be a set of indices  $\bar{\mathcal{K}}_{(m)} \subseteq \{1, \dots, n_k\}$  for which holds that

$$\forall \bar{k} \in \bar{\mathcal{K}}_{(m)} : \hat{V}_{\bar{k}m} + \hat{T}_{\bar{k}m} = \hat{p}_m \text{ for } \hat{g}_m > 0, \quad (\text{A.16a})$$

$$\forall \underline{k} \in \{1, \dots, n_k\} \setminus \bar{\mathcal{K}}_{(m)} : \hat{V}_{\underline{k}m} + \hat{T}_{\underline{k}m} < \hat{p}_m \text{ for } \hat{g}_m > 0. \quad (\text{A.16b})$$

From (A.13), (A.14), (A.16a) and (A.16b), it then follows that

$$1_{n_k}^T \hat{\gamma}_{(m)} = \sum_{\bar{k} \in \bar{\mathcal{K}}_{(m)}} \left[ \hat{\gamma}_{(m)} \right]_{\bar{k}} = \hat{\delta}_m. \quad (\text{A.17})$$

The stationarity conditions in  $V_{km}$  for  $1 \leq k \leq n_k$  and the complementary slackness condition of (3.5d) are given by

$$-1_q^T \hat{\chi}_{(m)}^{(k)} + \left[ \hat{\gamma}_{(m)} \right]_k = 0, \quad (\text{A.18})$$

$$\left[ \hat{\chi}_{(m)}^{(k)} \right]_i \left( \left[ F_{(m)}^{(k)} \right]_{(i,:)} (Ah(\hat{z}) + M\hat{w}) - \hat{V}_{km} \right) = 0, \quad (\text{A.19})$$

for  $1 \leq i \leq q$ , where  $\hat{\chi}_{(m)}^{(k)} \in \mathbb{R}_{\geq 0}^q$  is the optimal Lagrange multiplier of the  $(km)$ -th constraint in (3.5d). Now assume that  $\hat{V}_{\bar{k}m} > \max \left( F_{(m)}^{(\bar{k})} (Ah(\hat{z}) + M\hat{w}) \right)$ . It then follows from (A.19) that  $\hat{\chi}_{(m)}^{(\bar{k})} = 0$  and thus  $\left[ \hat{\gamma}_{(m)} \right]_{\bar{k}} = 0$  due to (A.18). Since then  $\sum_{\bar{k} \in \bar{\mathcal{K}}_{(m)}} \left[ \hat{\gamma}_{(m)} \right]_{\bar{k}} = 1_{n_k}^T \hat{\gamma}_{(m)} = \hat{\delta} = 0$  (see (A.17)) but we previously showed that  $\hat{\delta} = \sigma$  for  $\hat{g}_m > 0$  (see derivations for  $\hat{s}$ ), it follows by contradiction that

$$\hat{V}_{\bar{k}m} = \max \left( F_{(m)}^{(\bar{k})} (Ah(\hat{z}) + M\hat{w}) \right) \text{ for } \hat{g}_m > 0. \quad (\text{A.20})$$

**Value of  $\hat{T}$ :** The stationarity conditions in  $T_{km}$  for  $1 \leq k \leq n_k$ ,  $1 \leq m \leq n_s$ ,  $1 \leq i \leq q$  and the complementary slackness condition of (3.5e) are given by

$$-1_q^T \hat{\phi}_{(m)}^{(k)} + \left[ \hat{\gamma}_{(m)} \right]_k = 0, \quad (\text{A.21})$$

$$\left[ \hat{\phi}_{(m)}^{(k)} \right]_i \left( \left[ G_{(m)}^{(k)} \right]_{(i,:)} \rho(\hat{z}) - \hat{T}_{km} \right) = 0, \quad (\text{A.22})$$

## A Appendix

where  $\hat{\phi}_{(m)}^{(k)} \in \mathbb{R}_{\geq 0}^q$  is the optimal Lagrange multiplier of the  $(km)$ -th constraint in (3.5e). Similarly to  $\hat{V}$ , assume now that  $\hat{T}_{\bar{k}m} > \max \left( G_{(m)}^{(\bar{k})} \rho(\hat{z}) \right)$  with  $\bar{k} \in \bar{\mathcal{K}}_{(m)}$  as defined in (A.16) for  $\hat{g}_m > 0$ . It then follows from (A.22) that  $\hat{\phi}_{(m)}^{(\bar{k})} = 0$  and thus  $\left[ \hat{\gamma}_{(m)} \right]_{\bar{k}} = 0$  due to (A.21). Since then  $\sum_{\bar{k} \in \bar{\mathcal{K}}_{(m)}} \left[ \hat{\gamma}_{(m)} \right]_{\bar{k}} = 1_{n_k}^T \hat{\gamma}_{(m)} = \hat{\delta} = 0$  (see (A.17)) but we previously showed that  $\hat{\delta} = \sigma$  for  $\hat{g}_m > 0$  (see derivations for  $\hat{s}$ ), it follows by contradiction that

$$\hat{T}_{\bar{k}m} = \max \left( G_{(m)}^{(\bar{k})} \rho(\hat{z}) \right) \text{ for } \hat{g}_m > 0. \quad (\text{A.23})$$

**Collecting all optimizers:** We can now write  $\hat{s}_m$  for  $1 \leq m \leq n_s$  as

$\hat{g}_m > 0$  :

$$\begin{aligned} \hat{s}_m &\stackrel{(\text{A.12})}{=} \hat{g}_m = \hat{p}_m - b_m \\ &\stackrel{(\text{A.16a})}{=} \hat{V}_{\bar{k}m} + \hat{T}_{\bar{k}m} - b_m \\ &\stackrel{(\text{A.3}), (\text{A.20}), (\text{A.23})}{=} \max \left( F_{(m)}^{(\bar{k})} (Ah(\hat{z}) + M|h(\hat{z})|) \right) + \max \left( G_{(m)}^{(\bar{k})} \rho(\hat{z}) \right) - b_m \\ &\stackrel{(\text{A.16a})}{=} \max_{0 \leq k \leq n_k} \left\{ \max \left( F_{(m)}^{(k)} (Ah(\hat{z}) + M|h(\hat{z})|) \right) + \max \left( G_{(m)}^{(k)} \rho(\hat{z}) \right) \right\} - b_m, \end{aligned}$$

and thus

$$\hat{s}_m = \begin{cases} g_m(\hat{z}), & g_m(\hat{z}) > 0, \\ 0, & \text{otherwise} \end{cases}, \quad (\text{A.24})$$

or equivalently  $\hat{s} = \max(0, g(\hat{z}))$ , where  $g$  is defined in (3.4). Therefore, substitution of  $\hat{s}$ ,  $\hat{w}$ , and  $\hat{y}$  in (3.6) by (A.24), (A.3), and (A.7) yields  $J_{\text{smooth}}(\hat{z}, \hat{s}, \hat{w}, \hat{y}) = \tilde{J}(\hat{z})$ .

**Same minimum:** It remains to prove the first claim. Let  $(z^*, s^*)$  denote a global optimizer of (3.2). By design, it must hold that  $s^* = \max(0, g(z^*))$  since, if  $s^* > \max(0, g(z^*))$ , choosing  $(z^*, \max(0, g(z^*)))$  achieves a smaller minimum but  $(z^*, s^*)$  is assumed to be a global optimizer. Thus, the minimum  $\tilde{J}(z^*) = J(z^*, \max(0, g(z^*)))$  is attained. Previously, we showed that the minimum objective value at a critical point  $(\hat{z}, \hat{s}, \hat{w}, \hat{y}, \hat{V}, \hat{T}, \hat{p})$  of (3.5) is  $\tilde{J}(\hat{z})$ . We now prove the equivalence of (3.2) and (3.5) by



contradiction: Assume that  $\tilde{J}(z^*) < \tilde{J}(\hat{z})$ . Then  $(\bar{z}, \bar{s}, \bar{w}, \bar{y}, \bar{V}, \bar{T}, \bar{p})$  with

$$\begin{aligned}\bar{z} &= z^*, \\ \bar{s} &= \max(0, g(z^*)), \\ \bar{w} &= |h(z^*)|, \\ \bar{y} &= \max(\underline{R}\rho(z^*), \bar{R}\rho(z^*)), \\ \bar{V}_{km} &= \max\left(F_{(m)}^{(k)}(Ah(z^*) + M\bar{w})\right), \\ \bar{T}_{km} &= \max\left(F_{(m)}^{(k)}\rho(z^*)\right), \\ \bar{p}_m &= \max\left(\bar{V}_{(:,m)} + \bar{T}_{(:,m)}\right),\end{aligned}$$

for  $1 \leq k \leq n_k$  and  $1 \leq m \leq n_s$  is a feasible point of (3.5) with

$$J_{\text{smooth}}(\bar{z}, \bar{s}, \bar{w}, \bar{y}, \bar{V}, \bar{T}, \bar{p}) = \tilde{J}(z^*) < \tilde{J}(\hat{z}),$$

but  $(\hat{z}, \hat{s}, \hat{w}, \hat{r}, \hat{t})$  is a global optimizer by assumption; thus,  $\tilde{J}(z^*) \geq \tilde{J}(\hat{z})$  follows by contradiction. Similarly, now assume that  $\tilde{J}(z^*) > \tilde{J}(\hat{z})$ . Then  $(\hat{z}, \max(0, g(\hat{z})))$  is a feasible point of (3.2) with  $\tilde{J}(z^*) > \tilde{J}(\hat{z})$ , but  $(z^*, s^*)$  is a global optimizer by assumption; thus,  $\tilde{J}(z^*) = \tilde{J}(\hat{z})$  when combined with  $\tilde{J}(z^*) \geq \tilde{J}(\hat{z})$  as proven above, which concludes the proof.  $\square$

### A.1.2 Necessity of First-Order Optimality Conditions

*Proof of Prop. 3.2.* Let  $\hat{x} = [\hat{z}^T, \hat{s}^T, \hat{w}^T, \hat{y}^T, \hat{V}_{(\cdot)}^T, \hat{T}_{(\cdot)}^T, \hat{p}^T]^T \in \mathbb{R}^n$  denote a collected critical point of (3.5) with  $\tilde{g}(\hat{x}) \leq 0$ , where

$$\tilde{g}(\hat{x}) = \begin{bmatrix} \begin{bmatrix} +h(\hat{z}) \\ -h(\hat{z}) \\ \underline{R}\rho(\hat{z}) \\ \bar{R}\rho(\hat{z}) \end{bmatrix} - \begin{bmatrix} \hat{w} \\ \hat{w} \\ \hat{y} \\ \hat{y} \end{bmatrix} \\ \vdots \\ F_{(m)}^{(k)}(Ah(\hat{z}) + M\hat{w}) - 1_q \hat{V}_{km} \\ \vdots \\ G_{(m)}^{(k)}\rho(\hat{z}) - 1_q \hat{T}_{km} \\ \vdots \\ \hat{V}_{(:,m)} + \hat{T}_{(:,m)} - 1_{n_k} p_m \\ \vdots \\ \hat{p} - b - \hat{s} \\ -\hat{z} + \bar{z} \\ \hat{z} - \bar{z} \\ -\hat{s} \end{bmatrix},$$

## A Appendix

i.e.,  $\tilde{g}(\hat{x}) \in \mathbb{R}^m$  collects all constraints of (3.5). Since there are no equality constraints, the Mangasarian-Fromovitz constraint qualification (MFCQ) from Prop. 2.3 requires that there exists a direction

$$d = \left[ d^{(z)T}, d^{(s)T}, d^{(w)T}, d^{(y)T}, D_{(\cdot)}^{(V)T}, D_{(\cdot)}^{(T)T}, d^{(p)T} \right]^T,$$

with  $d^{(z)} \in \mathbb{R}^{n_z}$ ,  $d^{(s)} \in \mathbb{R}^{n_s}$ ,  $d^{(w)} \in \mathbb{R}^{n_w}$ ,  $D^{(V)} \in \mathbb{R}^{n_k \times n_s}$ ,  $D^{(T)} \in \mathbb{R}^{n_k \times n_s}$ , and  $d^{(p)} \in \mathbb{R}^{n_s}$ , such that  $J_{g_A}(\hat{x})d < 0$  for a given  $\hat{x}$ , where  $\mathcal{A} = \{a \in \{1, \dots, m\} \mid \tilde{g}_a(\hat{x}) = 0\}$  denotes the active constraint index set and  $J_{\tilde{g}_A}(\hat{x})$  is the Jacobian matrix of all constraints active at  $x = \hat{x}$ .

In the worst case, all constraints are active. That said, use the fact that – between  $z_i \leq \bar{z}_i$  and  $z_i \geq \underline{z}_i$  for  $1 \leq i \leq n_z$  – only one of the two constraints can be active at any given time since  $\underline{z}_i < \bar{z}_i$  by assumption. To make the presentation more compact, we write these two constraints as  $\left| z - \frac{1}{2}(\underline{z} + \bar{z}) \right| \leq \frac{1}{2}(\bar{z} - \underline{z})$ . The non-differentiability at  $\tilde{z} = \frac{1}{2}(\underline{z} + \bar{z})$  is of no consequence as both constraints are inactive at  $z = \tilde{z}$ , and its derivative is hence given by

$$\frac{d}{dz} \left( \left| z - \tilde{z} \right| - \frac{1}{2}(\bar{z} - \underline{z}) \right) = \text{diag}(\text{sign}(z - \tilde{z})).$$

Then we have

$$\begin{aligned} \max(J_{\tilde{g}_A}(\hat{x})d) &= \max \left( \begin{array}{c} \left[ \begin{array}{c} +J_h(\hat{z}) \\ -J_h(\hat{z}) \\ \underline{R}J_\rho(\hat{z}) \\ \bar{R}J_\rho(\hat{z}) \end{array} \right] d^{(z)} - \begin{bmatrix} d^{(w)} \\ d^{(w)} \\ d^{(y)} \\ d^{(y)} \end{bmatrix} \\ \max_{\substack{1 \leq k \leq n_k \\ 1 \leq m \leq n_s}} \left\{ \max \left( F_{(m)}^{(k)} \left( A J_h(\hat{z}) d^{(z)} + M d^{(w)} \right) \right) - D_{km}^{(V)} \right\} \\ \max_{\substack{1 \leq k \leq n_k \\ 1 \leq m \leq n_s}} \left\{ \max \left( G_{(m)}^{(k)} J_\rho(\hat{z}) d^{(z)} \right) - D_{km}^{(T)} \right\} \\ \max_{1 \leq m \leq n_s} \left\{ \max \left( D_{(\cdot, m)}^{(V)} - D_{(\cdot, m)}^{(T)} \right) - d_m^{(p)} \right\} \\ d^{(p)} - d^{(s)} \\ \text{diag}(\text{sign}(\hat{z} - \tilde{z})) d^{(z)} \\ -d^{(s)} \end{array} \right) \\ &= \max \left( \begin{array}{c} \left[ \begin{array}{c} \bar{h}(d^{(z)}) - d^{(w)} \\ \bar{r}(d^{(z)}) - d^{(y)} \\ \bar{F}_{(\cdot)}(d^{(z)}) - D_{(\cdot)}^{(V)} \\ \bar{G}_{(\cdot)}(d^{(z)}) - D_{(\cdot)}^{(T)} \\ \bar{d}(D^{(V)}, D^{(T)}) - d^{(p)} \\ d^{(p)} - d^{(s)} \\ \text{diag}(\text{sign}(\hat{z} - \tilde{z})) d^{(z)} \\ -d^{(s)} \end{array} \right] \end{array} \right), \end{aligned} \quad (\text{A.25})$$

where

$$\bar{h}(d^{(z)}) = \max\left(+J_h(\hat{z})d^{(z)}, -J_h(\hat{z})d^{(z)}\right), \quad (\text{A.26a})$$

$$\bar{r}(d^{(z)}) = \max\left(\underline{R}J_\rho(\hat{z})d^{(z)}, \bar{R}J_\rho(\hat{z})d^{(z)}\right), \quad (\text{A.26b})$$

$$\bar{F}_{km}(d^{(z)}, d^{(w)}) = \max\left(F_{(m)}^{(k)}\left(AJ_h(\hat{z})d^{(z)} + Md^{(w)}\right)\right), \quad (\text{A.26c})$$

$$\bar{G}_{km}(d^{(z)}) = \max\left(G_{(m)}^{(k)}J_\rho(\hat{z})d^{(z)}\right), \quad (\text{A.26d})$$

$$\bar{d}_m(D^{(V)}, D^{(T)}) = \max\left(D_{(:,m)}^{(V)} - D_{(:,m)}^{(T)}\right). \quad (\text{A.26e})$$

Since we assume that all functions are sufficiently smooth, their gradients are at least continuous, and since  $\hat{z}$  is bounded (see (3.5h)), the expressions in (A.26) are bounded for bounded  $d$  since they are also continuous. Choosing

$$\begin{aligned} \hat{d}^{(z)} &= -\epsilon \operatorname{sign}(\hat{z} - \tilde{z}), \\ \hat{d}^{(w)} &= \bar{h}(\hat{d}^{(z)}) + \epsilon, \\ \hat{d}^{(y)} &= \bar{r}(\hat{d}^{(z)}), \\ \hat{D}_{(\cdot)}^{(V)} &= \bar{F}_{(\cdot)}(\hat{d}^{(z)}, \hat{d}^{(w)}) + \epsilon, \\ \hat{D}_{(\cdot)}^{(T)} &= \bar{G}_{(\cdot)}(\hat{d}^{(z)}) + \epsilon, \\ \hat{d}^{(p)} &= \bar{d}(\hat{D}^{(V)}, \hat{D}^{(T)}) + \epsilon, \\ \hat{d}^{(s)} &= \max(0, \hat{d}^{(p)}) + \epsilon, \end{aligned}$$

with  $\epsilon \in \mathbb{R}_+$  yields

$$\max \left( \begin{array}{c} \bar{h}(\hat{d}^{(z)}) - \hat{d}^{(w)} \\ \bar{r}(\hat{d}^{(z)}) - \hat{d}^{(y)} \\ \bar{F}_{(\cdot)}(\hat{d}^{(z)}, \hat{d}^{(w)}) - \hat{D}_{(\cdot)}^{(V)} \\ \bar{G}_{(\cdot)}(\hat{d}^{(z)}) - \hat{D}_{(\cdot)}^{(T)} \\ \bar{d}(\hat{D}^{(V)}, \hat{D}^{(T)}) - \hat{d}^{(p)} \\ \hat{d}^{(p)} - \hat{d}^{(s)} \\ \operatorname{diag}(\operatorname{sign}(\hat{z} - \tilde{z}))\hat{d}^{(z)} \\ -\hat{d}^{(s)} \end{array} \right) < 0,$$

and thus there exists a  $\hat{d}$  for any  $\hat{x}$  such that  $J_{\hat{g}_A}(\hat{x})\hat{d} < 0$ , concluding the proof.  $\square$

## A.2 Complexity Analysis

**Proposition A.1** (Parameterized Reachable Set). *Let  $f : \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \mapsto \mathbb{R}^{n_x}$  denote the system dynamics and let  $n = \min(n_x, n_u)$ . Further, we assume that the number of dependent factors and dependent generators of the given initial set grows at most with*

## A Appendix

$O(na)$  and  $O(n)$ , and the number of elementary operations to evaluate  $f_i$  for  $1 \leq i \leq n_x$  and all its necessary derivatives grows at most with  $O(n)$ . The complexity of computing the parameterized reachable set using the conservative polynomialization approach (see Sec. 2.7.3) can then be bounded from above by  $O(n^5 + n^2a^2 + n^3a \log n)$ .

*Proof.* We simply follow the proof in [59, Sec. 4.1.4], where we make the necessary changes according to the made assumptions.

We discuss the complexity of all listed operations in [59, Eq. (4.18)] due to the different assumptions separately:

- Linear map:  $O(n^3)$  [59, Prop. 3.1.18]
- Minkowski sum:  $O(na)$  [59, Prop. 3.1.19]
- Exact sum: The merge operation has complexity  $O((na)^2)$  [59, Prop. 3.1.5] and the compact operation has complexity  $O(n^2a \log n)$  [59, Prop. 3.1.7], so that the overall complexity can be bounded from above by  $O(n^2a^2)$ .
- Cartesian product:  $O(1)$  [59, Prop. 3.1.22]
- Linear combination:  $O(n^2) + O(n^2) = O(n^2)$  [59, Prop. 3.1.26]
- Quadratic map [59, Prop. 3.1.26]

$$\begin{aligned} & O((na)^2) + O(n^2nn) + O(nnnn) + O(nn(n + na \log(nn))) \\ &= O(n^2a^2 + n^4 + n^3a \log n) \end{aligned}$$

- Reduce:  $O(n^2)$  [59, Tab. 3.5]
- Restructure: Not used to preserve dependencies
- Zonotope enclosure:  $O((na)n + nn) = O(an^2)$  [59, Prop. 3.1.14]
- Interval enclosure:  $O(n^3)$  [59, Tab. 3.3]

Thus, adding all the listed complexities above as in [59, Eq. (4.18)] yields

$$O(a^2n^2 + n^4 + an^3 \log n).$$

Since we assume a dimension of  $n$  and that the number of generators and number of elementary operations to evaluate  $f$  and its derivatives grows at most with  $O(n)$ , the upper bound on the complexity of computing the Lagrange remainder follows directly from [59, Sec. 4.1.4] as  $O(n^5)$ . The overall complexity thus follows, concluding the proof.  $\square$

**Proposition A.2** (Polynomial Zonotope Composition). *Given are the polynomial zonotopes  $\mathcal{A} = \langle \cdot, G, G_{\text{rest}}, E \rangle_{PZ} = \{a(\alpha)\}_\alpha$ , where  $G \in \mathbb{R}^{n \times m}$ ,  $G_{\text{rest}} \in \mathbb{R}^{n \times l}$ , and  $E \in \mathbb{N}_{\geq 0}^{n_a \times m}$ , and  $\mathcal{B} = \langle \cdot, G^{(1)}, 0, E^{(1)} \rangle_{PZ} = \{b(\beta)\}_\beta$ , where  $G^{(1)} \in \mathbb{R}^{n_a \times q}$  and  $E^{(1)} \in \mathbb{N}^{n_b \times q}$ . Further, let  $d = \max_{1 \leq k \leq m} \max(E_{(:,k)})$  be the maximum monomial exponent. An upper bound on the complexity of computing the polynomial zonotope  $\{a(b(\beta))\}_\beta = \langle \cdot, \hat{G}, G_{\text{rest}}, \hat{E} \rangle_{PZ}$  is given by*

$$O\left((d(n_a + n_b) + mn)q^d\right).$$

*Proof.* We first derive the complexity of computing  $\{b_i(\beta)^d\}_\beta$  for  $i \in \{1, \dots, n_a\}$  and maximum exponent  $d$  of  $a(\alpha)$ . Then, we extend this complexity result to the computation of  $\{a(b(\beta))\}_\beta$ . Note that the centers of  $\mathcal{A}$  and  $\mathcal{B}$  are ignored as they can be included in the generator and exponent matrices.

Let  $\{b_i(\beta)^k\}_\beta = \langle \cdot, G_{(i,:)}^{(k)}, 0, E^{(k)} \rangle_{PZ}$  for  $k \in \{0, \dots, d\}$  and any  $i \in \{1, \dots, n_a\}$  with  $G_{(i,:)}^{(k)} \in \mathbb{R}^{1 \times q^k}$  and  $E^{(k)} \in \mathbb{N}^{n_b \times q^k}$ . The computation of  $E^{(k)}$  has complexity  $O(n_b q^k)$  since  $E^{(k)} \in \mathbb{N}^{n_b \times q^k}$  (see [61, Prop. 12, Proof]). Similarly, the computation of  $G_{(i,:)}^{(k)} \in \mathbb{R}^{1 \times q^k}$  for  $i \in \{1, \dots, n_a\}$  has complexity  $O(q^k)$ . Because  $E^{(k)}$  is the same for all  $i$ , the complexity of computing  $\{b_i(\beta)^k\}_\beta$  for all  $i \in \{1, \dots, n_a\}$  and  $k \in \{0, \dots, d\}$  is  $O(d(n_a + n_b)q^d)$ .

Next, we form the  $m$  monomials of  $\{a(b(\beta))\}_\beta$ . Each monomial consists of at most  $n_a - 1$  multiplications of the already computed polynomial zonotopes  $\{b_i(\beta)^{k_i}\}_\beta$  for  $i \in \{1, \dots, n_a\}$  and corresponding exponent  $k_i \in \mathbb{N}$  with  $\sum_{k=1}^{n_a} k_i \leq d$ . The exponent matrix of the product of  $n_a$  polynomial zonotopes has at most dimension  $n_b \times q^d$  due to  $\sum_{k=1}^{n_a} k_i \leq d$  and thus the complexity of computing any monomial of  $\{a(b(\beta))\}_\beta$  is bounded from above by  $O(n_b q^d)$ . Finally, we multiply the  $k$ -th monomial, now represented as a polynomial zonotope with a generator matrix of dimension at most  $1 \times q^d$ , with the corresponding column  $G_{(:,k)} \in \mathbb{R}^n$  for all  $k \in \{1, \dots, m\}$  with complexity  $O(m(nq^d))$  and bring these  $m$  monomials into a single polynomial zonotope  $\{a(b(\beta))\}_\beta = \langle \cdot, \hat{G}, G_{\text{rest}}, \hat{E} \rangle_{PZ}$  by concatenation with complexity  $O(m(1))$ . Collecting all complexities then yields the desired result.  $\square$

### A.3 Derivatives of the LQR Feedback Matrix

In order to optimize over the gain matrix  $K(z)$ , we derive its Jacobian and Hessian matrices in App. A.3.1 and App. A.3.2, respectively.

To that end, let  $(A(P), B(P))$  denote a controllable LTI system with  $A(P) \in \mathbb{R}^{n_x \times n_x}$ ,  $B(P) \in \mathbb{R}^{n_x \times n_u}$ ,  $P \in [-1, 1]^{m \times n_u \times a}$ ,  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$ , and  $R \in \mathbb{S}_{++}^{n_u \times n_u}$ . Further, let  $z =$

## A Appendix

$\left[ P_{(\cdot)}^T, Q_{(\cdot)}^T, R_{(\cdot)}^T \right]^T \in \mathbb{R}^{n_z}$ , and we define  $n = \max(n_x, n_u)$ . For convenience, we restate that the optimal LQR feedback gain matrix is given as (we omit arguments from now on)

$$K = -R^{-1}B^T X, \quad (\text{A.27})$$

with

$$F = A^T X + XA - XBR^{-1}B^T X + Q = 0. \quad (\text{A.28})$$

For later convenience, we define here the closed-loop system matrix as  $A_{\text{cl}} = A + BK$ . Further, we introduce the commutation matrix  $M_D \in \{0, 1\}^{nm \times nm}$  for a matrix  $D \in \mathbb{R}^{n \times m}$ , which is implicitly defined by

$$M_D D_{(\cdot)} = \left( D^T \right)_{(\cdot)}.$$

According to Sec. 2.2, we can replace differentials with their corresponding differential tensors, and thus we derive the below expressions directly for these tensors. Due to the definition of the sums and products of these differential tensors, one can replace  $d\Box$  at any point with  $\text{d}\Box$  to recover the ‘‘differential’’ version. Let the differential tensors  $dA \in \mathbb{R}^{1 \times n_z \times n_x \times n_x}$ ,  $dB \in \mathbb{R}^{1 \times n_z \times n_x \times n_u}$ ,  $d^2 A \in \mathbb{R}^{n_z \times n_z \times n_x \times n_x}$ ,  $d^2 B \in \mathbb{R}^{n_z \times n_z \times n_x \times n_u}$ ,  $dQ \in \mathbb{R}^{1 \times n_z \times n_x \times n_x}$ , and  $dR \in \mathbb{R}^{1 \times n_z \times n_u \times n_u}$ , all with respect to  $z$ , be given. Lastly, let  $e = c_e n$  denote the number of elementary operations required to evaluate any element of  $A$ ,  $B$ , or any of its differential tensors, where  $c_e \in \mathbb{R}_+$ . Thus, for a given differential tensor  $M \in \mathbb{R}^{n_z \times n_z \times n \times n}$ , its evaluation complexity is bounded from above by  $O(en_z^2 n^2) = O(n_z^2 n^3)$ .

### A.3.1 Jacobian of Feedback Matrix

**Proposition A.3** (Jacobian Matrix of LQR Feedback Matrix). *Let  $(A(P), B(P))$  be a controllable LTI system with  $A(P) \in \mathbb{R}^{n_x \times n_x}$  and  $B(P) \in \mathbb{R}^{n_x \times n_u}$  for  $P \in [-1, 1]^{m_c n_u \times a}$ , where each element of  $A$  and  $B$  is at least once differentiable with respect to  $P$ . Further, denote with  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $R$  given LQR weighting matrices. The vectorized Jacobian matrix  $\frac{dK_{(\cdot)}}{dz}$  is given by*

$$\frac{dK_{(\cdot)}}{dz} = - \left( K^T \boxtimes R^{-1} \right) dR_{(\cdot)} - \left( X \boxtimes R^{-1} \right) M_B dB_{(\cdot)} - \left( I_n \boxtimes R^{-1} B^T \right) dX_{(\cdot)}, \quad (\text{A.29})$$

where

$$\begin{aligned} dX_{(\cdot)} = & - \left( A_{\text{cl}}^T \boxplus A_{\text{cl}}^T \right)^{-1} \left( \left( (X \boxtimes I_n) M_A + (I_n \boxtimes X) \right) dA_{(\cdot)} + \left( K^T \boxtimes K^T \right) dR_{(\cdot)} \right. \\ & \left. + \left( (X \boxtimes K^T) M_B + \left( K^T \boxtimes X \right) \right) dB_{(\cdot)} + dQ_{(\cdot)} \right), \end{aligned} \quad (\text{A.30})$$

with an upper bound on its computational complexity of  $O(n_x^{2\omega} + n^4 n_z)$ , where  $\omega \in \mathbb{R}_+$  is the exponent for the current complexity of matrix multiplication.

*Proof.* Applying the differential operator to (A.28) yields

$$\begin{aligned}
0 &= dA^T X + A^T dX + dXA + XdA + dXBK + XdBK + K^T dRK + K^T dB^T X \\
&\quad + K^T B^T dX + dQ \\
0 &= dX (A + BK) + (A + BK)^T dX + dA^T X + XdA + K^T dRK + XdBK \\
&\quad + K^T dB^T X + dQ \\
0 &= \left( (A_{\text{cl}}^T \otimes I_n) + (I_n \otimes A_{\text{cl}}^T) \right) dX_{(\cdot)} + \left( (X \otimes I_n) M_A + (I_n \otimes X) \right) dA_{(\cdot)} \\
&\quad + \left( K^T \otimes K^T \right) dR_{(\cdot)} + \left( (X \otimes K^T) M_B + (K^T \otimes X) \right) dB_{(\cdot)} + dQ_{(\cdot)},
\end{aligned}$$

where we used (A.27), and the vectorization follows from (2.1). Solving for  $dX_{(\cdot)}$  then yields (A.30). With

$$\begin{aligned}
dK &= d \left( -R^{-1} B^T X \right) \\
&= R^{-1} dR R^{-1} B^T X - R^{-1} dB^T X - R^{-1} B^T dX,
\end{aligned}$$

we finally arrive at (A.29) after vectorization, which concludes the first part of the proof.

Looking at (A.30), the complexity of computing  $(A_{\text{cl}}^T \boxplus A_{\text{cl}}^T)^{-1}$  is at most  $O((n_x^2)^\omega)$  which follows from  $(A_{\text{cl}}^T \boxplus A_{\text{cl}}^T) \in \mathbb{R}^{n_x^2 \times n_x^2}$  and the computational complexity of the matrix inverse being equal to that of the matrix multiplication [17, Prop. 16.6]. All remaining operations in (A.29) and (A.30) can then be computed by a fixed number of “matrix-vector” products of a  $(n^2 \times n^2)$ -dimensional matrix and a  $(1 \times n_z \times n^2 \times 1)$ -dimensional differential tensor, which has complexity  $O(n^4 n_z)$  since the matrix-vector product of a  $n^2$ -by- $n^2$  matrix and  $n^2$ -by-1 vector has complexity  $O((n^2)^2)$ , but we need to multiple a 1-by- $n_z$  vector with a scalar for each “inner” multiplication. Since forming the Kronecker product of two  $n \times n$ -dimensional matrices has complexity  $O(n^4)$ , the overall complexity of computing the Jacobian of the vectorized gain matrix with respect to  $z$  is at most  $O(n^3 n_z + n^4 + n_x^{2\omega} + n^4 n_z) = O(n_x^{2\omega} + n^4 n_z)$  when including the evaluation complexity  $O(n^3 n_z)$  of the precomputed differential tensors, which concludes the proof.  $\square$

### A.3.2 Hessian of Feedback Matrix

**Proposition A.4** (Hessian Matrix of LQR Feedback Matrix). *Let  $(A(P), B(P))$  be a controllable LTI system with  $A(P) \in \mathbb{R}^{n_x \times n_x}$  and  $B(P) \in \mathbb{R}^{n_x \times n_u}$  for  $P \in [-1, 1]^{m_c n_u \times a}$ , where each element of  $A$  and  $B$  is at least twice differentiable with respect to  $P$ . Further, denote with  $Q \in \mathbb{S}_{++}^{n_x \times n_x}$  and  $R$  given LQR weighting matrices. The Hessian matrix  $H_{K_{ij}}(z)$  of the  $i$ -th row and  $j$ -th column of  $K \in \mathbb{R}^{n_u \times n_x}$  with  $1 \leq i \leq n_u$  and  $1 \leq j \leq n_x$  is given by*

$$H_{K_{ij}}(z) = \frac{1}{2} \left( d^2 K_{ij} + (d^2 K_{ij})^T \right), \tag{A.31}$$

## A Appendix

with a computational complexity for all  $n_x n_u$  entries of at most  $O(n_x^{2\omega} + n^4 n_z^2)$ , where

$$\begin{aligned} d^2 K^T = & -d^2 X B R^{-1} + 2dX \left( -dB + B R^{-1} dR \right) R^{-1} \\ & + X \left( -d^2 B + 2 \left( dB - B R^{-1} dR \right) R^{-1} dR \right) R^{-1}, \end{aligned} \quad (\text{A.32})$$

$$d^2 X_{(\cdot)} = \left( A_{\text{cl}}^T \boxplus A_{\text{cl}}^T \right)^{-1} d^2 C_{(\cdot)}, \quad (\text{A.33})$$

$$d^2 C = - \left( d^2 M + d^2 M^T \right), \quad (\text{A.34})$$

$$\begin{aligned} d^2 M = & + X \left( d^2 A + d B R^{-1} \left( -2B^T dX - 2dR K - dB^T X \right) + d^2 B K \right) \\ & + dX \left( 2dA + 2dBK - 2B R^{-1} dR K - B R^{-1} B^T dX \right) \\ & - K^T dR R^{-1} dR K, \end{aligned} \quad (\text{A.35})$$

$dK$  and  $dX$  are given as in Prop. A.3, and  $\omega \in \mathbb{R}_+$  is the exponent for the current complexity of matrix multiplication.

*Proof.* Applying the differential operator to (A.28) yields (see proof of Prop. A.3)

$$\begin{aligned} dF = & dA^T X + A^T dX \\ & + dXA + XdA \\ & - \left( dX B R^{-1} B^T X + X d B R^{-1} B^T X - X B R^{-1} dR R^{-1} B^T X + X B R^{-1} dB^T X \right. \\ & \left. + X B R^{-1} B^T dX \right) \\ & + dQ \end{aligned}$$



and applying the differential operator again results in

$$\begin{aligned}
 d^2F = & + d^2A^T X + dA^T dX + dA^T dX + A^T d^2X \\
 & + d^2XA + dXdA + dXdA + Xd^2A \\
 & - ( - d^2XBK - dXdBK + dXBR^{-1}dRK + dXBR^{-1}dB^T X + dXBR^{-1}B^T dX \\
 & - dXdBK - Xd^2BK + XdBR^{-1}dRK + XdBR^{-1}dB^T X + XdBR^{-1}B^T dX \\
 & + dXBR^{-1}dRK + XdBR^{-1}dRK + 2K^T dRR^{-1}dRK + K^T dRR^{-1}dB^T X \\
 & + K^T dRR^{-1}B^T dX \\
 & + dXBR^{-1}dB^T X + XdBR^{-1}dB^T X + K^T dRR^{-1}dB^T X - K^T d^2B^T X \\
 & - K^T dB^T dX \\
 & + dXBR^{-1}B^T dX + XdBR^{-1}B^T dX + K^T dRR^{-1}B^T dX - K^T dB^T dX \\
 & - K^T B^T d^2X) \\
 d^2F = & + d^2A^T X + 2dA^T dX + A^T d^2X \\
 & + d^2XA + 2dXdA + Xd^2A \\
 & + d^2XBK + 2dXdBK - 2dXBR^{-1}dRK - 2dXBR^{-1}dB^T X - 2dXBR^{-1}B^T dX \\
 & + Xd^2BK - 2XdBR^{-1}dRK - 2XdBR^{-1}dB^T X - 2XdBR^{-1}B^T dX \\
 & - 2K^T dRR^{-1}dRK - 2K^T dRR^{-1}dB^T X - 2K^T dRR^{-1}B^T dX \\
 & + K^T d^2B^T X + 2K^T dB^T dX \\
 & + K^T B^T d^2X.
 \end{aligned}$$

Next, we define

$$d^2F = A^T d^2X + d^2XA + d^2XBK + K^T B^T d^2X - d^2C, \quad (\text{A.36})$$

where

$$\begin{aligned}
 -d^2C = & + d^2A^T X + 2dA^T dX \\
 & + 2dXdA + Xd^2A \\
 & + 2dXdBK - 2dXBR^{-1}dRK - 2dXBR^{-1}dB^T X - 2dXBR^{-1}B^T dX \\
 & + Xd^2BK - 2XdBR^{-1}dRK - 2XdBR^{-1}dB^T X - 2XdBR^{-1}B^T dX \\
 & - 2K^T dRR^{-1}dRK - 2K^T dRR^{-1}dB^T X - 2K^T dRR^{-1}B^T dX \\
 & + K^T d^2B^T X + 2K^T dB^T dX \\
 = & d^2M + d^2M^T,
 \end{aligned}$$

## A Appendix

collects all terms not containing  $d^2X$ , and where

$$\begin{aligned}
d^2M &= +Xd^2A + 2dXdA + 2dXdBK - 2dXBR^{-1}dRK - 2XdBR^{-1}B^TdX \\
&\quad - dXBR^{-1}B^TdX \\
&\quad + Xd^2BK - 2XdBR^{-1}dRK - XdBR^{-1}dB^TX - K^TdRR^{-1}dRK \\
&= X \left( d^2A - 2dBR^{-1}B^TdX + d^2BK - 2dBR^{-1}dRK - dBR^{-1}dB^TX \right) \\
&\quad + dX \left( 2dA + 2dBK - 2BR^{-1}dRK - BR^{-1}B^TdX \right) \\
&\quad - K^TdRR^{-1}dRK \\
&= +X \left( d^2A + dBR^{-1} \left( -2B^TdX - 2dRK - dB^TX \right) + d^2BK \right) \\
&\quad + dX \left( 2dA + 2dBK - 2BR^{-1}dRK - BR^{-1}B^TdX \right) \\
&\quad - K^TdRR^{-1}dRK.
\end{aligned}$$

Further, it holds that

$$A^Td^2X + d^2XA + d^2XBK + K^TB^Td^2X = A_{\text{cl}}^Td^2X + d^2XA_{\text{cl}}.$$

After vectorization of  $d^2F = 0$  from (A.36), it thus follows that

$$\left( A_{\text{cl}}^T \boxplus A_{\text{cl}}^T \right) d^2X_{(\cdot)} = d^2C_{(\cdot)},$$

and therefore

$$d^2X_{(\cdot)} = \left( A_{\text{cl}}^T \boxplus A_{\text{cl}}^T \right)^{-1} d^2C_{(\cdot)}.$$

Using the definition of  $K$  from (A.27) and applying the differential operator twice yields

$$\begin{aligned}
d^2K^T &= d \left( -dXBR^{-1} - XdBR^{-1} + XBR^{-1}dRR^{-1} \right) \\
&= -d^2XBR^{-1} - dXdBR^{-1} + dXBR^{-1}dRR^{-1} \\
&\quad - dXdBR^{-1} - Xd^2BR^{-1} + XdBR^{-1}dRR^{-1} \\
&\quad + dXBR^{-1}dRR^{-1} + XdBR^{-1}dRR^{-1} - 2XBR^{-1}dRR^{-1}dRR^{-1} \\
&= -d^2XBR^{-1} - 2dXdBR^{-1} + 2dXBR^{-1}dRR^{-1} \\
&\quad - Xd^2BR^{-1} + 2XdBR^{-1}dRR^{-1} \\
&\quad - 2XBR^{-1}dRR^{-1}dRR^{-1} \\
&= -d^2XBR^{-1} + 2dX \left( -dB + BR^{-1}dR \right) R^{-1} \\
&\quad + X \left( -d^2B + 2 \left( dB - BR^{-1}dR \right) R^{-1}dR \right) R^{-1}.
\end{aligned}$$

Evaluating these products and sums of tensors, we obtain  $d^2K^T \in \mathbb{R}^{n_z \times n_z \times n_x \times n_u}$ , which yields  $d^2K$  by transposition. Since

$$d^2K_{ij} = dz^T d^2K_{ij} dz, \quad 1 \leq i \leq n_u, \quad 1 \leq j \leq n_x,$$

### A.3 Derivatives of the LQR Feedback Matrix

the Hessian tensor  $H_K(z) \in \mathbb{R}^{n_z \times n_z \times n_u \times n_x}$  that contains the Hessian matrices of all  $n_u n_x$  entries of  $K \in \mathbb{R}^{n_u \times n_x}$  is given by

$$[H_K(z)]_{(:, :, i, j)} = H_{K_{ij}}(z) = \frac{1}{2} \left( d^2 K_{ij} + d^2 K_{ij}^T \right), \quad (\text{A.37})$$

for  $1 \leq i \leq n_u$  and  $1 \leq j \leq n_x$  due to Th. 2.5.

The computation of  $(A_{\text{cl}}^T \boxplus A_{\text{cl}}^T)^{-1}$  in (A.33) has complexity  $O(n_x^{2\omega})$  (see proof of Prop. A.3). Further, all remaining operations in (A.32) to (A.35) can be computed as a product of an  $(n_z \times n_z \times n \times n)$ -dimensional differential tensor and a  $(n \times n)$ -dimensional matrix, the tensor product between two  $(1 \times n_z \times n \times n)$ -dimensional differential tensors, and a “matrix-vector” tensor product of an  $(n_x^2 \times n_x^2)$ -dimensional matrix with a  $(n_z \times n_z \times n_x^2 \times 1)$ -dimensional differential tensor: The complexity of the first operation is  $O(n^\omega n_z^2)$  since normal matrix multiplication has complexity  $O(n^\omega)$  but here each scalar multiplication is replaced by the multiplication of an  $(n_z \times n_z)$ -dimensional matrix with a scalar. The complexity of the second operation is at most  $O(n^\omega n_z^2)$  since normal matrix multiplication has complexity  $O(n^\omega)$  but here each scalar multiplication is replaced by the outer product of two  $n_z$ -dimensional vectors. Finally, the complexity of the third operation is bounded from above by  $O((n^2)^2 n_z^2)$  since the matrix-vector product of a  $(n^2 \times n^2)$ -dimensional matrix and a  $n^2$ -dimensional matrix is  $O((n^2)^2)$ , but here each element is a  $(n_z \times n_z)$ -dimensional matrix and so each element multiplication of the “outer” matrix-vector product is the multiplication of a scalar with a  $(n_z \times n_z)$ -dimensional matrix. Thus, the combined complexity is then  $O(n^3 n_z^2 + n_x^{2\omega} + n^\omega n_z^2 + n^\omega n_z^2 + n^4 n_z^2) = O(n_x^{2\omega} + n^4 n_z^2)$  (computational complexity of  $dK$  and  $dX$  from Prop. A.3 can be thus be ignored) when including the evaluation complexity of all differential tensors of at most  $O(n^3 n_z^2)$ , which concludes the proof.  $\square$



# Bibliography

- [1] A. Alessio and A. Bemporad. “A survey on explicit model predictive control”. In: *Nonlinear Model Predictive Control*. Springer, 2009, pp. 345–369.
- [2] M. Althoff. “An introduction to CORA 2015”. In: *Workshop on Applied Verification for Continuous and Hybrid Systems*. EasyChair, 2015, pp. 120–151.
- [3] M. Althoff. “On computing the minkowski difference of zonotopes”. In: arXiv, 2015.
- [4] M. Althoff. “Reachability analysis and its application to the safety assessment of autonomous cars”. PhD thesis. Technische Universität München, 2010.
- [5] M. Althoff. “Reachability analysis of nonlinear systems using conservative polynomialization and non-convex sets”. In: *Conference on Hybrid Systems: Computation and Control*. ACM, 2013, pp. 173–182.
- [6] M. Althoff, G. Frehse, and A. Girard. “Set propagation techniques for reachability analysis”. In: *Annual Review of Control, Robotics, and Autonomous Systems* 4.1 (2021), pp. 369–395.
- [7] M. Althoff and B.H. Krogh. “Avoiding geometric intersection operations in reachability analysis of hybrid systems”. In: *Conference on Hybrid Systems: Computation and Control*. ACM, 2012.
- [8] M. Althoff, O. Stursberg, and M. Buss. “Computing reachable sets of hybrid systems using a combination of zonotopes and polytopes”. In: *Nonlinear Analysis: Hybrid Systems* 4.2 (2010), pp. 233–249.
- [9] M. Althoff, O. Stursberg, and M. Buss. “Reachability analysis of nonlinear systems with uncertain parameters using conservative linearization”. In: *Conference on Decision and Control*. IEEE, 2008, pp. 4042–4048.
- [10] J.A.E. Andersson et al. “CasADi: a software framework for nonlinear optimization and optimal control”. In: *Mathematical Programming Computation* 11.1 (2018), pp. 1–36.
- [11] H.E. Bell. “Gershgorin’s theorem and the zeros of polynomials”. In: *The American Mathematical Monthly* 72.3 (1965).
- [12] C. Belta and L. C. G. J. M. Habets. “Controlling a class of nonlinear systems on rectangles”. In: *Transactions on Automatic Control* 51.11 (2006), pp. 1749–1759.
- [13] H.G. Bock and K.J. Plitt. “A multiple shooting algorithm for direct solution of optimal control problems”. In: *IFAC Proceedings Volumes* 17.2 (1984), pp. 1603–1608.

## Bibliography

- [14] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge, United Kingdom: Cambridge university press, 2004.
- [15] S. Boyd et al. *Linear matrix inequalities in system and control theory*. Philadelphia, USA: SIAM, 1994.
- [16] J. van den Brand. “A deterministic linear program solver in current matrix multiplication time”. In: *Symposium on Discrete Algorithms*. ACM-SIAM, 2020, pp. 259–278.
- [17] P. Bürgisser, M. Clausen, and M.A. Shokrollahi. *Algebraic complexity theory*. Berlin, Germany: Springer, 1997.
- [18] C. Cartis, N.I.M. Gould, and P.L. Toint. “On the evaluation complexity of constrained nonlinear least-squares and general constrained nonlinear optimization using second-order methods”. In: *Journal on Numerical Analysis* 53.2 (2015), pp. 836–851.
- [19] M. Černý. “Goffin’s algorithm for zonotopes”. In: *Kybernetika* 48.5 (2012), pp. 890–906.
- [20] N. Chan and S. Mitra. “Verifying safety of an autonomous spacecraft rendezvous mission”. In: *Workshop on Applied Verification of Continuous and Hybrid Systems*. EasyChair, 2017, pp. 20–32.
- [21] X. Chen and S. Sankaranarayanan. “Decomposed reachability analysis for nonlinear systems”. In: *Real-Time Systems Symposium*. IEEE, 2016, pp. 13–24.
- [22] M.B. Cohen, Y.T. Lee, and Z. Song. “Solving linear programs in the current matrix multiplication time”. In: *Journal of the ACM* 68.1 (2021), pp. 1–39.
- [23] P. Collins et al. “Computing the evolution of hybrid systems using rigorous function calculus”. In: *IFAC Proceedings Volumes* 45.9 (2012), pp. 284–290.
- [24] T.H. Cormen et al. *Introduction to algorithms*. Cambridge, USA: MIT press, 2022.
- [25] S. Dasgupta, C.H. Papadimitriou, and U. Vazirani. *Algorithms*. New York, USA: McGraw-Hill Higher Education, 2006.
- [26] M. Diehl et al. “Fast direct multiple shooting algorithms for optimal robot control”. In: *Lecture Notes in Control and Information Sciences*. Springer, 2006, pp. 65–93.
- [27] A.L. Dontchev and R.T. Rockafellar. *Implicit functions and solution mappings*. New York, USA: Springer, 2014.
- [28] C. Durieu, É. Walter, and B. Polyak. “Multi-input multi-output ellipsoidal state bounding”. In: *Journal of Optimization Theory and Applications* 111.2 (2001), pp. 273–303.
- [29] C.H. Edwards. *Advanced calculus of several variables*. Academic Press, New York, USA, 2012.

- [30] H. J. Ferreau, H. G. Bock, and M. Diehl. “An online active set strategy to overcome the limitations of explicit MPC”. In: *Journal of Robust and Nonlinear Control* 18.8 (2008), pp. 816–830.
- [31] J.-A. Ferrez K. Fukuda and T.M. Liebling. “Solving the fixed rank convex quadratic maximization in binary variables by a parallel zonotope construction algorithm”. In: *European Journal of Operational Research* 166.1 (2005), pp. 35–50.
- [32] K. Fukuda. “From the zonotope construction to the Minkowski addition of convex polytopes”. In: *Journal of Symbolic Computation* 38.4 (2004), pp. 1261–1272.
- [33] A. A. Gaganov. “Computation complexity of the range of a polynomial in several variables”. In: *Cybernetics* 21.4 (1985), pp. 418–421.
- [34] V. Gaßmann and M. Althoff. “Implementation of ellipsoidal operations in CORA 2022”. In: *Workshop on Applied Verification of Continuous and Hybrid Systems*. Vol. 90. EasyChair, 2022, pp. 1–17.
- [35] V. Gaßmann and M. Althoff. “Polynomial controller synthesis of nonlinear systems with continuous state feedback using trust regions”. In: *Open Journal of Control Systems* 2 (2023), pp. 310–324.
- [36] V. Gaßmann and M. Althoff. “Scalable zonotope-ellipsoid conversions using the Euclidean zonotope norm”. In: *American Control Conference*. IEEE, 2020, pp. 4715–4721.
- [37] V. Gaßmann and M. Althoff. “Verified polynomial controller synthesis for disturbed nonlinear systems”. In: *IFAC-PapersOnLine* 54.5 (2021), pp. 85–90.
- [38] L. Geretti et al. “ARCH-COMP22 category report: Continuous and hybrid systems with nonlinear dynamics”. In: *Workshop on Applied Verification of Continuous and Hybrid Systems*. Vol. 90. EasyChair, 2022, pp. 86–112.
- [39] A. Girard. “Controller synthesis for safety and reachability via approximate bisimulation”. In: *Automatica* 48.5 (2012), pp. 947–953.
- [40] A. Girard. “Reachability of uncertain linear systems using zonotopes”. In: *Hybrid Systems: Computation and Control*. Springer, 2005, pp. 291–305.
- [41] A. Girard and C. Le Guernic. “Efficient reachability analysis for linear systems using support functions”. In: *IFAC Proceedings Volumes* 41.2 (2008), pp. 8966–8971.
- [42] A. Girard, G. Pola, and P. Tabuada. “Approximately bisimilar symbolic models for incrementally stable switched systems”. In: *Transactions on Automatic Control* 55.1 (2010), pp. 116–126.
- [43] F. De Groote et al. “Evaluation of direct collocation optimal control problem formulations for solving the muscle redundancy problem”. In: *Annals of Biomedical Engineering* 44.10 (2016), pp. 2922–2936.
- [44] F. Gruber and M. Althoff. “Scalable robust model predictive control for linear sampled-data systems”. In: *Conference on Decision and Control*. IEEE, 2019, pp. 438–444.

## Bibliography

- [45] C. Le Guernic. “Reachability analysis of hybrid systems with linear continuous dynamics”. PhD thesis. Université Joseph-Fourier - Grenoble I, 2010.
- [46] L. C. G. J. M. Habets, P. J. Collins, and J. H. Van Schuppen. “Reachability and control synthesis for piecewise-affine hybrid systems on simplices”. In: *Transactions on Automatic Control* 51.6 (2006), pp. 938–948.
- [47] A. Halder. “On the parameterized computation of minimum volume outer ellipsoid of Minkowski sum of ellipsoids”. In: *Conference on Decision and Control*. IEEE, 2018.
- [48] G.D. Halikias et al. “New bounds on the unconstrained quadratic integer programming problem”. In: *Journal of Global Optimization* 39.4 (2007), pp. 543–554.
- [49] B.C. Hall. *Lie groups, lie algebras, and representations*. Cham, Switzerland: Springer, 2015.
- [50] S.-P. Han and O.L. Mangasarian. “Exact penalty functions in nonlinear programming”. In: *Mathematical Programming* 17.1 (1979), pp. 251–269.
- [51] C. Helmberg. *Semidefinite programming for combinatorial optimization*. Berlin, Germany: ZIB, 2000.
- [52] D. Hess, M. Althoff, and T. Sattel. “Formal verification of maneuver automata for parameterized motion primitives”. In: *Conference on Intelligent Robots and Systems*. IEEE, 2014, pp. 1474–1481.
- [53] J.D. Hoffman and S. Frankel. *Numerical methods for engineers and scientists*. New York, USA: McGraw-Hill, 2018.
- [54] R.A. Horn and C.R. Johnson. *Matrix analysis*. Cambridge, United Kingdom: Cambridge University Press, 1985.
- [55] B. Houska, H.J. Ferreau, and M. Diehl. “ACADO Toolkit – an open-source framework for automatic control and dynamic optimization”. In: *Optimal Control Applications and Methods* 32.3 (2011), pp. 298–312.
- [56] R.E. Kalman. “Contributions to the theory of optimal control”. In: *Bol. soc. mat. mexicana* 5.2 (1960), pp. 102–119.
- [57] D.E. Kirk. *Optimal control theory: An introduction*. Mineola, NY: Dover Publications, 2004.
- [58] M. Kloetzer and C. Belta. “A fully automated framework for control of linear systems from temporal logic specifications”. In: *Transactions on Automatic Control* 53.1 (2008), pp. 287–297.
- [59] N. Kochdumper. “Extensions of polynomial zonotopes and their application to verification of cyber-physical systems”. PhD thesis. Technische Universität München, 2022.
- [60] N. Kochdumper and M. Althoff. “Representation of polytopes as polynomial zonotopes”. In: arXiv, 2019.



- [61] N. Kochdumper and M. Althoff. “Sparse polynomial zonotopes: A novel set representation for reachability analysis”. In: *Transactions on Automatic Control* 66.9 (2021), pp. 4043–4058.
- [62] N. Kochdumper, B. Schürmann, and M. Althoff. “Utilizing dependencies to obtain subsets of reachable sets”. In: *International Conference on Hybrid Systems: Computation and Control*. ACM, 2020, pp. 1–10.
- [63] N. Kochdumper et al. “AROC: a toolbox for automated reachset optimal controller synthesis”. In: *Conference on Hybrid Systems: Computation and Control*. ACM, 2021, pp. 1–6.
- [64] A. Kopetzki, B. Schürmann, and M. Althoff. “Methods for order reduction of zonotopes”. In: *Conference on Decision and Control*. IEEE, 2017, pp. 5626–5633.
- [65] A.B. Kurzhanski and P. Varaiya. “On ellipsoidal techniques for reachability analysis. Part I: External approximations”. In: *Optimization Methods and Software* 17.2 (2002), pp. 177–206.
- [66] A.A. Kurzhanskiy and P. Varaiya. “Ellipsoidal toolbox (ET)”. In: *Conference on Decision and Control*. IEEE, 2006, pp. 1498–1503.
- [67] A.A. Kurzhanskiy and P. Varaiya. “Reach set computation and control synthesis for discrete-time dynamical systems with disturbances”. In: *Automatica* 47.7 (2011), pp. 1414–1426.
- [68] G. Lafferriere, G.J. Pappas, and S. Yovine. “Symbolic reachability computation for families of linear vector fields”. In: *Journal of Symbolic Computation* 32.3 (2001), pp. 231–253.
- [69] A. Laub. “A schur method for solving algebraic riccati equations”. In: *Transactions on Automatic Control* 24.6 (1979), pp. 913–921.
- [70] A. J. Laub. *Matrix analysis for scientists and engineers*. Vol. 91. Philadelphia, USA: SIAM, 2005.
- [71] P. Leopardi. “Distributing points on the sphere”. PhD thesis. University of New South Wales, 2007.
- [72] D. Limon et al. “Robust tube-based MPC for tracking of constrained linear systems with additive disturbances”. In: *Journal of Process Control* 20.3 (2010), pp. 248–260.
- [73] J. Linhart. “Approximation of a ball by zonotopes using uniform distribution on the sphere”. In: *Archiv der Mathematik* 53.1 (1989), pp. 82–86.
- [74] J.R. Magnus and H. Neudecker. *Matrix differential calculus with applications in statistics and econometrics*. Hoboken, USA: John Wiley & Sons, 2019.
- [75] U. Malik et al. “On the gap between the quadratic integer programming problem and its semidefinite relaxation”. In: *Mathematical programming* 107.3 (2006), pp. 505–515.

## Bibliography

- [76] O.L. Mangasarian and S. Fromovitz. “The Fritz John necessary optimality conditions in the presence of equality and inequality constraints”. In: *Journal of Mathematical Analysis and Applications* 17.1 (1967), pp. 37–47.
- [77] S. Manzinger, M. Leibold, and M. Althoff. “Driving strategy selection for cooperative vehicles using maneuver templates”. In: *Intelligent Vehicles Symposium*. IEEE, 2017, pp. 647–654.
- [78] D. Q. Mayne et al. “Robust output feedback model predictive control of constrained linear systems”. In: *Automatica* 42.7 (2006), pp. 1217–1222.
- [79] D. Q. Mayne et al. “Robust output feedback model predictive control of constrained linear systems: Time varying case”. In: *Automatica* 45.9 (2009), pp. 2082–2087.
- [80] D.Q. Mayne et al. “Tube-based robust nonlinear model predictive control”. In: *Journal of Robust and Nonlinear Control* 21.11 (2011), pp. 1341–1353.
- [81] F. Messerer and M. Diehl. “An efficient algorithm for tube-based robust nonlinear optimal control with optimal linear feedback”. In: *Conference on Decision and Control*. IEEE, 2021, pp. 6714–6721.
- [82] M. Morari and J.H. Lee. “Model predictive control: Past, present and future”. In: *Computers & Chemical Engineering* 23.4-5 (1999), pp. 667–682.
- [83] M.E. Muller. “A note on a method for generating points uniformly on n-dimensional spheres”. In: *Communications of the ACM* 2.4 (1959), pp. 19–20.
- [84] K.G. Murty and S.N. Kabadi. “Some NP-complete problems in quadratic and nonlinear programming”. In: *Mathematical Programming* 39.2 (1987), pp. 117–129.
- [85] C.L. Navasca and A.J. Krener. “Solution of Hamilton Jacobi Bellman equations”. In: *Conference on Decision and Control*. IEEE, 2000, pp. 4442–4447.
- [86] C.H. Papadimitriou. “On the complexity of integer programming”. In: *Journal of the ACM* 28.4 (1981), pp. 765–768.
- [87] I. Papusha et al. “Automata theory meets approximate dynamic programming: Optimal control with temporal logic constraints”. In: *Conference on Decision and Control*. IEEE, 2016, pp. 434–440.
- [88] G. Pola, A. Girard, and P. Tabuada. “Approximately bisimilar symbolic models for nonlinear control systems”. In: *Automatica* 44.10 (2008), pp. 2508–2516.
- [89] G. Pola and P. Tabuada. “Symbolic models for nonlinear control systems: Alternating approximate bisimulations”. In: *Journal on Control and Optimization* 48.2 (2009), pp. 719–733.
- [90] G. Pola et al. “Symbolic models for nonlinear time-delay systems using approximate bisimulations”. In: *Systems & Control Letters* 59.6 (2010), pp. 365–373.
- [91] S. V. Rakovic et al. “Parameterized tube model predictive control”. In: *Transactions on Automatic Control* 57.11 (2012), pp. 2746–2761.

- [92] A.V. Rao. “A survey of numerical methods for optimal control”. In: *Advances in the Astronautical Sciences* 135.1 (2009), pp. 497–528.
- [93] J. B. Rawlings, D. Q. Mayne, and M. Diehl. *Model predictive control: Theory, computation, and design*. Wisconsin, WI: Nob Hill Publishing, 2009.
- [94] G. Reissig. “Computing abstractions of nonlinear systems”. In: *Transactions on Automatic Control* 56.11 (2011), pp. 2583–2598.
- [95] G. Reissig, A. Weber, and M. Rungger. “Feedback refinement relations for the synthesis of symbolic controllers”. In: *Transactions on Automatic Control* 62.4 (2017), pp. 1781–1796.
- [96] P. Riedinger, J. Daafouz, and C. Iung. “About solving hybrid optimal control problems”. In: *IMACS05* (2005).
- [97] S. Sadraddini and R. Tedrake. “Linear encodings for polytope containment problems”. In: *Conference on Decision and Control*. IEEE, 2019, pp. 4367–4372.
- [98] B. Schürmann. “Using reachability analysis in controller synthesis for safety-critical systems”. PhD thesis. Technische Universität München, 2022.
- [99] B. Schürmann and M. Althoff. “Guaranteeing constraints of disturbed nonlinear systems using set-based optimal control in generator space”. In: *IFAC-PapersOnLine* 50.1 (2017), pp. 11515–11522.
- [100] B. Schürmann and M. Althoff. “Optimal control of sets of solutions to formally guarantee constraints of disturbed linear systems”. In: *American Control Conference*. IEEE, 2017, pp. 2522–2529.
- [101] B. Schürmann and M. Althoff. “Optimizing sets of solutions for controlling constrained nonlinear systems”. In: *Transactions on Automatic Control* 66.3 (2021), pp. 981–994.
- [102] B. Schürmann, N. Kochdumper, and M. Althoff. “Reachset model predictive control for disturbed nonlinear systems”. In: *Conference on Decision and Control*. IEEE, 2018, pp. 3463–3470.
- [103] B. Sendov. *Hausdorff approximations*. Dordrecht, The Netherlands: Kluwer Academic Publishers, 1990.
- [104] O. von Stryk. “Numerical solution of optimal control problems by direct collocation”. In: *Optimal Control*. Birkhäuser, 1993, pp. 129–143.
- [105] O. von Stryk and R. Bulirsch. “Direct and indirect methods for trajectory optimization”. In: *Annals of Operations Research* 37.1 (1992), pp. 357–373.
- [106] P. Sun and R.M. Freund. “Computation of minimum-volume covering ellipsoids”. In: *Operations Research* 52.5 (2004), pp. 690–706.
- [107] P. Tabuada. *Verification and control of hybrid systems: A symbolic approach*. New York, USA: Springer, 2009.
- [108] H.R. Tiwary. “On computing the shadows and slices of polytopes”. In: arXiv, 2008.

## Bibliography

- [109] H.R. Tiwary. “On the hardness of computing intersection, union and Minkowski sum of polytopes”. In: *Discrete & Computational Geometry* 40.3 (2008), pp. 469–479.
- [110] L. Vandenberghe and S. Boyd. “Semidefinite programming”. In: *SIAM Review* 38.1 (1996), pp. 49–95.
- [111] A. Wächter and L.T. Biegler. “On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming”. In: *Mathematical Programming* 106.1 (2005), pp. 25–57.
- [112] M. Wetzlinger, N. Kochdumper, and M. Althoff. “Adaptive parameter tuning for reachability analysis of linear systems”. In: *Conference on Decision and Control*. IEEE, 2020, pp. 5145–5152.
- [113] M. Wetzlinger et al. “Adaptive reachability algorithms for nonlinear systems using abstraction error analysis”. In: *Nonlinear Analysis: Hybrid Systems* 46 (2022), p. 101252.
- [114] M. Wetzlinger et al. “Fully automated verification of linear systems using inner- and outer-approximations of reachable sets”. In: *Transactions on Automatic Control* (2023), pp. 1–16.
- [115] M. Wetzlinger et al. “Fully-automated verification of linear systems using reachability analysis with support functions”. In: *International Conference on Hybrid Systems: Computation and Control*. ACM, 2023.
- [116] E. M. Wolff and R. M. Murray. “Optimal control of nonlinear systems with temporal logic specifications”. In: *Robotics Research*. Springer, 2016, pp. 21–37.
- [117] A. Yershova and S.M. LaValle. “Deterministic sampling methods for spheres and  $SO(3)$ ”. In: *Conference on Robotics and Automation*. Vol. 4. IEEE, 2004, pp. 3974–3980.
- [118] M. Zamani et al. “Symbolic models for nonlinear control systems without stability assumptions”. In: *Transactions on Automatic Control* 57.7 (2012), pp. 1804–1809.
- [119] M.N. Zeilinger et al. “On real-time robust model predictive control”. In: *Automatica* 50.3 (2014), pp. 683–694.
- [120] G.M. Ziegler. *Lectures on polytopes*. New York, USA: Springer, 1995.