# ExaNIML

## **An Exascale Library for Numerically Inspired Machine Learning**

Severin Reiz<sup>§</sup>, Hasan Ashraf<sup>§</sup>, Tobias Neckel<sup>§</sup>, George Biros<sup>†</sup>, Hans-Joachim Bungartz<sup>§</sup> Technical University of Munich, <sup>†</sup> University of Texas at Austin, reiz@in.tum.de, ashraf@in.tum.de, neckel@in.tum.de, gbiros@ices.utexas.edu, bungartz@in.tum.de

## Introduction

### Motivation

- Significant gap in communities: Machine Learning (ML)  $\leftrightarrow$  high-performance computing (HPC)
- ML needs considerable computing power
- $\rightarrow$  we need adequate software!
- ExaNIML: Library with algorithms
- -with modern applications from ML community
- -with enough concurrency for next generation distributed computing systems

### Approach

- Methods from scientific computing domain
  - "Undervalued" near-linear complexity methods (fast multipole methods)<sup>1</sup>
  - Adaptive sparse grids to mitigate curse of dimensionality<sup>2</sup>
- **HPC:** Exploit potential of supercomputers
  - -Concurrency: Choose suitable algorithms for parallel computing
  - Extract computational bottlenecks as low-level drivers in C++ or Kokkos
- -Performance Portability in-light of the upcoming new GPU and CPU architectures

### **Classification with Kernel Methods**



#### Kernel Matrix

Occurs in many domains ...

- multi-class classification
- model-order reduction
- uncertainty quantification
- partial differential equations

**Example: Binary classification** 

### Ridge regression

- N data points  $x_i \in \mathbb{R}^d$  and N binary labels  $y_i$ •  $f(x) = \operatorname{sign}(\sum_{i=1}^{N} k(x, x_i)w_i) \rightarrow u = f(x_{test}) = K * w$
- Solving a linear system: kernel matrix Koften not stable, nearly singular. Solve  $K \to K + \lambda I$  instead

#### **Our approach: Kernel Matrix Approximation**

- Often K is a dense N-by-N matrix; this quadratic complexity often is the **computational bottleneck**
- To reach  $\mathcal{O}(N)$  algorithms it requires approximation
- For the majority of applications off-diagonal blocks of K admit good low-rank approximations

Many ML libraries offer Kernel methods: to our knowledge none of them offer kernel matrix approximation

### **Key Computational Bottlenecks**

### Geometry-oblivious Fast Multipole Method<sup>1</sup>

- Hierarchically off-diagonal low-rank
- Speeds up algebraic operations





### First scalability results

Multiplication<sup>1</sup>



### $\mathcal{O}(N)$ linear solver<sup>3</sup>



4-process distributed  $\mathcal{H}$ -Matrix compression. Mixed colored sections and factors are shared for dis-

Dependency graph for asynchronous task analysis tributed nodes





Time for compression (left) and multiplication Time for ULV-Factorization (left) and forward-(right) for a 6-d gaussian kernel matrix of 64Mby-64M

solve (right)

### **Dimensionality Reduction**



Reduce the dimensionality of dataset Manifold Learning Algorithms

- (Kernel) Principal Component Analysis
- Isomap algorithm
- Hessian local eigenmaps, ...

### Classification on Embedded Space

- Example forced to 2D manifold (plotting)
- Classification on lower dimensional manifold
- $\rightarrow$  Sparse grid classification<sup>4</sup>

### **Approximation with Sparse Grids**

### Sparse grids

- Reduce number of grid points
- One approach: Combi-Technique Combine Full Grids (red and blue, c.f. figure on right)
- Suitable for 5-20 dimensions

### Sparse Grids in Embedded Space

- 1. Manifold learning algorithm for coarse embedded space
- 2. **Fine** approximation in embedded space with **Sparse Grids**

Synergy between Point-based and Gridbased Methods









References	[2] HJ. Bungartz and M. Griebel, "Sparse grids," Acta numerica, vol. 13, pp. 147–269, 2004.
	[3] D. Y. Chenhan, S. Reiz, and G. Biros, "Distributed o (n) linear solver for dense symmetric hierarchical semi-separable matrices," in 2019 IEEE 13th International Symposium on Embedded Multicore/Many-core Systems-on-Chip (MCSoC), pp. 1–8, IEEE, 2019.
[1] C. D. Yu, S. Reiz, and G. Biros, "Distributed-memory hierarchical compression of dense SPD matrices," in Proceedings of the International Con ference for High Performance Computing, Networking, Storage, and Analysis, SC '18, (Piscataway, NJ, USA), pp. 15:1–15:15, IEEE Press, 2018.	- [4] B. Peherstorfer, D. Pflüger, and HJ. Bungartz, "Density estimation with adaptive sparse grids for large data sets," in <i>Proceedings of the 2014 SIAM international conference on data mining</i> , pp. 443–451, SIAM, 2014.