# ANNUAL REPORT 2024

# *Preface*



We hope you enjoy reading!

## Your Members of ZIB

# EXECUTIVE SUMMARY

2024 marked a year of dynamic advancement and distinguished recognition for the Zuse Institute Berlin (ZIB), underscoring its central mission to drive innovation in Application-Oriented Mathematics, High-Performance Computing (HPC), and Artificial Intelligence. Amid continuous global shifts and emerging technological landscapes, ZIB demonstrated remarkable resilience, scientific depth, and strategic foresight.

At the institutional level, ZIB successfully addressed an important gap in its scientific leadership. In its November 2024 meeting, ZIB's Board of Directors decided to entrust Prof. Andrea Walther (Humboldt-Universität zu Berlin) with the responsibilities of an additional Vice President at ZIB – an important step that strengthens ZIB's leadership in key areas of research and strategic development. Andrea will head the "Parallel and Distributed Computing" branch of ZIB including the NHR Center.

A major focus of the year was the successful preparation and submission of proposals for the continuation of key third-party funded projects that are central to ZIB's institutional framework. These include the Cluster of Excellence MATH+, the Research Campus MODAL, and the National High-Performance Computing Center (NHR) at ZIB. The effort required extensive reporting and comprehensive coordination with strategic partners and reflects ZIB's long-term commitment to research excellence and infrastructural development in these cornerstone initiatives. The first success in this challenging and labor-intensive process was the approval of a five-year extension for MODAL in December 2024.

A significant institutional milestone was celebrated with the 40th anniversary of ZIB. This landmark occasion brought together esteemed figures from science, politics, and industry, with Berlin's Senator for Higher Education and Research, Dr. Ina Czyborra, delivering a keynote at the celebratory ceremony. The associated scientific conference showcased ZIB's role as a nexus for forward-looking research and interdisciplinary collaboration.

ZIB was also proud to celebrate several prestigious awards bestowed upon its researchers, underlining the scientific excellence cultivated at the institute. Among them, Prof. Sebastian Pokutta and colleagues received the renowned Gödel Prize for their groundbreaking contributions in theoretical computer science. Prof. Hans-Christian Hege was honored with the IEEE Career Award, recognizing a lifetime of visionary contributions to visual computing. Moreover, the next generation of researchers at ZIB was celebrated through Stephanie Riedmüller's first-place achievement in the Heureka Student Award. This was complemented by ZIB academic placement seeing many of its students and researchers take up high-profile positions in academia and industry.

Significant strides were made in the expansion of ZIB's high-performance computing and data infrastructure. Major efforts were invested in upgrading and modernizing the computing and storage systems supporting the Tier-2 systems of the NHR Center and ZIB's own Tier-3 HPC and AI infrastructure. These enhancements not only ensure future scalability and reliability but also strengthen ZIB's ability to meet the growing demands of data-intensive and AI-driven re-

search. With a focus on sustainable architecture and interoperability, the upgrades support the continued evolution of ZIB as a national hub for high-end computing services. Exemplifying the success of these efforts, the NHR Center's Lise system earned international acclaim – securing third place worldwide in the IO500 competition at the ISC High Performance Conference 2024.

ZIB's commitment to Berlin as a vibrant hub for research and innovation remains steadfast. Through close partnerships with local universities, research institutions, and public stakeholders, ZIB actively contributes to the scientific and technological landscape of the region. A strategic roadmap, currently under development, outlines ZIB's future priorities in sustainable HPC, AI, and digital research infrastructure rooted in scientific excellence – further anchoring its role as a critical enabler for cutting-edge computational research and innovation in the Berlin-Brandenburg metropolitan area and beyond.

Scientific progress across a broad spectrum of application areas continued at a high pace. In this year's report, ten brief research highlights exemplify the broad impact of ZIB's method development and collaborations. In the realm of neuroscience, the article *"Neuro-Connectomics | Understanding the function of the brain"* explores the mapping of brain connectivity. Advances in AI are featured in *"Geometric deep learning | Pioneering new frontiers in AI"* and *"Neural networks guide mathematical discovery"*, both pushing the boundaries of intelligent systems. Complex systems are addressed in *"Network dynamics | Collective variables: The key to understanding dynamics on net-*works"*, which emphasizes effective reduction strategies. Model selection, reduction learning and efficiency take center stage in *"Active learning of surrogates"*, while real-world applications of AI are illustrated in *"Estimating canopy height at scale"* and *"Towards AI-accelerated molecular structure prediction pipelines on supercomputers"*. Energy and mobility challenges are tackled through *"Energy system transition | Fast, reliable and multi-perspective mathematical solution strategies"* and *"The art of charging electric buses"*. Finally, the importance of reproducible and robust research workflows is underscored in *"Improving data analysis workflows"*. Together, these highlights showcase ZIB's diverse and interdisciplinary impact across scientific domains.

These achievements are further contextualized by two in-depth interviews: one with Tim Conrad on data lakes and the Mathematical Research Data Initiative (MaRDI), and another with Carsten Schäuble on strengthening ZIB's HPC and AI infrastructure.

As reflected in this report, 2024 has been a year in which ZIB not only strengthened its foundational programs but also elevated its role as a thought leader in computational science. Looking ahead, ZIB remains dedicated to pushing the frontiers of interdisciplinary research, fostering innovation, and providing impactful services to the scientific community in Berlin and beyond. ⌣

<div align="center">

Christof Schütte      Sebastian Pokutta
(President)      (Vice President)
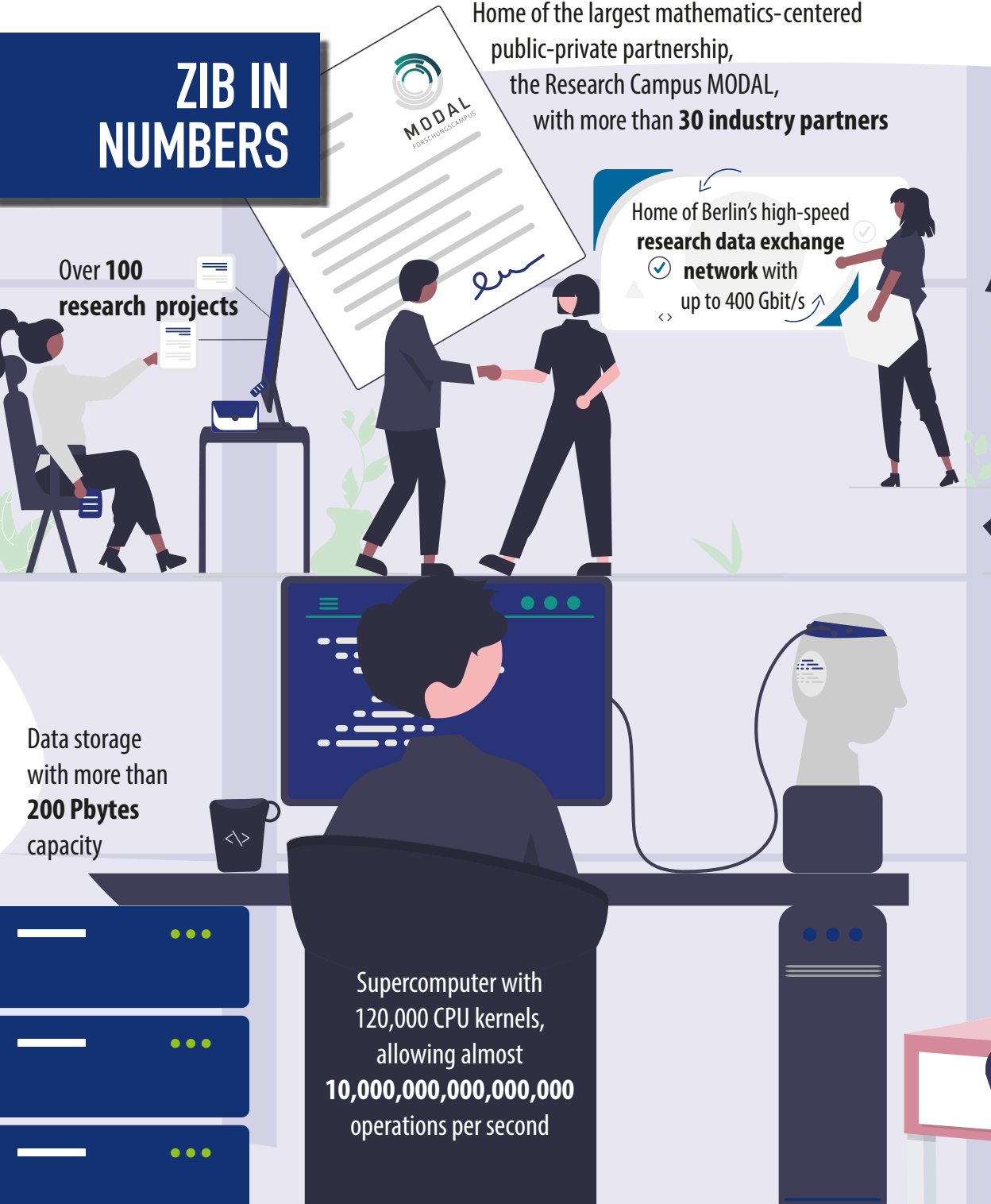
</div>

# ZIB IN NUMBERS

Home of the largest mathematics-centered public-private partnership, the Research Campus MODAL, with more than **30 industry partners**

MODAL
FORSCHUNGSCAMPUS

Home of Berlin's high-speed **research data exchange network** with up to 400 Gbit/s

Over **100 research projects**

Data storage with more than **200 Pbytes** capacity

Supercomputer with 120,000 CPU kernels, allowing almost **10,000,000,000,000,000** operations per second

Core institution of the Cluster of Excellence

Berlin Mathematics Research Center

MATH+

More than **5,000 visitors** per year

Advanced **AI computing infrastructure**

More than **300 researchers** and research service staff

# CONTENTS

Reconstruction of five cortical columns from the barrel cortex of a rat. This brain area is associated with the processing of tactile sensory information from whiskers at the rat's snout. The colors represent the different types of morphological cell of which the columns are composed.

# NEURO–
# CONNECTOMICS

# Understanding the function of the brain

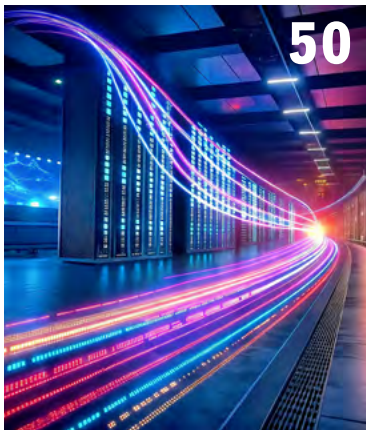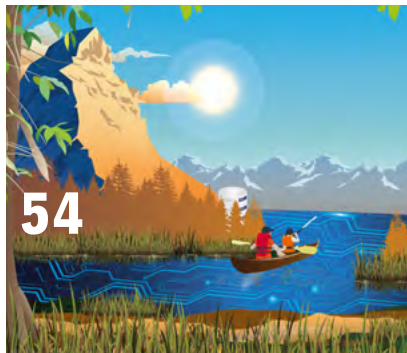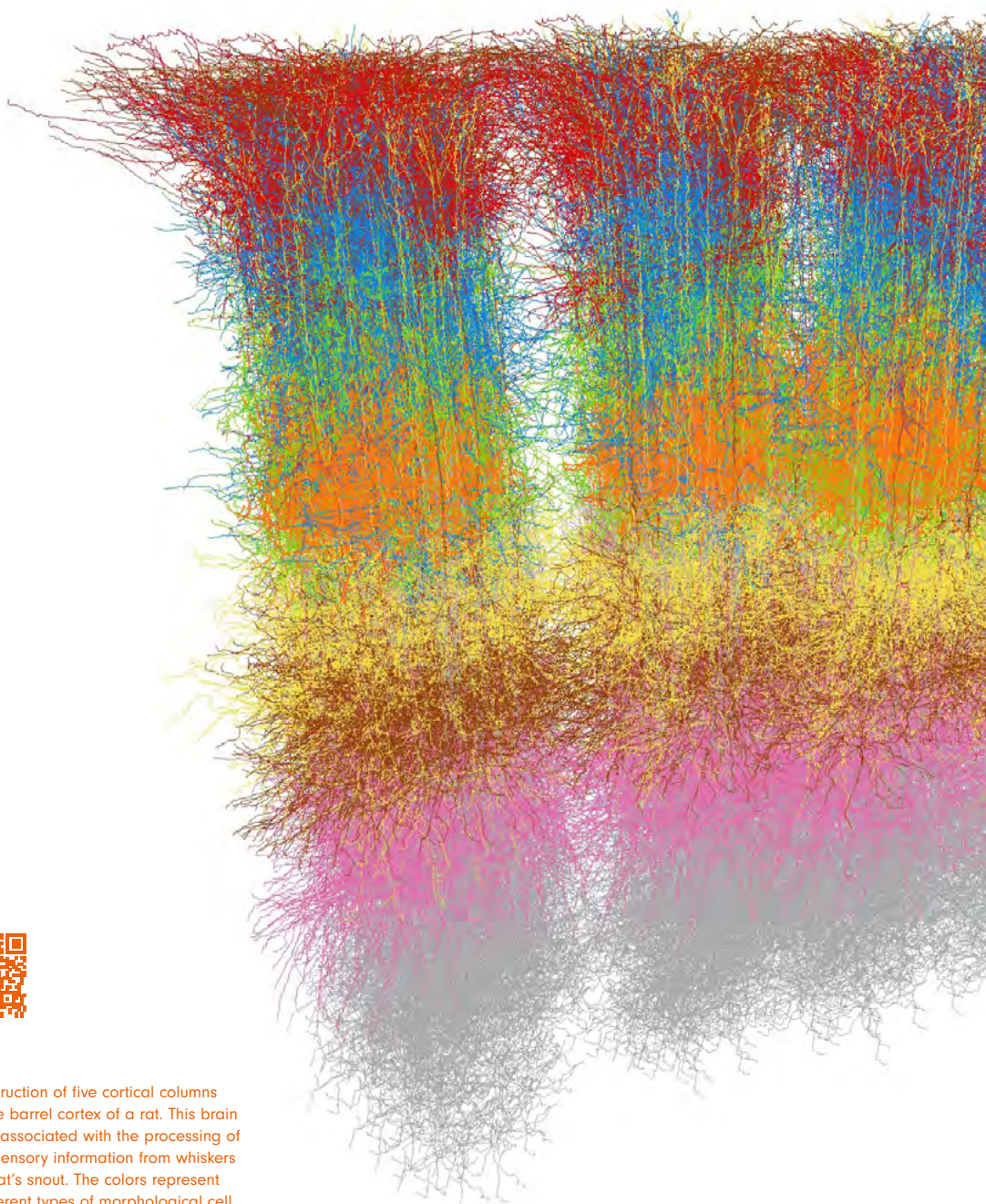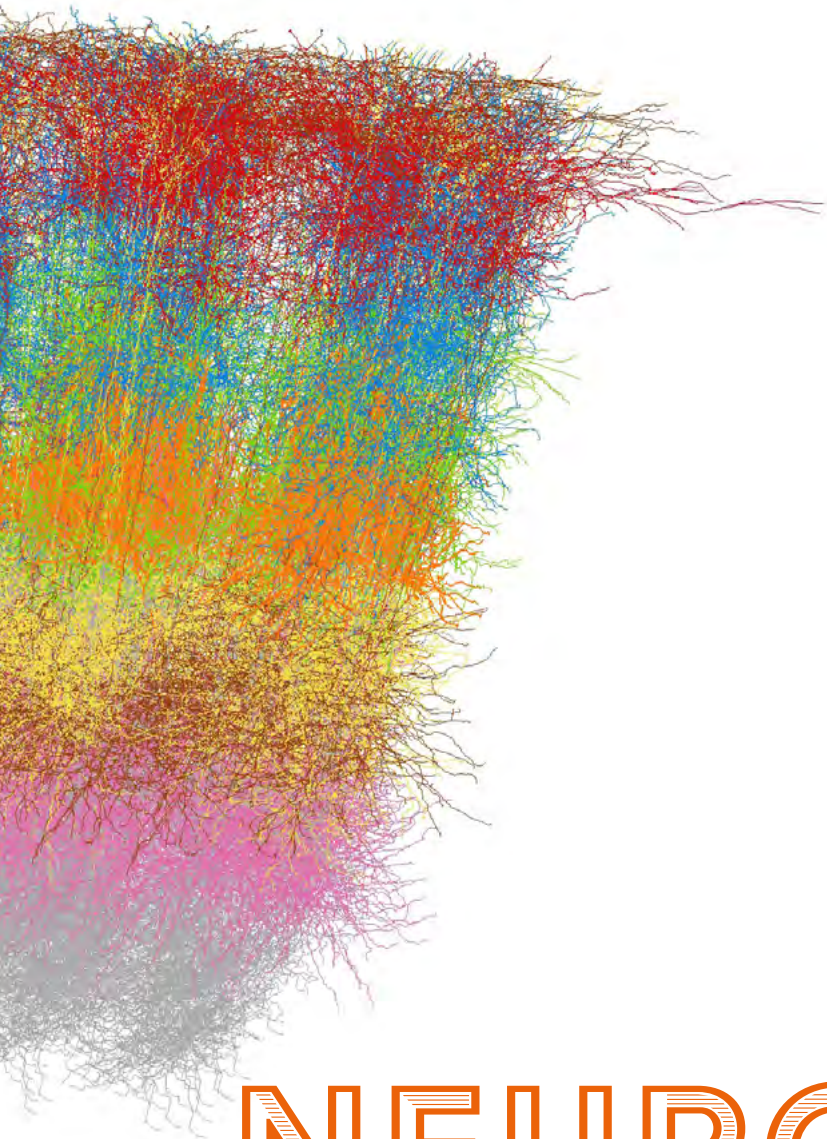Basic neurobiological research aims to investigate the structural composition, function, and development of biological neural networks (including the human nervous system). This research lays the foundation for a better understanding of the development of neurological diseases, such as Alzheimer's disease, and enables the development of targeted medications and therapies for mental disorders and neurological diseases.

Neurobiological research, like many of today's research fields, is extremely diverse and conducted on various levels. Although understanding the human brain is a primary goal, the possibilities of directly examining the brain at the cellular level are severely limited due to the non-invasive imaging techniques available. In smaller organisms, on the other hand, it is possible to examine the structure of the brain in much greater detail or even completely at the cellular level by reconstructing individual nerve cells (neurons) and their sy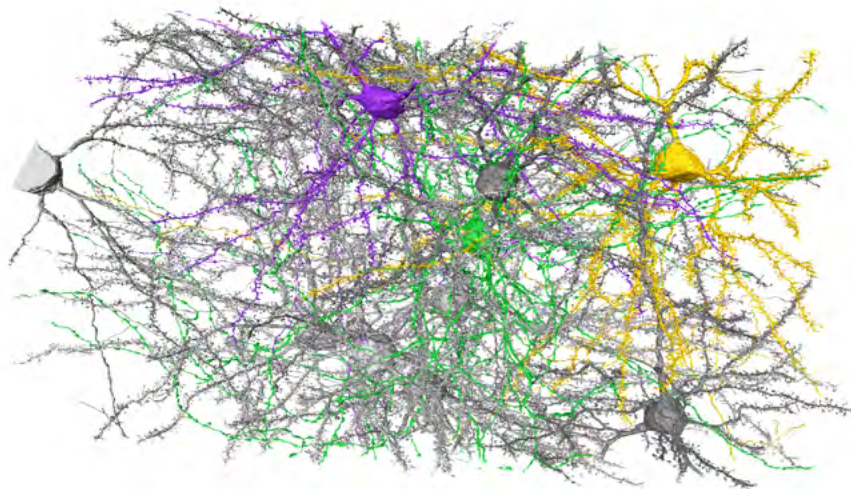naptic connections using modern imaging and mathematical image processing techniques. Basic neurobiological research is therefore usually conducted on various types of animal, such as fruit flies, zebrafish, mice or rats.

At the Zuse Institute Berlin, the "Visual Data Analysis" research group has been conducting intensive research with neurobiologists in various ways for many years. The goal in these projects is to support the work of neurobiologists through state-of-the-art techniques of mathematical image analysis and visualization. The following three projects highlight some aspects of these collaborations.

Within the DFG Research Unit "RobustCircuit", led by Prof. Hiesinger from the Free University of Berlin (FU), neurobiologists from FU and the University of Mainz are investigating the role of noise – in the sense of random fluctuations – in the development and function of the brain on various spatial and functional



Electron microscopy reconstruction of a human cortex tissue sample consisting of around 16,000 neurons. Pyramidal cells are shown in red, shaded by the cortical layer; interneurons are shown in blue. Data: Shapson-Coe et al. 2024; visualization: Amira.

scales. For this purpose, three-dimensional (3D) images and movies are acquired that allow one to see how the brain develops. Our task in this research unit is to develop automated techniques to analyze these data, for example, to extract statistical data about how the neurons grow, how they interact, and what molecular factors influence their growth. By automating these analysis tasks, we enable neurobiologists to analyze an unprecedented amount of data to derive robust statistical conclusions that would otherwise be infeasible.

In cooperation with Prof. Engert from Harvard University (Boston, USA), we are working on the creation of brain atlases and their use for behavioral studies as well as on the comparison of brains from different species and animal groups. For this, we have developed a web-based platform that allows the integration of the original image data, extracted structural data (brain regions and neurons), and functional data (neurons activated in behavioral experiments and gene expression data) into a common coordinate system. As a neurobiological model system, this project uses the zebrafish, whose larva is fully transparent for a few days after hatching and is therefore particularly suitable for capturing functional data. To complement the atlas, we have developed advanced visual analysis tools that enable neurobiologists, based on functional measurements taken during behavioral studies, to identify hypothetical neural circuits that are responsible for performing specific functions in response to certain visual, acoustic or tactile stimuli.

In a third project, in collaboration with Prof. Oberlaender from the MPI for Neurobiology of Behavior (caesar), Bonn, and Prof. Macke from the University of Tübingen, we developed analysis tools for the examination of connectome data, which represent the synaptic connections between neurons. Using simulation-based Bayesian inference, we developed a generative mathematical model that, with very few parameters, was able to reconstruct the structural properties of various connectomes. With this innovative technique, we can now for the first time generate network models directly from dense electron microscopy data that capture features of empirically observed connectomes from subcellular to network scales, and that, at the same time, enable the derivation of the synaptic specificity parameters that are necessary and sufficient to explain each of these connectivity features. With the help of our model, surprisingly, we identified strong similarities in the synaptic specificity parameters for the mouse visual cortex and human temporal cortex.

Daniel Baum

# GEOMETRIC DEEP LEARNING

# Pioneering new

Over the past decade, artificial intelligence and data science have undergone a transformative shift, largely propelled by advancements in deep learning techniques. High-dimensional learning tasks, such as protein structure prediction and natural language processing, have proven feasible with sufficient data, although the number of data samples needed to achieve reliable results generally increases exponentially with dimension – a phenomenon known as the *curse of dimensionality*. These breakthroughs can be attributed to the fact that most relevant tasks are not generic but have inherent regularities due to the underlying low-dimensionality and structure of the physical world. Revealing such regularities through geometric concepts not only facilitates the study of state-of-the-art architectures but also provides a blueprint for deriving novel ones that adhere to prior (physical) knowledge. In this area of *geometric deep learning*, neural designs targeting graph-structured data represent a prominent example. Unlike images, graphs do not have a regular structure: nodes may have a varying number of neighbors and are not sampled from a regular grid. Furthermore, a key property of graphs is

that the ordering of nodes is usually assumed to be arbitrary, and thus any function acting on graphs should not depend on it. Geometrically speaking, the set of all reorderings, called the permutation group, makes up the symmetries of graphs, and we are interested in neural networks for which the output is either invariant for global, graph-wise prediction or equivariant for node-wise ones, meaning that the ordering in the output is tied to the ordering of the input.

Graph neural networks have found widespread applications and have in particular become the de facto standard for transduction in deep learning. Other than inductive learning, which tries to infer a general model from labeled examples in order to predict labels of unseen ones, transductive approaches learn labels simultaneously on training and test data. Transduction therefore avoids solving a more general problem as an intermediate step and thus faces a potentially simpler problem as compared to inductive learning. At ZIB, we derived novel transductive learning approaches for morphometric grading of disease states. In particular, we conditioned graph convolutional net-

# frontiers in AI

works on a population graph, whose nodes represented subject-specific shapes and whose edge weights encoded similarities between subjects. Based on this graph, the learning task could be formulated as a semi-supervised node classification problem, where labels are only given for nodes corresponding to subjects from the training set. The resulting grading systems could be successfully applied to the detection of Alzheimer's disease from hippocampi shapes and the scoring of osteophytes in knee bones. We further adapted the transductive approach to derive a retrieval system for cultural heritage objects that learns object embeddings such that pair-wise distances encode task-specific similarities.

Existing graph neural networks adhere to the geometry of the input domain, that is, the graph; however, they assume that signals take values in vector spaces. In many cases – for example, when dealing with shape data – the signals belong to non-trivial, geometric spaces on their own, and it is only consequential to ask for models that use the geometry of the signal space as well. We constructed a novel graph convolutional layer based on a manifold-valued graph diffusion equation and derived node-wise multilayer perceptrons, both of which are equivariant not only with respect to node permutations but also isometries of the feature manifold. These filters are not only promising for manifold input signals but also allow for geometric representation learning. In particular, heterogeneous degree distributions and strong clustering in complex networks can often be explained by assuming an underlying hierarchy that is well captured in hyperbolic space. Indeed, our geometric filters enabled novel hyperbolic embedding strategies and learning of hierarchical structures that led to substantial improvements in full-graph classification. We further derived an equivariant temporal convolutional network that enables gesture recognition from skeletal data. These successes provide a strong impetus for further investigation of geometric deep learning techniques and create pathways for exploring their potential for open challenges in artificial intelligence and its applications.
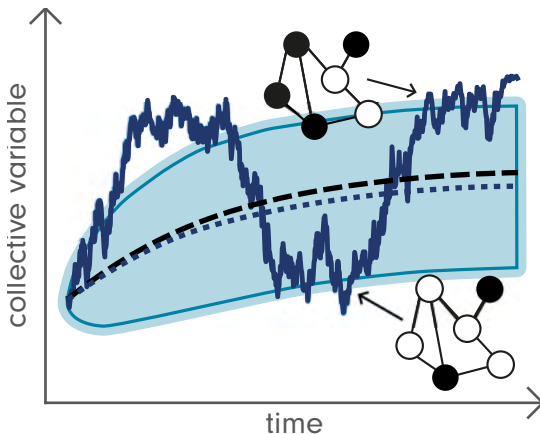
Christoph v. Tycowicz

# Collective variables: The key to understanding dynamics on networks

Networks of interacting agents are widely used to model dynamic social phenomena such as the spreading of a disease, the diffusion of a political opinion within a society, or the adoption of sustainable behaviors and technologies in response to environmental challenges. In such networks, nodes represent individual agents and edges stand for some type of social interaction.

Each node has a state that evolves over time depending on the states of neighboring nodes. Often, stochastic effects are included to account for uncertainty in the dynamics and for the variability of agent behavior. These types of spreading processes lie at the heart of numerous open problems in a wide range of disciplines, such as understanding social collective behavior, assessing systemic risk in financial systems, or controlling modern power grids.

Although individual agent behavior may be intuitive and easy to understand, the emergent macroscopic dynamics resulting from their interactions are often complex and unpredictable. Small changes in the environment or agent behavior can lead to qualitatively different outcomes on a global scale, making it challenging to devise macroscopic models without prior knowledge of the system's behavior. Even with simple interaction rules, analyzing emergent behavior analytically is often difficult. Hence, one usually resorts to numerical simulations for evaluating such systems, but these become computationally intensive or even infeasible for large networks. To balance detail with feasibility, reduced-order models are sought to capture central phenomena while enabling more nuanced insights and better computational efficiency.



**Figure 1: The solution of the mean-field equation (the black dashed line) closely approximates the collective dynamics given by the shares of nodes in certain classes.**

At ZIB, we tackle this challenge by identifying low-dimensional representations that capture the most significant properties of a system's dynamics. We develop constructive methods to find *collective variables* (CVs) that project the system into a lower-dimensional space, filtering out unnecessary detail while retaining essential information about the system's behavior. Good CVs reduce dimensionality while also offering insights into macroscopic dynamics by highlighting the most relevant features.
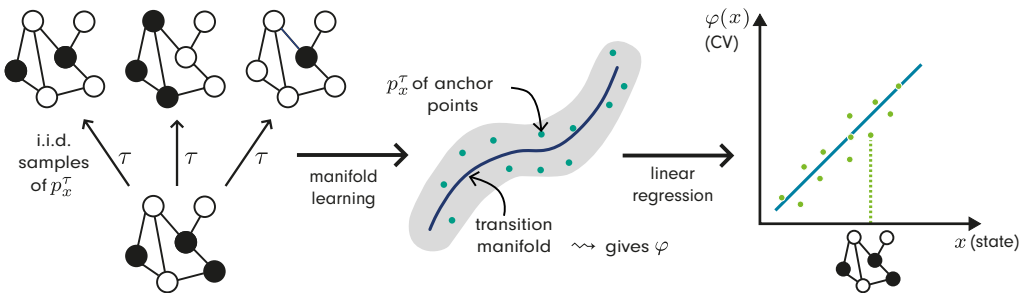
Our research not only enables the identification of collective variables but also facilitates the derivation of a reduced macroscopic model – evolution equations for the macroscopic state defined by the CVs – that accurately approximates the low-dimensional projection of the original dynamics. While the reduced model often remains stochastic, our work demonstrates that, in certain cases, the random actions of many agents in the network effectively cancel out each other, leading to approximately deterministic macroscopic dynamics.

In collaboration with FU Berlin and the Potsdam Institute for Climate Impact Research (PIK), we derived conditions under which Markovian discrete-state systems on networks converge to a mean-field limit. This theory states that, for particular networks, the shares of nodes in certain classes align with the solution of a mean-field ordinary differential equation in the large population limit (Figure 1). Moreover, we proposed a method to algorithmically learn interpretable CVs for spreading processes on networks without requiring prior expert knowledge of a network's topology or dynamics. The approach involves sampling network states, running multiple brief simulations from these states, and extracting optimal CVs by applying manifold learning techniques to the set of transition densities (Figure 2) – a method developed with input from ZIB researchers, known as the transition manifold approach. The learned CVs are then extended to unseen data using total-variation-regularized linear regression. These CVs are interpretable because the inferred parameters reveal the role and importance of specific network features.

The methods and results developed at ZIB have been successfully applied to spreading dynamics across various network types, including Erdős–Rényi random graphs, stochastic block models, random regular graphs, ring-shaped networks, and scale-free networks, demonstrating their flexibility and effectiveness. Our methods enable the exploration of spreading processes on networks, capturing high-impact phenomena like epidemic thresholds and network fragmentation. By incorporating real-world complexities, they offer valuable tools for analyzing and comprehending complex social dynamics.

<div align="right">Marvin Lücke, Stefanie Winkelmann</div>



**Figure 2: Samples of the network dynamics are used to learn a low-dimensional transition manifold. A regression step allows the associated collective variables to be interpreted.**
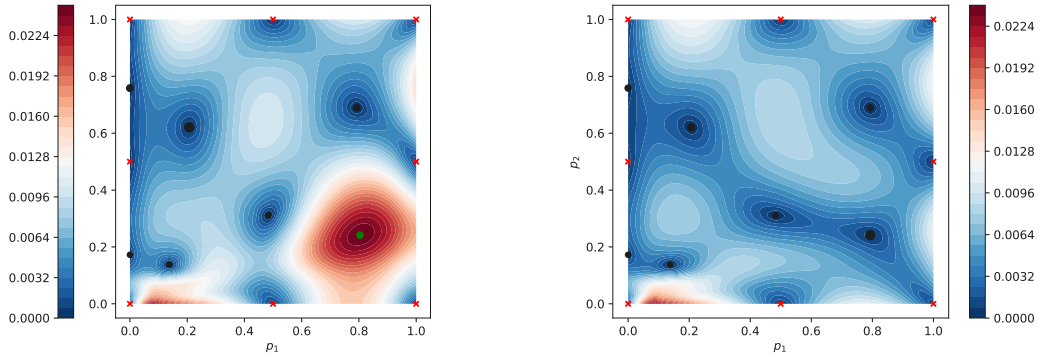
# ACTIVE LEARNING OF SURROGATES

Figure 1: Adaptive design for a two-parameter identification problem on $[0, 1]^2$. The given design (left) is improved to the next one (right) by including an evaluation point where the estimated error (color-coded) is largest. The size of the points shows the simulation accuracy for that evaluation point.

Non-destructive testing and quality control in technology and civil engineering, as well as the continuous updating of digital twins, require the repeated solution of parameter identification problems to estimate the state of the system from measurements. These systems are often described by forward models in the form of parametrized systems of partial differential equations (PDEs). These include Maxwell's equations in optical metrology applications, solid mechanics in bridge surveillance, and the heat equation for time-of-death estimation in forensic medicine.

Estimating the state parameters requires many solutions of these PDE systems in an optimization method for computing maximum posterior point estimates or in Markov-Chain Monte Carlo methods for sampling the posterior probability distribution. Each individual solution may require significant computational resources, often making the parameter estimation procedure too slow for real-time applications.

Replacing the PDE solver with a fast surrogate model that maps parameters directly to the measurable outputs can bypass expensive simulations. Surrogate models using Gaussian process regression, neural networks, sparse grids, or polynomial interpolation are
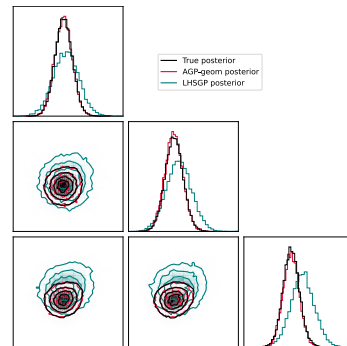


Figure 2: Marginals of a sampled posterior of three parameters, obtained from an adaptively constructed surrogate model (red) and from a surrogate based on a latin hypercube design (green) with comparable computational effort.

trained offline using data consisting of parameter–measurement pairs. Often, a large number of training data points are required, making the construction of surrogate models expensive. This cost is typically reduced by selecting the most informative parameters for acquiring training data. When computing training data via PDE simulations, there is an additional design choice: the accuracy with which to solve the PDE numerically.

# Design of computer experiments

In collaboration with scientists from Friedrich-Alexander-Universität (FAU), Fraunhofer IISB, and Bundesanstalt für Materialforschung und -Prüfung (BAM), and supported by Bundesministerium für Bildung und Forschung (BMBF) and Deutsche Forschungsgemeinschaft (DFG), researchers at ZIB have developed adaptive algorithms for selecting both parameter positions and simulation accuracies for computing training data. The method aims to select a design D (a particular choice of parameter positions and corresponding simulation accuracies) minimizing the error introduced by the surrogate approximation for a given computational budget. It combines a posteriori error estimators E(D) for the surrogate model, as provided explicitly by Gaussian process regression surrogates, with a priori estimates of the computational effort W(D) required for solving the PDE up to the requested accuracies. In a greedy fashion, starting from a given design D, we seek an improved design D′ that can be realized with a small increment ΔW of computational budget – that is, we solve the design of computer experiments problem

$$\min_{D'} E(D') \text{ subject to } W(D') \leqslant W(D) + \Delta W,$$

(see Figure 1).

Depending on how the surrogate model is used – either for computing point estimates of the parameters or for sampling the posterior density – the introduced error propagates into different quantities of interest. Thus, the notion of error depends on the application. We have derived error estimators for two cases: first, for the deviation of reconstructed parameters relative to the uncertainty inherent due to measurement errors, and second, for the Kullback–Leibler divergence between the true posterior density and the posterior obtained when using a surrogate model (see Figure 2).

# Significant performance improvements

Compared to a priori designs such as factorial designs, latin hypercubes or low-discrepancy sequences, and even compared to adaptively selected parameter positions with fixed simulation accuracy, the additional adaptive choice of accuracies significantly improves the efficiency of surrogate model creation. In both simple test examples and actual PDE inverse problems, the required computational effort for training data simulation has been reduced by a factor of ten or more (see Figure 3). This is particularly true if gradients of the forward model can be cheaply evaluated. ⌡
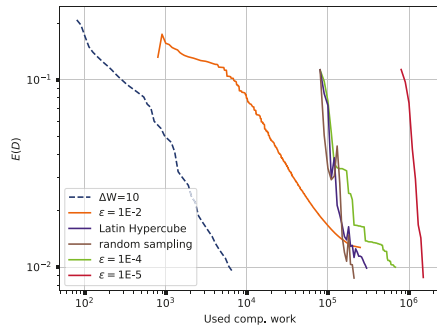
Martin Weiser



Figure 3: Expected parameter identification error over computational work for setting up the surrogate model. The position and accuracy of the adaptive design (ΔW = 10) is far more efficient than using a fixed low precision (ε = 1E − 2), which in addition does not reach the desired error level of $10^{-2}$, or latin hypercube or random designs, or fixed high precision.

# NEURAL NETWORKS
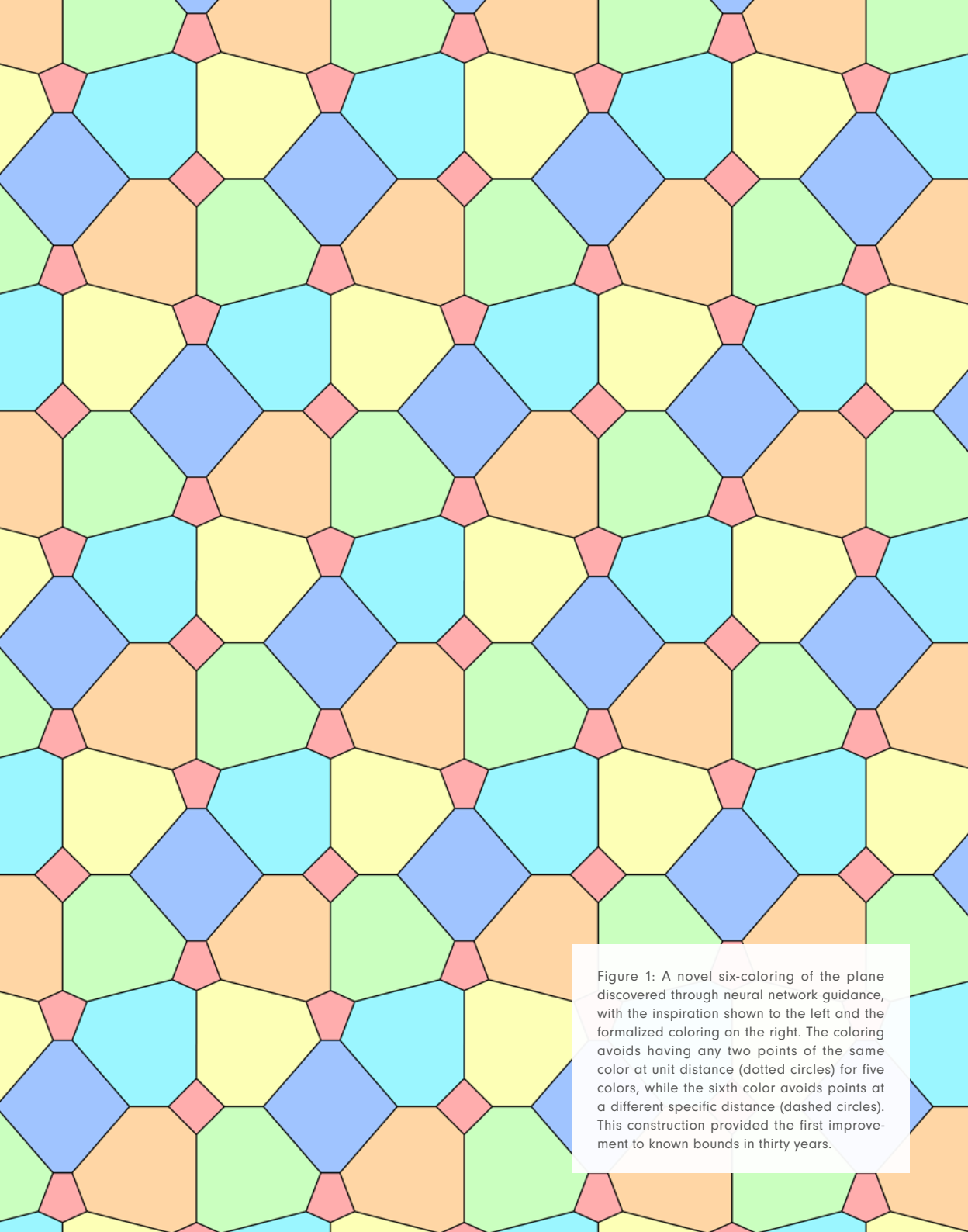
# GUIDE

# MATHEMATICAL

# DISCOVERY

Figure 1: A novel six-coloring of the plane discovered through neural network guidance, with the inspiration shown to the left and the formalized coloring on the right. The coloring avoids having any two points of the same color at unit distance (dotted circles) for five colors, while the sixth color avoids points at a different specific distance (dashed circles). This construction provided the first improvement to known bounds in thirty years.

The interface between theoretical mathematics and practical computation continues to evolve rapidly. While computers have become increasingly sophisticated tools for mathematical exploration and verification, discovering new mathematical constructions remains a delicate interplay between human insight and computational guidance. This is particularly challenging for problems that combine discrete choices with continuous aspects. At the Zuse Institute, we have developed a novel approach using neural networks to guide mathematical intuition in precisely such settings, leading to the first improvement in thirty years for a variant of the famous Hadwiger–Nelson problem in geometric graph theory.
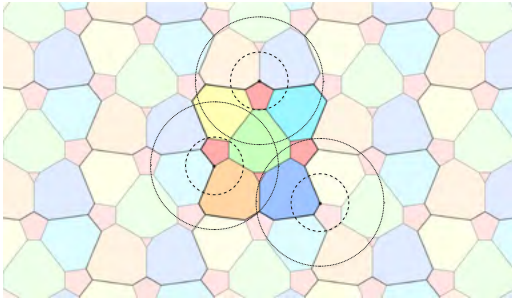


Figure 1: A novel six-coloring of the plane discovered through neural network guidance. The coloring avoids having any two points of the same color at unit distance (dotted circles) for five colors, while the sixth color avoids points at a different specific distance (dashed circles). This construction provided the first improvement to known bounds in thirty years.
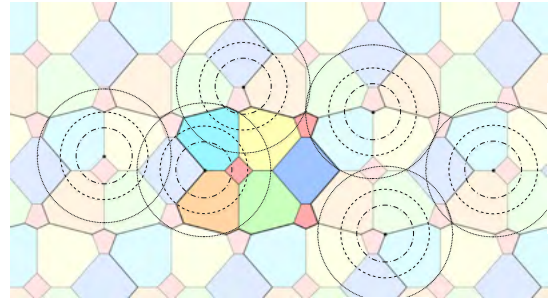


Figure 2: Our second novel six-coloring construction, which extends the range of realizable distances in the sixth color to 0.657. The dotted circles indicate unit distances that must be avoided by five colors, while the dashed circles show the larger distance that must be avoided by the sixth color.
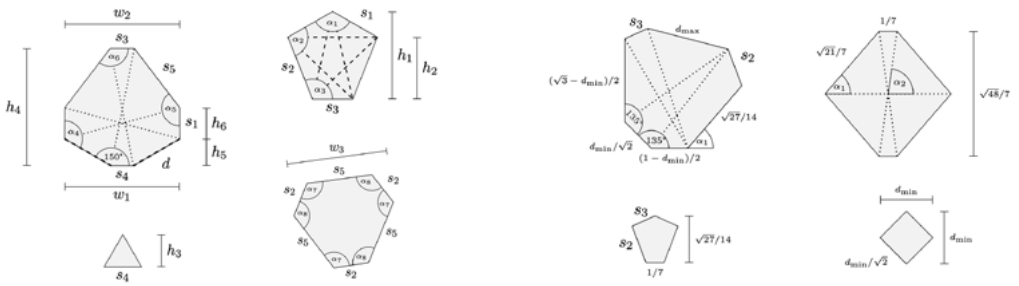


Figure 3: Decomposition of our two novel colorings into their fundamental building blocks. These regular patterns emerged naturally from the neural network's continuous optimization process.

The Hadwiger–Nelson problem, first posed in 1950, asks a deceptively simple question: how many colors are needed to paint a plane so that no two points at a distance of exactly 1 from each other have the same color? More formally, we want to know the chromatic number of the plane $\chi(\mathbb{R}^2)$, that is the smallest integer c for which there exists a function $f : \mathbb{R}^2 \rightarrow \{1, \ldots, c\}$ such that $f(x) \neq f(y)$ for any $x, y \in \mathbb{R}^2$ satisfying $\|x - y\|_2 = 1$.

Despite its elementary statement, this problem has resisted resolution for over 70 years. We know that at least five colors are needed, famously proved by De Grey in 2018, and that seven colors are always sufficient, shown through an explicit construction in 1950. One natural variant asks whether six colors might suffice if we allow the sixth color to avoid a different distance than one.

The key idea was to reformulate this geometric coloring problem as a continuous optimization task that neural networks could tackle effectively. Rather than searching for a discrete coloring function $f$, we trained networks to output probability distributions $p_x \in [0, 1]^c$ over c colors at each point $x$, minimizing the expected occurrence of forbidden patterns through a differentiable loss function $\mathcal{L}(p) = \mathbb{E}_{\|x-y\|=1}[p_x \cdot p_y]$.

The results exceeded our expectations. The neural networks consistently suggested patterns that, after careful mathematical analysis and formalization, led to two novel six-colorings of the plane. These constructions significantly expanded the known range of distances that could be avoided in the sixth color, extending it from the previously known interval of [0.415, 0.447] to [0.354, 0.657]. This represents the first improvement to these bounds since the work of Hoffman and Soifer in 1993.

Beyond these specific results, our work showcases a promising direction for AI-assisted mathematics. Rather than attempting fully automated proofs, we use neural networks to explore possible constructions and suggest patterns that human mathematicians can then analyze rigorously. This computer-guided intuition proved particularly valuable in our case, where the final constructions involved intricate geometric patterns that would have been difficult to discover through traditional methods.

Looking ahead, the framework we developed is not limited to the Hadwiger–Nelson problem. Our approach has already shown promise in exploring related problems, including variants involving monochromatic triangles with prescribed side lengths and colorings of three-dimensional space. More broadly, similar techniques could be applied to other mathematical domains where continuous relaxations are possible, such as using graphons for problems in extremal graph theory. The success of neural networks in suggesting novel mathematical constructions, as demonstrated by recent work across various domains, highlights the growing synergy between machine learning and pure mathematics. This synthesis of computational exploration and mathematical rigor represents a promising direction for tackling other long-standing open problems in mathematics. ⌣

Konrad Mundinger, Max Zimmer,
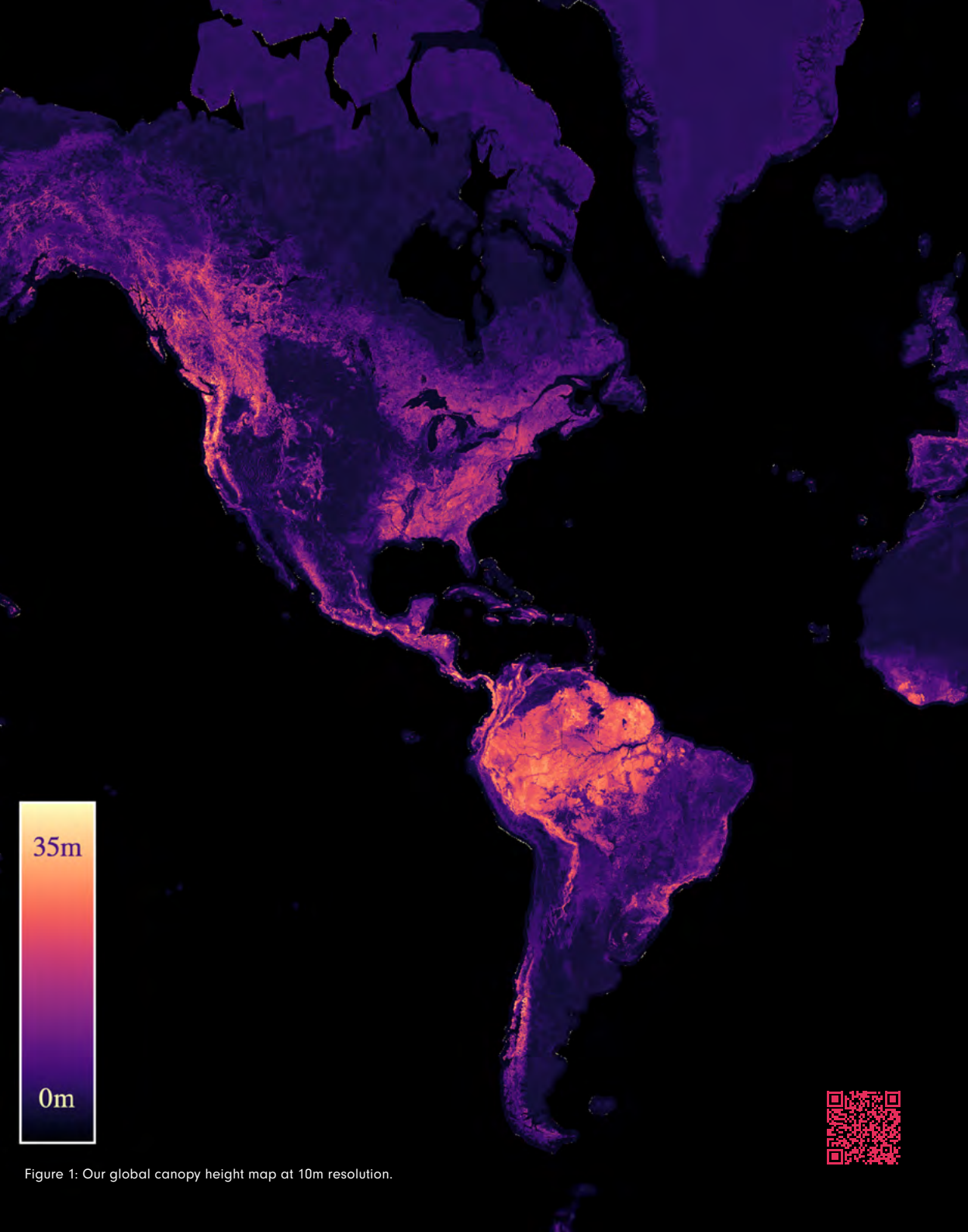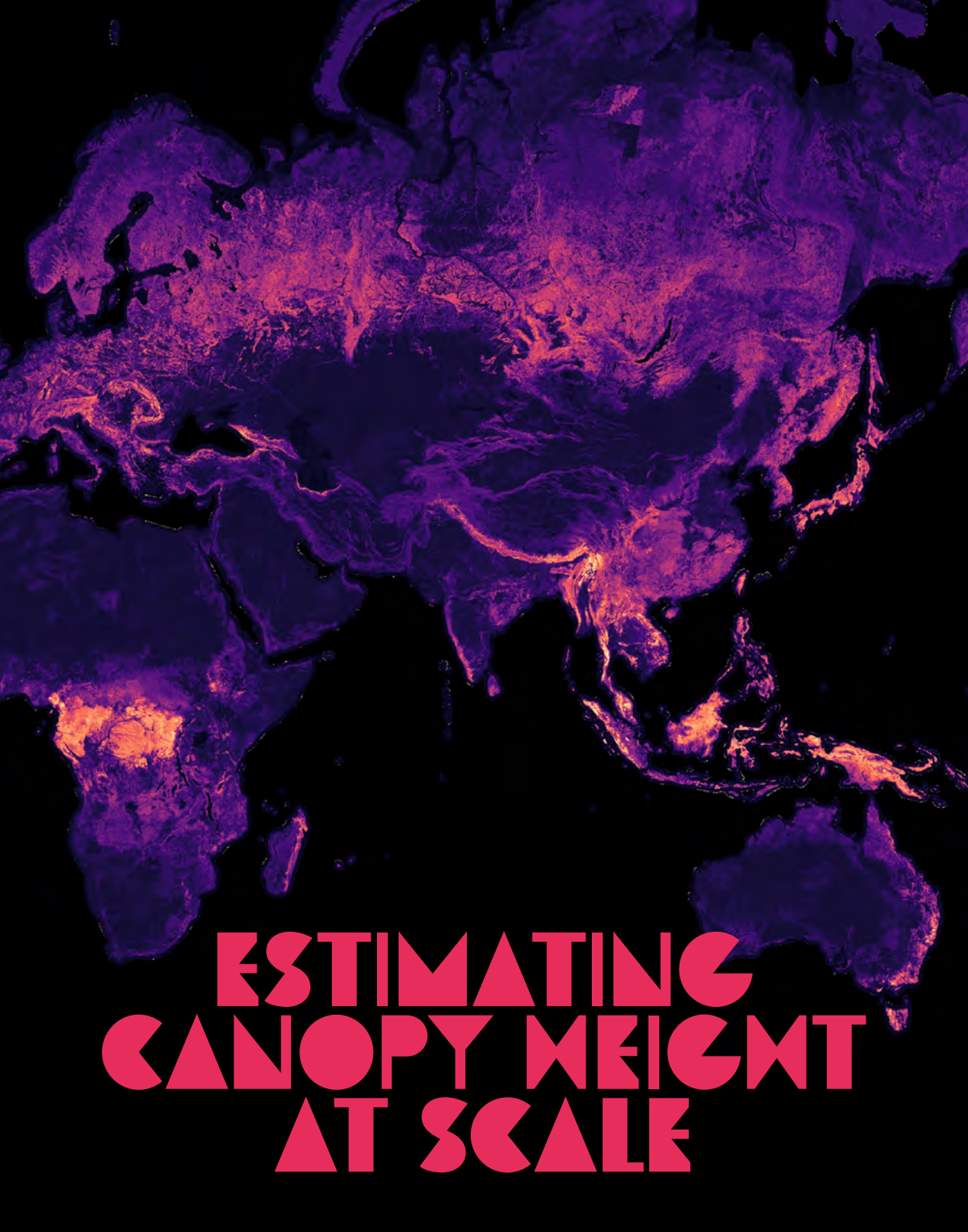Christoph Spiegel, and Sebastian Pokutta

Figure 1: Our global canopy height map at 10m resolution.

# ESTIMATING CANOPY HEIGHT AT SCALE

Forests face increasing threats from climate change and land use changes like deforestation. Effective forest management and conservation are vital for climate adaptation and mitigation, aligning with the UN's Sustainable Development Goals (SDGs), the Bonn Challenge on forests, and the Glasgow Declaration on halting forest loss. Forest conservation efforts, including regeneration, afforestation, and reforestation, are essential for climate mitigation under the Paris Agreement. Precise and up-to-date high-resolution information about the structure, health, and carbon stocks of forests in the world is critical for promoting appropriate measures to tackle forest loss and to effectively combat climate change. However, traditional methods of assessing forest carbon stocks through field inventories are labor-intensive and lack the granularity needed for regional and local monitoring, failing to capture sudden losses from events like fires or insect attacks. Further, many countries also lack comprehensive inventories and use basic methods for estimating forest carbon changes.
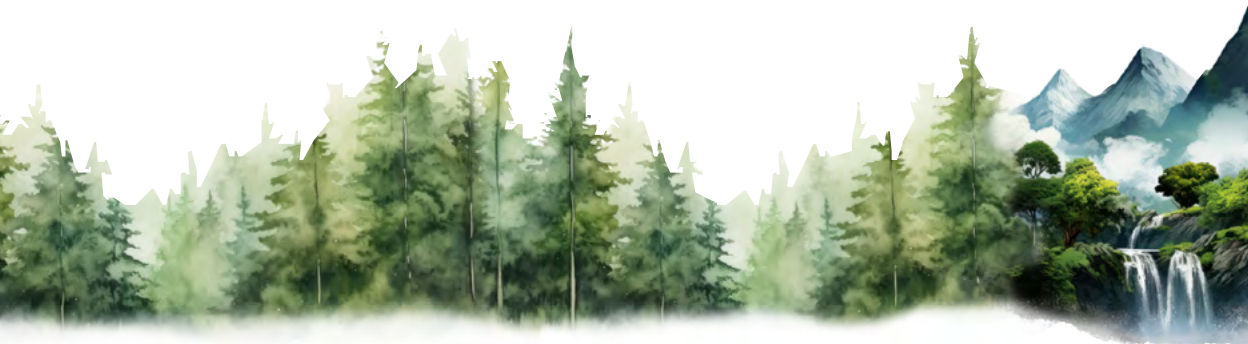
To address this gap, the AI4Forest project, a collaboration between ZIB, LSCE Paris, University of Münster, ENS Paris, and TUM, aims to enhance the accuracy of forest monitoring through deep learning techniques by leveraging publicly available satellite data from optical, radar, and GEDI spaceborne sensors. One important application is the estimation of one of the Essential Climate Variables, namely the above-ground biomass (AGB), using tree canopy height as a proxy. While GEDI provides only point measurements of canopy height, a comprehensive global height map is currently lacking.

We were able to provide such a global canopy height map at 10m resolution by training a neural network that accurately predicts canopy heights from satellite imagery across large areas. Using four channels from Sentinel-1 and ten channels from Sentinel-2, we achieved a mean absolute error of 2.43m in canopy height predictions. Sparse GEDI point measurements served as the ground-truth for our model. Figure 2 gives a schematic overview of our workflow. Our map significantly improves upon existing global-scale maps and can be used as a reference to derive more accurate biomass carbon stocks, helping policymakers to make more accurately informed decisions to guide forest management efforts.

We identified both the potential to further improve the resolution and accuracy of our canopy height map, as well as the limitation of our map in only showing estimates for the year 2020, without capturing temporal dynamics, which would allow for a better monitoring of forest health and disturbances. We have addressed both issues.

Our novel approach, which leverages a full 12-month time series of Sentinel-2 imagery rather than using a single aggregated composite, yields substantial performance gains by allowing the model to capture seasonal patterns and exploit geolocation shifts in satellite measurements. We present the first 10m resolution temporal canopy height map of the European continent for the period 2019–2022. This new map is significantly better than our previous one for 2020
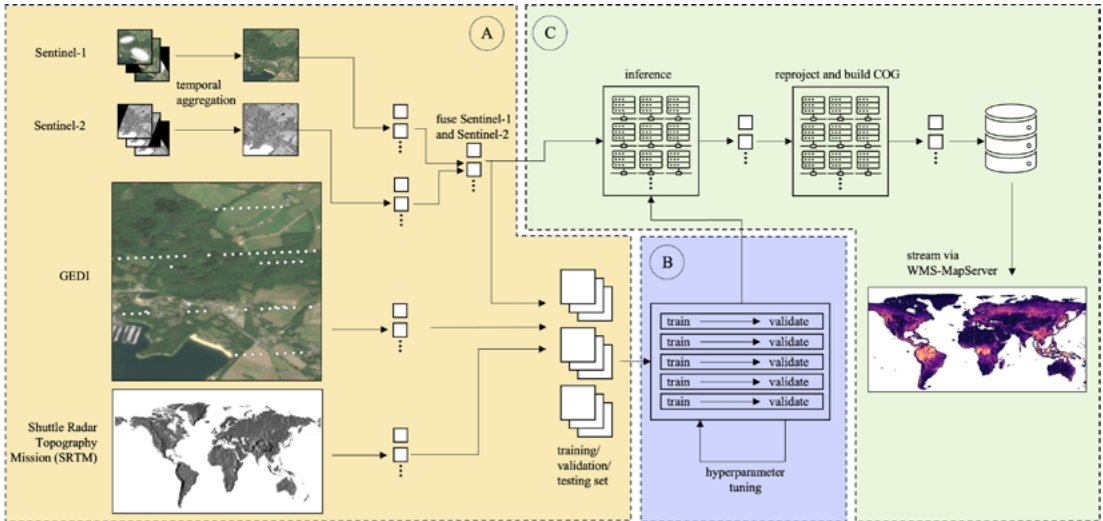
Figure 2: The pipeline for estimating canopy heights from satellite imagery.

and now also captures temporal dynamics, providing more precise estimates than previous studies.

A huge challenge we faced was the substantial amount of training data required. To cover monthly imagery from 2017 to 2023 for approximately 20,000 landmass tiles, we needed around 1200 TB of Sentinel-2 data. Despite these challenges, our pipeline and the resulting temporal height map are publicly available, enabling comprehensive large-scale monitoring of forests and facilitating future research and ecological analyses. A further direction of our work is to extend our map to the entire world. The current temporal map can be explored at https://europetreemap. projects.earthengine.app/view/ temporalcanopyheight.

<div align="right">

Jan Pauls, Max Zimmer,
Berkant Turan,
Sebastian Pokutta et al.

</div>

# FAST, RELIABLE AND MULTI-PERSPECTIVE MATHEMATICAL SOLUTION STRATEGIES

# Energy system transition under uncertainty

Planning a transition to a decarbonized energy system presents notable mathematical challenges. These arise from balancing ecological sustainability and economic feasibility while managing uncertainties in energy demand, weather-dependent renewables, and geopolitical price fluctuations. Despite these uncertainties, reliable investment decisions need to be made, requiring optimization models with millions of variables that must be frequently optimized under slight variations. We collaborate with industry partners and solver companies to address these challenges, integrating real-world insights with state-of-the-art solver technology to develop advanced modeling and optimization techniques.

# Solution strategies with first-order methods

Efficient solutions require leveraging massively parallel high-performance computing (HPC) and GPU technology. Solving energy-related mixed-integer problems (MIPs) is computationally demanding, including already computationally expensive linear programming (LP) relaxations. At ZIB, we push computational boundaries by integrating decomposition techniques with interior-point methods in our solver, PIPS-IPM++, harnessing HPC for high-accuracy LP solutions. For MIPs, first-order methods (FOMs) provide a faster, lower-accuracy alternative by iteratively applying gradient information via matrix products, a process well suited for GPUs. We embed these LP techniques into a heuristic framework called fix-and-propagate (FP), enabling efficient optimization of large-scale energy system models.

# Decision support under uncertainty

Long-term energy investments must consider uncertainties that cannot be fully captured by a single optimization outcome. We integrate high-performance LP solvers into robust optimization frameworks to facilitate sensitivity and robustness analyses. Using a modeling-to-generate-alternatives approach, we provide stakeholders with multiple near-optimal solutions, which offer a broad catalog of investment options close to the economically optimal set. Given the large scale and numerical intricacy of energy system models, efficient re-optimization is crucial. Our warm-starting techniques significantly reduce computation time from several hours on parallel machines to just minutes on a single computer, enhancing decision-making efficiency.
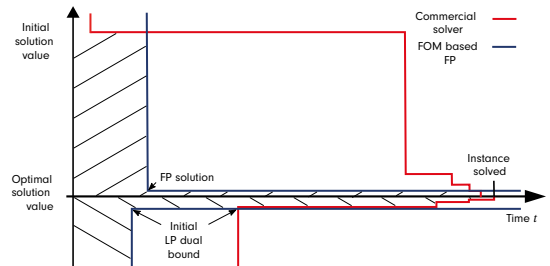


Figure 1: Comparison of the specialized FOM-based FP.

# Conflicting targets in energy system optimization

When dealing with conflicting targets for large MIPs, computing all Pareto-optimal solutions (that is, such that no objective can be improved without worsening another) is often impossible and not beneficial to the decision-maker. Instead, we focus on generating a representative subset that balances solution diversity with computational efficiency while considering both convex and non-convex regions. For example, when optimizing Berlin's district heating network, we investigate variants of the classical ε-constraint algorithm for three objectives ($f_1$, $f_2$, $f_3$). Through lexicographic optimization, we restrict the search region projected onto the $f_2$/$f_3$-hyperplane to a rectangular shape defined by including $f_1$-optimal (usually cost-optimal) solutions. Infeasible subregions can be eliminated early by employing the appropriate order of the subproblems. Depending on the generation of the ε-constraints – grid-wise or dynamically – we can either compute well-distributed or clustered solution sets to produce valuable insights for the decision-maker.



Figure 2: Range of installable capacities within 2% of the optimal system cost for Berlin (blue) and Brandenburg (orange), showing the optimal solution (x) and a selected alternative (x).

# Advancing energy-system planning for the Berlin energy system

At ZIB, we enhance large-scale energy-system planning by leveraging HPC, FOMs, and heuristic frameworks. Validated on the Berlin-Brandenburg multi-sectoral energy system design and the Berlin district heating system, our approach integrates robust optimization in decision-support tools, enabling faster re-optimization and providing stakeholders with diverse, near-optimal solutions for more informed decision-making – even in uncertain times. ⌡

Nils-Christian Kempke, Niels Lindner,
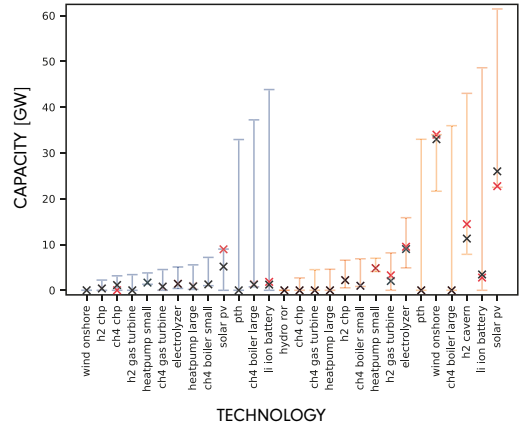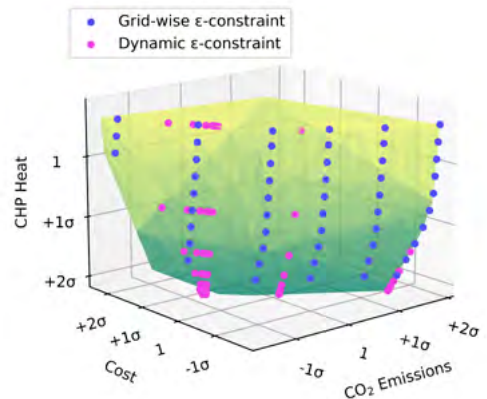Stephanie Riedmüller, Janina Zittel



Figure 3: Comparison of two variants of the ε-constraint algorithm in generating the Pareto front for a tri-objective unit commitment problem, with the feasible region (green) depicted in the near-optimal subregion.

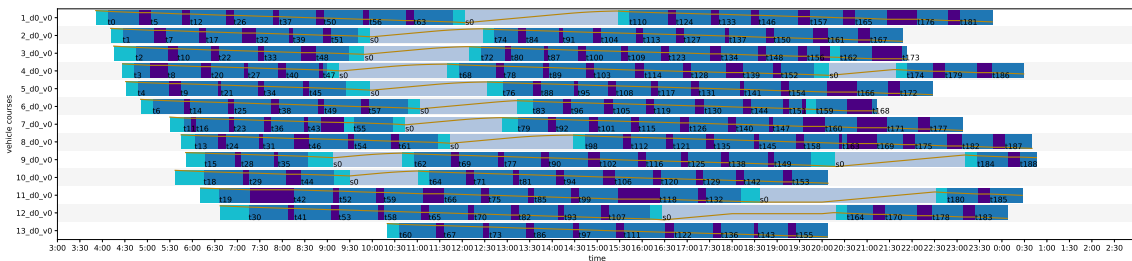# The art of
## charging
# electric buses

Figure 1: Gantt chart of a vehicle schedule (passenger trips in dark blue) with charging events (gray) for a small test scenario.

According to the German Clean Road Vehicle Procurement Act, in force since 2021, 45% of newly purchased public transport vehicles must be powered by alternative fuels, and 65% from 2026. In line with this legal requirement, many local transport companies are currently converting their bus fleets to battery-powered electric drive systems. Large cities such as Berlin, Hamburg, and Munich are working to fully electrify their bus fleets by 2030.

Electric buses (e-buses) are still more expensive than their counterparts with combustion engines and their purchase requires investments in specialized charging infrastructure, which for cost reasons is being built preferably in the form of slow chargers in bus depots and additionally as fast chargers at selected terminals. The comparatively shorter range of electric buses, especially in winter and summer, when heating and air conditioning consume significant amounts of additional energy, often requires charging en route. Charging times and detours that may occur hence need to be planned. It is also crucial to bear in mind that energy prices can fluctuate widely during the day and that the capacity of the local electricity grid may be subject to dynamic constraints, which is particularly important for night-time depot charging of the

entire fleet. Load peaks due to simultaneous charging should be prevented as far as possible. It is therefore imperative not only to optimize vehicle deployment through vehicle rotations, but also to carefully select the charging events, which requires integrated bus circulation and charging planning (see Figure 1).

We developed a novel method at the Zuse Institute Berlin in the MobilityLab of the MODAL research campus together with our industry partner IVU Traffic Technologies AG. The algorithm enables the calculation of cost-minimal bus schedules with integrated optimization of the battery's state of charge (SoC), while complying with capacity conditions at individual and combined charging points. The treatment of mixed e-bus and diesel fleets, which is important in a transition phase, is also possible.

Bus batteries are charged like conventional batteries using the so-called CC-CV (Constant Current - Constant Voltage) scheme. The SoC increases almost linearly up to about 80% under a constant charging rate, while for the last 20% the maximum permissible charging power then decreases rapidly. However, this last fifth of the battery capacity is important in order to make the best possible use of the performance of
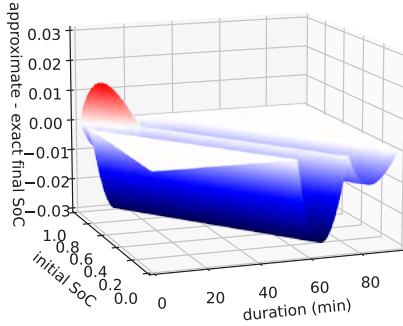
Figure 2: Charging paradox: Approximation error of the charging increment due to a linear spline approximation that underestimates the charging curve. The SoC is scaled to the interval [0,1]. Blue values indicate an underestimation, and red ones an overestimation.
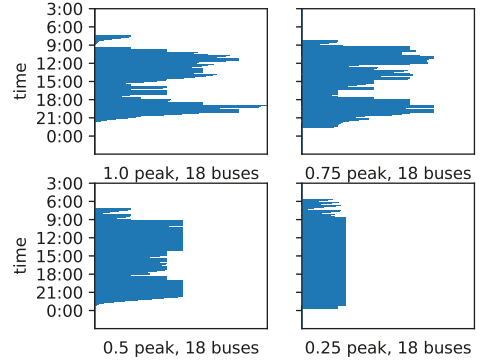


Figure 3: Reducing peak loads by distributing charging events more evenly throughout the day while maintaining the same fleet size, known as peak shaving. The total power consumption in the depot over the planning horizon for various maximum peak loads is shown.

expensive electric buses. This charging process is generally modeled in terms of a so-called charge curve $\zeta$, where $\zeta(t)$ gives the final SoC of an initially empty battery that has been charged for a period of time ($t$). Pre-existing planning methods are based on a linear or piecewise linear approximation of this function. One of the main innovations of our approach is the insight that this approach turns out to be inadequate for two reasons.

On the one hand, the aforementioned conditions result in local transport operators implementing an active charging management, which dynamically throttles the charging power depending on the current electricity price and the grid load. It can therefore not be assumed that charging is always carried out at maximum power, but that the selection of the charging curve itself is a variable. But instead of one charging curve, then all possible charging curves would have to be taken into account.

On the other hand, it is – at first sight paradoxically – possible that an underestimating approximation of the charging curve results in an overestimation of the state of charge, as shown in Figure 2.
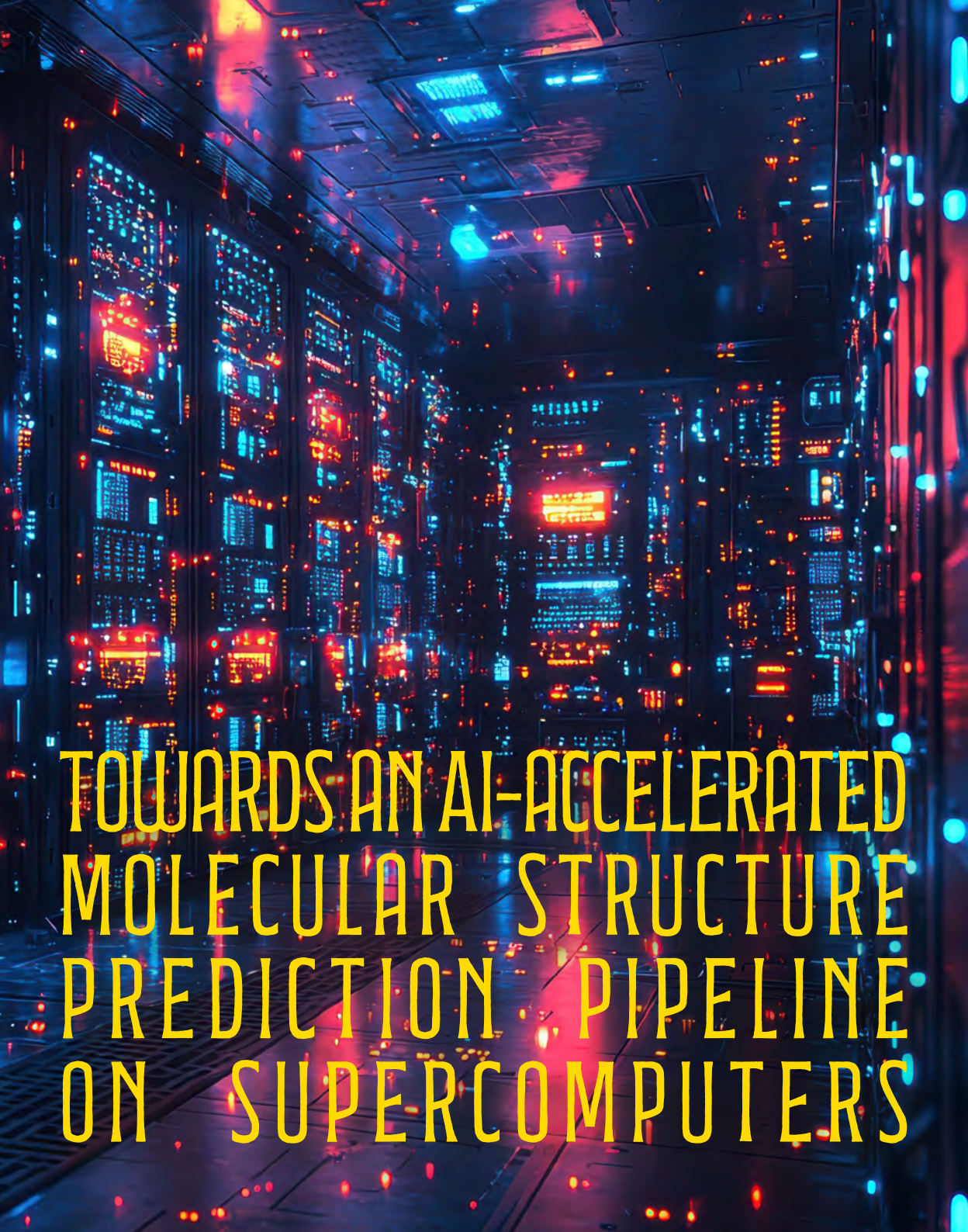
Indeed, batteries usually have an initial residual state of charge $y$ and we are interested in the charge increment $\Delta\zeta(y,t)=\zeta(\zeta^{-1}(y)+t)-y$, which indicates how much the charging event increases the SoC. But this function depends mainly on the charging rate, which can be overestimated by piecewise linear approximations of the charging curve. Instead of the charging curve, we approximate the charge increment function with a piecewise linear underestimation and are able to show that the SoC is never overestimated. In addition, approximations involving only three sections with time increments of five minutes already lead to demonstrably negligible errors in the per mille range.

Our charging model also makes it easier to take into account electricity costs and grid capacities that depend on the time of day. In all cases we considered, it was possible to reduce the peak grid load of an unrestricted reference scenario by 50% to 75% without incurring an additional vehicle demand by optimizing the cycles and the charging in an integrated manner; see Figure 3 for an example. ⌣
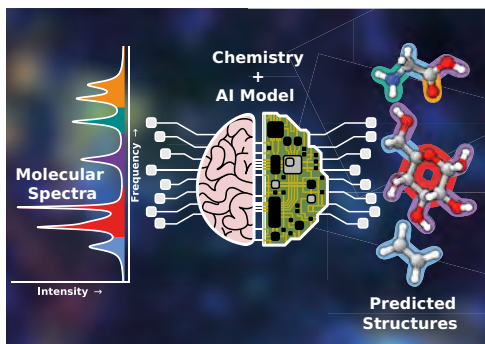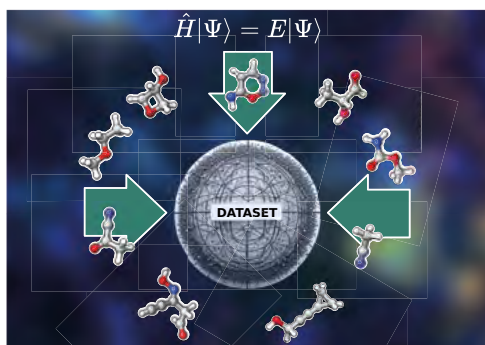
Ralf Borndörfer & Fabian Löbel

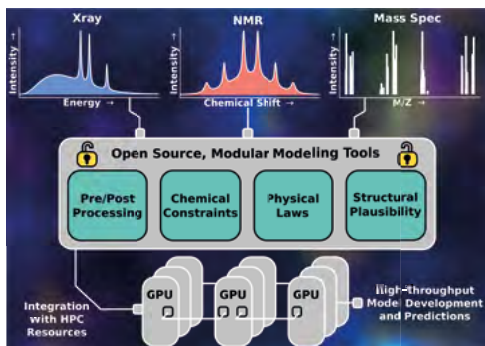# TOWARDS AN AI-ACCELERATED MOLECULAR STRUCTURE PREDICTION PIPELINE ON SUPERCOMPUTERS

**AI-accelerated molecular structure prediction pipeline:** Methods for inverse molecular structure prediction directly from spectroscopy measurements are developed by leveraging expert knowledge in chemistry together with the modeling power of AI on supercomputers.



**Datasets for ML training:** The initial dataset, composed of hundreds of thousands of small molecules using high-level quantum-chemical calculations.



**The software pipeline:** The resulting code base will provide HPC users with an open-source, modular toolbox that can be used with different kinds of spectroscopic measurements. The software packages will be easily deployable on HPC resources so that users can easily train new models or perform high-throughput structural predictions on modern hardware.

Predicting molecular structures from experimental data is a fundamental challenge in molecular sciences, crucial for applications such as drug discovery, materials science, and sustainable energy solutions. Traditional structure prediction methods are both computationally expensive and time-intensive. Advances in artificial intelligence (AI) and machine learning (ML) now offer an alternative, enabling more efficient and accurate structure predictions based on experimental and computational data.

Our journey towards an AI-accelerated pipeline for structure prediction began with developing machine learning methods to explore the high-dimensional space of molecular configurations and identify energy barriers along potential reaction paths.

## Machine learning approach for energy barrier estimation

Energy barriers determine different molecular conformations and are crucial for understanding molecular properties. We developed an ML framework to estimate reaction energy barriers without explicit and expensive transition state calculations. Using a dataset of over 11,000 reactions, this approach employs kernel ridge regression (KRR) to predict energy barriers based on molecular descriptors, including Coulomb matrices, bond distances, atomic charges, and electronic properties such as electronegativity and hardness. By encoding these characteristics, the model captures key structural and energetic trends that influence reaction feasibility.

A key advantage of this ML approach is its computational efficiency, enabling rapid screening of chemical spaces to identify feasible molecular structures. AI-driven models trained on reaction energy barriers refine and enhance structure determination processes. Furthermore, continuous integration of experimental and computational data enhances predictive accuracy, ensuring adaptability to evolving trends in molecular chemistry.

## Computational screening and quantum chemistry approach

A computational study validated and refined ML-based energy barrier predictions. Advanced quantum chemistry techniques were applied to analyze specific molecular interactions and assess the feasibility of transformations. The results confirmed that while some reactions proceed smoothly, others require additional energy due to structural and electronic differences.

This study illustrates how ML-generated predictions can be verified and improved using quantum chemical methods, leading to a more systematic and reliable approach to molecular structure prediction. By analyzing a wide range of possible reactions, the study helps refine search strategies for viable molecular transformations and enhances AI-driven structure determination. Due to the high accuracy achieved with the ML model, expensive computations can be eliminated.

## Recovering important protein configurations with coarse-grain models

Predicting stable protein structures plays a pivotal role in drug discovery and therapeutic advancements. Traditionally, molecular dynamics simulations are used to explore, discover, and validate configurations, but for large proteins, this process is computationally expensive. Coarse-grain models simplify these simulations but often lack chemical accuracy or detect only limited stable states.

By using ML to replace traditional simulations with faster but still accurate coarse-grain models, these problems have recently been circumvented. ML coarse-grain models reduce complex simulations to a few representative atoms while preserving structural variety. Recent ML-driven coarse-grain models not only enhance molecular configuration prediction across diverse proteins but also establish a pathway for versatile, chemically accurate, and efficient structure prediction models.

## AI-accelerated molecular structure prediction pipeline on supercomputers

Building on prior work, we advance AI-based molecular structure prediction. Our methodological design and implementation target state-of-the-art supercomputer resources at NHR@ZIB, combining computationally demanding molecular simulations with AI/ML-based predictive modeling to analyze patterns and make predictions from large datasets.

We are currently developing a structure prediction method for molecules based on computed spectra, advancing the direct prediction of molecular structures from infrared (IR) spectroscopy data – a crucial frontier in molecular science. Since IR spectra are challenging to analyze manually, we train AI models on quantum chemistry datasets to enable highly accurate predictions.

## Outlook

Computational chemistry, ML, and supercomputers provide a robust framework for structure prediction. The ML approach accelerates energy barrier estimation and broadens the scope of possible molecular transformations, while quantum chemical studies validate and refine these predictions. By integrating these methodologies with experimental data analysis, particularly IR spectroscopy, we are developing state-of-the-art AI-driven models capable of accurately predicting molecular structures. This approach advances molecular science with broad scientific and industrial applications.
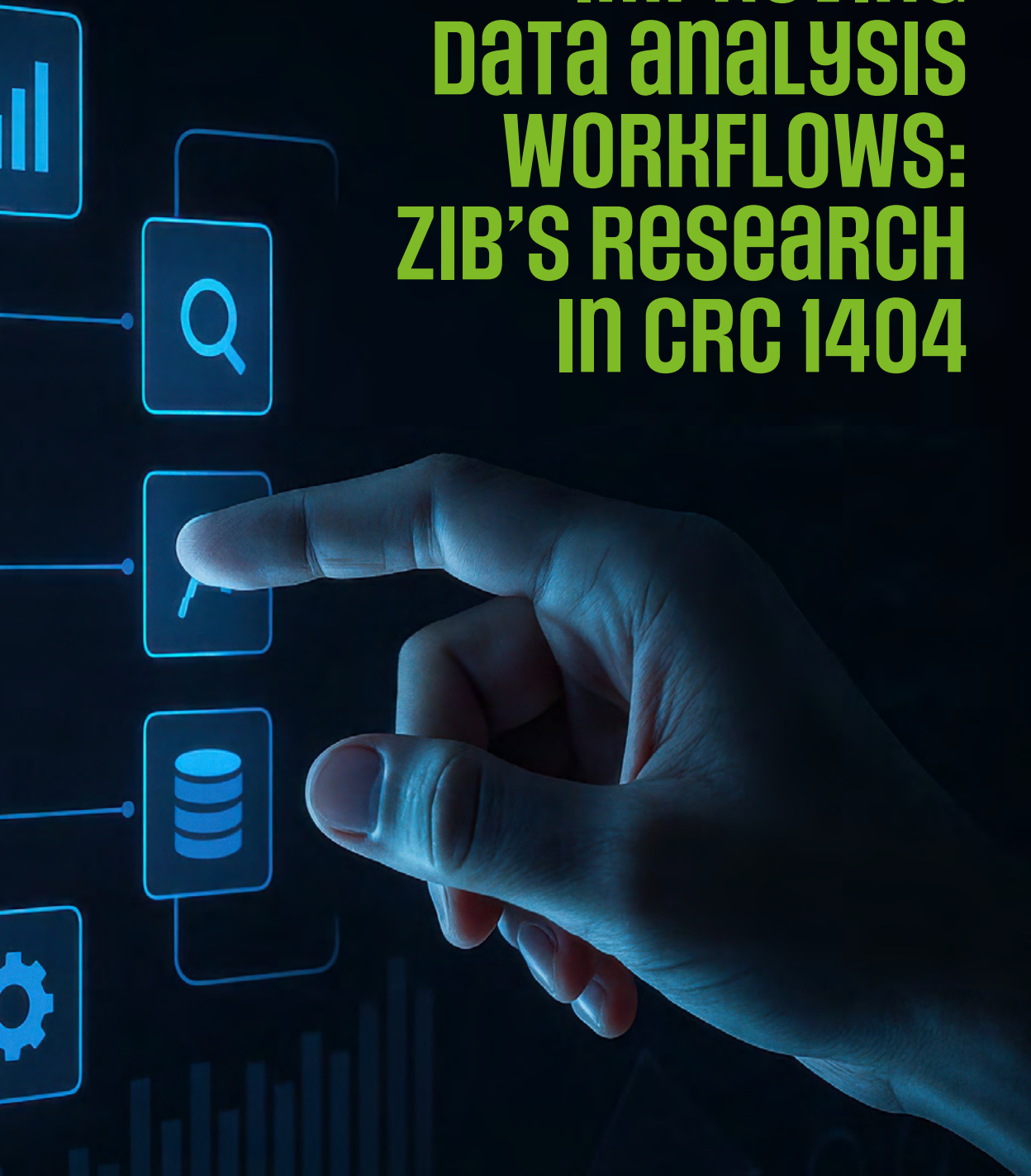
Anita Ragyanszki, Nicholas Charron

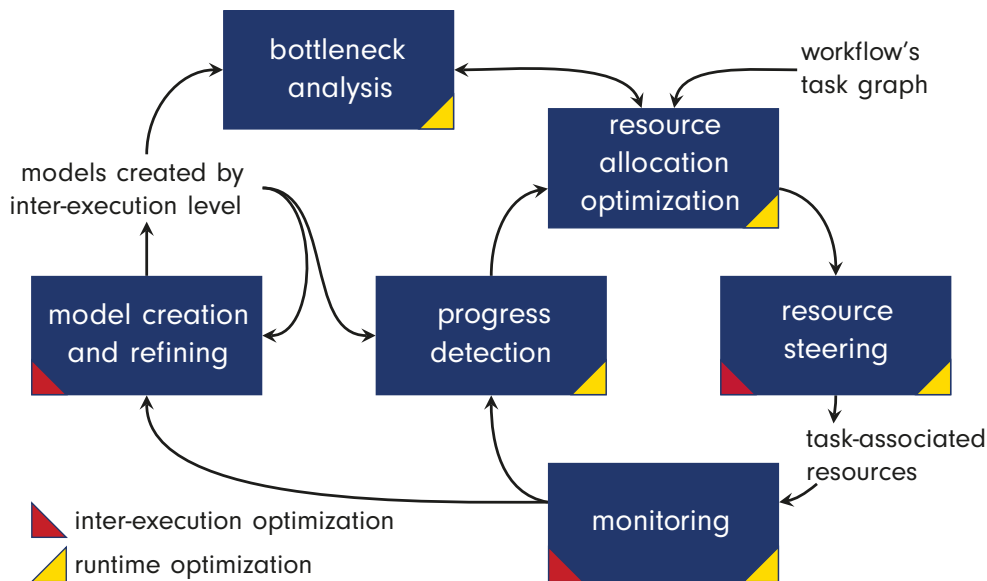# IMPROVING DATA ANALYSIS WORKFLOWS: ZIB'S RESEARCH IN CRC 1404

Fig. 1: Inter- and intra-workflow optimization cycle.

ZIB is the partner of the DFG-funded collaborative research center "FONDA – Foundations of Workflows for Large-Scale Scientific Data Analysis" (CRC 1404) led by Humboldt-Universität zu Berlin. FONDA aims to enhance human productivity in data analysis workflows (DAWs). By now, almost all scientific disciplines use DAWs to process ever-increasing amounts of data. But portability, availability, and dependability of scientific workflows are limited and means to increase them are necessary. In addition, improving productivity is crucial not only for computer science but also across various fields. Consequently, many FONDA subprojects include both computer scientists and experts from other natural sciences. This interdisciplinary approach fosters a deeper understanding of real-world workflows, ensuring practical applicability.

We participate in subproject B4, focused on optimizing the execution of DAWs. By treating workflow tasks as black boxes, we developed a generic optimization architecture (Figure 1). Through task execution monitoring,

we create models that represent each task's requirements using mathematical functions. These models enable rapid testing of different resource allocations, orders of magnitude faster than state-of-the-art workflow simulations based on discrete event simulation. Our 'tests' perform bottleneck analysis by not only predicting the makespan of tasks but also their actual resource needs and the dependencies affecting execution (Figure 2). Analyzing chains of tasks is also possible by using one task's output as an input for another, allowing comprehensive bottleneck analysis of entire workflow executions. Based on this fast bottleneck analysis, heuristic algorithms can optimize resource allocation before and during execution, enhancing accuracy and efficiency. To facilitate just-in-time optimization, we determine the current progress of the workflow and executing tasks through live monitoring data.

To address our monitoring requirements, we explored various methods and developed a custom solution for efficient yet comprehensive I/O monitoring. We discov-
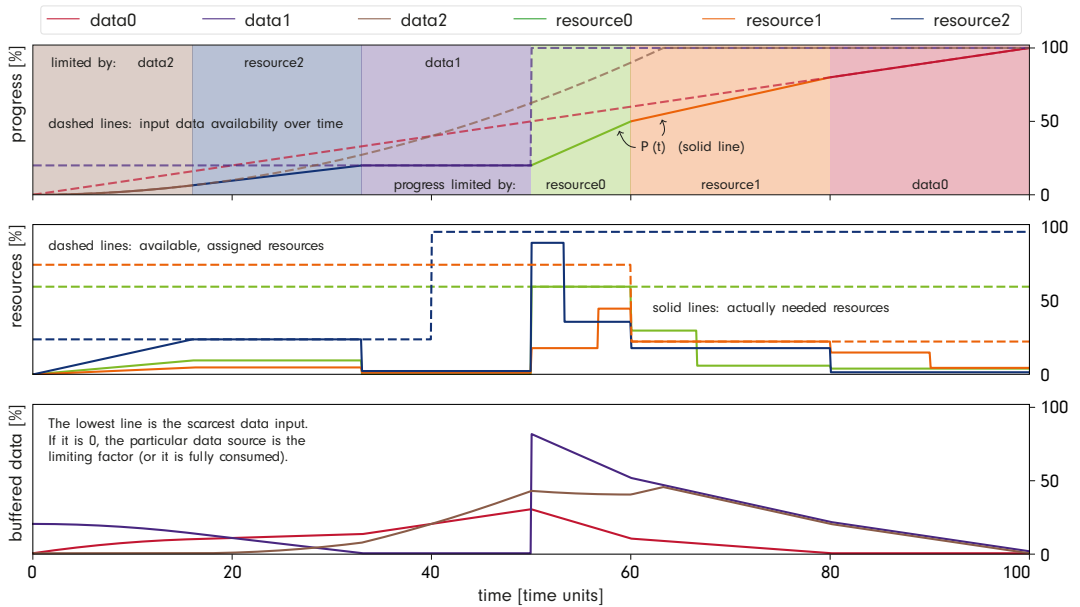
Fig. 2: Detailed bottleneck and resource analysis using BottleMod.

ered that it is challenging to track low-level metrics like I/O requests and correlate them with higher-level concepts such as workflow tasks. For appropriate performance prediction of I/O-heavy workloads, we developed I/O models that consider the caching behavior of the Linux kernel for different I/O methods that significantly improve the prediction accuracy compared to event-based simulation. With our IOSIG plug-in for GCC, we allow pragma annotations to the source code expressing the input/output (I/O) characteristics for certain I/O streams and then decide during execution to which devices the input/output should best be re redirected.

Research for the first phase of FONDA officially began in July 2020 and concluded four years later. We presented and published our results at international scientific conferences, won a 'best poster award' at eScience'24 and Masoud Jami, né Gholami, successfully finished and received his PhD in computer science contributing to the topics from Humboldt-Universität

zu Berlin in 2024. We applied for a second phase of FONDA and our subproject, which now receives funding following a thorough review by DFG in February 2024. The second phase of FONDA started in July 2024 and will also span four years. In this phase, we will shift our focus from executing a single workflow within one data center to running workflows concurrently across multiple data centers. This new scenario introduces challenges such as typically slower connections between different data centers, making data placement significantly more important. Additionally, the concurrent execution of multiple workflows naturally increases the complexity of optimization approaches. While we can build on our previous work from the first phase, these new factors present unique hurdles that will require innovative solutions.

Joel Witzke, Florian Schintke

# Advancing Tier-3 HPC and AI infrastructure

## Interview with Carsten Schäuble, Head of IT and Data Services at ZIB.



**Editorial Team:** Could you please give us an overview of ZIB's Tier-3 HPC infrastructure?

**Carsten Schäuble:** There are two different HPC installations at ZIB: the Tier-2 NHR system and ZIB's Tier-3 AI and HPC infrastructure. The Tier-3 HPC infrastructure at ZIB supports diverse scientific computing tasks, offering scalable computing power for simulations, data analytics, and AI applications. Our system balances traditional HPC workloads with the latest AI research. In recent years, we have focused on enhancing flexibility and responsiveness to the evolving needs of interdisciplinary research.

**Editorial Team:** What are the latest AI-related additions?

**Carsten Schäuble:** We have expanded our AI infrastructure with high-performance GPU clusters, optimized storage solutions, and AI-specific software frameworks. These enhancements improve scalability and efficiency for deep learning models. We also provide container-based tools for streamlined development and deployment.

**Editorial Team:** How does this affect research at ZIB?

**Carsten Schäuble:** The infrastructure enables faster AI model training, benefiting fields like computational biology and materials science. Researchers can now process complex computations with greater speed and accuracy, allowing for more rapid iterations and deeper exploration of scientific questions.

**Editorial Team:** How does ZIB ensure its infrastructure remains cutting-edge?

**Carsten Schäuble:** We regularly upgrade hardware and software, collaborate with industry leaders, and integrate emerging technologies based on researchers' evolving needs. Feedback from our user community is essential in shaping our roadmap and ensuring long-term relevance.

**Editorial Team:** What role does energy efficiency play?

**Carsten Schäuble:** Energy efficiency is crucial. We implement advanced cooling, power-efficient hardware, and workload optimization to minimize energy consumption, while exploring renewable energy options. It is a key part of our infrastructure strategy and long-term sustainability efforts.

**Editorial Team:** How do Kubernetes or similar solutions fit into ZIB's infrastructure?

**Carsten Schäuble:** Kubernetes and similar container orchestration platforms play a significant role in managing workloads efficiently. At ZIB, we use Kubernetes for scalable deployment of AI models and HPC applications, ensuring better resource utilization, automated workload balancing, and seamless integration with cloud and hybrid computing environments.

**Editorial Team:** How does ZIB handle large-scale data management?

**Carsten Schäuble:** ZIB is linked to Berlin's research network with more than 200 GBit/s, allowing for ultra-fast data transfers. Our infrastructure is designed to handle truly petabyte-scale datasets, ensuring researchers can store, process, and analyze massive amounts of data efficiently.

**Editorial Team:** What's next for ZIB's HPC and AI infrastructure?

**Carsten Schäuble:** Our future plans include expanding AI capabilities with next-gen GPUs, integrating hybrid computing models, and strengthening academic and industry partnerships to drive innovation. We are also preparing for future workloads in quantum-inspired simulation and large-scale graph analytics.

**Editorial Team:** How is user training supported at ZIB?

**Carsten Schäuble:** We offer targeted workshops, documentation, and personalized consulting to support users of all experience levels, especially as part of the excellent training and support program of the National Alliance of High-Performance Computing (NHR). This ensures efficient access to resources and accelerates research productivity.

**Editorial Team:** Thank you for your insights.

**Carsten Schäuble:** My pleasure. We are committed to advancing scientific computing at ZIB.

# Developing data capabilities for AI at ZIB

## Interview with ZIB researcher Tim Conrad

**Editorial Team:** Mr Conrad, thank you for taking the time to talk to us. Can you briefly explain what the Data Lake at the Zuse Institute Berlin (ZIB) is and what role it plays?

**Tim Conrad:** My pleasure. The Data Lake at ZIB is a central data infrastructure – or repository – which can be used to store, manage, and share research data of (almost) any size. The interesting part is that all kinds of data can be stored, whether it is structured (think of tables) or unstructured (think of images or videos) or in-between (for example logging data). A big advantage of the lake is that you can keep track of various versions of your datasets – as you might know already from Git repositories. And you can even create multiple branches of the same dataset to enable different team members to work on the same data without changing it for the rest of the team. Overall, I believe that this allows researchers at ZIB and beyond to work with and exchange data more efficiently.

**Editorial Team:** How exactly can the Data Lake be used in research?

**Tim Conrad:** One key aspect is the support for data-intensive research projects. Our scientists can store very large datasets from simulations or experiments and share them with others – which might be more complicated if you used standard file system-based directories. However, scientists whose datasets are not so large can also benefit from the Data Lake, in particular when they want to use modern analysis methods such as artificial intelligence and machine learning which can be performed more efficiently using a data lake.

**Editorial Team:** Could you elaborate on how AI and machine learning benefit from the Data Lake?

**Tim Conrad:** Absolutely. Machine learning and AI models require high-quality data for effective training. The Data Lake's integrated metadata and search functions help researchers quickly find and reuse relevant datasets without duplication. This is especially valuable in interdisciplinary projects, where different team members may need access to specific datasets while applying distinct methodological approaches.

**Editorial Team:** That sounds like an efficient way to manage and use research data. Can you give an example of a specific project at ZIB that is already benefiting from the Data Lake?

**Tim Conrad:** One such project is MaRDI, the Mathematical Research Data Initiative, where ZIB is one of the project partners. MaRDI focuses on the systematic collection, management, and use of mathematical research data. A key component of this initiative is the development of standards and technologies for data management in mathematics, facilitating the exchange of mathematical models and research findings. By leveraging the Data Lake, MaRDI provides a structured environment where researchers can store and access high-quality, reusable datasets, particularly when these are too large for standard data repositories.

**Editorial Team:** How do the Data Lake at ZIB and MaRDI complement each other?

**Tim Conrad:** MaRDI directly benefits from the Data Lake as it provides a reliable infrastructure for storing and managing mathematical data. At the same time, MaRDI contributes important insights into the development of specific solutions for managing complex mathematical datasets. The close integration of both systems ensures the sustainable use of mathematical research data and enables more efficient collaboration within the mathematical community.

**Editorial Team:** That's really impressive. How do you see the future of the Data Lake at ZIB?

**Tim Conrad:** The future lies in expanding its capabilities. We're working on integrating more advanced search and indexing functions to make data discovery even easier. Another exciting direction is enabling real-time data analysis within the lake, so researchers can process and analyze data directly without having to move large datasets. Finally, we want to enhance interoperability with other national and international research infrastructures to facilitate broader scientific collaboration.

**Editorial Team:** Thank you for these insights, Mr Conrad. It's clear that the Data Lake at ZIB can indeed become an important resource in supporting data-driven research.

**Tim Conrad:** Thank you, it was a pleasure discussing this!

# IMPRINT