# Mathematisches Forschungsinstitut Oberwolfach

# Computation and Learning in High Dimensions

Organized by
Markus Bachmayr, Aachen
Albert Cohen, Paris
Angela Kunoth, Köln
Olga Mula, Eindhoven

17 August – 22 August 2025

ABSTRACT. Computations with functions depending on a large numbers of variables are at the core of many problems in science and engineering. They arise naturally in physical models described by partial differential equations (PDEs) depending on many parameters, in purely data-driven tasks such as optimization and machine learning, and in hybrid contexts combining physical models with data.

Traditionally, dealing with such high dimensionality was avoided by the use of simplified models. With the availability of more computational power and the development of sophisticated approximation schemes and algorithms, however, such tasks in high dimensions are increasingly treated directly on the basis of general mathematical principles.

The naive use of classical approximation methods for such problems typically leads to computational costs that scale exponentially with respect to the dimension, an effect known as the *curse of dimensionality*. To make computations tractable, *nonlinear strategies* that leverage in more subtle ways inherent properties of the problem are inevitably required. In recent years, many new and diverse approaches have emerged from different fields. Shaping up the theoretical foundations for the analysis and development of these approaches requires new interactions between approximation theory, numerical analysis, probability theory, mathematical and statistical learning theory, and optimization. This workshop aimed to deepen the mathematical foundations of the underlying numerical concepts that drive this new evolution of computational methods, and to promote the exchange of ideas arising from various disciplines about how to treat high-dimensional problems.

# Introduction by the Organizers

**Scientific context and challenges:** Problems involving high-dimensional spaces come up in many different applications, and pose a challenge for accurate and reliable numerical solution concepts. Such problems arise in:

- Physics-Driven approaches: PDE models depending on large numbers of variables arise naturally in the context of probability theory (such as Kolmogorov equations of random processes), in quantum physics (such as the electronic Schrödinger equation in molecular models), or more generally, in particle and kinetic equations.
- Data-Driven approaches: Another relevant situation where high-dimensional computations arise concern problems depending on many parameters. An important example in this category is machine learning, where one aims to learn functions of very large numbers of variables (such as classifiers that take an image as their input) from large sets of data.
- Hybrid approaches: In inverse problem tasks combining classical PDE models with measurement data, one needs to handle functions on high-dimensional spaces such as posterior distributions in Bayesian formulations. PDEs involving large numbers of parameters also arise in this context and represent yet another example of high-dimensional objects.

Mathematical methods that can successfully address such problems need to exploit in a subtle way their structural features to capture higher-dimensional objects by a tractable amount of information. A common feature of these methods is that they produce approximations with a high degree of nonlinearity, ranging from sparse and adaptive basis expansions as well as low-rank tensor representations to compositional approximations as in neural networks.

Not every one of these approaches is equally suited to every problem. In addition, often conflicts between theoretical accuracy and practical feasibility occur. Indeed, a higher degree of expressive power and nonlinearity – in particular in the case of neural networks which provide a large degree of universality – frequently leads to substantial difficulties in the numerical computation of corresponding approximations. In high-dimensional problems, the practical numerical challenges are frequently so apparent that the mere existence of efficient approximations is often much less relevant than the existence of practically feasible algorithms for computing them.

Next we briefly discuss the current state of the art and challenges concerning theory and numerical algorithms for the abovely mentioned three families of high-dimensional problems (physics-driven, data-driven, and hybrid problems). In the field of numerical analysis of PDEs, spatially adaptive methods using wavelet expansions or finite element bases have advanced the frontiers of computability for certain PDE classes, and feasible algorithms with optimal scaling costs are known. Although standard adaptive paradigms using local spatial refinements are not tractable in high dimensions, similar optimality results have been obtained for certain high-dimensional PDEs using anisotropic sparse expansions or low-rank

tensor representations. Such results remain, however, an exception so far: for many other types of nonlinear approximations, such as deep neural network approximations of PDE solutions, complexity or even convergence guarantees for numerical schemes remain out of reach. One can draw a similar landscape regarding model order reduction of parametrized PDEs: although solid theoretical foundations exist for the approximation of parametric elliptic problems, a cohesive theory and efficient methods to address other types of PDE classes remains an open problem for very active research.

The challenges faced by the above-discussed forward PDE computations are inherited and even amplified in inverse state or parameter estimation problems. One reason is that observational data is often scarce due to prohibitive cost or severe obstructions to acquiring it. PDE data is also a scarce commodity due to the complexity of solving forward problems of complex physical processes. Current efforts are devoted to understand the role of model reduction, sparse recovery, efficient sampling, and other high-dimensional methods in this context.

Finally, purely data-driven problems arising in deep learning pose, to some extent, complementary challenges to data assimilation. Here, one is confronted with extremely large data-sets, and handling all the information appropriately becomes a challenge in itself. In addition, numerous question arise concerning how to leverage the approximation power of neural networks with practical numerical algorithms. In particular, the role of deep and highly over-parametrized architectures is an important direction of research to understand generalization properties. Many new and diverse ideas are being explored to shed light on these questions.

**Goals of the workshop:** The field of high-dimensional approximation is rapidly evolving, with a diversity of vibrant developments, and ideas coming from various fields that have traditionally been considered as disjoint areas of knowledge. Building a cohesive, overarching theory requires new interactions between approximation theory, numerical analysis, probability theory, mathematical and statistical learning theory, and optimization. The ambition of this workshop was to enhance interactions between these fields in order to deepen the mathematical foundations of the underlying numerical concepts that drive the new evolution of computational methods for high-dimensional problems. We proposed to gather leading experts with an interdisciplinary track record in combining these mathematical topics and having made significant recent contribution to one of the above three main problem classes (physics-driven, data-driven, hybrid problems). Of course, we also invited promising researchers at early career stages.

The event helped to promote the exchange of ideas arising from various disciplines about how to treat high-dimensional problems. In particular, given certain conceptual similarities that occur in a variety of application domains, we encountered a wealth of synergies and cross-fertilization. These concepts are in our opinion not only relevant for the development of efficient solution methods for large-scale and inherently high-dimensional problems but also for the formulation of rigorous mathematical models for quantifying the extraction of essential information from complex objects in many dimensions.

Specific examples of topics that were addressed in our workshop were:

- Sparse methods for parametric PDEs and PDEs with stochastic data
- Optimal transport in high dimensions
- Multilevel and high-dimensional meshless methods
- Incorporating anisotropy in analysis, estimation, compression and encoding
- Linear and nonlinear reduced modeling for forward and inverse problems
- Data assimilation and Bayesian inverse problems
- Convergence theory and analysis for model reduction and low-rank methods
- Theoretical and numerical aspects of sparse approximation and compressed sensing
- Design and analysis of estimators in high dimensional machine learning
- High-dimensional approximation using low-rank tensor structures
- Understanding the gap between analysis and practical performance of approximations by deep neural networks
- Performance and stability of optimization techniques in machine learning

For instance, a previous workshop 'Computation and Learning in High Dimensions' held in August 2021 (organized by two of us with R. DeVore and W. Dahmen) brought together some of the leading researchers in deep learning to interact with the approximation, numerical analysis, and computational harmonic analysis communities, see the Oberwolfach Report 36/2021.

Among the several recognizable outcomes of this and previous of our workshops were:

(i) a wide range of new results quantifying the performance of approximation when using deep neural networks,

(ii) fast online computational algorithms based on adaptive partition for mathematical learning,

(iii) injection of the notion of sparsity into stochastic models to identify computational paradigms that are more efficient than Monte Carlo techniques,

(iv) a coherent theory to explain why techniques like sparse representation and reduced modeling work and how they can be improved.

**Workshop: Computation and Learning in High Dimensions**

**Table of Contents**

# Abstracts

## A short, non exhaustive, and biased historical account on a cycle of MFO Workshops

### Albert Cohen

This talk was an attempt to draw a big picture on the main topics that influenced the cycle of nine workshops started in 2004,
("Multiscale Methods and Wavelets",

"Computation and Learning in High Dimensions").
A predecessor was the 1992 workshop on "Wavelets" that preceeded the growth in importance of non-linear (best $n$-term) approximation in the 1990's, with applications to data compression, statistical estimation, adaptive numerical simulation.

This led to the concept of sparsity which dominated these workshops in the years 2005–2015 with the emergence of compressed sensing.

Breaking the curse of dimensionality in learning and computation has motivated new forms of sparse approximation in the years 2010–2025, as well as the developments of non-linear reduced models. DNN and related tools can be embbeded in this general framework and seem to play an increasing role in numerical modelling since 2020.

## Funny things may happen when using NNs to solve PDEs

### Mark A. Peletier

(joint work with Daan Bon, Benjamin Caris, and Olga Mula)

Define the shallow neural network $\mathcal{U}_n(\theta)(x) = \sum_{i=1}^n a_i \varphi(x - b_i)$ for some smooth $\varphi \in C^\infty(\mathbb{R})$, with $\theta = (a_i, b_i)_{i=1}^n \in \mathbb{R}^{2n}$. In this talk we studied properties of the set $\mathcal{M}_n := \{\mathcal{U}_n(\theta) : \theta \in \mathbb{R}^{2n}\} \subset C^\infty(\mathbb{R})$ generated by such NNs. We showed that $\mathcal{M}_n$ is a "singular manifold" in $L^2(\mathbb{R})$.

We discussed downsides of existing approaches to using the set $\mathcal{M}_n$ as candidate solutions for e.g. the Allen-Cahn equation; we showed how these existing approaches lead to singularities in the evolution (one of the "funny things").

Using the fact that the Allen-Cahn equation is the $L^2$-gradient flow of some functional $\mathcal{E}$, we postulate that the correct way of using $\mathcal{M}_n$ is to consider the metric-space gradient flow of $\mathcal{E}$ in the metric space $(\mathcal{M}_n, d)$, where the metric $d$ is the ambient $L^2$-metric, $d(u, v) = \|u - v\|_{L^2}$. We showed that (under conditions) a minimizing-movement scheme for this setup converges as $\Delta t \to 0$ to a curve of maximal slope of $\mathcal{E}$ on the completion $\overline{(\mathcal{M}_n, d)}$. We also showed with an example that parametrizations may be discontinuous in time for such curves of maximal slope, and that in some cases this can not be avoided.

**Optimization-Free Diffusion Model - A Perturbation Theory Approach**

Mathias Oster

(joint work with Yuehaw Khoo, Yifan Peng)

Diffusion models have emerged as a powerful framework in generative modeling, typically relying on optimizing neural networks to estimate the score function via forward SDE simulations. In this work, we propose an alternative method that is both optimization-free and forward SDE-free. By expanding the score function in a sparse set of eigenbasis of the backward Kolmogorov operator associated with the diffusion process, we reformulate score estimation as the solution to a linear system, avoiding iterative optimization and time-dependent sample generation. We analyze the approximation error using perturbation theory and demonstrate the effectiveness of our method on high-dimensional Boltzmann distributions and real-world datasets.

References

[1] Y. Khoo and M. Oster and Y. Peng, *Optimization-Free Diffusion Model – A Perturbation Theory Approach*, arXiv preprint, url=https://arxiv.org/abs/2505.23652 (2025).

**On the expressivity of deep Heaviside networks**

Johannes Schmidt-Hieber

(joint work with Insung Kong, Juntong Chen, Sophie Langer)

The first models of an artificial neural network used the Heaviside activation function

$$x \mapsto \mathbf{1}(x \geq 0).$$

The highlight of this early literature is the Perceptron Convergence Theorem proving that one Heaviside neuron equipped with a simple update rule will perfectly classify two linearly separable classes after a finite number of update steps.

As the Heaviside activation function is non-differentiable, it had to be replaced by smoothed versions once the backpropagation algorithm was introduced. Modern deep network architectures mainly rely on the ReLU and increasingly on the SiLU activation function.

There is renewed interest in Heaviside networks. The straight-through estimator [1] provides a scalable method to train Heaviside networks. Moreover, Hopfield networks are based on the Heaviside networks and have seen renewed interest due to their connection with memorization in transformers [3]. Finally, the success of DeepSeek requires to understand the quantization of the activation function in a neural network.

While there has been an extensive body of literature studying approximation and generalization for ReLU networks, results for Heaviside networks are missing.

Deep Heaviside networks can be shown to have poor approximation properties. The approximation rates can be characterized by the width of the first hidden layer and do not make full use of the depth in the network. As each unit only

outputs one bit, the hidden layers are information bottlenecks constraining the availability of input information in the deeper layers.

However, one can equip deep Heaviside networks with either skip connections or linearly activated neurons and thereby overcome the information bottleneck. In both cases, one can then derive nearly matching upper and lower bounds for the Vapnik-Chervonenkis (VC) dimension and the worst case approximation error for Hölder balls.

Based on this, one can derive moreover generalization guarantees for the empirical risk minimizer computed over deep Heaviside networks.

All details can be found in [2].

### References

[1] Y. Bengio, N. Léonard, A. Courville *Estimating or propagating gradients through stochastic neurons for conditional computation* ArXiv:1308.3432
[2] I. Kong, J. Chen, S. Langer, J. Schmidt-Hieber *On the expressivity of deep Heaviside networks*, ArXiv:2505.00110
[3] H. Ramsauer et al. *Hopfield networks is all you need*, International Conference on Learning Representations (2021).

## A Space-Time Adaptive Low-Rank Method for High-Dimensional Parabolic PDEs

Manfred Faldum

(joint work with Markus Bachmayr)

In this work we investigate second-order parabolic differential equations on high-dimensional product domains $\Omega = \Omega_1 \times \cdots \times \Omega_d$ on a time interval $[0, T]$. As a model problem we consider the heat equation

$$
\begin{aligned}
\partial_t u - \Delta u = g & \quad \text{in } (0, T) \times \Omega, \\
u(0, \cdot) = h & \quad \text{in } \Omega, \\
u = 0 & \quad \text{on } (0, T) \times \partial\Omega.
\end{aligned}
$$

on the high-dimensional unit cube $\Omega = (0, 1)^d$. The goal of this work is to construct a method that has guaranteed convergence with rigorous deterministic space-time error bounds that are computable even when $d \gg 1$ and has estimates on the computational complexity that shows that we circumvent the curse of dimensionality, that is, exponential growth of the computational complexity with respect to the number of spatial dimensions.

We use the space-time variational formulation from [1] with trial space $\mathcal{X} = L_2(0, T; H_0^1(\Omega)) \cap H^1(0, T; (H_0^1(\Omega))')$ and test space $\mathcal{Y} = L_2(0, T; H_0^1(\Omega)) \times L_2(\Omega)$. The bilinear form and right hand side are given by

(1)
$$b(u, (v, w)) = \int_0^T \langle \partial_t u, v \rangle_{H^{-1}, H^1} + \langle \nabla u, \nabla v \rangle_{L_2} \, \mathrm{d}t + \langle u(0), w \rangle_{L_2},$$
$$f((v, w)) = \int_0^T \langle g, v \rangle_{L_2} \, \mathrm{d}t + \langle h, w \rangle_{L_2}.$$

Our approach for approximating the solution of the corresponding variational formulation combines a sparse wavelet expansion in time with a low-rank approximation in the spatial variables. For a temporal basis $\{\theta_{\hat{\nu}}\}_{\hat{\nu} \in \mathcal{I}}$ and the one dimensional spatial basis $\{\psi_{\hat{\nu}}\}_{\hat{\nu} \in \mathcal{J}}$ on $(0, 1)$, we aim for an approximation

$$u(t, x_1, \ldots, x_d) \approx \sum_{\mu \in \Lambda_t \subset \mathcal{I}} \theta_\mu(t) \sum_{(\nu_1, \ldots, \nu_d) \in \Lambda_\mu} \mathbf{u}_{\mu, \nu_1, \ldots, \nu_d} \, d_{\mu, \nu}^{\mathcal{X}} \psi_{\nu_1}(x_1) \cdots \psi_{\nu_d}(x_d)$$

with a finite index set $\Lambda_t$. Additionally, each time index $\mu$ can have a distinct spatial index $\Lambda_\mu$ with product structure $\Lambda_\mu = \Lambda_\mu^1 \times \cdots \times \Lambda_\mu^d \subset \mathcal{J} \times \cdots \times \mathcal{J}$. This product structure allows for the use of hierarchical low-rank tensor representations for each $\mathbf{u}_\mu = (\mathbf{u}_{\mu, \nu_1, \ldots, \nu_d})_{\nu \in \Lambda_\mu}$. The index sets $\Lambda_t$ and $\Lambda_\mu$ for $\mu \in \mathcal{I}$ can be chosen adaptively.

As one dimensional spatial basis functions as well as temporal basis functions, we choose wavelet Riesz bases. Then, $\{\Psi_{\mu, \nu}\}_{(\mu, \nu) \in \mathcal{I} \times \mathcal{J}^d}$ with

(2)
$$\Psi_{\mu, \nu} = \frac{\theta_\mu \otimes \psi_{\nu_1} \otimes \cdots \otimes \psi_{\nu_d}}{\sqrt{\|\psi_{\nu_1} \otimes \cdots \otimes \psi_{\nu_d}\|_{H^1}^2 + \|\theta_\mu\|_{H^1}^2 \|\psi_{\nu_1} \otimes \cdots \otimes \psi_{\nu_d}\|_{H^{-1}}^2}}$$

is a Riesz basis of the trial space X. A Riesz basis of $\mathcal{Y}$ can be constructed analogously. Based on the Riesz bases property, we have

$$\Big\| \sum_{(\mu, \nu) \in \mathcal{K}} \mathbf{v}_{\mu, \nu} \Psi_{\mu, \nu} \Big\|_{\mathcal{X}} \simeq \|\mathbf{v}\|_{\ell_2},$$

which guarantees an error of the same magnitude if we approximate $\mathbf{u} \in \ell_2$ in the basis expansion instead of $u$ in the space-time space $\mathcal{X}$.

The corresponding sequence $\mathbf{u}$ is given by the solution of the bi-infinite matrix vector equation $\mathbf{B}\mathbf{u} = \mathbf{f}$, where the matrix $\mathbf{B}$ as well as the right hand side $\mathbf{f}$ are given by using the Riesz bases of $\mathcal{X}$ and $\mathcal{Y}$ in (1).

To solve this bi-infinite matrix vector equation, we use the approximate Richardson iteration [2] applied to the normal equation $\mathbf{B}^\top \mathbf{B} \mathbf{u} = \mathbf{B}^\top \mathbf{f}$. In this procedure two reduction operators are applied to the iterates. The tensor recompression keeps the arising ranks of the low-rank approximations quasi- optimal for a given error tolerance. This routine is based on truncation of the hierarchical singular value decomposition of the low-rank representations. The second routine, basis coarsening, keeps the lower-dimensional support quasi- optimal for a given error tolerance. It is based on the spatio-temporal contractions [3], which measure the

influence of the lower-dimensional basis functions separately for each time basis index.

For the analysis of the method we propose a new approximation class for the temporal operator which is necessary due to the interaction between hierarchical tensor formats of different time indices. One of the main challenges is the fact that the parabolic operator is an isomorphism with respect to spaces not endowed with a cross norm. Hence, the scaling factor in (2) is not separable. Therefore, as in [2], we use a method for preconditioning operators in low-rank format by exponential sum approximations.

The method is shown to converge and satisfy similar complexity bounds as the existing adaptive low-rank method for elliptic problems [2, 3], establishing its suitability for parabolic problems on high-dimensional spatial domains. The construction also yields computable rigorous a posteriori error bounds for the total error depending on the activated basis functions and ranks in the approximation.

For more details as well as numerical results for the heat equation in high dimensions, demonstrating the practical efficiency, we refer to [3].

## References

[1] C. Schwab. and R. Stevenson, *Space-Time Adaptive Wavelet Methods for Parabolic Evolution Problems*, Mathematics of Computation **78** (2009), 1293–1318.

[2] M. Bachmayr and W. Dahmen, *Adaptive low-rank methods: problems on Sobolev spaces*, SIAM Journal on Numerical Analysis **54** (2016), 744–796.

[3] M. Bachmayr and M. Faldum, *A space-time adaptive low-rank method for high-dimensional parabolic partial differential equations*, Journal of Complexity **82** (2024).

## Scalable Bayesian Optimization via Online Gaussian Processes

Marcel Neugebauer

**Abstract.** Bayesian optimization is a state-of-the-art method for optimizing black box functions. It typically assumes that the unknown function is a sample path of a Gaussian process, which serves as surrogate model. As observations are collected, the belief is getting updated, enabling both prediction and uncertainty quantification. An acquisition function then guides the selection of new evaluation points by leveraging the posterior belief to balance exploration of uncertain regions and exploitation of promising areas. However, standard algorithms recompute the entire Gaussian process with each new observation or hyperparameter update, limiting scalability for large datasets.

To overcome this drawback, we employ a low-rank approximation of the Gaussian process kernel matrix that enables both the incorporation of new observations and online hyperparameter learning. This leads to online Gaussian processes and scalable optimization.

## 1. Introduction

Bayesian optimization is a procedure that aims to discover

$$\boldsymbol{x}^* \in \arg\max_{\boldsymbol{x} \in [0,1]^d} f(\boldsymbol{x}),$$

where $f : [0,1]^d \to \mathbb{R}$ is a black box function. This means that the underlying mapping of $f$ is unknown, yet we can observe values

$$y(\boldsymbol{x}) = f(\boldsymbol{x}) + \varepsilon$$

typically at high cost, where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ with $\sigma^2 \geq 0$. Notably, the lack of access to gradients of $f$ precludes the use of gradient-based methods. Black box functions are pervasive across various fields such as hyperparameter tuning, robotics, asset allocation, advertising and drug discovery. Bayesian optimization seeks to optimize such functions by executing the following iterative process:

1. Start with observations $(\boldsymbol{x}_1, y_1), ..., (\boldsymbol{x}_N, y_N)$.
2. Utilize the observations to construct a surrogate model for $f$.
3. Employ an acquisition function to determine $\boldsymbol{x}_{N+1}$ and evaluate $y_{N+1}$.
4. Add $(\boldsymbol{x}_{N+1}, y_{N+1})$ to the observations and repeat the process.

## 2. Gaussian Processes as Surrogate Models

The Gaussian process (GP) stands out as a preferred surrogate model for $f$ due to its widespread adoption. Modelling $f$ as a GP, denoted $f \sim \mathcal{GP}(0, k)$, we assume that for any finite number of points $\boldsymbol{z}_1, ..., \boldsymbol{z}_n \in [0,1]^d$, we have

$$\begin{bmatrix} f(\boldsymbol{z}_1) \\ \vdots \\ f(\boldsymbol{z}_n) \end{bmatrix} \sim \mathcal{N} \left( \boldsymbol{0}, \begin{bmatrix} k(\boldsymbol{z}_1, \boldsymbol{z}_1) & \cdots & k(\boldsymbol{z}_1, \boldsymbol{z}_n) \\ \vdots & \ddots & \vdots \\ k(\boldsymbol{z}_n, \boldsymbol{z}_1) & \cdots & k(\boldsymbol{z}_n, \boldsymbol{z}_n) \end{bmatrix} \right),$$

where $k : [0,1]^d \times [0,1]^d \to \mathbb{R}$ is a symmetric and positive definite kernel. A commonly used example is the Matérn 5/2 kernel, defined as

$$k_{s^2, \ell}(\boldsymbol{x}, \boldsymbol{x}') := s^2 \left( 1 + \frac{\sqrt{5}\, \|\boldsymbol{x} - \boldsymbol{x}'\|_2}{\ell} + \frac{5\|\boldsymbol{x} - \boldsymbol{x}'\|_2^2}{3\ell^2} \right) \exp\left( -\frac{\sqrt{5}\, \|\boldsymbol{x} - \boldsymbol{x}'\|_2}{\ell} \right),$$

where $s^2 > 0$ and $\ell > 0$ are hyperparameters. These are typically estimated from data by maximizing the log marginal likelihood. Given $f \sim \mathcal{GP}(0, k)$ and data $\mathcal{D}_N := \{(\boldsymbol{x}_1, y_1), ..., (\boldsymbol{x}_N, y_N)\}$, the posterior process is $f \,|\, \mathcal{D}_N \sim \mathcal{GP}(m', k')$ for

(1)
$$m'(\boldsymbol{x}) := \boldsymbol{k}(\boldsymbol{x})^T \left( \boldsymbol{K} + \sigma^2 \boldsymbol{I}_N \right)^{-1} \boldsymbol{y},$$
$$k'(\boldsymbol{x}, \boldsymbol{x}') := k(\boldsymbol{x}, \boldsymbol{x}') - \boldsymbol{k}(\boldsymbol{x})^T \left( \boldsymbol{K} + \sigma^2 \boldsymbol{I}_N \right)^{-1} \boldsymbol{k}(\boldsymbol{x}'),$$

where $\boldsymbol{K} := [k(\boldsymbol{x}_i, \boldsymbol{x}_j)]_{1 \leq i,j \leq N}$, $\boldsymbol{k}(\boldsymbol{x}) := [k(\boldsymbol{x}, \boldsymbol{x}_1), ..., k(\boldsymbol{x}, \boldsymbol{x}_N)]^T$, $\boldsymbol{y} := [y_1, ..., y_N]^T$. This posterior allows for inference about the black box function $f$ and motivates using $\mathcal{GP}(m', k')$ as the surrogate model when data $\mathcal{D}_N$ is available.

To compute the posterior process, the Cholesky decomposition of $\boldsymbol{K} + \sigma^2 \boldsymbol{I}_N$ is typically employed. This algorithm has a computational complexity of $\mathcal{O}(N^3)$

for a single posterior GP, which becomes expensive when $N$ is large. Another context in which an issue arises is Bayesian optimization, where many GPs must be computed repeatedly and computing each from scratch is inefficient. If Bayesian optimization runs for $T$ iterations and performs $R$ hyperparameter updates per GP, then standard algorithms incur a complexity of

$$\mathcal{O}\left(\sum_{N=1}^{T} RN^3\right) = \mathcal{O}\left(R\frac{T^2(T+1)^2}{4}\right) = \mathcal{O}(RT^4).$$

To overcome this drawback, we use an approximation that supports both the incorporation of new observations and online hyperparameter learning. This enables the computation of GPs in an online fashion.

Let $\boldsymbol{\theta}$ denote the set of kernel hyperparameters. For $M \ll N$, we consider a kernel approximation

$$(2) \qquad k_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{x}') \approx \sum_{i,j=1}^{M} \phi_i(\boldsymbol{x})\lambda_{i,j}(\boldsymbol{\theta})\phi_j(\boldsymbol{x}'),$$

where $\Lambda := [\lambda_{i,j}(\boldsymbol{\theta})]_{1 \leq i,j \leq M}$ is a real, symmetric, and positive definite matrix and $\boldsymbol{\Phi} := [\phi_i(\boldsymbol{x}_j)]_{1 \leq i \leq M, 1 \leq j \leq N}$ is a real matrix of rank $M$. Note that $\phi_i$ depends only on the points, while $\lambda_{i,j}$ depends solely on the hyperparameters. This separation of variables enables online computation of GPs. Here are three examples on how such a kernel approximation can be realized:

- Hilbert space methods for low-rank GP regression, see, e.g., [6]
- Kernel interpolation for scalable structured GPs, see, e.g., [5]
- Quadrature Fourier features, see, e.g., [3]

In the following, I employ the techniques from [6]. Inserting the kernel approximation (2) into the GP formulas (1) leads to

$$(3) \qquad \begin{aligned} m'_{\boldsymbol{\theta}}(\boldsymbol{x}) &\approx \phi(\boldsymbol{x})^T(\boldsymbol{\Phi\Phi}^T + \sigma^2\boldsymbol{\Lambda})^{-1}\boldsymbol{\Phi y}, \\ k'_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{x}') &\approx \sigma^2\phi(\boldsymbol{x})^T(\boldsymbol{\Phi\Phi}^T + \sigma^2\boldsymbol{\Lambda})^{-1}\phi(\boldsymbol{x}'), \end{aligned}$$

where $\phi(\boldsymbol{x}) := [\phi_1(\boldsymbol{x}), ..., \phi_M(\boldsymbol{x})]^T$. If the kernel approximation (2) converges uniformly, then both the posterior mean and kernel approximations (3) also converge uniformly.

These approximations allow computing a GP from scratch in $\mathcal{O}(NM^2)$ time by using the Cholesky decomposition $\boldsymbol{L}\boldsymbol{L}^T = \boldsymbol{\Phi\Phi}^T + \sigma^2\boldsymbol{\Lambda}^{-1}$. When a new observation is incorporated, the Cholesky decomposition can be modified via a rank-1 update in $\mathcal{O}(M^2)$ time, for example using Givens rotations. This allows adding a new observations in $\mathcal{O}(M^2)$ time.

Let $\theta_1 \in \boldsymbol{\theta}$ be a kernel hyperparameter. Using the kernel approximation (2), we approximate the log marginal likelihood as

$$\log p(\boldsymbol{y}\,|X) \approx -\frac{1}{2\sigma^2}(\boldsymbol{y}^T\boldsymbol{y} - \boldsymbol{y}^T\boldsymbol{\Phi}^T(\boldsymbol{L}\boldsymbol{L}^T)^{-1}\boldsymbol{\Phi}\boldsymbol{y})$$

$$-\frac{1}{2}\left(\log|\boldsymbol{\Lambda}| + 2\sum_{i=1}^{M}\log L_{ii} + (N-M)\log\sigma^2\right) - \frac{N}{2}\log 2\pi$$

where $X := \{\boldsymbol{x}_1, ..., \boldsymbol{x}_N\}$. We approximate the derivative with respect to $\theta_1$ via

$$\frac{\partial \log p(\boldsymbol{y}\,|X)}{\partial \theta_1} \approx -\frac{1}{2}\left(\operatorname{tr}\left(\boldsymbol{\Lambda}^{-1}\frac{\partial\boldsymbol{\Lambda}}{\partial\theta_1}\right) - \sigma^2\operatorname{tr}\left(\boldsymbol{\Lambda}^{-1}(\boldsymbol{L}\boldsymbol{L}^T)^{-1}\boldsymbol{\Lambda}^{-1}\frac{\partial\boldsymbol{\Lambda}}{\partial\theta_1}\right)\right.$$

$$\left. - \boldsymbol{y}^T\boldsymbol{\Phi}^T(\boldsymbol{L}\boldsymbol{L}^T)^{-1}\boldsymbol{\Lambda}^{-1}\frac{\partial\boldsymbol{\Lambda}}{\partial\theta_1}\boldsymbol{\Lambda}^{-1}(\boldsymbol{L}\boldsymbol{L}^T)^{-1}\boldsymbol{\Phi}\boldsymbol{y}\right)$$

and a similar approximation can be derived for the derivative with respect to $\sigma^2$. Since $\boldsymbol{\Phi}$ does not depend on $\boldsymbol{\theta}$, it does not need to be recomputed and a single hyperparameter update can be performed in $\mathcal{O}(M^3)$ time.

Using this approximate GP in Bayesian optimization provides a computational advantage. For $T$ iterations and $R$ hyperparameter updates per GP, the computational cost decreases from $\mathcal{O}(RT^4)$ to $\mathcal{O}(RM^3T)$.

## 3. Thompson Sampling for Acquisition

Thompson sampling, described in Section 7.9 of reference [2], selects the next point by drawing a sample path $a : [0,1]^d \to \mathbb{R}$ of the posterior $\mathcal{GP}(m', k')$ and choosing

$$\boldsymbol{x}_{\text{next}} \in \underset{\boldsymbol{x}\in[0,1]^d}{\arg\max}\, a(\boldsymbol{x}).$$

To compute a sample path, we use the approximate GP described by (3). Let

$$\boldsymbol{\omega} \sim \mathcal{N}\left((\boldsymbol{\Phi}\boldsymbol{\Phi}^T + \sigma^2\boldsymbol{\Lambda}^{-1})^{-1}\boldsymbol{\Phi}\boldsymbol{y}, \sigma^2(\boldsymbol{\Phi}\boldsymbol{\Phi}^T) + \sigma^2\boldsymbol{\Lambda}^{-1})^{-1}\right),$$

then $g(\boldsymbol{x}) := \boldsymbol{\phi}(\boldsymbol{x})^T\boldsymbol{\omega}$ is a sample path of the approximate GP. If $\boldsymbol{\phi}(\boldsymbol{x})$ is continuously differentiable, the Lipschitz constant $L$ of $g$ can be bounded by

$$L \leq \|\boldsymbol{\omega}\|_2 \sup_{\boldsymbol{x}\in[0,1]^d}\|D\boldsymbol{\phi}(\boldsymbol{x})\|_2.$$

Thus, $g$ can be optimized using its gradient and Lipschitz bound.

I presented a numerical example using the $2D$-Ackley function with the kernel $k_{s^2,\ell}$ and initial hyperparameters $\sigma^2 = 1$, $\ell = 1$ and $s^2 = 1$. The results support the theoretical complexity estimates. Higher-dimensional variants based on product kernels are examined in the present setting in [4].

## References

[1] M. Balakirsky, K. Li and S. Mak, *Trigonometric Quadrature Fourier Features for Scalable Gaussian Process Regression*, Proceedings of Machine Learning Research **238** (2024), 3484–3492.

[2] R. Garnett, *Bayesian Optimization*, Cambridge University Press (2023).

[3] A. Krause and M. Mutny, *Efficient High Dimensional Bayesian Optimization with Additivity and Quadrature Fourier Features*, Advances in Neural Information Processing Systems **31** (2018).

[4] M. Neugebauer, *Scalable Bayesian Optimization via Online Gaussian Processes (working title)*, doctoral dissertation, in preparation.

[5] H. Nickisch and A. Wilson, *Kernel Interpolation for Scalable Structured Gaussian Processes (KISS-GP)*, Proceedings of Machine Learning Research **37** (2015), 1775–1784.

[6] A. Solin and S. Sarkka, *Hilbert space methods for reduced-rank Gaussian process regression*, Statistics and Computing **30** (2020), 419–446.

## Iterative thresholding low-rank time integration

Matthieu Dolbeault

(joint work with Markus Bachmayr, Polina Sachsenmaier)

High-dimensional evolution equations pose multiple challenges in terms of numerical approximation. First, one needs a suitable representation of the solution with respect to space variables, avoiding the curse of dimensionality; low-rank tensor methods allow to reduce the complexity of the numerical approximation from exponential to linear in the dimension. Second, this representation needs to evolve over time, in order to adapt to the dynamics of the solution; rank-adaptive strategies benefit from their great flexibility, at the expense of a more delicate balance between accuracy and computational complexity.

We consider a numerical scheme based on low-rank matrix decompositions, viewed as a preliminary step towards tensor formats. The time approximation relies on high order collocation and allows us to use large time steps, making our approach an intermediate between dynamical low-rank approximation [1, 2] and space-time methods [3, 4]; it is closely related to high order BUG integrators from [5, 6]. Our analysis focuses on the prototypical Schrödinger equation

$$i\partial_t u + \Delta u = V u, \qquad t \in [0, T],$$

with initial data $u_0 : \mathbb{T}^d \to \mathbb{C}$, under the assumption that $u_0$ and $V$ are sufficiently smooth and of finite rank. Given a singular value decomposition of the solution

$$u(t, x) = \sum_{i \geq 1} \sigma_i(t) u_i^{(1)}(t, x_1) u_i^{(2)}(t, x_2), \qquad x = (x_1, x_2) \in \mathbb{T}^d = \mathbb{T}^{d_1} \times \mathbb{T}^{d_2},$$

the best rank-$r$ approximation of $u$ in $L^2$ is obtained by keeping only the first $r$ terms, and attains an $L^2$ error

$$\left\| u(t) - \sum_{i=1}^r \sigma_i(t) u_i^{(1)}(t, x_1) u_i^{(2)}(t, x_2) \right\|_{L^2(\mathbb{T}^d)}^2 = \sum_{i>r} |\sigma_i(t)|^2.$$

Ideally, one would like to construct an approximation of $u$ of fixed accuracy $\eta > 0$ with optimal ranks

$$r(t) = \min \left\{ r \geq 0, \quad \sum_{i>r} |\sigma_i|^2 \leq \eta^2 \right\}.$$

However, this objective is out of reach for time-stepping methods, since the error from the previous time steps makes it impossible to distinguish between singular values that are just above or just below the threshold $\alpha(t) = \sigma_{r(t)}(t)$. Instead, we consider a relaxed notion of rank: for $\alpha = \inf_{t \in [0,T]} \alpha(t)$, define

$$\tilde{r} = \sup_{t \in [0,T]} r(t) + \sum_{i > r} \frac{\sigma_i(t)^2}{\alpha^2} = \sup_{t \in [0,T]} \frac{\|u - S_\alpha u\|_{L^2(\mathbb{T}^d)}^2}{\alpha^2},$$

where $S_\alpha$ is the soft thresholding operator

$$S_\alpha u = \sum_{i \geq 1} \max(\sigma_i - \alpha, 0)\, u_i^{(1)} \otimes u_i^{(2)}.$$

Note that if the soft thresholding was replaced by hard thresholding (defined similarly to $S_\alpha$, but with $\max(\sigma_i - \alpha, 0)$ replaced by $\sigma_i \mathbb{1}_{\sigma_i > \alpha}$), we would simply recover $\sup_{t \in [0,T]} r(t)$.

With this, we can state our main theorem. Let $J \in \mathbb{N}$ be the desired order of the scheme, and assume that $u \in C([0,T], H^{2J}(\mathbb{T}^d))$ and $V \in H^{4J}(\mathbb{T}^d)$. For a time step $h > 0$, denote by $\tilde{u}_n$ the scheme at time $nh$, with $n \leq N := \lfloor T/h \rfloor$.

**Theorem** ([7])**.** *The proposed scheme achieves the global error bound*

$$\sup_{1 \leq n \leq N} \|\tilde{u}_n - u(nh)\| \lesssim \exp(cT) \left( \eta + h^{J+1} + h^{2J}\, T \right).$$

*The ranks of the scheme are bounded by*

$$\mathrm{rank}(\tilde{u}_n) \lesssim \frac{\tilde{r}}{h^3},$$

*and the intermediate ranks are at most twice as large.*

The algorithm is inspired by the iterative soft thresholding strategy from [8], combined with a Picard iteration and a refined analysis of the rank reduction from a final recompression at the end of the iteration.

The above result is the first to control the global error by an exponential in time, while maintaining near-optimal ranks. A naive approach would instead yield an exponential in the number of time steps, which is significantly worse for small time steps $h$. In our case, we take $h$ such that $h^{J+1} + h^{2J}$ is of order $\eta$, and the only issue with the limit $h \to 0$ is the factor $1/h^3$ in the rank bound. This factor could be improved to $1/h^2$ by increasing the number of fixed point iterations, but it is not clear if it can be improved further. Extensions to other equations, for instance parabolic problems, are also the subject of ongoing work.

REFERENCES

[1] O. Koch and C. Lubich, *Dynamical low-rank approximation*, SIAM Journal on Matrix Analysis and Applications **29** (2007), 434–454.

[2] G. Ceruti, J. Kusch and C. Lubich, *A rank-adaptive robust integrator for dynamical low-rank approximation*, BIT Numerical Mathematics **62** (2022), 1149–1174.

[3] M. Bachmayr and M. Faldum, *A space-time adaptive low-rank method for high-dimensional parabolic partial differential equations*, arXiv:2302.01658 (2023).

[4] M.-S. Dupuy, V. Ehrlacher and C. Guillot, *A space-time variational formulation for the many-body electronic Schrödinger evolution equation*, arXiv:2405.18094 (2024).

[5] S. Li, Y. Jiang and Y. Cheng, *High-order implicit low-rank method with spectral deferred correction for matrix differential equations*, arXiv:2412.09400 (2024).

[6] F. Nobile and S. Riffaud, *Robust high-order low-rank BUG integrators based on explicit Runge-Kutta methods*, arXiv:2502.07040 (2025).

[7] M. Bachmayr, M. Dolbeault and P. Sachsenmaier, *Iterative thresholding low-rank time integration*, arXiv:2507.15848 (2025).

[8] M. Bachmayr and R. Schneider, *Iterative methods based on soft thresholding of hierarchical tensors*, Foundations of Computational Mathematics **17** (2017), 1037–1083.

## Direct interpolative construction of quantized tensor trains

MICHAEL LINDSEY

(joint work with Maxime Müller)

Quantized tensor trains (QTTs) have recently emerged as a framework for the numerical discretization of continuous functions, with the potential for widespread applications in numerical analysis, including rank-structured solvers and preconditioners based on "quantum-inspired" algorithms such as DMRG.

We advance the theory and practice of QTT approximation from the point of view of multiscale polynomial interpolation. This perspective clarifies why QTT ranks decay with increasing depth, quantitatively controls QTT rank in terms of smoothness of the target function, accommodates the presence of sharp features through a generalized construction, and motivates new fast algorithms for the construction of QTTs with performance guarantees.

Finally, we leverage the perspective of multiscale interpolation to offer the first direct construction of the fast Fourier transform (FFT) as a QTT operator, equipped with a priori compression guarantees.

## Parametric regularity plays a crucial role for lattice QMC methods – from DNNs to precision oncology

FRANCES Y. KUO

(joint work with Alexander D. Gilbert, Alexander Keller, Dirk Nuyens, Graham Pash, Ian H. Sloan, Karen E. Willcox)

Quasi-Monte Carlo (QMC) methods have been successful for high dimensional integration, multivariate function approximation, density estimation, etc., in many application areas including uncertainty quantification problems driven by PDEs with random coefficients. Lattice QMC points can be tailor-constructed to the applications, to achieve a convergence rate close to $\mathcal{O}(N^{-1})$ or better, with the implied constant independent of the dimension $s$ in appropriately weighted space settings. The success of QMC relies on the underlying smoothness and dimension structure of the target function: the "parametric regularity" — the mixed partial derivatives of the function with respect to the parametric variables. This is the underlying theme of this talk.

Part 1 of this talk is based on [1] where we obtain explicit *parametric regularity bounds* for a standard feed-forward Deep Neural Network (DNN), as well as a periodic variant with a built-in sine layer. These bounds depend on the network parameters and the choice of the smooth activation function. By imposing restrictions on the network parameters to match the regularity features of the target function, we prove that DNNs with $N$ *tailor-constructed lattice training points* can achieve the generalization error ($L_2$ approximation error) bound $\texttt{tol} + \mathcal{O}(N^{-r/2})$, where $\texttt{tol} \in (0,1)$ is the tolerance achieved by the training error in practice, and $r$ characterises the decaying importance of the input variables in the target function. In our numerical experiments, we restrict the network parameters during training by adding a *tailored regularization* term, and we show that for an algebraic equation mimicking the parametric PDE problems the DNNs trained with tailored regularization perform significantly better.

Part 2 of this talk is based on [2] where we consider a class of *parametric semi-linear parabolic PDEs* used to model tumor growth and treatment, capturing infiltration of the tumor into surrounding healthy tissue, proliferation of the existing tumor, and patient response to chemo- and radiotherapies. Considerable inter-patient variability, inherent heterogeneity of the disease, sparse and noisy data collection, and model inadequacy all contribute to significant uncertainty in the model parameters. It is crucial that these uncertainties can be efficiently propagated through the model to compute quantities of interest (QoIs), which in turn may be used to inform clinical decisions. We show that QMC methods can be successful in computing expectations of meaningful QoIs. Well-posedness results and *parametric regularity bounds* are developed for the model at hand and used to show a theoretical error bound for the case of uniform random fields. The theoretical linear error rate $\mathcal{O}(N^{-1})$ is verified numerically, demonstrating the superiority of the method over standard Monte Carlo. Encouraging computational results are also provided for lognormal random fields, prompting further theoretical development.

## References

[1] A. Keller, F. Y. Kuo, D. Nuyens, I. H. Sloan, *Regularity and tailored regularization of Deep Neural Networks, with application to parametric PDEs in uncertainty quantification*, arXiv: 2502.12496 (2025)

[2] A. D. Gilbert, F Y. Kuo, Dirk Nuyens, Graham Pash, Ian H. Sloan, Karen E. Willcox, *Quasi-Monte Carlo for uncertainty quantification of tumor growth and treatment, modeled by a parametric semilinear parabolic reaction-diffusion PDE*, in preparation (2025)

# Optimal Solvers for Infinite-Dimensional Sparse Approximations in Adaptive Stochastic Galerkin Finite Element Methods

Henrik Eisenmann

(joint work with Markus Bachmayr, Martin Eigel, Igor Voulis)

For a class of ellliptic diffusion problems on a domain $D$, we aim to approximate the mapping from diffusion coefficients $a = a(y)$ to the corresponding solutions $u(y)$ satisfying

$$-\nabla_x \cdot (a(y)\nabla_x u(y)) = f.$$

We treat here the case of infinitely many parameters, which is common when random fields are represented in series expansions. We show convergence with uniform rate of an adaptive stochastic Galerkin method, and given an expansion of the random field of the form

$$(1) \qquad a(y) = f\Big( \sum_{j \in \mathbb{N}} \sum_{k \in \mathbb{N}} y_{j,k} \theta_{j,k}(x) \Big)$$

with functions $\theta_{j,k}$ having multilevel structure, and an analytic function $f$, it is shown to produce quasi-optimal approximations with almost optimal computational costs.

The solution map is well approximated by a series

$$\tilde{u}(x,y) = \sum_{\nu \in F} u_\nu(x) L_\nu(y)$$

with function valued coefficients $u_\nu$ and product Legendre polynomials $L_\nu$. For this expansion to have a quasi-optimal number of degrees of freedom, it is essential to allow each function $u_\nu$ to be approximated in a different discrete space $V_\nu \subset H_0^1(D)$.

We utilize finite element frames when estimating the residual to circumvent problems arising from jump discontinuities of the residual on an overlay of many different meshes. This allows to show the saturation property, that is, error reduction by a uniform factor in each step. For optimality, we show a stability property of finite element subframes connected to conforming triangulations.

With this machinery at hand, the main result is an adaptive stochastic Galerkin method for the parametric problem, that produces approximate solutions $u^k$ to $u$ achieving the following:

- Uniform error reduction in each step:

$$\|u^{k+1} - u\|_{L_2(Y,H_0^1(D))} \leq \delta \|u^k - u\|_{L_2(Y,H_0^1(D))}.$$

- Quasi-optimal approximations to $u$: The total number of triangles $N_k$ required for $u_k$ scale optimally, that is if $u$ is in the approximation class $\mathcal{A}^s$ with

$$|u|_{\mathcal{A}^s} = \sup_{N>0} N^s \min_{|\mathbb{T}| \leq N} \min_{v \in \mathcal{V}(\mathbb{T})} \|u - v\|_{L_2(Y,H_0^1(D))}$$

then

$$N_k - N_0 \leq C \|u^k - u\|_{L_2(Y,H_0^1(D))}^{-\frac{1}{s}} |u|_{\mathcal{A}^s}.$$

- Optimal computational complexity in the affine case up to logarithmic factor: If $a(y)$ has an affine linear representation, then the costs of computing $u_k$ is bounded by a fixed multiple of

$$\left(1 + |\log \|u - u^k\|_{L_2(Y,H_0^1(D))}| + \log |u|_{\mathcal{A}^s}\right)^3 \|u - u^k\|_{L_2(Y,H_0^1(D))}^{-\frac{1}{s}} |u|_{\mathcal{A}^s}^{\frac{1}{s}}.$$

- Almost optimal computational complexity in the non-affine case: If $a$ can be expanded in the form of (1) with $f$ having a sufficiently large holomorphic extension, then there is $r \in \mathbb{N}$ such that for any $s' < s$ the costs of computing $u_k$ is bounded by a fixed multiple of

$$\left(1 + |\log \|u - u^k\|_{L_2(Y,H_0^1(D))}| + \log |u|_{\mathcal{A}^s}\right)^r \|u - u^k\|_{L_2(Y,H_0^1(D))}^{-\frac{1}{s'}} |u|_{\mathcal{A}^s}^{\frac{1}{s'}}.$$

REFERENCES

[1] M. Bachmayr, H. Eisenmann, I. Voulis, *Adaptive stochastic Galerkin finite element methods: Optimality and non-affine coefficients*, arXiv:2503.18704 (2025).
[2] M. Bachmayr, M. Eigel, H. Eisenmann, I. Voulis, *A convergent adaptive finite element stochastic Galerkin method based on multilevel expansions of random fields*, arXiv:2403.13770 (2025).

## The Polytope division method

EVIE NIELEN

In this talk, we introduced Configuration Optimization Problems (COPs). These problems involve minimizing a loss function over a set of discrete points $\eta \subset P$. Examples of these problems can be found in areas like Model Order Reduction, Active Learning, and Optimal Experimental Design. While exact solutions are often incomputable, heuristic solutions can be found via the weak Greedy Sampling Method (wGSM), particularly in low-dimensional cases. wGSM recursively updates $\eta$ by computing an error estimate over a discrete sample set $S \subset P$. However, as the dimensionality grows, the sample size suffers from the curse of dimensionality.

To address this, we discuss the Polytope Division Method (PDM), a scalable greedy-type approach that adaptively partitions the parameter space and targets regions of high loss. PDM achieves linear scaling with problem dimensionality and offers a first step towards a solution approach for high-dimensional COPs, however we also discuss two downsides of PDM: Finding a proof of convergence and explorability. To resolve these issues, we expand upon this method by introducing the randomized Polytope Division Method (r-PDM). Next, we make the connection between greedy-type methods and birthing processes. We describe the greedy algorithm as a stochastic process and introduce a transition kernel $\lambda$, and error function $G$. Under appropriate assumptions of this transition kernel and error function, we formulate a theorem describing the convergence in probability, including convergent rates. Lastly, we compare r-PDM with a stochastic process with uniform rates on a small interpolation example. Following the results of

the theorem, both methods convergence in expectation, but r-PDM has a lower variance.

## Nonlinear model-order reduction via optimal transport for electronic structure calculations
### Geneviève Dusson
(joint work with Maxime Dalery, Virginie Ehrlacher, Alexei Lozinski)

Parametric partial differential equations exhibiting transport-dominated behavior pose a significant challenge for classical linear model-order reduction (MOR) techniques, which often fail to capture essential solution features such as translation or deformation. In this talk, I presented a nonlinear MOR strategy based on optimal transport for the computation of the ground state of the electronic Schrödinger equation parametrized by nuclei positions. The reduced solutions are constructed as modified Wasserstein barycenters of selected high-fidelity solutions, which are, given a collection of probability measures $\{\mu_i\}_{i=1}^n$, defined as the minimizer of the weighted sum of squared distances:

$$\arg\min_{\mu} \sum_{i=1}^{n} \lambda_i \, d^2(\mu, \mu_i),$$

where for $i \in \{1, \ldots, n\}$, $\lambda_i \geq 0$, $\sum_{i=1}^{n} \lambda_i = 1$, and $d$ is a chosen modified Wasserstein distance. This notion therefore provides a natural way to interpolate between solutions that are probability distributions, allowing us to exploit the geometry of the solution manifold for varying positions of the nuclei more effectively.

In the first part of the talk we presented the contribution [1] where we analyze a toy eigenvalue problem mimicking the electronic Schrödinger equation in one dimension, involving a fixed nuclear configuration and a single electron, for which analytical solutions are available. More precisely we consider parameters $\mathbf{r} := (r_1, \ldots, r_M) \in \mathbb{R}^M$ and $\mathbf{z} := (z_1, \ldots, z_M) \in (\mathbb{R}_+^*)^M$ for $M \in \mathbb{N}^*$, and we are interested in the lowest eigenvalue $E_{\mathbf{r},\mathbf{z}} \in \mathbb{R}$ and the corresponding (strictly positive) eigenstate $u_{\mathbf{r},\mathbf{z}} \in H^1(\mathbb{R})$ satisfying

$$(1) \qquad -\frac{1}{2} u_{\mathbf{r},\mathbf{z}}'' + \left( - \sum_{m=1}^{M} z_m \delta_{r_m} \right) u_{\mathbf{r},\mathbf{z}} = E_{\mathbf{r},\mathbf{z}} u_{\mathbf{r},\mathbf{z}}.$$

We rescale the eigenvector to associate it to a probability distribution. We propose a greedy algorithm to select the best snapshots; the selection is based on the computation of the best barycenter for a mixture Wasserstein distance [3, 4], well suited to the solutions of the considered equation, which are mixtures of Slater distributions. The provided numerical results are promising, exhibiting a very fast decay of the error with respect to the number of selected snapshots, both during the offline (selection) phase and the online phase.

In the second part of the talk, we turned to a more challenging problem where we aim at approximating the pair density from the electronic density, as presented

in [2]. Given a ground-state wavefunction $\Psi_{\mathbf{R}}$ for a nuclear configuration $\mathbf{R}$, the electronic density is defined (up to a scalar factor) as

$$\rho_{\mathbf{R}}(x) = \int_{\mathbb{R}^{d(N-1)}} |\Psi_{\mathbf{R}}(x, x_2, \ldots, x_N)|^2$$

and the pair density is defined (also up to a scalar factor) as

$$\tau_{\mathbf{R}}(x, y) = \int_{\mathbb{R}^{d(N-2)}} |\Psi_{\mathbf{R}}(x, y, x_3, \ldots, x_N)|^2,$$

thus the electronic density is the marginal of the pair density. For tractability we still limit ourselves to 1D particles ($d = 1$), but we consider several electrons. When developing a similar greedy algorithm using a selection based on Wasserstein barycenters either on the density or the pair density, we encounter a marginal inconsistency: the Wasserstein barycenter between marginals is in general not equal to the marginal of the corresponding Wasserstein barycenter. Due to this, we proposed modified Wasserstein barycenters, exactly satisfying given marginal constraints. These modified barycenters have so far been developed for Gaussian distributions and Gaussian mixture distributions. They are defined as the solution to an optimization problem which can be analytically solved for Gaussian distributions, and requires a post-processing step for Gaussian mixtures. This enables efficient approximations leveraging the more easily accessible knowledge to the marginals of probability distributions, compared to the full distributions. This shows better approximation results with distributions that are rotating or translating at a varying speed. This is expected to be particularly useful in high-dimensional settings where the access to full solutions is limited or expensive.

REFERENCES

[1] M. Dalery, G. Dusson, V. Ehrlacher, and A. Lozinski, *Nonlinear reduced basis using mixture Wasserstein barycenters: application to an eigenvalue problem inspired from quantum chemistry,* arXiv preprint arXiv:2307.15423, (2023).
[2] M. Dalery, G. Dusson, and V. Ehrlacher, *Marginal-preserving modified Wasserstein barycenters for Gaussian distributions and Gaussian mixtures,* hal preprint hal-04696783v2, (2024).
[3] J. Delon, and A. Desolneux, *A Wasserstein-type distance in the space of Gaussian mixture models*, SIAM Journal on Imaging Sciences, 13(2), 936-970 (2020).
[4] G. Dusson, V. Ehrlacher, and N. Nouaime, *A Wasserstein-type metric for generic mixture models, including location-scatter and group invariant measures,* arXiv preprint arXiv:2301.07963, (2023).

## Inverse optimal transport and related problems

CLARICE POON

(joint work with Francisco Andrade, Gabriel Peyré)

*Overview.* Estimating parameters from samples of an optimal probability distribution is essential in applications ranging from socio-economic modeling to biological system analysis. In these settings, the probability distribution arises as the solution to an optimization problem that captures either static interactions among

agents or the dynamic evolution of a system over time. We introduce a general methodology based on a new class of loss functions, called sharpened Fenchel-Young losses, which measure the sub-optimality gap of the optimization problem over the space of probability measures. We provide explicit stability guarantees for two relevant settings in the context of optimal transport: the first is inverse optimal transport and the second is inverse gradient flows. This is based on the two papers [2] and [1].

*Inverse optimal transport (iOT).* The entropic optimal transport problem is as follows: Let $\epsilon > 0$ be a fixed entropic regularization parameter. Given two probability measures $\alpha \in \mathcal{P}(\mathcal{X})$ and $\beta \in \mathcal{P}(\mathcal{Y})$, a cost function $c \in \mathcal{C}(\mathcal{X} \times \mathcal{Y})$, find $\pi(c) \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})$ that solves

$$\pi(c) = \mathrm{argmin}_{\pi \in \mathcal{U}(\alpha, \beta)} \int c(x, y) d\pi(x, y) + \epsilon \mathrm{KL}(\pi | \alpha \otimes \beta).$$

where $\mathcal{U}(\alpha, \beta)$ denotes the set of probability measures on $\mathcal{X} \times \mathcal{Y}$ that have marginals $\alpha, \beta$. The *inverse problem* is to recover the cost function $c$ given $n$ samples $(x_i, y_i) \overset{iid}{\sim} \pi(c)$. Note that these samples also give access to the empirical marginals $\hat{\alpha}^n = \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}$ and $\hat{\beta}^n = \frac{1}{n} \sum_{i=1}^{n} \delta_{y_i}$. These kinds of problems were introduced by Galichon in a series of works, see for example [3].

*Inverse gradient flow (iJKO).* Suppose one observes samples iid samples of probability distributions $\rho_k$ for $k = 1, 2, \ldots$, where

$$\rho_{k+1} = \mathrm{argmin}_{\rho \in \mathcal{P}(\mathcal{X})} \mathcal{F}(\rho) + \frac{1}{2\tau} W_2^2(\rho, \rho_k)$$

where $W_2^2$ is the (entropy regularized) Wasserstein distance with Euclidean metric. The *inverse problem* is to recover the functional $\mathcal{F} : \mathcal{P}(\mathcal{X}) \to \mathbb{R}$. This is (on a very formal level) the so-called Jordan Kinderlehre Otto discretization of the PDE $\mathrm{div}(\mu_t \nabla \delta F(\mu_t)) + \partial_t \mu_t = 0$ with $\mu_{k\tau} \approx \rho_k$ when $\tau$ is small. We will call this the iJKO problem and for simplicity, consider the case where we have observations of two snapshots $\rho_0$ and $\rho_1$ One particular example of interest is where $\mathcal{F}(\rho) = \int V(x) d\rho(x)$ and in this context, we are tasked with recovering the potential function $V$ from iid samples of $\rho_k$. Such problems are of particular interest for understanding cell population dynamics in single-cell genomics.

*Fenchel-Young losses for inverse optimization.* One can view the two inverse problems above under the common umbrella of inverse optimization: Recover the parameter $c$ from noisy/sampled observations of an optimization solution $\pi(c) = \mathrm{argmin}_c \langle c, \pi \rangle + \Omega(\pi)$. In the context of iOT, $\Omega(\rho) := \mathrm{KL}(\pi | \alpha \otimes \beta) + \iota_{\mathcal{U}(\alpha, \beta)}$ while in the case of iJKO problem, $\Omega(\rho) = W_2^2(\rho, \rho_0)$. Given observation $\hat{\pi}$ and a discrepancy $D : \mathcal{P}(\mathcal{Z}) \times \mathcal{P}(\mathcal{Z}) \to [0, \infty]$ with $D(\rho, \rho) = 0$, the sharpened Fenchel-Young loss is

$$L(c, \hat{\pi}, \Omega, D) := \langle c, \hat{\pi} \rangle + \Omega(\hat{\pi}) - \inf_{\pi} \{ \langle c, \pi \rangle + \Omega(\pi) + D(\pi, \hat{\pi}) \}.$$

As a function of $c$, this loss satisfies the following three properties: For all $c$, $L(c, \hat{\pi}) \geq 0$ and $L(c, \hat{\pi}) = 0$ if $\hat{\pi} = \pi(c)$; It is differentiable if the inner problem

over $\pi$ has a unique solution; it is convex in $c$ since the infimum over affine functions is concave. In terms of relationship to well-established functions: if $D \equiv 0$, this is known as the Fenchel-Young loss and the non-negativity property is simply the Fenchel-Young inequality; if $D$ is the Bregman distance induced by $\Omega$, then this is the Fitzpatrick function. In practice, we parameterize $c$ in a linear manner $c_\theta = \theta^\top \phi := \sum_{j=1}^S \theta_j \phi_j$ for some basis $\{\phi_j\}_j$, $\hat{\pi}$ corresponds to an empirical measure (from sampled data), and $\Omega$ is only given approximately as $\hat{\Omega}$ since it often incorporates empirical data. In the following, since we are interested in minimizing over $c$, we drop the $\Omega(\hat{\pi})$ term when writing the loss.

For iOT, the sampled loss given data $\hat{\pi}^n$ is

$$J_n(\theta) = \langle \theta^\top \phi, \hat{\pi}^n \rangle - \inf_{\pi \in \mathcal{U}(\hat{\alpha}^n, \hat{\beta}^n)} \left\{ \langle \theta^\top \phi, \pi \rangle + \epsilon \mathrm{KL}(\pi | \hat{\alpha}^n \otimes \hat{\beta}^n) \right\}.$$

where $\hat{\alpha}^n$ and $\hat{\beta}^n$ are the marginals of $\hat{\pi}^n$, where we have taken $D \equiv 0$.

For iJKO, the sampled loss given empirical data $\hat{\rho}_0, \hat{\rho}_1$ is

$$J_n(\theta) = \langle \theta^\top \phi, \hat{\rho}_1 \rangle - \inf_{\alpha \in \mathcal{P}(\mathcal{X})} \left\{ \langle \theta^\top \phi, \alpha \rangle + W_{2,\epsilon}^2(\alpha, \hat{\rho}_0) + r\mathrm{KL}(\alpha | \hat{\rho}_1) \right\},$$

where we have taken $D(\alpha, \rho) = r\mathrm{KL}(\alpha | \hat{\rho}_1)$ for some $r > 0$.

Finally, due to the noisy data, we consider the regularized problem

$$\tag{1} \min_{\theta \in \mathbb{R}^S} \lambda R(\theta) + J_n(\theta),$$

for some (convex lower semi-continuous) regularizer $R$ with parameter $\lambda > 0$, which is often taken as the $\ell_1$ norm (to enforce sparsity) or nuclear norm (to enforce low-rankness).

*Main results.* We state our main results for the iOT setting [2]. Similar results can be derive for the iJKO problem [1].

*Theorem 1 on sample complexity. Fix the entropy regularization parameter $\epsilon > 0$. Let $\pi^\star$ be the entropic OT plan associated with cost $c^\star = (\theta^\star)^\top \phi$, and let $\alpha^\star, \beta^\star$ be its marginals. Assume that $\alpha^\star, \beta^\star$ are compactly supported, and the cost parameterization $\phi$ is such that its centered version is nondegenerate: define $\bar{\phi}(x,y) = \phi(x,y) - \int \phi(x,y)d\alpha^\star(x) - - \int \phi(x,y)d\beta^\star(y) + \int \phi(x,y)d\alpha^\star(x)d\beta^\star(y)$, and assume that*

$$\left( \mathbb{E}_{\alpha^\star \otimes \beta^\star}[\bar{\phi}_i(x,y)\bar{\phi}_j(x,y)] \right)_{i,j}$$

*is invertible. Then, the iOT loss $J$ defined with the full data $\pi^\star$ is locally strongly convex, locally Lipschitz smooth and is twice differentiable. Moreover, for all $t > 0$, with probability at least $1 - e^{-t}$, the minimizer $\hat{\theta}_n^\lambda$ to the sampled problem (1) is unique and satisfies*

$$\tag{2} \|\hat{\theta}_n^\lambda - \theta^\star\|_2 = \mathcal{O}\left( \sqrt{\frac{m_\alpha m_\beta (\log(n) + t)}{n}} \right) + \mathcal{O}(\lambda).$$

Let us make some remarks on the theorem: By convex duality on the inner problem can be written as a supremum over functions $f \in \mathcal{C}(\mathcal{X})$ and $f \in \mathcal{C}(\mathcal{Y})$,

leading to the alternative formulation

$$J_n(\theta) = \inf_{f,g} \langle \theta^\top \phi - (f \oplus g), \hat{\pi}^n \rangle + \epsilon \int \exp\left( \frac{f(x) + g(y) - \theta^\top \phi(x,y)}{\epsilon} \right) d\hat{\alpha}^n(x) d\hat{\beta}^n(y).$$

Due to the term $\theta^\top \phi - (f \oplus g)$, the recovered cost is invariant to addition by functions of the form $u(x) + v(y)$ – this is why we impose the assumption on linear independence of the centered parameterization functions $\bar{\phi}$ to ensure uniqueness of the minimizer. The main efforts in the proof is devoted to establishing local curvature properties of the loss and using concentration of measure results to establish the high probability estimates under sampled data.

One often takes $R(\theta) = \|\theta\|_1$ to enforce sparsity of the solution. Via a dual certificate/source condition, we can guarantee that the recovered solution $\hat{\theta}_\lambda^n$ has the same support as the underlying ground truth $\theta^\star$: First recall that the iOT loss $J$ with full data $\pi^\star$ is twice differentiable. Let $H := \nabla^2 J(\theta^\star)$ and define the certificate $z^\star := \mathrm{argmin}\{\langle z, (H^\star)^{-1} z \rangle : z \in \partial \|\theta^\star\|_1\}$. It is said to be *nondegenerate* if $z^\star$ is in the relative interior of $\partial \|\theta^\star\|_1$). We have the following result:

*Theorem 2 on sparsistency. Consider the setting of Theorem 1. Suppose that the certificate $z^*$ is non-degenerate. Let $\hat{\theta}$ minimize (1) with $R(\theta) = \lambda \|\theta\|_1$. Let $\delta > 0$. Then, for all sufficiently small regularization parameters $\lambda$ and sufficiently many number of samples $n$ with $\lambda \geq C\sqrt{\log(n/\delta)}/\sqrt{n}$, with probability at least $1 - \delta$, the minimizer $\hat{\theta}$ has the same support as $\theta^\star$.*

For simple settings (such as sampling from Gaussians), the non-degeneracy condition can be checked numerically and we carry out such a numerical investigation in [2]. Similar investigations are carried out for the iJKO setting in [1].

## References

[1] Andrade, Peyré, Poon, *Learning from Samples: Inverse Problems over measures via Sharpened Fenchel-Young Losses*, https://arxiv.org/abs/2505.07124

[2] Andrade, Peyré, Poon. *Sparsistency for inverse optimal transport*. In The Twelfth International Conference on Learning Representations.

[3] Dupuy, Galichon, Sun. *Estimating matching affinity matrices under low-rank constraints*. Information and Inference: A Journal of the IMA, 8(4):677–689, 2019.

## A super-resolution approach to classification

Hrushikesh Mhaskar

(joint work with Ryan O'Dowd)

The traditional approach in machine learning is to treat the classification problem as a problem of function approximation. This creates a gap between the theory, where one requires the target function to be smooth, and the practice, where the class boundaries may be non-smooth, even touching each other, or the distribution might be supported on a set of measure 0 in the ambient space. We propose a novel paradigm, where we consider each class $k$ appearing with a probability distribution $\mu_k$, $k = 1, \ldots, K$. The data consists of samples $\{x_j\}$ drawn from an

unknown convex combination of these distributions. We then use our localized kernels developed for solving super-resolution signal separation problems [3, 4, 1] to separate the supports of the measures $\mu_k$. With these supports identified accurately, one can then seek the label of one of the points in each of the supports, and extend it to the entire support. In this "cautious active learning" manner [2] one can solve the problem with a small if not minimal number of labels queried at judiciously chosen points $x_j$.

REFERENCES

[1] C. K. Chui and H. N. Mhaskar. *Signal decomposition and analysis via extraction of frequencies.* Applied and Computational Harmonic Analysis, 40(1):97–136, 2016.
[2] A. Cloninger and H. N. Mhaskar. *Cautious active clustering.* Applied and Computational Harmonic Analysis, 54:44–74, 2021.
[3] H. N. Mhaskar and J. Prestin. *On the detection of singularities of a periodic function.* Advances in Computational Mathematics, 12(2-3):95–131, 2000.
[4] H. N. Mhaskar and J. Prestin. *On local smoothness classes of periodic functions.* Journal of Fourier Analysis and Applications, 11(3):353–373, 2005.

## The Reduced Basis Method for problems of fractional order
### Karsten Urban

A *parameterized partial differential equation (PPDE)* is often formulated as follows: for any parameter $\mu \in \mathcal{P}$ out of a compact parameter set $\mathcal{P} \subset \mathbb{R}^P$, one seeks a solution $u(\mu) \in X$ ($X$ being a Hilbert space, the *trial* space) such that $a(u(\mu), v; \mu) = f(v; \mu)$ for all $v \in Y$ ($Y$ being a Hilbert space, the *test* space), where $a : X \times Y \times \mathcal{P} \to \mathbb{R}$ is a bounded bilinear form and $f \in Y'$ being a given right-hand side ($Y'$ denotes the dual space of $Y$). The *Reduced Basis Method (RBM)* aims at realizing an extremely efficient numerical approximation method in cases where the PPDE has to be evaluated for many parameters, in realtime or on devices with restricted memory or CPU.

The starting point is a detailed (or high-fidelity) discretization in terms of finite-dimensional subspaces $X^N \subset X$, $Y^N \subset Y$ of (usually large) dimension $N \in \mathbb{N}$ ($Y^N$ might also depend on $\mu$), such that the detailed approximation $u^N(\mu)$ is as close to $u(\mu)$ as desired (at the price of increasing numerical cost, of course). Then, certain samples $\mu^{(i)} \in \mathcal{P}$, $i = 1, ..., n \ll N$ are selected in an offline training phase (usually by a greedy algorithm) and the detailed solution method is used to compute *snapshots* $\xi^{(i)} := u^N(\mu^{(i)})$, $i = 1, ..., n$, which in turn are used as a basis for the reduced trial space $X_n := \text{span}\{\xi^{(i)} : i = 1, ..., n\}$. The reduced test space $Y_n(\mu)$ might be either fixed or formed e.g. by supremizers, [4].

The goal of the RBM is realize a good approximation to the *solution manifold* defined as $\mathcal{F} := \{u^N(\mu) : \mu \in \mathcal{P}\}$. The RB solution $u_n(\mu) \in X_n$ is defined as the (Petrov-)Galerkin projection onto $X_n$ along $Y_n(\mu)$. This allows for best approximation statement, i.e.,

$$(1) \qquad \|u^N(\mu) - u_n(\mu)\|_X \leq \frac{C_\mu}{\beta_\mu} \inf_{w_n \in X_n} \|u^N(\mu) - w_n\|_X,$$

where $C_\mu$ is the continuity and $\beta_\mu$ is the inf-sup constant of the bilinear form $a(\cdot, \cdot; \mu)$. This justifies that the benchmark for the RBM approximation is the *Kolmogorov n-width* defined by

$$(2) \qquad d_n(\mathcal{F}) := \inf_{\substack{X_n \subset X \\ \dim(X_n)=n}} \sup_{\mu \in \mathcal{P}} \inf_{w_n \in X_n} \|u^N(\mu) - w_n\|_X.$$

It was proven in [3, 5] that

$$(3) \qquad d_n(\mathcal{F}) \leq C \exp(-c\, n^{1/Q^a}),$$

with constants $0 < c, C < \infty$, if $a(\cdot, \cdot; \mu)$ is bounded, inf-sup stable and *affine*, i.e.,

$$(4) \qquad a(u, v; \mu) = \sum_{q=1}^{Q^a} \vartheta_q^a(\mu)\, a_q(u, v) \quad \forall \mu \in \mathcal{P}, u \in X, v \in Y$$

and $X_n$, $Y_n$ are inf-sup-stable (uniformly for all parameters $\mu$ and dimensions $n$). On the other hand, it is also known from [1, 3] that $d_n(\mathcal{F}) \geq n^{-1/2}$ for certain first order linear transport problems. In that light, we address two questions:

- How does the RBM and $d_n(\mathcal{F})$ behave for PPDEs of order $s \in (1, 2)$?
- Is there additional gain for the RBM as those problems are non-local?

In order to answer these questions, we consider the parameterized fractional order operator

$$(5) \qquad \mathcal{A}_s(\mu)\, u := {}_0\mathcal{D}_x^{s/2}(d(\mu)\, {}_0\mathcal{D}_x^{s/2}u) + r(\mu)u = f(\mu) \quad \text{on } \Omega := (0, 1),$$

where ${}_0\mathcal{D}_x^\beta$ denotes the left-sided Riemann-Liouville fractional derivative of order $\beta \in (0, \frac{1}{2})$, $d(\mu) \in L_\infty(\Omega)$ such that $d(x; \mu) \geq d_0 > 0$ for $x \in \Omega$ a.e., $r(\mu) \in L_\infty(\Omega)$ and $f(\mu) \in L_2(\Omega)$. Using the right-sided Riemann-Liouville fractional derivative ${}_x\mathcal{D}_1^\beta$, we define the bilinear form

$$a(u, v; \mu) := -(d(\mu){}_0\mathcal{D}_x^{s/2}u, {}_x\mathcal{D}_1^{s/2}v)_{L_2(\Omega)} + (r(\mu)\, u, v)_{L_2(\Omega)},$$

and the variational formulation of (5) amounts determining $u \in \tilde{H}^{s/2}(\Omega)$ such that

$$(6) \qquad a(u(\mu), v; \mu) = (f(\mu), v)_{L_2(\Omega)} \quad \text{for all } v \in \tilde{H}^{s/2}(\Omega),$$

where $\tilde{H}^\beta(\Omega)$ denotes the space of functions, whose zero extension to $\mathbb{R}$ is in $H^\beta(\mathbb{R})$. The bilinear form is bounded, i.e., $a(u, v; \mu) \leq C_{d,r}(\mu)\|u\|_{\tilde{H}^{s/2}(\Omega)}\|v\|_{\tilde{H}^{s/2}(\Omega)}$ for all $u, v \in \tilde{H}^{s/2}(\Omega)$ with $C_{d,r}(\mu) := 2\left(\|d(\mu)\|_{L_\infty(\Omega)} + \|r(\mu)\|_{L_\infty(\Omega)}\right)$. Moreover, defining the *average* $\nu(d)$ and the *range* $\rho(d)$ of a function $d : \Omega \to \mathbb{R}$ as

$$\nu(d) := \tfrac{1}{2}\left(\operatorname*{ess\,sup}_{x \in \Omega} d(x) + \operatorname*{ess\,inf}_{x \in \Omega} d(x)\right), \quad \rho(d) := \tfrac{1}{2}\left(\operatorname*{ess\,sup}_{x \in \Omega} d(x) - \operatorname*{ess\,inf}_{x \in \Omega} d(x)\right),$$

setting $\underline{r} := \operatorname*{ess\,inf}_{x \in \Omega} r(x)$ as well as $\gamma_{s,d}(\mu) := \nu(d(\mu))|\cos\left(s\frac{\pi}{2}\right)| - \rho(d(\mu))$ we have to assume that

$$(7) \qquad c_{s,d,r}(\mu) := \gamma_{s,d}(\mu)\tfrac{1}{4}\Gamma(\tfrac{s}{2} + 1)^2 + \underline{r} \geq 0.$$

Under this condition, we show that the bilinear form $a(\cdot, \cdot; \mu)$ is coercive, i.e., $a(u, u; \mu) \geq \alpha_{s,d}(\mu)\|u\|_{\tilde{H}^{s/2}(\Omega)}^2$ for $u \in \tilde{H}^{s/2}(\Omega)$, where $\alpha_{s,d}(\mu) = \gamma_{s,d}(\mu)\,\Gamma(s/2 +$

$1)^4/8$ and that there exists a unique solution of the variational problem (6) such that $\|u(\mu)\|_{\tilde{H}^{s/2}(\Omega)} \leq \frac{1}{\alpha_{s,d}(\mu)}\|f(\mu)\|_{L_2(\Omega)}$.

Discretizing the variational problem by continuous piecewise linear finite elements with mesh size $h$, the error of the discrete solution $u_h(\mu)$ can be bounded as $\|u(\mu) - u_h(\mu)\|_{\tilde{H}^{s/2}(\Omega)} \leq Ch^\beta \|f(\mu)\|_{L_2(\Omega)}$ for any $\beta \in [0, \frac{s}{2} - \frac{1}{2})$ and some positive $C > 0$. For the resulting stiffness matrix $A_h(\mu)$, there exists a constant $0 < c \neq c(h)$ such that $\kappa(A_h) \leq c \frac{h^{-s}}{\alpha_{s,d}(\mu)}$.

For the application of the RBM, we have to assume that the operator $\mathcal{A}_s(\mu)$ in (5) has affine components, i.e.,

$$c(\mu) = \sum_{q=1}^{Q^c} \vartheta_q^c(\mu)\, c_q, \qquad c \in \{d, r, f\},$$

where $\vartheta_q^c : \mathcal{P} \to \mathbb{R}$ and $d_q, r_q \in L_\infty(\Omega)$, $f_q \in L_2(\Omega)$. Then, the above quoted result from [3] yields $d_n(\mathcal{F}) \leq C \exp(-c\, n^{1/(Q^d + Q^r)})$, which, however, does not show the dependency on the order $s$. Instead, we show that for $d(\mu) \equiv 1$ (the general case is technically more involved, but gives similar bounds)

$$d_n(\mathcal{P}) \leq C \exp\left\{ -c_\Omega\, \frac{\alpha_s n}{|\mathcal{P}|}\right\},$$

where $\alpha_s = |\cos\left(s\frac{\pi}{2}\right)|\, \Gamma(s/2+1)^4/8$ is the coercivity constant for $d \equiv 1$ and $c_\Omega$ is constant only depending on $\Omega$.

We also show numerical results concerning the convergence of the finite element discretization, the conditioning of the stiffness matrix, the decay of the RBM error and the speedup of the RBM. The talk is based upon [2], where details can be found.

## REFERENCES

[1] F. Arbes, C. Greif and K.. Urban, *The Kolmogorov N-width for linear transport: Exact representation and the influence of the data*, Adv. Comput. Math. **51**, no. 13 (2025).
[2] R. Aylwin, G. Oruc and K.. Urban, *Fractional differential equations: non-constant coefficients, simulation and model reduction*, Ulm University, preprint (2025).
[3] M. Ohlberger and S.. Rave, *Reduced Basis Methods: Success, Limitations and Future Challenges*, Proceedings of the Conference Algoritmy 2016, 1–12.
[4] G. Rozza and K. Veroy, *On the stability of the reduced basis method for Stokes equations in parametrized domains*, Computer Methods in Applied Mechanics and Engineering **196**, no. 7, 1244–1260 (2007).
[5] K.. Urban, *The Reduced Basis Method in Space and Time: Challenges, Limits and Perspectives*, in *Model Order Reduction and Applications*, M. Falcone and G. Rozza (eds.), Cetraro 2021, C.I.M.E. Series, Springer 2023, 1–73.

# Singular Perturbations, Min-Max Optimization, and Accuracy Control

Wolfgang Dahmen

(joint work with Zhu Wang)

Trying to recover a (complex) "physical state of interest" from experimental observations/measurements is a ubiquitous task in science and technology. Since due to pysical constraints or acquisition cost, such data, as a sole source of information, are usually far from sufficient for an acccurate recovery so that one has to make use of further prior information. This could be obtained from governing physical laws that are typically given in terms of systems of partial differential equations (PDEs) involving "uncertain" problem data such as coefficient fields, equations of states, initial conditions or source terms - in brief *parameters* the PDE depends on. Related inverse tasks, like state- or parameter-estimation, require an efficient exploration of the "solution manifold", comprised of all solutions obtained when traversing the parameter domain. Viewing the solutions of the parameter dependent PDEs as functions of spatio-temporal and parametric variables, the recovery tasks is intrinsically high-dimensional. This has stirred interest in leveraging machine learning concepts to mitigate complexity challenges, due to the Curse of Dimensionality. The central theme in this talk is a "learning-based" construction of accurate and yet efficient surrogates for the underlying *parameter-to-solution* (PtS) map under the following provisions: (i) a rigorous accuracy quantification with respect to "model-compliant" norms; (ii) "supervised learning" in terms of regression should be based on *residual-type* training losses; (iii) estimation quality should be as robust as possible for desirable large parameter ranges that may cause near-degeneracies in the family of PDEs.

(i) is important to avoid aggravating the level of ill-posedness when using the surrogate in the context of an inverse problem. For instance, imposing "excess regularity" narrows the range of the solution operator. (ii) is to avoid the need to compute a large number of high-fidelity solutions as training data for regression over hypothesis classes of high expressivity (such as deep neural network). Since the nonlinearity of such hypothesis classes renders entailed optimization problem non-convex, one encounters an inherent uncertainty in optimization success. Therefore, a priori expressivity results for the hypothesis class are of little (if any) help. This has motivated the notion of *variationally correct residual losses* which roughly means that the loss itself provides up to uniform proportionality constants lower and upper bounds for the deviation of an estimator from the exact solution with respect to a model-compliant norm [1]. Such norms are actually dictated by appropriate *stable variational formulations* for the underlying PDE family. Specifically, we show how variational correctness of the residual loss is based then on the fact that errors in the trial norm are uniformly proportional to the residual, measured in the dual test-norm. Since a dual norm of a Hilbert space involves the supremization over that Hilbert space, it can generally not be evaluated exactly (unless the Hilbert space is a product of $L_2$-spaces and hence self-dual). We discuss several strategies for constructing equivalent computable quantities [1, 4]. The

common starting point of these techniques is to first transform the PDE into an equivalent system of first order PDEs which increases the flexibility of identifying appropriate stable variational formulations, including formulations with self-dual test spaces, see also [8]. In particular, we highlight the fact that nevertheless the proportionality between errors and residuals - hence variational correctness - may degrade significantly when the parameter range includes singularly perturbed models (e.g. high-contrast diffusion models, [4]), which concerns issue (iii). We therefore discuss, as a potential remedy, the role of so called *ultra-weak* formulations in combination with *optimal test-norms* which are well-defined as long as the operator, induced by a perhaps more conventional formulation, is bijective (regardless of a perhaps prohibitively large condition in the standard formulation), see e.g. [5, 7]. Endowing the test space with the optimal test-norm, can be viewed as "preconditioning" the operator equation on the infinite-dimensional level, namely the induced operator becomes an isometry between the trial space and the dual test space. Thus, the residual in this dual test-norm is *equal* to the error in the trial-norm, even in the regime of nearly degenerating coefficients. The price is that now the learning problem boils down to solving a *min-max optimization* problem, reflecting a non-trivial dual norm. This option has indeed been explored previously (see e.g. [3, 2, 9]) but it is fair to say that the employed ad-hoc approaches of alternating gradient ascent and descent steps for a given hypothesis class and an associated "adversarial test-class" has significant robustness issues. As a potential remedy we propose the following strategy: First we reformulate the the min-max problem for the quotients, defining the dual norm, as an equivalent *affine-quadratic saddle point problem*. This exploits the fact that the supremizer is the Riesz-lift of a symmetric coercive (in this sense elliptic) variational problem which is characterized as the minimizer of an affine-quadratic energy functional. As a next step we formulate a fictitious *primal-dual proximal* scheme in Hilbert space whose convergence, due to the very benign affine quadratic structure of the Lagrangian, can be quantified. A slight issue is to properly account for the different norms imposed by the underlying PDE model. Understanding the convergence properties of such an idealized scheme, allows one then to determine tolerances within which the execution of the involved proximal maps can be inexact while still warranting convergence to the exact limit. Again drawing on the Hilbert space formulation, we derive a-posteriori bounds that are used to check whether the approximate proximal map meets the permissible perturbation tolerance. If this is the case we obtain a convergent scheme with certified error bounds. In this sense we obtain a rigorous "conditional convergence result". Specifically, we obtain a certifiable solution of the initial "global" min-max problem, provided that we can solve the much simpler minimization problems, posed by the proximal maps, albeit over nonlinear hypothesis classes. The fact, that the exact minimizer gets increasingly closer to the initial guess is a further favorable aspect regarding practical realizations. Nevertheless, meeting the tolerances will in general require gradually enlarging the hypothesis classes, as a further constituent of this paradigm. For that purpose

we plan to employ such expansion strategies from [6], developed in the context of minimizing convex Hilbert space energies over neural network hypothesis classes.

## References

[1] M. Bachmayr, W. Dahmen, M. Oster, Variationally correct neural residual regression for parametric PDEs: on the viability of controlled accuracy, http://arxiv.org/abs/2405.20065, to appear in IMA Journal of Numerical Analysis

[2] G.Bao, X.Ye, Y.Za, H.Zhu, Numerical solution of inverse problems by weak adversarial networks, Inverse Problems, **36** (2020), doi: org/10.1088/1361-6420/abb447.

[3] F. Chen, J. Huang, C. Wang, H. Yang, Friedrichs learning: weak solutions of partial differential equations via deep learning, arXiv: 2012.08023, 2021.

[4] P. Cortés Castillo, W. Dahmen, J. Gopalakrishnan, DPG loss functions for learning parameter-to-solution maps by neural networks, June 23, 2025, http://arxiv.org/abs/2506.18773.

[5] W. Dahmen, C.Huang, C.Schwab, G.Welper, Adaptive Petrov-Galerkin methods for first order transport equations, SIAM J. Numer. Anal., **50** (5) (2012), 2420–2445.

[6] W. Dahmen, W. Li, Y. Teng, Z. Wang, Expansive Natural Neural Gradient Flows for Energy Minimization, arXiv: 2507.13475 [math.NA], July 2025.

[7] L. Demkowicz, J. Gopalakrishnan, The Discontinuous Petrov-Galerkin Method, Acta Numerica, Cambridge University Press, 2025.

[8] J. Opschoor, P. Petersen, C. Schwab, First Order System Least Squares Neural Networks, arXiv:2409.20264v1 [math.NA], 30 Sep 2024.

[9] Y. Zang, G. Bao, X. Ye, H. Zhou, Weak adversarial networks for high-dimensional partial differential equations, Journal of Computational Physics, **411** (2020), 109409, https://doi.org/10.1016/j.jcp.2020.109409,

# Nonlinear manifold approximation using compositional polynomial networks

Anthony Nouy

(joint work with Antoine Bensalah, Joel Soffo)

We consider the problem of approximating a subset $M$ of a normed space $X$ by a low-dimensional manifold $M_n$, using samples from $M$. We propose a nonlinear approximation method where $M_n$ is defined as the range of a smooth nonlinear decoder $D$ defined on $\mathbb{R}^n$ with values in a possibly high-dimensional linear space $X_N$, and a linear encoder $E$ which associates to an element from $M$ its coefficients $E(u)$ on a basis of a $n$-dimensional subspace $X_n \subset X_N$, where $X_n$ and $X_N$ are optimal or near to optimal linear spaces, depending on the selected error measure. The linearity of the encoder allows to easily obtain the parameters $E(u)$ associated with a given element $u$ in $M$. The proposed decoder is a polynomial map from $\mathbb{R}^n$ to $X_N$ which is obtained by a tree-structured composition of polynomial maps, estimated sequentially from samples in $M$. Rigorous error and stability analyses are provided, as well as an adaptive strategy for constructing a decoder that guarantees an approximation of the set $M$ with controlled mean-squared or wort-case errors, and a controlled stability (Lipschitz continuity) of the encoder and decoder pair.

Then we consider the problem of approximating online a new element $u \in M$ by an element of the manifold $M_n$, either with an operator learning point of view when $M$ is the image of some operator defined on a parameter set, or by solving a best approximation problem with natural gradient schemes, using new and adaptive information.

REFERENCES

[1] A. Bensalah, A. Nouy, and J. Soffo, *Nonlinear manifold approximation using compositional polynomial networks.* arXiv e-prints arXiv:2502.05088, Feb. 2025.
[2] R. Gruhlke, A. Nouy, and P. Trunschke, *Optimal sampling for stochastic and natural gradient descent.* arXiv e-prints arXiv:2402.03113, 2024.

## On the computational and statistical complexity of predicting non-linear dynamical systems

RICHARD NICKL

We discuss recent progress in our understanding of Bayesian inference methods for parameters or states of time evolution phenomena modelled by non-linear partial differential equations (PDEs) such as Navier-Stokes, McKean-Vlasov, and reaction-diffusion systems. We will show that posteriors can deliver consistent solutions in the 'informative' large data/small noise limit, discuss probabilistic approximations to the fluctuations of such posterior measures in infinite dimensions, and how such results can be used to show that the non-convex problem of computation of the associated 'filtering' distributions are polynomial time problems. Relevant references are [1, 2, 3].

REFERENCES

[1] Dimitri Konen and Richard Nickl, *Data assimilation with the 2D Navier-Stokes equations: Optimal Gaussian asymptotics for the posterior measure*, ArXiv preprint (2025).
[2] R. Nickl and Sven Wang, *On polynomial-time computation of high-dimensional posterior measures by Langevin-type algorithms*, Journal of the European Mathematical Society (2024).
[3] R. Nickl, *Bayesian Non-linear Statistical Inverse Problems*, European Mathematical Society Press, (2023).

## Kernel methods in high dimensions with lengthscale informed sparse grids

ARETHA TECKENTRUP

Kernel methods, in the form of radial basis function interpolation and Gaussian process regression, have proved successful as a tool for various tasks in approximation and inference. This talk focussed on presenting recent advances in the design and numerical analysis of kernel methods in the context of modern applications, which typically involve very high dimensions. To combat this challenge, we

introduce anisotropic kernels and lengthscale informed sparse grids that allow for accurate reconstruction and efficient compuatation in this setting.

## Stable Nonlinear Dynamical Approximation with Dynamical Sampling

Daan Bon

(joint work with Benjamin Caris, Olga Mula)

We present a nonlinear dynamical approximation method for time-dependent Partial Differential Equations (PDEs). The approach makes use of parametrized decoder functions, e.g. (shallow) neural networks. The parameters of these functions are evolved in time by finding a curve of parameters that induces a curve of decoders, such that the time derivative of the decoders matches up with the considered PDE evolution. This is ensured by projecting the right hand side of the PDE evolution onto the span of the partial derivatives of the decoder w.r.t. the parameters, an approach known as the Dirac-Frenkel principle.

These projections are w.r.t. an ambient Hilbert space, and thus usually require the need to perform some integration over a spatial domain. This can be costly if the domain is high dimensional, as many quadrature points are needed, or if the PDE solution has a local and moving support. We therefore propose an approach that instead makes use only of a restricted amount of local information of the functions, through certain linear functionals (such as a small Gaussian average). These functionals are then evolved in time together with the approximation, so that we do not require to cover the full spatial domain with quadrature points. This evolution of the linear functionals is driven by the optimization of a stability constant, which is directly part of an error bound between the true solution and our proposed approximation.

We show several applications of the method in both low and high dimensions, and in particular show that it works well even if the PDE solution has a local and moving support, for which quadrature methods would have been costly.

## Neural and Spectral Operator surrogates on Gaussian spaces

Christoph Schwab

(joint work with C. Marcati, M. Maric and J. Zech)

We prove expression rate bounds of finite-parametric, spectral and neural surrogates for holomorphic maps between separable spaces. In the encoder-approximator-decode framework, with Riesz-basis encoders, and frame decoders, we prove expression rate bounds for two classes of finite-parametric surrogates: i) spectral surrogates obtained by $N$-term trunations of Wiener-Hermite polynomial chaos expansions and ii) neural surrogates obtained by approximation of parametric maps with several families of neural networks: ReLU, clipped ReLU and RePU activations are considered. We work under abstract hypotheses on weighted summability of encoded inputs on sequence spaces, and with Gaussian measures charging the set

of admissible operator inputs. Operator approximation rates are in mean-square and in hilbertian gaussian Sobolev spaces.

Work is based on and generalizing [1], [2], using novel sparsity results of parametric, holomorphic maps from [3].

REFERENCES

[1] L. Herrmann, Ch. Schwab, J. Zech, Neural and spectral operator surrogates: construction and expression rate bounds, Advances in Computational Mathematics **50** (4), 72, (2024)
[2] D. Dũng, V.K. Nguyen, Ch. Schwab, J. Zech, Analyticity and sparsity in uncertainty quantification for PDEs with Gaussian random field inputs, Springer LNM 2334 (2023), Springer Verlag
[3] C. Marcati, Ch. Schwab, J. Zech, Sparsity for Infinite-Parametric Holomorphic Functions on Gaussian Spaces, arXiv:2504.21639

## Deep neural network analysis made easy

MARIO ULLRICH

(joint work with Cornelia Schneider, Jan Vybiral)

Recently, Daubechies, DeVore, Foucart, Hanin, and Petrova introduced a system of piece-wise linear functions, which can be easily reproduced by artificial neural networks with the ReLU activation function and which form a Riesz basis of $L_2([0,1])$. This work was generalized by two of the authors to the multivariate setting. We show that this system serves as a Riesz basis also for Sobolev spaces $W^s([0,1]^d)$ and Barron classes $\mathbb{B}^s([0,1]^d)$ with smoothness $0 < s < 1$. We apply this fact to re-prove some recent results on the approximation of functions from these classes by deep neural networks. Our proof method avoids using local approximations and allows us to track also the implicit constants as well as to show that we can avoid the curse of dimension. Moreover, we also study how well one can approximate Sobolev and Barron functions by ANNs if only function values are known.

REFERENCES

[1] I. Daubechies and R. DeVore and S. Foucart and B. Hanin and G. Petrova, *Nonlinear approximation and (deep) ReLU networks*, Constr. Approx, **55** (2022), 127–172.
[2] C. Schneider, M. Ullrich and J. Vybiral, *Nonlocal techniques for the analysis of deep ReLU neural network approximations*, arXiv:2504.04847.

# Wavelet compressed, modified Hilbert transform in the space-time discretization of the heat equation

Helmut Harbrecht

(joint work with Christoph Schwab and Marco Zank)

## 1. Introduction

The efficient numerical solution of initial-boundary value problems is central in computational science and engineering. Accordingly, numerical methods have been developed to a high degree of sophistication and maturity. Foremost among these are time-stepping schemes, which are motivated by the causality of the physical phenomena modelled by the equations. They discretize the evolution equation via sequential numerical solution of a sequence of spatial problems. In recent years, however, especially motivated by applications from numerical optimal control, so-called *space-time methods* have emerged: these methods aim at the "one-shot" solution of the initial-boundary value problem as a well-posed operator equation on a space-time cylinder. We present here such a space-time method for the efficient solution of linear, parabolic initial boundary value problems.

## 2. Problem formulation

Without loss of generality, we set $I := (0,1)$. We intend to find the function $u(x,t)$ satisfying the following linear, parabolic initial-boundary value problem

(1)
$$\partial_t u - \operatorname{div}(A\nabla u) = f \quad \text{in } Q = \Omega \times I,$$
$$u = 0 \quad \text{on } \Gamma_{\mathrm{D}} \times \overline{I},$$
$$\frac{\partial u}{\partial n} = 0 \quad \text{on } \Gamma_{\mathrm{N}} \times \overline{I},$$
$$u(\cdot, 0) = 0 \quad \text{in } \Omega$$

with given right-hand side $f$ and, for simplicity, homogeneous Dirichlet boundary conditions on $\Gamma_{\mathrm{D}} \subset \partial\Omega$, homogeneous Neumann boundary conditions on $\Gamma_{\mathrm{N}} \subset \partial\Omega$ and homogeneous initial conditions. Here, $\Omega \subset \mathbb{R}^n$, $n \in \{2,3\}$, is a bounded domain with Lipschitz boundary and $A \in [L^\infty(\Omega)]^{n \times n}$ is a uniformly elliptic diffusion matrix.

## 3. Variational formulation

We introduce the temporal spaces by

$$H^1_{0,}(I) = \{z \in H^1(I)|\ z(0) = 0\}, \quad H^1_{,0}(I) = \{z \in H^1(I)|\ z(1) = 0\}$$

and define the respective fractional-order Sobolev spaces by interpolation

$$H^s_{0,}(I) = [H^1_{0,}(I), L^2(I)]_s, \quad H^s_{,0}(I) = [H^1_{,0}(I), L^2(I)]_s$$

for $s \in (0,1)$. We moreover denote the space that consists of all functions from $H^1(\Omega)$, fulfilling homogeneous Dirichlet conditions on $\Gamma_\mathrm{D}$, by $H^1_{\Gamma_\mathrm{D}}(\Omega)$. With these spaces at hand, we introduce the following intersection spaces

$$H^{1,1/2}_{\Gamma_\mathrm{D};0,}(Q) = \big(H^1_{\Gamma_\mathrm{D}}(\Omega) \otimes L^2(I)\big) \cap \big(L^2(\Omega) \otimes H^{1/2}_{0,}(I)\big),$$

$$H^{1,1/2}_{\Gamma_\mathrm{D};,0}(Q) = \big(H^1_{\Gamma_\mathrm{D}}(\Omega) \otimes L^2(I)\big) \cap \big(L^2(\Omega) \otimes H^{1/2}_{,0}(I)\big),$$

equipped with the sum norm, and the duality pairing

$$\langle \cdot, \cdot \rangle_Q : \big[H^{1,1/2}_{\Gamma_\mathrm{D};,0}(Q)\big]' \times H^{1,1/2}_{\Gamma_\mathrm{D};,0}(Q) \to \mathbb{R}$$

as continuous extension of the $L^2(Q)$ inner product. Then, the space-time variational formulation of (1) reads: Seek $u \in H^{1,1/2}_{0;0,}(Q)$ such that

(2)    $\forall w \in H^{1,1/2}_{0;,0}(Q) : \quad b(u,w) := \langle \partial_t u, w \rangle_Q + \langle \nabla_x u, \nabla_x w \rangle_{[L^2(Q)]^n} = \langle f, w \rangle_Q.$

It is uniquely solvable and induces an isomorphism

$$\partial_t - \Delta \in \mathcal{L}_\mathrm{iso}\big(H^{1,1/2}_{0;0,}(Q), \big[H^{1,1/2}_{0;,0}(Q)\big]'\big),$$

compare [5].

## 4. Modified Hilbert transform

For a given function $z \in L^2(I)$ with Fourier coefficients

$$z_k = \sqrt{2} \int_0^1 z(t) \sin\left(\left(\frac{\pi}{2} + k\pi\right)t\right) \mathrm{d}t$$

and series representation

$$z(t) = \sum_{k=0}^\infty z_k \sqrt{2} \sin\left(\left(\frac{\pi}{2} + k\pi\right)t\right), \quad t \in I,$$

the modified Hilbert transform is defined by the series

$$(\mathcal{H}_T z)(t) = \sum_{k=0}^\infty z_k \sqrt{2} \cos\left(\left(\frac{\pi}{2} + k\pi\right)t\right), \quad t \in I.$$

It defines an isometry

$$\mathcal{H}_T \colon H^s_{0,}(I) \to H^s_{,0}(I)$$

for all $s \in [0,1]$, see e.g. [5] for the details. This property is the key feature for the space-time formulation, which we will use in the following. Namely, instead of (2), we shall from now on consider the variational formulation: Seek $u \in H^{1,1/2}_{0;0,}(Q)$ such that

(3)                    $\forall w \in H^{1,1/2}_{0;0,}(Q) : \quad b(u, \mathcal{H}_T w) = \langle f, \mathcal{H}_T w \rangle_Q.$

This formulation is unconditionally stable, satisfying the stability estimate

$$\|u\|_{L^2(\Omega) \otimes H^{1/2}_{0,}(I)} \leq \|f\|_{[L^2(\Omega) \otimes H^{1/2}_{,0}(I)]'}$$

provided that $f \in [L^2(\Omega) \otimes H^{1/2}_{,0}(I)]'$.

## 5. Discretization

A discretization by product spaces of Lagrangian finite elements in space and splines in time yields a linear system of equations of the form

$$(4) \qquad \left( \mathbf{A}_j^t \otimes \mathbf{M}_j^x + \mathbf{M}_j^t \otimes \mathbf{A}_j^x \right) \mathbf{u}_j = \mathbf{f}_j.$$

Here, the temporal system matrices are given given by

$$\mathbf{A}_j^t = \left[ \langle \partial_t \phi_{j,k'}^t, \mathcal{H}_T \phi_{j,k}^t \rangle_{L^2(0,1)} \right]_{k,k'}, \quad \mathbf{M}_j^t = \left[ \langle \phi_{j,k'}^t, \mathcal{H}_T \phi_{j,k}^t \rangle_{L^2(0,1)} \right]_{k,k'}$$

and the spatial system matrices are given by

$$\mathbf{M}_j^x = \left[ \langle \phi_{j,k'}^x, \phi_{j,k}^x \rangle_{L^2(\Omega)} \right]_{k,k'}, \quad \mathbf{A}_j^x = \left[ \langle \nabla_x \phi_{j,k'}^x, \nabla_x \phi_{j,k}^x \rangle_{L^2(\Omega)^n} \right]_{k,k'},$$

while $\mathbf{f}_j$ is the related right-hand side. The above matrices are positive definite and the matrices $\mathbf{A}_j^t$, $\mathbf{A}_j^x$, $\mathbf{M}_j^x$ are symmetric, whereas $\mathbf{M}_j^t$ is nonsymmetric. However, since the modified Hilbert transform is a nonlocal operator, we shall apply spline wavelets in order to be able to compress the matrices $\mathbf{A}_j^t$ and $\mathbf{M}_j^t$ in accordance with [2]. Then, also the graph-algorithm based computation of their inverses becomes possible, see [3]. This enables to apply the BPX-type preconditioner

$$\sum_{\ell=0}^{j} \left( \mathbf{I}_j^t \otimes \mathbf{I}_\ell^j \right) \left( \left( \mathbf{A}_j^t + (h_\ell^x)^{-2} \mathbf{M}_j^t \right)^{-1} \otimes \mathbf{I}_\ell^x \right) \left( \mathbf{I}_j^t \otimes \mathbf{I}_j^\ell \right)$$

for the efficient iterative solution of (4), compare [1]. Finally, due to the use of a multilevel basis in time, a *sparse-tensor product discretization* between space and time becomes easily realizable. In this case, the temporal coordinate is basically for free and the complexity corresponds, up to a poly-logarithmic factor, to the discretization of a purely spatial partial differential equation. We refer to [4] for all the details.

## References

[1] J. Bramble, J. Pasciak, and J. Xu, *Parallel multilevel preconditioners*, Math. Comput. **55** (1990), 1–22.

[2] W. Dahmen, H. Harbrecht, and R. Schneider, *Compression techniques for boundary integral equations — asymptotically optimal complexity estimates*, SIAM J. Numer. Anal. **43** (2006), no. 6, 2251–2271.

[3] H. Harbrecht and M.D. Multerer, *A fast direct solver for nonlocal operators in wavelet coordinates.* J. Comput. Phys. **428** (2021), 110056.

[4] H. Harbrecht, C. Schwab, and M. Zank, *Wavelet compressed, modified Hilbert transform in the space-time discretization of the heat equation*, IMA J. Numer. Anal. (2025), to appear.

[5] M. Zank, *Inf-sup stable space-time methods for time-dependent partial differential equations*, Monographic Series TU Graz: Computation in Engineering and Science, vol. 36, TU Graz, Austria, 2020.

## Solving the (Multi-)Electronic Schrödinger Equation with Deep Learning

PHILIPP GROHS

I presented some recent results from [1, 2].

### REFERENCES

[1] M. Scherbela, L. Gerard, P. Grohs, Towards a transferable fermionic neural wavefunction for molecules. Nature Communications. 2024 Jan 2;15(1):120. Epub 2024 Jan 2. `doi: 10.1038/s41467-023-44216-9`
[2] M. Scherbela, N. Gao, P. Grohs and S. Günnemann, Accurate Ab-initio Neural-network Solutions to Large-Scale Electronic Structure Problems, 2025, `//arxiv.org/abs/2504.06087`

## Kernel interpolation on generalized sparse grids

MICHAEL MULTERER

(joint work with Michael Griebel, Helmut Harbrecht)

Let $\mathcal{H} = \bigotimes_{i=1}^{m} \mathcal{H}^{(i)}$ be a tensor product Hilbert space of functions $f \colon \boldsymbol{\Omega} \to \mathbb{R}$ defined on the product region $\boldsymbol{\Omega} = \Omega_1 \times \cdots \times \Omega_m$ with the unidirectional regions $\Omega_i \subset \mathbb{R}^{d_i}$. Each unidirectional space $\mathcal{H}^{(i)}$ is assumed to be a reproducing kernel Hilbert space on $\Omega_i \subset \mathbb{R}^{d_i}$ with reproducing kernel $\kappa_i$, $i = 1, \ldots, m$. Therefore, the space $\mathcal{H}$ is a reproducing kernel Hilbert space itself with product kernel $\boldsymbol{\kappa}(\boldsymbol{x}, \boldsymbol{y}) = \prod_{i=1}^{m} \kappa_i(x_i, y_i)$ defined on the product region $\boldsymbol{\Omega} = \Omega_1 \times \cdots \times \Omega_m$.

Given the unidirectional sets of *data sites* $X^{(i)} \subset \Omega_i$, $i = 1, \ldots, m$, we define the product grid $\boldsymbol{X} := \{[x_1, \ldots, x_m] : x_i \in X^{(i)}\}$. Associated to $\boldsymbol{X}$, we introduce the subspace of kernel translates $\mathcal{H}_{\boldsymbol{X}} := \operatorname{span}\{\boldsymbol{\kappa}(\boldsymbol{x}, \cdot) : \boldsymbol{x} \in \boldsymbol{X}\}$. By the reproducing property, the $\mathcal{H}$-orthogonal projection $f_{\boldsymbol{X}}$ of any function $f \in \mathcal{H}$ onto $\mathcal{H}_{\boldsymbol{X}}$ coincides with the interpolant $f_{\boldsymbol{X}}(\boldsymbol{x}_i) := \sum_{j=1}^{|\boldsymbol{X}|} \alpha_j \boldsymbol{\kappa}(\boldsymbol{x}_j, \boldsymbol{x}_i) = f(\boldsymbol{x}_i)$ for all $\boldsymbol{x}_i \in \boldsymbol{X}$. Introducing the data vector $\boldsymbol{f} := [f(\boldsymbol{x}_i)]_i$, the coefficient vector of the interpolant can be computed by solving the linear system

$$(\boldsymbol{K}^{(1)} \otimes \cdots \otimes \boldsymbol{K}^{(m)})\boldsymbol{\alpha} = \boldsymbol{f} \tag{1}$$

with the unidirectional kernel matrices $\boldsymbol{K}^{(i)} := [\kappa_i(x, y)]_{x,y \in X^{(i)}}$. Our goal is to approximately solve the linear system (1) by using generalized sparse grids and to have control on the resulting approximation error.

For the sparse grid construction, we start from nested sequences $X_0^{(i)} \subset X_1^{(i)} \subset \cdots \subset X^{(i)} \subset \Omega_i$ with decreasing *fill distance* $h_j^{(i)} := \sup_{x \in \Omega_i} \min_{y \in X_j^{(i)}} \|x - y\|_2 \sim 2^{-j}$ for $i = 1, \ldots, m$. Associated to the sequence of points, there is a sequence of subspaces $\mathcal{H}_0^{(i)} \subset \mathcal{H}_1^{(i)} \subset \cdots \subset \mathcal{H}^{(i)}$, $\mathcal{H}_j^{(i)} := \operatorname{span}\{\kappa_i(x, \cdot) : x \in X_j^{(i)}\}$ of kernel translates for each unidirectional space. Assuming $\mathcal{H}^{(i)} \cong H^{s_i}(\Omega_i)$ for $s_i > d_i/2$, there hold the univariate error estimates

$$\left\|f - P_j^{(i)} f\right\|_{H^{t_i}(\Omega_i)} \lesssim \left(h_j^{(i)}\right)^{t'_i - t_i} \|f\|_{H^{t'_i}(\Omega_i)}, \quad 0 \le t_i \le s_i \le t'_i \le 2s_i, \; f \in H^{t'_i}(\Omega_i),$$

see [1, 6, 7]. Herein $P_j^{(i)} f = f_{X_j^{(i)}}$ denotes the $\mathcal{H}^{(i)}$-orthogonal projection onto the subspace $\mathcal{H}_j^{(i)}$. Similarly, the *detail projection* $Q_j^{(i)} := P_j^{(i)} - P_{j-1}^{(i)}$, $P_{-1}^{(i)} := 0$, satisfies

$$\left\| Q_j^{(i)} f \right\|_{H^{t_i}(\Omega_i)} \lesssim \left( h_j^{(i)} \right)^{t_i' - t_i} \| f \|_{H^{t_i'}(\Omega_i)}, \quad 0 \le t_i \le s_i \le t_i' \le 2s_i, \ f \in H^{t_i'}(\Omega_i).$$

Given the detail projections, we define the *sparse grid projection*

$$\widehat{\boldsymbol{P}}_J^{\boldsymbol{w}} : \mathcal{H} \to \widehat{\mathcal{H}}_J^{\boldsymbol{w}}, \quad \widehat{\boldsymbol{P}}_J^{\boldsymbol{w}} f = \sum_{\boldsymbol{j}^\mathsf{T} \boldsymbol{w} \le J} \boldsymbol{Q}_{\boldsymbol{j}} f := \sum_{\boldsymbol{j}^\mathsf{T} \boldsymbol{w} \le J} \left( Q_{j_1}^{(1)} \otimes \cdots \otimes Q_{j_m}^{(m)} \right) f.$$

**Theorem 1.** *Let* $f \in \boldsymbol{H}^{\boldsymbol{t}'}(\boldsymbol{\Omega})$ *and let* $\boldsymbol{X}$ *be quasi-uniform, i.e., the fill-distance scales like the separation radius. Denote by* $N := \dim \widehat{\mathcal{H}}_J^{\boldsymbol{w}}$ *the number of degrees of freedom in the sparse tensor product space* $\widehat{\mathcal{H}}_J^{\boldsymbol{w}}$ *and set*

$$(2) \qquad \beta := \frac{\min\{r_1/w_1, \dots, r_m/w_m\}}{\max\{d_1/w_1, \dots, d_m/w_m\}}, \quad r_i := t_i' - t_i.$$

*Assume that the minimum in the numerator is attained* $P \in \mathbb{N}$ *times and the maximum in the denominator is attained* $R \in \mathbb{N}$ *times. Then, the sparse grid kernel interpolant in* $\widehat{\mathcal{H}}_J^{\boldsymbol{w}}$ *satisfies the error estimate*

$$\left\| (I - \widehat{\boldsymbol{P}}_J^{\boldsymbol{w}}) f \right\|_{\boldsymbol{H}^{\boldsymbol{t}}(\boldsymbol{\Omega})} \lesssim N^{-\beta} (\log N)^{(P-1) + \beta(R-1)} \| f \|_{\boldsymbol{H}^{\boldsymbol{t}'}(\boldsymbol{\Omega})}.$$

For all $\boldsymbol{w} > \boldsymbol{0}$, there holds $\beta \le \beta^\star := \min\{r_1/d_1, \dots, r_m/d_m\}$, see [2]. The maximum rate $\beta^\star$ is always achievable. More precisely, the maximum rate $\beta = \beta^\star$ is attained for all $\boldsymbol{w} > \boldsymbol{0}$ such that $r_\ell/r_i \le w_\ell/w_i \le d_\ell/d_i$, $i = 1, 2, \dots, m$, where, without loss of generality, $r_\ell/d_\ell = \beta^\star$. Several canonical choices for the weight are possible:

- *Equilibration of accuracy*: Set $w_i \sim r_i$ such that
$$2^{-Jr_1/w_1} = 2^{-Jr_2/w_2} = \cdots = 2^{-Jr_m/w_m}.$$

- *Equilibration of degrees of freedom*: Set $w_i \sim d_i$ such that
$$2^{Jd_1/w_1} = 2^{Jd_2/w_2} = \cdots = 2^{Jd_m/w_m}.$$

- *Equilibration of cost-benefit rate*: Set $w_i \sim (d_i + 2s_i)$ such that
$$2^{j_1(d_1 + r_1)} \cdot 2^{j_2(d_2 + r_2)} \cdots 2^{j_m(d_m + r_m)} = 2^{J \cdot const.}, \quad \boldsymbol{j}^\mathsf{T} \boldsymbol{w} = J.$$

Due to the Galerkin orthogonality, the detail projections satisfy

$$(\boldsymbol{Q}_{\boldsymbol{j}} u, \boldsymbol{Q}_{\boldsymbol{j}'} v)_{\mathcal{H}} = 0 \quad \text{for } \boldsymbol{j} \ne \boldsymbol{j}' \text{ and any } u, v \in \mathcal{H}.$$

As a consequence, the sparse grid combination technique is exact, see, e.g., [5]. Let $\boldsymbol{P}_{\boldsymbol{j}} := \sum_{\boldsymbol{\ell} \le \boldsymbol{j}} \boldsymbol{Q}_{\boldsymbol{\ell}}$ be the tensor product projection and define the *combination technique index set* $\mathcal{J}_J^{\boldsymbol{w}} := \{ \boldsymbol{j} \in \mathbb{N}_0^m : J - |\boldsymbol{w}| < \boldsymbol{j}^\mathsf{T} \boldsymbol{w} \le J \}$. Then, there holds

$$\widehat{\boldsymbol{P}}_J^{\boldsymbol{w}} f = \sum_{\boldsymbol{j} \in \mathcal{J}_J^{\boldsymbol{w}}} c_{\boldsymbol{j}}^{\boldsymbol{w}} \boldsymbol{P}_{\boldsymbol{j}} f, \quad \text{where } c_{\boldsymbol{j}}^{\boldsymbol{w}} := \sum_{\substack{\boldsymbol{j}' \in \{0,1\}^m \\ (\boldsymbol{j} + \boldsymbol{j}')^\mathsf{T} \boldsymbol{w} \le J}} (-1)^{|\boldsymbol{j}'|}.$$

As a consequence, we only need to solve the smaller tensor product problems

$$\boldsymbol{K_j\alpha_j} := \big(\boldsymbol{K}_{j_1}^{(1)} \otimes \cdots \otimes \boldsymbol{K}_{j_m}^{(m)}\big)\boldsymbol{\alpha_j} = \boldsymbol{f_j}, \quad \boldsymbol{j} \in \mathcal{J}_J^{\boldsymbol{w}}.$$

For computational efficiency, the unidirectional kernel matrices are compressed using *samplets*, see [3]. Samplets are discrete localized signed measures with vanishing moments. More precisely, let $\mathcal{X}_0' \subset \mathcal{X}_1' \subset \cdots \subset \mathcal{X}_J' := \mathcal{X}' := \mathrm{span}\{\delta_x : x \in X\}$ denote a multiresolution analysis. We equip $\mathcal{X}$ with the topology of $\mathbb{R}^{|X|}$ and, since $\mathcal{X}_j' \subset \mathcal{X}_{j+1}'$, we can orthogonally decompose $\mathcal{X}_{j+1}' = \mathcal{X}_j'\oplus\mathcal{S}_j'$. For $\mathcal{S}_j'$, we introduce the orthonormal bases $\{\sigma_{j,k}\}_k$. Recursively applying the decomposition yields a *samplet basis* for $\mathcal{X}_J'$. For data compression, we may construct samplets with vanishing moments $(\sigma_{j,k}, p)_\Omega = 0$ for all polynomials $p$ of degree smaller than or equal to $q$.

Since $\mathcal{X}_J' \subset \mathcal{H}'$, samplets induce a multiresolution basis for $\mathcal{H}_X$ via the embedding $\sigma_{j,k} = \sum_\ell u_{j,k,\ell}\delta_{x_\ell} \mapsto \psi_{j,k} = \sum_\ell u_{j,k,\ell}\kappa(x_\ell,\cdot)$. The kernel matrix in samplet coordinates satisfies $[\langle\psi_{j,k},\psi_{j',k'}\rangle_\mathcal{H}]_{j,j',k,k'} = \boldsymbol{TKT}^\mathsf{T}$, $\boldsymbol{T} := [u_{j,k,\ell}]_{(j,k),\ell} \in \mathbb{R}^{|X|\times|X|}$. The vanishing moment property can be employed for the compression of kernel matrices, given that the kernel is *asymptotically smooth* according to

$$\frac{\partial^{|\boldsymbol{\alpha}|+|\boldsymbol{\beta}|}}{\partial x^{\boldsymbol{\alpha}}\partial y^{\boldsymbol{\beta}}}\kappa(x,y) \le c_\kappa \frac{(|\boldsymbol{\alpha}|+|\boldsymbol{\beta}|)!}{\rho^{|\boldsymbol{\alpha}|+|\boldsymbol{\beta}|}\|x-y\|_2^{|\boldsymbol{\alpha}|+|\boldsymbol{\beta}|}}, \quad c_\kappa, \rho > 0.$$

Let the data set $X$ be quasi-uniform and let $\mathcal{T}$ be a hierarchical cluster tree for $X$. Then, setting all coefficients of the transformed kernel matrix $\boldsymbol{K}^\Sigma := \boldsymbol{TKT}^\mathsf{T}$ to zero which satisfy $\mathrm{dist}(\tau,\tau') \ge \eta\max\{\mathrm{diam}(\tau),\mathrm{diam}(\tau')\}$, $\eta > 0$, $\tau,\tau' \in \mathcal{T}$ results in a consistency error of $\|\boldsymbol{K}^\Sigma - \boldsymbol{K}_\eta^\Sigma\|_F/\|\boldsymbol{K}^\Sigma\|_F \lesssim (\eta\rho/d)^{-2(q+1)}$. The compressed matrix $\boldsymbol{K}_\eta^\Sigma$ only contains $\mathcal{O}(|X|\log|X|)$ entries and can be computed with cost $\mathcal{O}(|X|\log|X|)$, see [3] and the references therein. Afterwards, the linear system can efficiently be solved using a sparse direct solver, see [4].
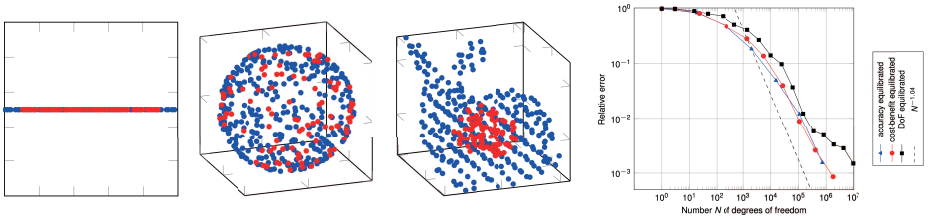


FIGURE 1. Data sites in $1+2+3$ dimensions and convergence of the different sparse grid approximations.

To illustrate the approach, we consider an example in $1+2+3$ dimensions using nested subsets of uniformly random points in $[0,1]$ and $\mathbb{S}^2$ and nested subsets from a volume mesh of the Stanford bunny, see the left hand side of Figure 1. The kernel on each region is given by the Matérn-$\big(\frac{25}{16} - \frac{d}{2}\big)$ kernel, $d = 1,2,3$. Therefore, the energy space is $H^{25/16}(\Omega_i)$ and the expected rate of convergence, when the error is measured in $L^2$, is $25/8$ for $m = 1,2,3$. As data, we consider $f \equiv 1$, which is

not contained in the ansatz space, and evaluate the sparse grid interpolant at 100 random points in each region, see again the left hand side of Figure 1. For the three discussed sparse grid constructions, the error almost perfectly decays with the expected rate $\beta = 25/24 \approx 1.04$, see the right hand side of Figure 1.

## References

[1] M. Griebel, H. Harbrecht, and M. Multerer. Kernel interpolation on sparse grids. *arXiv:2505.12282*, 2025.

[2] M. Griebel and H. Harbrecht. A note on the construction of $L$-fold sparse tensor product spaces. *Constr. Approx.*, 38(2):235–251, 2013.

[3] H. Harbrecht and M. Multerer. Samplets: Wavelet concepts for scattered data. In R. DeVore and A. Kunoth, editors, *Multiscale, Nonlinear and Adaptive Approximation II*, pages 299–326. Springer, Berlin-Heidelberg, 2024.

[4] H. Harbrecht, M. Multerer, O. Schenk, and C. Schwab. Multiresolution kernel matrix algebra. *Numer. Math.*, 156:1085–1114, 2024.

[5] H. Harbrecht, M. Peters, and M. Siebenmorgen. Combination technique based $k$-th moment analysis of elliptic problems with random diffusion. *J. Comput. Phys.*, 252:128–141, 2013.

[6] R. Schaback. Superconvergence of kernel-based interpolation. *J. Approx. Theory*, 235:1–19, 2018.

[7] I. Sloan and V. Kaarnioja. Doubling the rate: Improved error bounds for orthogonal projection with application to numerical analysis. *BIT Numer. Math.*, 65:10, 2025.

## Scalable Sequential Exponential Design for Bayesian Inverse Problems

Robert Scheichl

(joint work with Karina Koval, Roland Herzog, Tiangang Cui)

We propose a novel approach for sequential optimal experimental design (SOED) for Bayesian inverse problems involving expensive models with large-dimensional unknown parameters. The focus of this work is on designs that maximize the expected information gain (EIG) from prior to posterior which is a computationally challenging task in the non-Gaussian setting. This challenge is amplified in SOED, as the incremented expected information gain (iEIG) must be approximal multiple times in distinct stages, with both prior and posterior distributions often being intractable. To address this, we derive a derivative-based upper bound for the iEIG, which not only guides design placement but also enables the construction of projections onto likelihood-informed subspaces, facilitating parameter dimension rediction. By combining this approach with conditional measure transport maps for the sequence of posterior, we develop a unified framework for SOED, together with amortized inference, scalable to high- and infinite-dimensional problems. Numerical experiments for two inverse problems governed by partial differential equations (PDEs) demonstrate the effectiveness of designs that maximize our proposed upper bound.

# On the Lipschitz constant of random ReLU neural networks

## Felix Voigtlaender

(joint work with Sjoerd Dirksen, Patrick Finke, Paul Geuchen, Dominik Stöger)

Deep Learning, the application of machine learning techniques based on deep neural networks, has over the last decade lead to breakthrough results in diverse areas, including image classification [8] and natural language processing [10]. Up to now, however, there is no comprehensive theoretical explanation for this success, and there are still many limitations remaining [11]. Particularly, it has been empirically shown that trained neural networks are often susceptible to so-called *adversarial perturbations* [7, 5], where a small, often imperceptible, perturbation to the input can cause a significant change to the network output (leading e.g. to a misclassification).

In my talk at the MFO, I presented recent results from the paper [4], which makes a first step towards understanding this non-robustness of neural networks, by proving sharp bounds for the *Lipschitz constants* of neural networks with randomly chosen weights and biases. The paper concentrates on networks with the so-called *ReLU activation function* [8], given by $\varrho(x) = x_+ = \max\{0, x\}$. A ReLU neural network with $d$-dimensional input and $L$ hidden layers, consisting of $N$ neurons each, is given by

$$\Phi : \mathbb{R}^d \to \mathbb{R}, \quad \Phi = V^{(L+1)} \circ (\varrho \circ V^{(L)}) \circ \cdots \circ (\varrho \circ V^{(1)}),$$

where the ReLU is applied componentwise to vectors, and where the affine-linear maps $V^{(\ell)}z = W^{(\ell)}x + b^{(\ell)}$ are determined by the weight matrices $W^{(\ell)}$ and the bias vectors $b^{(\ell)}$, which we consider to be chosen at random, via a variant of the *He initialization* proposed in [8]. Specifically, we assume that all the entries of all the $W^{(\ell)}, b^{(\ell)}$ are jointly independent and normally distributed, with

$$W^{(1)} \in \mathbb{R}^{N \times d}, \quad W_{i,j}^{(1)} \sim \mathcal{N}(0, \tfrac{2}{N}),$$
$$W^{(\ell)} \in \mathbb{R}^{N \times N}, \quad W_{i,j}^{(\ell)} \sim \mathcal{N}(0, \tfrac{2}{N}), \quad 2 \leq \ell \leq L$$
$$W^{(L+1)} \in \mathbb{R}^{1 \times N}, \quad W_{i,j}^{(N+1)} \sim \mathcal{N}(0, 1).$$

This choice ensures — in the *zero bias case*, i.e., if all bias vectors $b^{(\ell)}$ are chosen as zero — that the network is isometric in expectation, meaning

$$\mathbb{E}\big[\|\Phi(x)\|_{\ell^2}^2\big] = \|x\|_{\ell^2}^2 \qquad \forall\, x \in \mathbb{R}^d,$$

and in fact this holds not only for the full network $\Phi$, but for the output of each layer of the network.

Studying such random neural networks is mainly important since they are used as the starting point for the training process, and because there are results showing that the weights don't move much during training under certain assumptions [1], so that one might hope that the final trained network behaves similarly to the randomly initialized network, at least in some respects. Moreover, it has been empirically shown that randomly initialized networks can perform well in many tasks, as long as the last layer is trained [9]. Finally, with respect to the *relative*

performance of different network architectures (number and size of the layers, etc.), such "random, untrained networks" tend to behave similarly to the fully trained networks [9].

The following results, which characterize the $\ell^p$ Lipschitz constant

$$\mathrm{Lip}_{\ell^p}(\Phi) := \sup_{x,y \in \mathbb{R}^d, x \neq y} \frac{|\Phi(x) - \Phi(y)|}{\|x - y\|_{\ell^p}}$$

of a random ReLU network $\Phi$, are informal versions of the fully precise results in [4]. We here focus on the *zero-bias case*, in which all bias vectors $b^{(\ell)}$ vanish.

The first result characterizes the $\ell^q$-norm of the gradient $\nabla\Phi(x_0)$ at a single point $x_0 \neq 0$. It essentially shows that this quantity behaves like the $\ell^q$-norm of a $d$-dimensional random vector with independent, standard normal entries. The proof heavily relies on a "decoupling technique" developed in [2, Lemma 2.1].

**Theorem** ([4]). Consider the zero-bias case, and let $x_0 \in \mathbb{R}^d \neq \{0\}$ be fixed. Suppose that $d \gg 1$ and $N \gg L^3$. Then $\Phi$ is almost surely differentiable at $x_0$, and for a suitable absolute constant $c > 0$, the following hold:

- For $q \leq 2$, we have

$$\mathbb{P}\Big( \|\nabla\Phi(x_0)\|_{\ell^q} \asymp d^{1/q} \Big) \geq 1 - \exp\big( -c \min\{d, N/L^2\} \big).$$

- For $q \geq \ln(d)$, we have

$$\mathbb{P}\Big( \|\nabla\Phi(x_0)\|_{\ell^q} \asymp \sqrt{\ln(d)} \Big) \geq 1 - \exp\big( -c \min\{\ln(d), N/L^2\} \big).$$

If $p, q \in [1, \infty]$ are conjugate, then $\|\nabla\Phi(x_0)\|_{\ell^q}$ is a lower bound for $\mathrm{Lip}_{\ell^p}(\Phi)$. Thus, the above theorem immediately implies the following corollary.

**Corollary** ([4]). Under the same assumptions as in the above theorem, we have

$$\mathrm{Lip}_{\ell^p}(\Phi) \geq \|\nabla\Phi(x_0)\|_{\ell^q} \gtrsim d^{1 - \frac{1}{p}} \quad \text{for all } p \geq 2, \quad \text{with high probability.}$$

This lower bound for the Lipschitz constant for the case $p \geq 2$ in fact turns out to be sharp up to log factors.

**Theorem** ([4]). Consider the zero-bias case, and let $\frac{N}{\ln^2(1+N)} \gg d \cdot L^3$. Then, for any $p \geq 2$ and a suitable absolute constant $c > 0$,

$$\mathbb{P}\Big( \mathrm{Lip}_{\ell^p}(\Phi) \lesssim d^{1 - \frac{1}{p}} \cdot \sqrt{\ln(N/d)} \Big) \geq 1 - \exp\big( -c \cdot d \cdot \ln(N/d) \big).$$

The theorem is derived by first establishing the case $p = 2$, and then applying the estimate $\| \cdot \|_{\ell^2} \leq d^{\frac{1}{2} - \frac{1}{p}} \cdot \| \cdot \|_{\ell^p}$, to extend the bound to $p \geq 2$. In the complementary regime $p \leq 2$, using the trivial estimate $\| \cdot \|_{\ell^2} \leq \| \cdot \|_{\ell^p}$, one obtains the upper bound $\mathrm{Lip}_{\ell^p}(\Phi) \lesssim \sqrt{d} \cdot \sqrt{\ln(N/d)}$ with high probability. Due to using the rather naive bound $\| \cdot \|_{\ell^2} \leq \| \cdot \|_{\ell^p}$, one might think that the resulting bound is probably loose. However, as the following result shows, this is in fact not the case, at least if one considers the depth $L$ to be fixed.

**Theorem** ([4]). Assume that $d \gg 1$ and $\frac{N}{\ln^2(N+1)} \gg d \cdot (CL)^{CL}$, and consider the zero-bias case. Then, for any $p \geq 1$, we have

$$\text{Lip}_{\ell^p}(\Phi) \geq \text{Lip}_{\ell^1}(\Phi) \gtrsim \sqrt{d}/\sqrt{L} \quad \text{with probability at least } 1 - \exp\left(-c \cdot \frac{d}{L}\right).$$

A somewhat surprising consequence of these bounds is that in the regime $p \geq 2$, it holds that $\text{Lip}_{\ell^p}(\Phi) \asymp \|\nabla\Phi(x_0)\|_{\ell^q}$, up to a log factor in $N/d$. In contrast, for the case $p < 2$ and $d/(L\ln(d)) \gg 1$, it holds that $\text{Lip}_{\ell^p}(\Phi) \gg \|\nabla\Phi(x_0)\|_{\ell^q}$. Intuitively, the reason for this is that while each individual gradient has relatively small ($\ell^q$) norm with high probability, due to the large number of possible gradients, the maximal norm over all possible gradients is much larger.

**Previous work.** In addition to the zero-bias case considered above, [4] also derives similar results for the case that the distribution of the biases is symmetric and "sufficiently well-behaved". Similar bounds to the ones discussed above, but only for the case $p = 2$ and for the zero-bias setting already appeared in [2]. Moreover, the paper [6] shows that for the case of *shallow* ReLU networks (that is, $L = 1$), the factor $\sqrt{\ln(N/d)}$ in the upper bound can be omitted for the case $p = 2$. This is extended in [4] to obtain matching upper and lower bounds for the shallow case, for arbitrary biases and arbitrary $p \in [1, \infty]$.

## REFERENCES

[1] Z. Allen-Zhu, Y. Li, Z. Song, *A convergence theory for deep learning via over-parameterization*, International conference on machine learning (ICLR), 2019.

[2] P. Bartlett, S. Bubeck, Y. Cherapanamjeri, *Adversarial examples in multi-layer random ReLU networks*, Advances in Neural Information Processing Systems **34** (2021).

[3] S. Buchanan, D. Gilboa, J. Wright, *Deep Networks and the Multiple Manifold Problem*, International Conference on Learning Representations (ICLR), 2021, arXiv:2008.11245.

[4] S. Dirksen, P. Finke, P. Geuchen, D. Stöger, F. Voigtlaender, *Near-optimal estimates for the $\ell^p$-Lipschitz constants of deep random ReLU neural networks*, arXiv preprint, arXiv:2506.19695.

[5] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, D. Song, *Robust Physical-World Attacks on Deep Learning Visual Classification*, Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), (2018).

[6] P. Geuchen, D. Stöger, T. Telaar, F. Voigtlaender, *Upper and lower bounds for the Lipschitz constant of random neural networks*, Information and Inference: A Journal of the IMA **14** (2025), iaaf009.

[7] I.J. Goodfellow, J. Shlens, C. Szegedy, *Explaining and harnessing adversarial examples*, International Conference on Learning Representations (ICLR), 2015, arXiv:1412.6572.

[8] K. He, X. Zhang, S. Ren, J. Sun, *Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification*, Proceedings of the IEEE international conference on computer vision (2015).

[9] A. M. Saxe, P. W. Koh, Z. Chen, M. Bhand, B. Suresh, and A. Y. Ng, *On random weights and unsupervised feature learning*, Proceedings of the 28th International Conference on Machine Learning (ICML-11).

[10] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, *Attention is all you need*, Advances in neural information processing systems **30** (2017), arXiv:1706.03762.

[11] C. Zhang, S. Bengio, M. Hardt, B. Recht, O. Vinyals, *Understanding deep learning (still) requires rethinking generalization*, Communications of the ACM **64** (2021).

## Stochastic gradient with least-squares control variates for smooth stochastic optimization problems

Fabio Nobile

(joint work with Matteo Raviola, Nathan Schaeffer)

The stochastic gradient descent (SGD) method is a widely used approach for solving stochastic optimization problems, but its convergence is typically slow. Existing variance reduction techniques, such as SAGA, improve convergence by leveraging stored gradient information; however, they are restricted to settings where the objective functional is a finite sum, and their performance degrades when the number of terms in the sum is large. In this work, we propose a novel approach which is well suited when the objective is given by an expectation over random variables with a continuous probability distribution. Our method constructs a control variate by fitting a linear model to past gradient evaluations using weighted discrete least-squares, effectively reducing variance while preserving computational efficiency. We establish theoretical sublinear convergence guarantees and demonstrate the method's effectiveness through numerical experiments on random PDE-constrained optimization problems.

References

[1] F. Nobile, M Raviola, N. Schaeffer, *Stochastic gradient with least-squares control variates*, (2025) available as arXiv:2507.20981.
[2] M.C. Martin; F. Nobile *PDE-Constrained Optimal Control Problems with Uncertain Parameters using SAGA*, SIAM/ASA J. Uncertainty Quantification **9(3)** (2021), 979–1012. DOI : 10.1137/18M1224076.
[3] A. Defazio, F. Bach, S. Lacoste-Julien, *Saga: A fast incremental gradient method with support for non-strongly convex composite objectives*, In *Advances in Neural Information Processing Systems*, eds Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence and K. Q. Weinberger, **27** (2014), 1646–1654.

## Computational Math with neural networks is hard

Michael Feischl

(joint work with Fabian Zehetgruber)

In [6], we show that under some widely believed assumptions, there are no higher-order algorithms for basic tasks in computational mathematics such as: Computing integrals with neural network integrands, computing solutions of a Poisson equation with neural network source term, and computing the matrix-vector product with a neural network encoded matrix. We show that this is already true for very simple feed-forward networks with at least three hidden layers, bounded weights, bounded realization, and sparse connectivity, even if the algorithms are allowed to access the weights of the network. We demonstrate sharpness of our results by providing fast quadrature algorithms for one-layer networks and giving numerical evidence that quasi-Monte Carlo methods achieve the best possible order of convergence for quadrature with neural networks.

Neural networks are excellent surrogates for (high-dimensional) functions and perform at least as good as virtually all currently used specialized (high-dimensional) approximation methods such as polynomials, rational approximation, sparse-grids, tensor networks, . . . . Prominent examples of these qualities are given in, e.g., [7]. Strong results are also available for more involved applications such as solving partial differential equations [3, 11] and inverse problems [1]. This even includes problems that are hard for classical approximation methods, such as high-dimensional problems, fractals or stochastic processes.

Thus, two natural questions arise: First, can we efficiently find those networks (for a recent approach to tackle this question, see [5]), and second, if we found them, can we efficiently do computations with them. After all, computing a surrogate is usually done with the intention of using it in another algorithm. In this work, we consider the latter question and derive the following result: Under the widely believed *Strong Exponential Time Hypothesis* (SETH), we show for three fundamental tasks from computational mathematics, that even with full knowledge of the neural network representation of the surrogate (including the weights), no higher-order algorithms exist for the tasks. We give a short overview of these tasks in the following.

*Quadrature:* We particularly see quadrature in the sense

$$\Phi \mapsto \int_\Omega \mathcal{R}_\Phi(x)\, dx$$

for a neural network $\Phi$, its realization function $\mathcal{R}_\Phi$, and a given domain $\Omega \subseteq \mathbb{R}^d$ as a fundamental task. This algorithm is used as a basic building block in countless algorithms, and even in the training of neural networks itself. E.g., for the training of PINNs [4], one usually has to approximate an integral type norm in order to evaluate the loss function, for Variational Monte Carlo (see, e.g., [2]) the same is true for a scalar product.

Thus, the interesting question is whether there exist higher-order quadrature algorithms that do not impose smoothness on the neural network. We show that this is not the case, at least under the assumption of the SETH. We demonstrate experimentally, that quasi-Monte Carlo methods achieve the best possible order of convergence, even for non-smooth neural network integrands.

*Solving PDEs:* A similar question arises in the approximation of PDE solutions. It is well-known that smooth maps can be approximated very well with neural networks, which is the foundation of many operator learning approaches. We refer to the overview articles [9] and the references therein and to [8] for expression rate bounds.

However, we show that smoothness is really fundamental here. Even for the much simpler linear problem of computing the solution $u_f$ of $-\Delta u_f = f$ with Dirichlet boundary conditions, we show that no higher-order algorithms exist if the right-hand side is represented by a neural network. This means that no algorithm can efficiently approximate the map

$$\Phi \mapsto u_\Phi \quad \text{with } -\Delta u_\Phi = \mathcal{R}_\Phi \text{ and } u_\Phi = 0 \text{ on } \partial\Omega.$$

*Matrix-vector multiplication:* Non-linear representations of high-dimensional objects have gained significant interest particularly in the context of low-rank tensor representations, see, e.g., [10] are encoded in tensor formats. While these formats come with very efficient arithmetic, we show that similarly efficient algorithms cannot exist for objects that are encoded with neural networks. To that end, we consider large matrices as one of the simplest objects that can be used to store high-dimensional data. Concretely, we consider matrices $M_\Phi \in \mathbb{R}^{2^d \times 2^d}$ defined by

$$(M_\Phi)_{ij} := \mathcal{R}_\Phi(b(i)_1, \ldots, b(i)_d, b(j)_1, \ldots, b(j)_d),$$

where $b(i)$ is the binary representation of $i$. We show that even simple matrix-vector products with such matrices cannot be computed with higher-order accuracy.

## References

[1] Jens Berg and Kaj Nyström. Neural networks as smooth priors for inverse problems for pdes. *Journal of Computational Mathematics and Data Science*, 1:100008, 2021.

[2] Giuseppe Carleo and Matthias Troyer. Solving the quantum many-body problem with artificial neural networks. *Science*, 355(6325):602–606, 2017.

[3] Weinan E and Bing Yu. The deep Ritz method: a deep learning-based numerical algorithm for solving variational problems. *Commun. Math. Stat.*, 6(1):1–12, 2018.

[4] Weinan E and Bing Yu. The Deep Ritz Method: A Deep Learning-Based Numerical Algorithm for Solving Variational Problems. *Communications in Mathematics and Statistics*, 6(1):1–12, March 2018.

[5] Michael Feischl, Alexander Rieder, and Fabian Zehetgruber. Towards optimal hierarchical training of neural networks. *arXiv preprint arXiv:2407.02242*, 2024.

[6] Michael Feischl and Fabian Zehetgruber. Computational mathematics with neural networks is hard, 2025. *arXiv preprint arXiv:2505.17751*, 2025.

[7] Philipp Grohs, Fabian Hornung, Arnulf Jentzen, and Philippe von Wurstemberger. A proof that artificial neural networks overcome the curse of dimensionality in the numerical approximation of Black-Scholes partial differential equations. *Mem. Amer. Math. Soc.*, 284(1410):v+93, 2023.

[8] Lukas Herrmann, Christoph Schwab, and Jakob Zech. Neural and spectral operator surrogates: unified construction and expression rate bounds. *Adv. Comput. Math.*, 50(4):Paper No. 72, 43, 2024.

[9] Nikola B. Kovachki, Samuel Lanthaler, and Andrew M. Stuart. Operator learning: Algorithms and analysis, 2024.

[10] I. V. Oseledets. Approximation of $2^d \times 2^d$ matrices using tensor decomposition. *SIAM J. Matrix Anal. Appl.*, 31(4):2130–2145, 2009/10.

[11] Yaohua Zang, Gang Bao, Xiaojing Ye, and Haomin Zhou. Weak adversarial networks for high-dimensional partial differential equations. *J. Comput. Phys.*, 411:109409, 14, 2020.

*Reporter: Janina Tikko*

# Participants

**Prof. Dr. Markus Bachmayr**
Institut für Geometrie und Praktische
Mathematik
RWTH Aachen
Templergraben 55
52062 Aachen
GERMANY


**Prof. Dr. Peter Binev**
Department of Mathematics
University of South Carolina
Columbia, SC 29208
UNITED STATES


**Daan Bon**
Dept. of Mathematics & Computer
Science
Eindhoven University of Technology
5600 MB Eindhoven
NETHERLANDS


**Max Brockmann**
Mathematisches Institut
Universität zu Köln
Weyertal 86-90
50931 Köln
GERMANY


**Prof. Dr. Albert Cohen**
Laboratoire Jacques-Louis Lions
Sorbonne Université
4, Place Jussieu
75005 Paris Cedex
FRANCE


**Prof. Dr. Wolfgang Dahmen**
Department of Mathematics
University of South Carolina
1523 Greene Street
Columbia, SC 29208
UNITED STATES


**Dr. Matthieu Dolbeault**
Institut für Geometrie und Praktische
Mathematik
RWTH Aachen
Templergraben 55
52062 Aachen
GERMANY


**Dr. Geneviève Dusson**
Faculté des Sciences et Techniques
Laboratoire Mathématiques de Besancon
Université de Franche-Comte
16, route de Gray
25030 Besançon Cedex
FRANCE


**Prof. Dr. Virginie Ehrlacher**
CERMICS - ENPC
Bât. Coriolis B 312
Cité Descartes, Champs-sur-Marne
6 et 8 Avenue Blaise Pascal
77455 Marne-la-Vallée Cedex 2
FRANCE


**Henrik Eisenmann**
Institut für Geometrie und
Praktische Mathematik
RWTH Aachen
Templergraben 55
52062 Aachen
GERMANY


**Manfred Faldum**
Institut für Mathematik
RWTH Aachen
Templergraben 55
52056 Aachen
GERMANY

**Prof. Dr. Michael Feischl**
Institute for Analysis and Scientific
Computing
TU Wien
Wiedner Hauptstraße 8-10
1040 Wien
AUSTRIA

**Isabella Carla Gonnella**
SISSA
International School for Advanced
Studies
Via Bonomea 265
34136 Trieste
ITALY

**Margarita Granzow**
Mathematisches Institut
Universität zu Köln
50931 Köln
GERMANY

**Prof. Dr. Lars Grasedyck**
Institut für Geometrie und
Praktische Mathematik
RWTH Aachen
Templergraben 55
52062 Aachen
GERMANY

**Prof. Dr. Philipp Grohs**
Fakultät für Mathematik
Universität Wien
Oskar Morgenstern Platz 1
1090 Wien
AUSTRIA

**Diane Guignard**
Department of Mathematics & Statistics
University of Ottawa
Ottawa ON K1N 6N5
CANADA

**Prof. Dr. Helmut Harbrecht**
Departement Mathematik und
Informatik
Universität Basel
Spiegelgasse 1
4051 Basel
SWITZERLAND

**Prof. Dr. Angela Kunoth**
Department Mathematik/Informatik
Universität zu Köln
Weyertal 86-90
50931 Köln
GERMANY

**Prof. Dr. Frances Y. Kuo**
School of Mathematics and Statistics
The University of New South Wales
Sydney NSW 2052
AUSTRALIA

**Prof. Dr. Michael Lindsey**
Department of Mathematics
University of California, Berkeley
970 Evans Hall
Berkeley CA 94720-3840
UNITED STATES

**Prof. Dr. Hrushikesh N. Mhaskar**
Institute of Mathematical Sciences
Claremont Graduate University
1232 N. Dartmouth Avenue
Claremont, CA 91711
UNITED STATES

**Prof. Dr. Olga Mula**
TU Eindhoven
P.O. Box 513
5600 Eindhoven
NETHERLANDS

**Prof. Dr. Michael Multerer**
Istituto Eulero
Università della Svizzera Italiana
Via la Santa 1
6962 Lugano
SWITZERLAND

**Marcel Neugebauer**
Mathematisches Institut
Universität zu Köln
50923 Köln
GERMANY

**Prof. Dr. Richard Nickl**
Centre for Mathematical Sciences
Wilberforce Road
Cambridge CB3 0WB
UNITED KINGDOM

**Evie Nielen**
Dept. of Mathematics & Computer
Science
Eindhoven University of Technology
5600 MB Eindhoven
NETHERLANDS

**Prof. Dr. Fabio Nobile**
CSQI - MATH
École Polytechnique Fédérale de
Lausanne
P.O. Box Station 8
1015 Lausanne
SWITZERLAND

**Prof. Dr. Anthony Nouy**
Centrale Nantes, Nantes Université
1, rue de la Noe
P.O. Box 92101
44321 Nantes Cedex 3
FRANCE

**Prof. Dr. Dirk Nuyens**
Departement Computerwetenschappen
KU Leuven
Celestijnenlaan 200A
3001 Leuven
BELGIUM

**Prof. Dr. Christoph Ortner**
Department of Mathematics
University of British Columbia
121-1984 Mathematics Road
Vancouver BC V6T 1Z2
CANADA

**Dr. Mathias Oster**
Institut für Geometrie und Praktische
Mathematik
RWTH Aachen
Templergraben 55
52062 Aachen
GERMANY

**Prof. Dr. Peter Oswald**
Institut für Numerische Simulation
Universität Bonn
Friedrich-Hirzebruch-Allee 7
53115 Bonn
GERMANY

**Prof. Dr. Mark A. Peletier**
Department of Mathematics and
Computer Science
Eindhoven University of Technology
MetaForum 5.062
P.O. Box 513
5600 MB Eindhoven
NETHERLANDS

**Maximilian Penka**
Department of Mathematics
TUM School of Computation,
Information and Technology
Technische Universität München
Boltzmannstraße 3
85748 Garching bei München
GERMANY

**Dr. Clarice Poon**
Mathematics institute, university of
Warwick
Coventry Cv47ay
UNITED KINGDOM

**Prof. Dr. Robert Scheichl**
Institut für Mathematik
Universität Heidelberg
Im Neuenheimer Feld 205
69120 Heidelberg
GERMANY

**Prof. Dr. Johannes
Schmidt-Hieber**
Department of Applied Mathematics
University of Twente
Drienerlolaan 5
7522 NB Enschede
NETHERLANDS

**Prof. Dr. Reinhold Schneider**
Fakultät II - Institut für Mathematik
Technische Universität Berlin
Sekr. MA 5 - 3
Straße des 17. Juni 136
10623 Berlin
GERMANY

**Prof. Dr. Christoph Schwab**
Seminar für Angewandte Mathematik
ETH Zurich
ETH Zentrum, HG G 57.1
Rämistrasse 101
8092 Zürich
SWITZERLAND

**Prof. Dr. Jonathan W. Siegel**
Department of Mathematics
Texas A & M University
College Station, TX 77843-3368
UNITED STATES

**Agustin Somacal**
Ecole Centrale Nantes
1, rue de la Noe
P.O. Box 92101
44321 Nantes Cedex 3
FRANCE

**Dr. Aretha Teckentrup**
School of Mathematics
University of Edinburgh
James Clerk Maxwell Building
Edinburgh EH9 3FD
UNITED KINGDOM

**Janina Tikko**
Mathematisches Institut
Universität zu Köln
50931 Köln
GERMANY

**Prof. Dr. Elisabeth Ullmann**
Department Mathematik
Technische Universität München
Boltzmannstraße 3
85748 Garching bei München
GERMANY

**Dr. Mario Ullrich**
Institut für Analysis
Johannes-Kepler-Universität
Altenberger Strasse 69
4040 Linz
AUSTRIA

**Prof. Dr. Karsten Urban**
Institut für Numerische Mathematik
Universität Ulm
Helmholtzstraße 20
89081 Ulm
GERMANY

**Prof. Dr. Felix Voigtlaender**
Katholische Universität
Eichstätt-Ingolstadt
Mathematisches Institut für maschinelles
Lernen und Data Science (MIDS)
Auf der Schanz 49
85049 Ingolstadt
GERMANY

**Tolunay Yilmaz**
Fakultät für Mathematik
Universität Wien
Oskar-Morgenstern-Platz 1
1090 Wien
AUSTRIA