

# Dynamic Car-Passenger Matching of Online and Reservation Requests

## **Marvin Erdmann**

BMW AG  
Mobility Services  
Parkring 19  
85748 Garching, Germany  
Email: marvin.erdmann@bmw.de

## **Florian Dandl**

Bundeswehr University Munich  
Institute for Intelligent Transportation Systems  
Werner-Heisenberg-Weg 39  
85779 Neubiberg, Germany  
Email: florian.dandl@unibw.de  
ORCID: 0000-0003-3706-9725

## **Bernd Kaltenhäuser**

Baden-Wuerttemberg Cooperative State University  
Department of Technical Management  
Friedrich-Ebert-Str. 30  
78054 Villingen-Schwenningen, Germany  
Email: bernd.kaltenhaeuser@dhbw-vs.de  
ORCID: 0000-0001-9786-6465

## **Klaus Bogenberger**

Bundeswehr University Munich  
Institute for Intelligent Transportation Systems  
Werner-Heisenberg-Weg 39  
85779 Neubiberg, Germany  
Email: klaus.bogenberger@unibw.de

Word count: 7111 words text + 1 tables x 250 words (each) = 7361 words

Number figures: 5

Submission Date: August 1, 2019

## 1 ABSTRACT

2 Dynamic Car-Passenger Matching is a variant of the Dial-a-Ride Problem (DaRP) and a crucial part of  
3 operating an On-Demand Mobility (ODM) service. The ODM fleet management algorithm has to be able to  
4 find close-to-optimal solutions quick and reliable in order to provide a high quality service to the users of  
5 the service. In this work, an approach is presented that combines customer requests for immediate pick-ups  
6 (online) and reservations for pick-ups in the near future. This approach is based on the method of Global  
7 Optimization with Time Windows (GOTW), which is able to benefit from both a very fast accept/reject  
8 response with an estimated pick-up time window generated by a greedy Nearest Neighbor Policy (NNP)  
9 heuristic as well as the optimization potential of a periodically executed Tabu Search. This work also intro-  
10 duces List-Based Assignments (LBAs) as a replacement of the initial NNP in order to shorten computational  
11 times of optimization without losing any feasible solution. The solutions provided by the algorithm using  
12 this new approach generate up to 67.7% less empty mileage compared to the solutions found by NNP. This  
13 advantage outweighs the slightly higher number of rejections and longer waiting times for most of the eval-  
14 uated scenarios. Especially in settings considering reservations, the new methodology clearly outperforms  
15 the benchmark algorithm by up to around 36% in terms of the overall quality of solution.

16 **Keywords:** Dynamic Car-Passenger Matching, On-Demand Mobility, Global Optimization with Time  
17 Windows, List-Based Assignments, Reservations

## 18 INTRODUCTION

19 In today's urban traffic travelers face many problems and inconveniences, among others congestion and the  
20 lack of space due to many parking cars. On-Demand Mobility (ODM) is an increasingly popular concept  
21 that tries to tackle these challenges by allowing people to profit from individual urban mobility without  
22 owning a car. The enhanced use of shared mobility services and the increased utilization of vehicles in cities  
23 would lead to less space occupied by parking cars without losing the benefits of always available mobility.

24 To operate a fleet with many cars and customers, a management algorithm is necessary to match service  
25 requests and vehicles of the fleet dynamically. It must be able to quickly find a reliable and time efficient  
26 solution for the whole system. Such a problem is called Dial-a-Ride-Problem (DaRP)(1). As a general-  
27 ization of the Traveling Salesman Problem it is NP-hard, which means finding optimal solutions requires  
28 exponentially rising computational time with increasing problem size. Due to the dynamism of the problem  
29 and the necessity of short response time to customers of the service, heuristics are often used to circumvent  
30 these long computational times. An often used method is called Nearest Neighbor Policy (NNP), which  
31 matches cars to new customers immediately based on the respective arrival times. Instead of local optimiza-  
32 tion (NNP), more sophisticated approaches like metaheuristics aim to optimize the whole fleet, but avoid  
33 exponential computational times by smart exploration of the solution space. Osman and Kelly (2) defined  
34 a metaheuristic as '[...] an iterative generation process which guides a subordinate heuristic by combining  
35 intelligently different concepts for exploring and exploiting the search spaces using learning strategies to  
36 structure information in order to find efficiently near-optimal solutions'.

37 The significance of metaheuristics to solve the DaRP for the ODM use case led to an increasing interest  
38 in this research area recently. The static version of the problem has been discussed in (3) by Cordeau and  
39 Laporte in which they use a metaheuristic procedure called Tabu Search (TS). Brandão used TS for the  
40 heterogeneous fixed fleet vehicle routing problem in (4). In (5), Prodhon and Prins present and compare the  
41 most important heuristics and metaheuristics for the Vehicle Routing Problem, in which they emphasize the  
42 good performances of TS algorithms. Pandi et al. (6) used a GPU-accelerated TS algorithm to solve the  
43 DaRP, which produced comparable results up to 30 times faster than a single-core CPU-based TS algorithm.

44 The models used in simulations to test the algorithms vary significantly, since there are countless  
45 specifics an ODM service may offer to its customers. Alonso-Mora et al. (7) considered the ride-sharing  
46 use case for immediate-pick-up requests using constraint optimization. Hyland and Mahmassani (8) focus  
47 on the ride-hailing use case. They compare six assignment strategies, four of which assign requests only  
48 once and could therefore provide information rather quickly to the customers. In (9), Sheridan et al. propose

1 a dynamic nearest neighbor policy to find good solutions for highly dynamic problems in a very fast way.  
 2 Though, none of the mentioned papers considers both reservations and online requests as possible ways for  
 3 customers to request a pick-up.

4 However, the incorporation of reservations into the dynamic optimization of online requests is an impor-  
 5 tant task for mobility services offering both options to their customers, like Uber (10) and Lyft (11). Short  
 6 accept/reject response times are as important for reservations as for online requests. Therefore, this initial  
 7 decision should be made by a fast heuristic and the optimization should only take place periodically instead  
 8 of each time a new request occurs. Hence, this work is based on an approach called Global Optimization  
 9 with Time Windows (12) which combines the quick response times of NNP with the optimization potential  
 10 of TS. The new approach considers reservations as well as online request to find nearly optimal solutions  
 11 for the combined problem. It focuses on the so called ride-hailing use case in which at most one customer  
 12 occupies a car.

13 The objective of this work is to present a method to combine ODM requests for immediate pick-ups  
 14 (online) with reservation requests in a close-to-optimal way using TS. To guarantee a fast response to the  
 15 customers, a heuristic based on the idea of a NNP is used to determine if a request is accepted or not. The  
 16 benefits of applying the presented methods in a model which uses the New York Taxi Data Set (13) will be  
 17 evaluated and compared to results found using only NNP, as it is an often used standard for the problem in  
 18 real world applications.

## 19 PROBLEM FORMULATION

20 The DaRP has been formulated in different variants since it first was established. In this work, a very  
 21 common formulation is used with constraints found in e.g. (7), (8) and (14).

22 A solution for the DaRP is considered to be optimal if both the fleet costs and the dissatisfaction of the  
 23 customers are minimal over a given period of time. As fleet costs are mainly produced by the movement and  
 24 corresponding fuel consumption of vehicles in the fleet, the overall distance between all visited locations  
 25 in the solution should be as short as possible. More specifically, the empty mileages between drop-off  
 26 locations and pick-up locations should be minimized, as the demanded routes from a customer's pick-up  
 27 location to his or her destination are independent of the assignments and should therefore not be part of the  
 28 optimization.

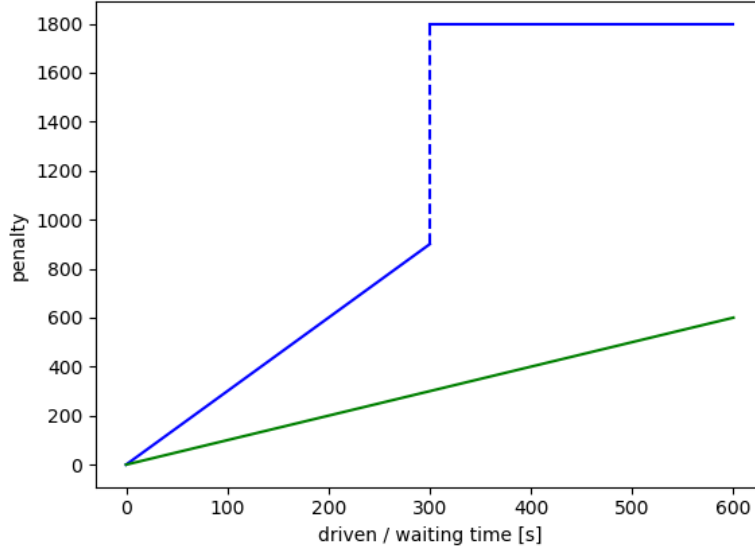
29 In the dynamic version of the DaRP, new online customer requests occur eventually. These requests are  
 30 not known ahead of time, so the operator can only solve the DaRP with a subset of active requests in order to  
 31 find solutions shortly after customers apply for the service. Hence, this optimization of new requests should  
 32 take place periodically after a certain time window  $t_{\text{period}}$  in which new requests are accepted or rejected. On  
 33 the one hand,  $t_{\text{period}}$  - defined by the fleet operator - has to be short enough to guarantee quick response times  
 34 to the customers. On the other hand, the longer  $t_{\text{period}}$  is chosen to last, the more requests are optimized at  
 35 once, increasing the optimization potential. So, the choice of this parameter is critical.

36 Each time an optimization takes place, all eligible customers  $N$  with pick-up locations  $i \in [1, \dots, N]$   
 37 and drop-off locations  $j \in [N + 1, \dots, 2N]$  and an ODM fleet with  $M$  vehicles at their next idle locations  
 38  $m \in [2N + 1, \dots, 2N + M]$  are considered. The set of pick-up and drop-off locations of the customers is  
 39 referred to as  $P$  and  $D$  respectively. The next idle location of a vehicle is either its actual position (if it is  
 40 idle) or the drop-off location of the last request it is matched to. The set of all next idle locations is defined  
 41 as  $C$ .  $V = P \cup D \cup C$  is the set of vertices of an undirected graph  $G = (V, A)$ , with  $A$  the set of arcs  
 42 containing all feasible pairs of locations  $(a, b)$ , with  $a, b \in V$ .

43 If in a solution a vehicle moves from  $a$  to  $b$ , the decision variable  $x_{a,b}^m = 1$ , otherwise  $x_{a,b}^m = 0$ . Every  
 44 arc  $(a, b)$  is weighted with a cost  $c_{a,b}$  proportional to the distance between  $a$  and  $b$ .

45 In this work, the objective function consists of two terms, both of which ought to be minimized.

$$\min (f_{\text{obj}}) = \min_{x_{a,b}^m} \left( \sum_{m \in M} \sum_{a \in D \cup C} \sum_{b \in P} c_{a,b} \times x_{a,b}^m + \sum_{i \in N} d_i (x_{a,b}^m) \right), \quad (1)$$



**FIGURE 1 Penalty model: customer dissatisfaction  $t_w \times s$  if  $t_w \leq t_{mw}$  or  $t_{mw} \times s \times r$  otherwise (blue); cars causing penalty equal to empty driven time in seconds (green).**

1 The first term represents the empty driven mileage of all cars in the fleet. Each empty car is defined to  
 2 produce a penalty of  $\frac{1}{s}$  when on its way to pick up a customer.

3 The second term is the sum of all customers' dissatisfaction. The dissatisfaction  $d_i$  of each customer  
 4  $i \in N$  depends on the waiting time  $t_w$  of this customer. With increasing waiting time,  $d_i$  rises linearly with  
 5 slope  $s$ . In this work  $s = \frac{3}{\text{sec}}$ , indicating the priority of customer satisfaction compared to fleet costs. If no  
 6 car would be able to pick up a customer before his or her maximum waiting time  $t_{mw}$ , the algorithm rejects  
 7 the request immediately. If a request is rejected,  $d_i$  is set to a constant value, which is  $r$  times higher than the  
 8 dissatisfaction at  $t_{mw}$ .  $r$  is set to the value of 2, reflecting the expectation of customers not using the service  
 9 in the future anymore when being rejected once. The penalty model is outlined in figure 1.

10 While optimizing the solution, the following constraints must be satisfied at every time:

$$\sum_{i \in C \cup D} \sum_{j \in P} x_{i,j}^m - x_{j,j+N}^m = 0 \quad \forall m \in M \quad (2)$$

11

$$\sum_{m \in M} x_{i,j}^m \leq 1 \quad \forall j \in P \quad (3)$$

12

$$t_{i,p} - t_{i,d} < 0 \quad \forall i \in N \quad (4)$$

13 Equation 2 ensures that if a customer  $j \in N$  is picked up, the drop-off follows immediately and is  
 14 done by the same car (pairing constraint). So, there is at most one customer in a car at each time (capacity  
 15 constraint). In Equation 3, the constraint is formulated that every customer is served at most once. The  
 16 precedence constraint is stated in equation 4 by only allowing pick-up times  $t_{i,p}$  earlier than drop-off times  
 17  $t_{i,d}$  for each customer  $i \in N$  respectively.

18 Additionally, reservations are considered to be another type of service customers might request. They  
 19 can be ordered some time before the desired pick-up time, which allows customers to plan trips in advance.  
 20 Reservations always have to be served in time in the proposed model, which means cars need to arrive at the  
 21 pick-up location no later than the requested pick-up time. If the operator cannot guarantee the service, the  
 22 reservation request is rejected.

## 1 METHODOLOGY

2 After a short introduction to the Nearest Neighbor Policy, the Tabu Search procedure and Global Opti-  
3 mization with Time Windows, the methods of List-Based Assignments and the incorporation of reservation  
4 requests in online request optimization are explained in this section.

### 6 Nearest Neighbor Policy

7 The Nearest Neighbor Policy (NNP) is a simple heuristic designed to find feasible solutions of the DaRP  
8 very fast. Every time a new online request occurs, the distances from each of the vehicles' next idle locations  
9 to the pick-up location of the new customer are calculated and the earliest possible pick-up times are derived.  
10 If the maximum waiting time of the customer is not surpassed, the vehicle that is able to pick up the new  
11 customer the earliest is matched permanently with that request, otherwise the request is rejected. The request  
12 is immediately added to the route of the car and it moves directly to the pick-up location if it is idle.

13 The benefit of using the NNP is its ability to produce feasible solutions without any considerable delay  
14 as the computational time necessary to find solutions with it scales linearly with the fleet size. As it only  
15 considers the very next request, it tends to fail finding optimal solutions for the whole system. Nonetheless,  
16 in real world applications those solutions are used to operate ODM fleets. Hence, this heuristic is used as a  
17 generator of initial solutions and as the benchmark for the proposed algorithm.

### 19 Tabu Search

20 The concept of Tabu Search (TS) is a variant of the so called Local Search technique, which in general  
21 improves initial solutions by applying modifications (or moves) to it iteratively. The basic idea of the TS  
22 metaheuristic is to avoid local minima of the objective function value by allowing moves in a local search  
23 that produce worse solutions while forbidding already visited (tabu) solutions.

24 In this work, a preliminary solution is generated with the NNP following the process explained in the  
25 previous section. At the end of each optimization period after  $t_{\text{period}}$ , this becomes the initial solution of  
26 the optimization. Three types of local moves are then applied to it. The swap move switches requests  
27  $(R_1, \dots, R_N)$  within the route of one car  $([R_1, R_2] \rightarrow [R_2, R_1])$ , the shift move changes the position of a  
28 request from one car's route to the end of another car's route  $([R_1, R_2], [R_3, R_4] \rightarrow [R_1], [R_3, R_4, R_2])$ .  
29 The third move is called interchange and exchanges the positions of two requests in two different cars'  
30 routes  $([R_1, R_2], [R_3, R_4] \rightarrow [R_1, R_4], [R_3, R_2])$ . These moves are executed for all cars and pairs of cars  
31 of the fleet and the resulting solution is tested to be feasible. All feasible solutions are then defined as the  
32 neighborhood of the initial solution.

33 The solutions in the neighborhood of the initial solution are sorted by their value of the objective func-  
34 tion. Beginning with the best solution, the so called Tabu List is checked, which indicates if a solution was  
35 already found recently. If so, the next solution is checked until a maximum number of solutions was checked  
36 and the search terminates. If a new solution was found before the search terminates, this solution becomes  
37 the initial solution of the next iteration. If this solution is better than the Best Solution Found (BSF) in this  
38 optimization step, it also becomes the BSF. The next iteration of optimization starts by applying local moves  
39 on the new solution. This procedure is repeated until a) the search terminates because it cannot find a solution  
40 that was not recently found or b) the search does not produce a BSF for a certain amount of iterations  $I_{\text{noBSF}}$   
41 or c) a maximum number of iterations  $I_{\text{max}}$  had been performed or d) the time limit set by the optimization  
42 period's length is reached. After the TS terminated, the last BSF is the resulting solution of the optimization.

### 44 Global Optimization with Time Windows

45 Global Optimization with Time Windows (GOTW) is a method introduced in (12) which is able to benefit  
46 from both the response speed of the NNP as well as the optimization potential of the TS.

47 At the beginning of each optimization period, an empty list  $L_{\text{opt}}$  is set up, which will contain all requests  
48 subject to optimization at the end of the optimization period. After the simulation of all cars' movements  
49 according to their currently matched requests, at every simulation step new requests that are added to the

1 problem are preliminary matched to a car or rejected, following the NNP. Every accepted request is added  
 2 to  $L_{opt}$  and if the preliminary matched car is idle, it moves to the pick-up location starting with the next  
 3 simulation step. A customer is not allowed to be rejected after his or her request was initially accepted and  
 4 is not allowed to be re-scheduled in a way he or she would have to wait longer than the maximum waiting  
 5 time, defined by the fleet operator. By using the NNP as an initial decision maker for accepting or rejecting  
 6 a request, it is possible to send a response to the customer very quick, either a rejection or a latest pick-up  
 7 time.

8 After adding new requests to the problem, it is checked if preliminary matched requests' pick-up loca-  
 9 tions are already reached by the respective cars. If so, those requests are matched permanently to the car,  
 10 added to the actual solution and deleted from  $L_{opt}$ .

11 This simulation circle is executed until the end of the current optimization period. At this point, all  
 12 requests in  $L_{opt}$  are subject to a global optimization using TS as described in the previous section. The  
 13 found solution is considered as permanent and the customer could now be informed about the exact pick-up  
 14 modalities.

15 Another optimization period begins by setting up a new  $L_{opt}$ . This procedure is performed periodically  
 16 until the end of simulation.

### 17 **List-Based Assignments**

18 The idea of List-Based Assignments (LBAs) is to improve the initial solution by enhancing the simple NNP  
 19 heuristic while decreasing the size of the solution space needed to be searched by the TS metaheuristic,  
 20 without eliminating feasible solutions. A short example is presented to illustrate the approach.

21 As described in section *Nearest Neighbor Policy*, the NNP searches for the vehicle that is able to pick  
 22 up a customer  $i \in N$  fastest, once a new request is submitted. The other cars' earliest pick-up times for  
 23 that particular request are not stored. When another request  $j \in N$  is submitted, this car is considered to be  
 24 occupied until the drop-off time of customer  $i$ , to guarantee a pick-up in time. If no other car is able to pick  
 25 up customer  $j$ , this request is rejected.

26 With the method of LBAs though, all cars that are able to pick up a customer  $i \in N$  before he or she  
 27 had to wait longer than the maximum waiting time  $t_{mw}$  are added to a list  $L_{LBA,i}$ . This list is sorted by the  
 28 earliest possible pick-up time and the fastest car is matched preliminary with the request. Thus, initially the  
 29 same preliminary match is made for customer  $i$  as it would have been done following the NNP.

30 However, when the aforementioned request  $j \in N$  occurs, the preliminary matched car is not considered  
 31 to be occupied until it would drop off customer  $i$ . If it is able to pick up customer  $j$  before the respective  
 32 maximum waiting time and if there is another car in  $L_{LBA,i}$  which had not been preliminary matched to  
 33 another request in the meantime and is still able to pick up customer  $i$  in time, the request of customer  $j$  can  
 34 be accepted by preliminary matching it to the car, that was first sent to request  $i$ .

35 By doing so, more online requests can be accepted in a respective optimization period compared to only  
 36 using NNP without losing the guarantee for accepted customers to be picked up in a time shorter than the  
 37 maximum waiting time and without considerable additional computational effort.

38 Furthermore, as only vehicles in the list  $L_{LBA,i}$  are able to pick up customer  $i$ , the TS can neglect solu-  
 39 tions, in which  $i$  is matched to a vehicle that is not in  $L_{LBA,i}$ . This decreases the number of possible solutions  
 40 tremendously, which translates to a much faster optimization. As this version of the DaRP is very dynamic,  
 41 this advantage is crucial for the success of the proposed algorithm in real world applications.

### 42 **Optimization of Reservations**

43 As customer  $i \in N$  requesting a ride at an exact point in time in the future  $t_{res,i}$  wants to be informed about  
 44 his or her acceptance or rejection as quick as possible, the initial decision is made by the NNP as for online  
 45 requests. In contrast to online requests, matched requests are not added to the actual route of the respective  
 46 car  $m \in M$ , but instead stored in another list, denoted as  $L_{res,m}$ .

47 As explained in section *Problem Formulation*, requests in  $L_{res,m}$  are not part of optimization and cars

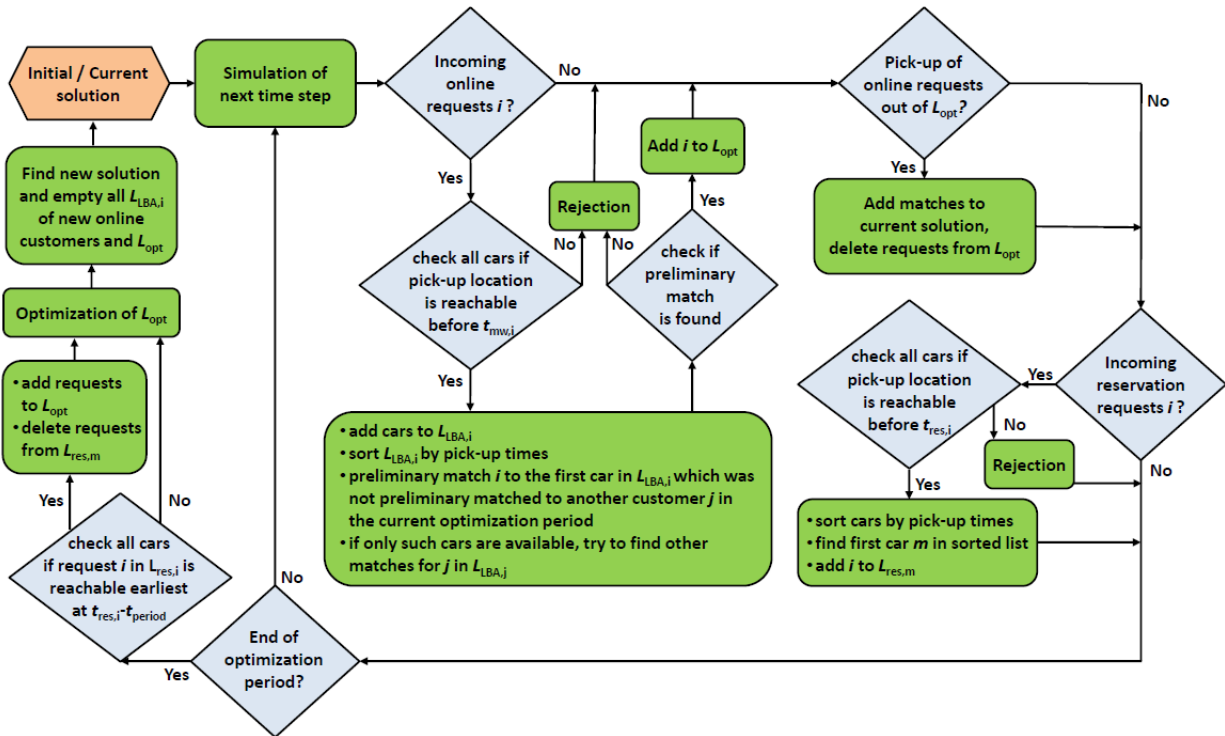


FIGURE 2 Concept of GOTW with LBAs, considering both online requests and reservations.

1 are not moving to the pick-up location of the next request on the list in the model proposed in this work.  
 2 Instead, at the end of each optimization period with length  $t_{\text{period}}$  the list of reservations  $L_{\text{res},m}$  is checked  
 3 for every car  $m \in M$ : if a car  $m$  has to start moving to the pick-up location of its next reservation request  
 4  $i$  before the next optimization period in order to arrive there in time, the request is deleted from  $L_{\text{res},m}$   
 5 and added to the route of the car. In the proposed algorithm it becomes subject to optimization (in the  
 6 benchmarking algorithm using only NNP, car  $m$  has to perform the service). If request  $i$  is scheduled to be  
 7 served by another car  $m'$  after the optimization and this car is not needed to start immediately to pick up  
 8 the customer in time, the request is deleted from that car's route and added to  $L_{\text{res},m'}$ . Therefore, car  $m'$  can  
 9 potentially serve future online requests or reservations  $i \in N$  with earlier reservation times  $t_{\text{res},j} < t_{\text{res},i}$ , as  
 10 long as the pick-up of customer  $i$  at  $t_{\text{res},i}$  is still guaranteed.

11 This procedure maximizes the time cars are able to perform other tasks before picking up customers  
 12 that requested the service in advance. It also produces significantly less empty driven miles, since it allows  
 13 considering the spatial distribution of cars closer to the actual pick-up time. By optimizing reservation  
 14 requests only when they become urgent a given time budget for optimization can be used more efficiently  
 15 since the TS can focus on the solution space involving more relevant assignments of the near future, resulting  
 16 in shorter computational times.

17 In figure 2 the principle of the GOTW with LBAs is depicted, considering online requests as well as  
 18 reservations.

## 19 CASE STUDY

20 In this section, the simulation settings will be explained in detail before the presentation, evaluation and  
 21 discussion of the results.

## 22 Simulation Settings

23 In this work, the open source New York Taxi Data (13) is used to simulate the demand. Because computa-  
 24 tional time for solving the problem increases with growing problem size, only 5% of all requests are taken  
 25

**TABLE 1 Summary of optimization and simulation settings**

TS local moves	swaps, shifts, interchanges
$I_{\text{noBSF}}$	10 iterations
$I_{\text{max}}$	100 iterations
penalty for empty driven mileage $s$	1 per second
waiting time penalty factor $s$	3 per second
rejection penalty factor $r$	2
online customer maximum waiting time $t_{\text{mw}}$	300 seconds
operator acceptance/rejection time	immediate
optimization period	60 seconds
reservation announcement time	10 – 60 minutes
reservation customer maximum delay	0 seconds
reservation percentages	0%, 50%, 100%
average number of customers per hour	1000 customers
fleet sizes	20, 40, ..., 300 cars
warm up period	60 minutes
simulation period	60 minutes
simulation time step	1 second

1 into account in this work in order to generate a meaningful number of runs to evaluate the statistical effect  
2 of the proposed methods. In each run, the selection of requests taken into account is randomly sampled.  
3 The simulations start at 6 p.m. and last until 7 p.m. of each day from Monday, June 20th 2016 to Thursday,  
4 June 23rd 2016, where around 1000 requests per simulated hour must be handled. The time of simulation  
5 is chosen to test the algorithms during the evening demand peaks in weekdays. This is the last full week  
6 the data source provides GPS information on pick-up and drop-off locations of the requests. Before each  
7 simulation, a one-hour warm up phase is run to ensure a realistic distribution of vehicles in service at the  
8 start of the simulation.

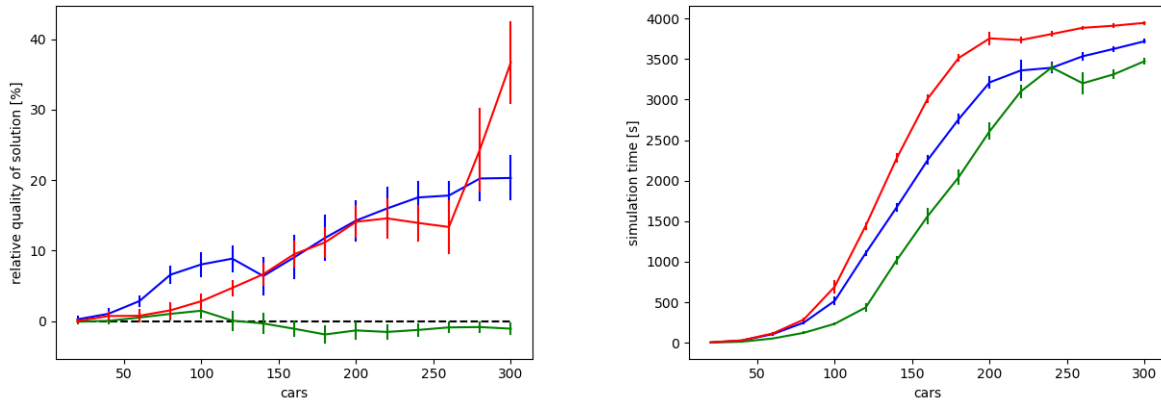
9 The positions of cars as well as pick-up and drop-off locations are projected on a simplified grid, mea-  
10 suring distances between two points on it by adding horizontal and vertical differences in position (Man-  
11 hattan distance). The grid space is set to ten meters and vehicles in the simulation are able to move this  
12 distance in one simulation step which is equivalent to one second of simulated time. This leads to a constant  
13 simulated velocity of 36 kilometers per hour.

14 To compare the results of the proposed algorithm and the NNP, each simulation is initialized twice  
15 with the same set of cars and customers, one following the NNP only, as described in section *Nearest*  
16 *Neighbor Policy*, the other one using the methods explained in sections *Tabu Search*, *Global Optimization*  
17 *with Time Windows*, *List-Based Assignments* and *Optimization of Reservations* to find optimized solutions  
18 at each optimization step. To find the benefits of both the list-based assignments (only for online requests,  
19 as described in section *List-Based Assignments*) and the incorporation of reservations (only for scheduled  
20 requests, as described in section *Optimization of Reservations*), simulations are run using three reservation  
21 percentages: 0%, 50% and 100%.

22 The optimization periods are 60 seconds long which is considered as a reasonable amount of time a  
23 customer is willing to wait for the exact information when he or she will be picked up. The same applies for  
24 the maximum waiting time of online customers, which is set to 300 seconds. Reservations can be made 10  
25 to 60 minutes in advance. Simulations including reservations last until every request becomes permanently  
26 matched and therefore contributes to the objective function value. The key performance indicators (KPIs)  
27 are measured at the end of each one hour simulation time window.

28 The number of available ODM fleet vehicles varies from 20 to 300 in 15 isometric steps to evaluate





(a) Qualities of solutions found with GOTW compared to NNP, representing relative values of the objective function in percent. (b) Computational times for Tabu Search optimization in seconds.

**FIGURE 3** Qualities of solutions and computational times for optimization for scenarios with 0% (green), 50% (blue) and 100% (red) reservations using GOTW (solid) and NNP (dashed).

1 the behavior of the algorithms with different ratios of demand (number of requests) and supply (number of  
2 cars).

3 Each set up is simulated ten times for each of the chosen simulation dates, leading to a total of 40 runs  
4 per setting and an overall number of 3600 simulations taken into account in this work. These simulations  
5 are run on a twelve-core Intel Xeon E5-2687W 3.0GHz processor. An overview of the scenarios is given in  
6 table 1.

7 In addition to that evaluation, to quantify the benefit of LBAs in terms of computational times, two sets  
8 of simulations with the same settings are initialized without the termination criterion (d) in the Tabu Search,  
9 that terminates the search after the optimization period's length. One set of simulations uses LBAs to cut  
10 the solution space at each optimization step, the other set searches the solution space defined by all cars of  
11 the fleet. The simulations consider scenarios with 50% of all requests being reservations.

12

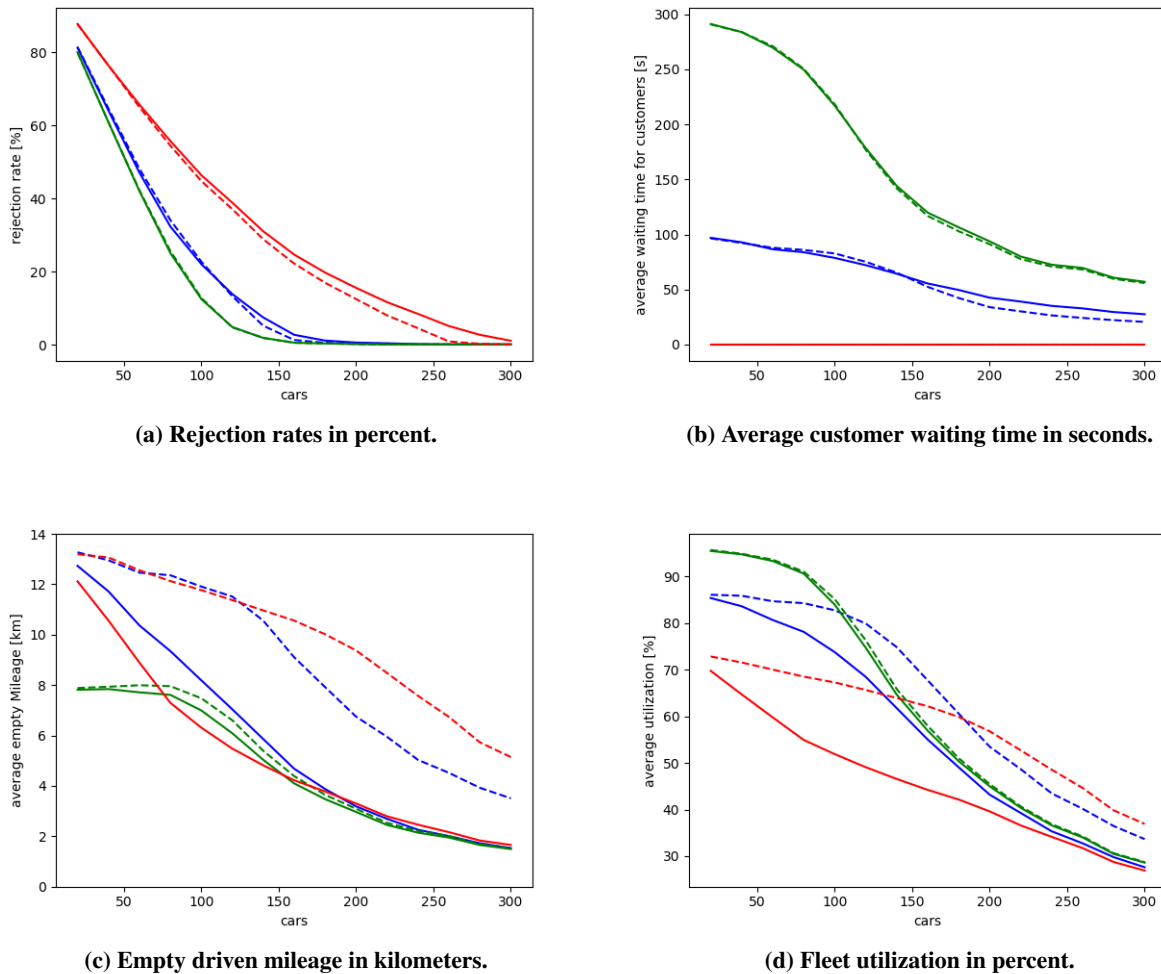
### 13 Results

14 The solution qualities as well as the computational times for the simulations using optimization are presented  
15 in figure 3 to evaluate the overall performances of both methods. Error bars indicate the 95%-confidence  
16 interval of the respective measurement. The quality of a solution (QoS) found with GOTW relative to  
17 solutions generated with NNP is defined as

$$\text{QoS} = \begin{cases} 1 - \frac{f_{\text{obj,GOTW}}}{f_{\text{obj,NNP}}} & \text{if } f_{\text{obj,GOTW}} > f_{\text{obj,NNP}} \\ \frac{f_{\text{obj,NNP}}}{f_{\text{obj,GOTW}}} - 1 & \text{otherwise.} \end{cases} \quad (5)$$

18 Additionally the average values for rejection percentages, average waiting times of customers, total  
19 empty mileages per car and fleet utilization are compared for both approaches, as shown in figure 4. To  
20 avoid confusion in the graphs, error bars are hidden in this figure.

21 In the case of no reservations, the overall impact of the new approach is the smallest for all ratios of  
22 available cars and customers. The rejection percentages are very similar over all simulated numbers of  
23 cars and fall below one percent in scenarios with more than 140 cars. The differences in driven and empty  
24 mileages are small in general, but the GOTW algorithm produces less empty mileage in every setting. In



**FIGURE 4 KPIs for scenarios with 0% (green), 50% (blue) and 100% (red) reservations using GOTW (solid) and NNP (dashed).**

1 scenarios with fleet sizes short of meeting the demand, the difference in produced empty mileage peaks,  
 2 leading to slightly better overall solutions. With 100 cars, GOTW produces solution qualities which are  
 3 1.49% higher than those generated with NNP. On the other hand, GOTW is outperformed by the NNP in  
 4 scenarios in which the rejection percentage is smaller than one percent. The worst performance is measured  
 5 with 180 cars, where GOTW is 1.87% worse than NNP on average. This result is caused by the difference  
 6 of around 2 seconds in average waiting time per accepted customer, which outweighs the slight benefit of  
 7 saved empty mileage for fleet sizes sufficient to hold rejection rates under one percent.

8 When reservations are considered though, the overall solution qualities produced with GOTW are gen-  
 9 erally higher than solutions generated with NNP. If only reservations are requested (100%-reservation case),  
 10 the objective function value is the sum of the penalties induced by rejections and empty miles. Each cus-  
 11 tomer is either picked up at the requested pick-up time or is rejected and therefore produces no penalty due  
 12 to individual waiting time. Compared to NNP, solutions found with GOTW imply generally less driven  
 13 mileage per car and specifically less empty mileage. With increasing number of vehicles in the fleet, the  
 14 empty mileage produced with solutions found by GOTW decreases much faster than in NNP-solutions, up  
 15 to the maximum simulated fleet sizes with 300 cars. In these scenarios GOTW produces 1.66 km average

1 empty mileage, compared to 5.14 km with NNP, a relative improvement of 67.7%. At the same time, the re-  
2 jection rate in simulations using the GOTW approach is considerably higher compared to the NNP method if  
3 only reservations are considered. This behavior becomes more and more obvious with increasing fleet sizes  
4 up to the point when the demand is met, in the case of GOTW at around 300 vehicles. The empty mileage  
5 becomes increasingly dominant with more cars in the fleet, because the number of rejections decreases.  
6 Since the difference in empty mileage produced in the GOTW and NNP solutions also increases, the results  
7 found with GOTW become more and more superior to the benchmark up to maximum fleet size of 300 cars  
8 with an improvement of 36.69%. A local minimum in the solution quality is found only when the difference  
9 in rejections peaks at around 260 cars, where in solutions with NNP around 0.91% of all requests had been  
10 rejected, while with GOTW 5.11% of the customers were rejected.

11 In the mixed case with one half of all requests reservations, the other half customers which want to be  
12 picked up as soon as possible, the behaviors of both aforementioned cases are recognizable in the resulting  
13 KPIs. The produced empty mileage becomes more dominant with increasing number of cars, as it is the  
14 case when only reservations are considered. As the GOTW approach produces solutions with less empty  
15 mileage than in NNP's solutions and the difference in this aspect grows with increasing fleet size, the overall  
16 solutions found with GOTW tend to become better the more cars are part of the fleet. Also, this trend is  
17 only retarded when the difference in rejections reaches its peak shortly before the demand is met, as it was  
18 observed before in the 100%-reservation case. Though, the number of cars necessary to meet the demand is  
19 much more similar to the one found in the use case considering only online requests (approximately 140),  
20 with around 180 vehicles instead of between 260 and 300. The differences in average waiting times of online  
21 customers are even higher than in the 0%-reservation case for fleet sizes meeting the demand. Still, these  
22 differences are not as critical for the objective function values as the empty mileage is the dominant term  
23 in this genre. In scenarios with fleet sizes not matching the demand, the average customer waiting times  
24 produced in solutions with GOTW tend to be shorter than in those found with NNP.

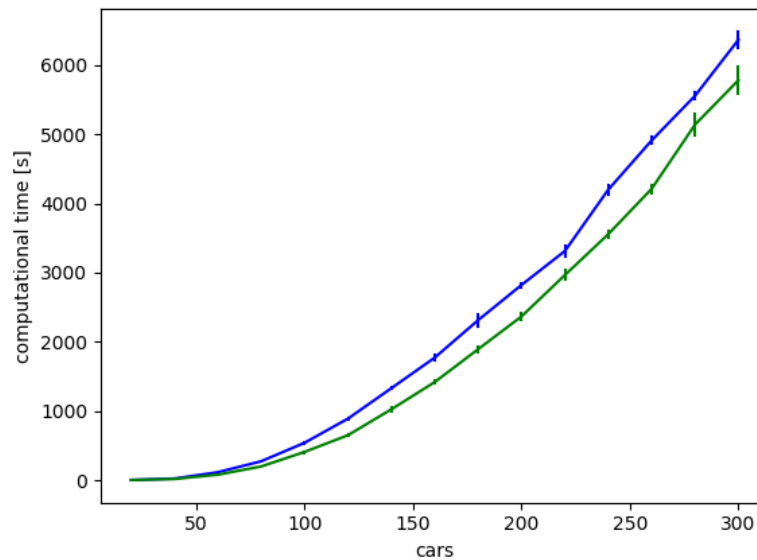
25 The computational times for simulations rise exponentially as expected due to the kind of the optimiza-  
26 tion problem. The converging behavior for bigger fleet sizes is caused by the termination condition (d) of  
27 the Tabu Search (see section *Tabu Search*) which terminates the search after a wall-clock time equal to the  
28 optimization period's length. The higher the reservation percentage of an evaluated set up at a certain fleet  
29 size, the longer the computational time becomes. That behavior can be explained by the method of taking  
30 reservations into account in optimization: requests might be part of optimization more often than once until  
31 they are matched permanently. Online requests on the other hand are matched just once and thereafter are  
32 not part of optimization, leading to lower average numbers of requests per optimization and therefore shorter  
33 computational times in total.

34 In figure 5, the computational times taken for the sets of simulations using Tabu Search optimization not  
35 limited by the termination criterion (d) are compared. It is apparent that using LBAs reduces the computa-  
36 tional effort significantly. With rising fleet size and therefore more solutions to evaluate in the optimization,  
37 the total difference in computational time taken to find solutions increases. For same fleet sizes and settings,  
38 solutions produced with LBAs are found between 7.5% (280 cars) and 31.4% (60 cars) faster compared to  
39 the equivalent algorithm without LBAs, with an average of 20.0%.

## 41 Discussion and Future Work

42 The results shown in the previous section give indications for strengths and weaknesses of the proposed  
43 algorithm.

44 A clear disadvantage is the performance in scenarios without reservations. In these scenarios, the  
45 proposed LBAs were expected to improve the solution quality by increasing the number of accepted requests  
46 compared to the simple NNP algorithm. However, due to longer average waiting times for customers the  
47 objective function values produced with LBAs and periodic optimization are worse than those found with  
48 NNP. One explanation for this behavior is the nature of the dynamic DaRP. At each optimization step, the  
49 goal of the optimization algorithm is to find the best car-passengers matches at this point in time. The



**FIGURE 5** Computational times without time-out termination criterion in TS with (green) and without (blue) LBAs.

1 future vehicle distribution changes according to origin-destination patterns of served demand and the made  
 2 assignments. While an optimal assignment at one point in time is clearly beneficial for the currently known  
 3 requests, the assignments might be worse considering not yet revealed future requests. That behavior was  
 4 described in detail by Dandl et al. in (15).

5 A similar explanation holds for the higher number of rejections in the cases with reservations. In these  
 6 scenarios, cars are preliminary matched with requests for reservations as long as the cars do not need to start  
 7 moving to pick up the customer in time. In that case, an optimization takes place which might find another  
 8 car which is closer to the pick-up location. This procedure is repeated until no closer car is found. Then the  
 9 match becomes permanent and the car starts moving to the pick-up location. As for the online scenario, the  
 10 results can be explained by the dynamic problem combined with the centralized distribution of requests in  
 11 the evaluated data. Cars which are able to pick up customers faster than other cars are more often close to the  
 12 center, since that is where the most new requests occur. That means that those cars tend to be matched with  
 13 more and more reservation requests until they do not have any capacity left to accept new ones. The NNP  
 14 does not search for closer cars when a car needs to start moving in order to pick up a customer in time. This  
 15 leads to much longer ways to the pick-up locations - represented by the significantly higher empty mileages  
 16 in these solutions - but it also decreases the density of requests in the schedule of cars near the center which  
 17 are therefore able to serve new requests more often.

18 The benefits of the proposed methods are clearly found in scenarios considering reservations. Especially  
 19 if the fleet size is chosen to be sufficient to meet the demand and the rejection rate is very low, the approach  
 20 presented in this work outperforms the NNP due to the empty mileage saved.

21 Using LBAs in order to avoid considering infeasible solutions in the optimization causes a signifi-  
 22 cant improvement in terms of computational run time necessary to find solutions with the Tabu Search  
 23 metaheuristic. This advantage allows more iterations per optimization period, potentially leading to better  
 24 solutions and decreases the time necessary to find these solutions, resulting in shorter response times to  
 25 customers, a crucial feature in highly dynamic problems.

26 Future work should further enhance the degree of realism increasing the problem scale to the actual  
 27 demand and the actual sizes of ODM fleets in cities like New York City as well as using real maps, traffic

1 information and routing instead of a simplified grid. Another emphasis could be the improvement of the  
2 search procedure itself to find better solutions in deeper solution spaces as well as the combination of the  
3 proposed model with a repositioning algorithm, which would increase the solution qualities by reducing the  
4 number of rejections and shortening individual customer waiting times.

## 5 **CONCLUSION**

6 This work's objective was to introduce an approach which is able to combine online requests of customers,  
7 which want to be picked up as soon as possible by a vehicle of an ODM fleet and reservation requests,  
8 in which customers demand a pick-up at a certain point in the future by cars of the same fleet. The pro-  
9 posed method is based on the concept of Global Optimization with Time Windows (GOTW). It combines  
10 the benefits of a quick initial decision, determining whether a request is accepted or rejected by a simple  
11 insertion heuristic based on the Nearest Neighbor Policy (NNP), and the optimization potential of a Tabu  
12 Search metaheuristic, which optimizes the matching of new requests and vehicles in the fleet periodically.

13 The presented methodology of List-Based Assignments (LBAs) is able to reduce the computational  
14 time of the Tabu Search procedure by 20% on average due to exclusion of infeasible solutions before the  
15 actual optimization takes place.

16 The proposed method to take requests for reservation into account in the optimization of the car-  
17 passenger matches generates less empty mileage driven by the ODM fleet than the NNP (up to 67.7%  
18 in scenarios with 300 cars and 100% reservations). This benefit outweighs the higher rejection rates in so-  
19 lutions found with the proposed algorithm compared to the benchmark solutions following the NNP by up  
20 to 36.69%.

21 Upcoming work should evaluate the shortcomings of the presented approaches in scenarios with no  
22 reservations, where the average customer waiting time is higher than in the benchmark solutions. Also,  
23 repositioning of idle cars to decrease average waiting times and rejection rates should be considered.

## 24 **AUTHOR CONTRIBUTION**

25 Marvin Erdmann, Florian Dandl, Bernd Kaltenhäuser and Klaus Bogenberger designed the research; Mar-  
26 vin Erdmann implemented the simulations and data analysis. Marvin Erdmann, Florian Dandl and Bernd  
27 Kaltenhäuser interpreted the results and prepared the manuscript. All authors reviewed the results and ap-  
28 proved the final version of the manuscript.

1 **REFERENCES**

- 2 [1] Psaraftis, H. N. A Dynamic Programming Solution to the Single Vehicle Many-to-Many Immediate  
3 Request Dial-a-Ride Problem. *Transportation Science*, 1980. 14: 130–154.
- 4 [2] Osman, I. H., and J. P. Kelly. Meta-Heuristics: An Overview. In *Meta-Heuristics* (I. H. Osman and J.  
5 P. Kelly, eds.), Springer US, Boston, 1996, pp. 1–21.
- 6 [3] Cordeau, J.-F., and G. Laporte. A tabu search heuristic for the static multi-vehicle dial-a-ride problem.  
7 *Transportation Research Part B*, 2003. 37: 579–594.
- 8 [4] Brandão, J. A tabu search algorithm for the heterogeneous fixed fleet vehicle routing problem. *Com-  
9 puters & Operations Research*, 2011. 38: 140–151.
- 10 [5] Prodhon, C., and C. Prins. Metaheuristics for Vehicle Routing Problems. In *Metaheuristics* (P. Siarry,  
11 ed.), Springer International Publishing Switzerland, Basel, 2016, pp. 407–437.
- 12 [6] Pandi, R. R., S. G. Ho, S. C. Nagavarapu, T. Tripathy, and J. Dauwels. GPU-Accelerated Tabu Search  
13 Algorithm for Dial-A-Ride Problem. *Proceedings of 21st International IEEE Conference on Intelligent  
14 Transportation Systems*, 2018. 2519–2524.
- 15 [7] Alonso-Mora, J., S. Samaranyake, A. Wallar, E. Frazzoli, and D. Rus. On-demand high-capacity ride-  
16 sharing via dynamic trip-vehicle assignment. *Proceedings of the National Academy of Sciences*, 2017.  
17 114: 462–467.
- 18 [8] Hyland, M., and H. S. Mahmassani. Dynamic autonomous vehicle fleet operations: Optimization-  
19 based strategies to assign AVs to immediate traveler demand requests. *Transportation Research Part  
20 C*, 2018. 92: 278–297.
- 21 [9] Sheridan, P. K., E. Gluck, Q. Guan, T. Pickles, B. Balçioğlu, and B. Benhabib. The Dynamic Nearest  
22 Neighbor Policy for the Multi-Vehicle Pick-Up and Delivery Problem. *Transportation Research Part  
23 A*, 2013. 49: 178–194.
- 24 [10] Uber. *Riding with Uber - Schedule a Ride*. Uber Technologies Inc., San Francisco. <https://www.uber.com/ride/how-uber-works/scheduled-rides>. Accessed June 26, 2019.
- 25  
26 [11] Lyft. *Scheduled Rides for Passenger*. Lyft Inc., San Francisco. [https://help.lyft.com/hc/en-  
27 us/articles/115013078668-Scheduled-rides-for-passengers](https://help.lyft.com/hc/en-us/articles/115013078668-Scheduled-rides-for-passengers). Accessed June 26, 2019.
- 28 [12] Erdmann, M., F. Dandl, and K. Bogenberger. Dynamic Car-Passenger Matching based on Tabu Search  
29 using Global Optimization with Time Windows. Presented at the 8th International Conference on Mod-  
30 eling, Simulation and Applied Optimization, Manama, Bahrain, 2019.
- 31 [13] Taxicab Passenger Enhancement Program. *2016 Yellow Taxi Trip Data*. NYC Taxi and Limou-  
32 sine Commision, New York City. [https://data.cityofnewyork.us/Transportation/2016-Yellow-Taxi-  
33 Trip-Data/k67s-dv2t](https://data.cityofnewyork.us/Transportation/2016-Yellow-Taxi-Trip-Data/k67s-dv2t). Accessed May 15, 2019.
- 34 [14] Dandl, F., and K. Bogenberger. Booking Processes in Autonomous Carsharing and Taxi Systems.  
35 *Proceedings of 7th Transportation Research Arena*, Vienna, Austria, 2018.
- 36 [15] Dandl, F., M. Hyland, K. Bogenberger, and H. S. Mahmassani. Evaluating the impact of spatio-  
37 temporal demand forecast aggregation on the operational performance of shared autonomous mobility  
38 fleets. *Transportation*, 2019. 114: 462–484