# Research Statement

In the era of increasingly heterogeneous High Performance Computing (HPC) environments with diverse architectures, managing and optimizing application performance is essential yet challenging. Modern computing systems comprise numerous components, each with its own configuration parameters, which often conflict with each other. For instance, my research shows that the same parameter setting on a supercomputer can achieve optimal execution time but also cause $4x$ run-to-run variation [1], meaning a week-long vaccine simulation could take a month in the worst case if the appropriate configuration parameters are not used. The complexity of understanding these relationships to make informed decisions underscores the need for building a data-driven performance model for automated decision making.

Recognizing these needs, my research objective is to design innovative data-driven performance modeling techniques to automate decision-making in these dynamic and heterogeneous environments. My research approach includes developing parallel and scalable data science solutions by fusing advancements in AI/ML with my expertise in HPC. My past research efforts address important research questions in many areas including performance characterization, comparison, bottleneck detection, reproducibility, and mitigation.

## Area 1: Represent Data to Facilitate Accurate Downstream Modeling; *PI: $770K; DOE Early CAREER.*
The first step of any Machine Learning (ML)-based modeling approach is to transform data from their native and multi-modal forms (e.g., graph, tree, table, image) to a numeric vector representation called "embeddings". My DOE Early Career award develops one such innovative way to represent multi-modal performance data that enables ML models to exploit both sample-sample and feature-feature relationships. In particular, my work imagines performance samples as graph nodes, and constructs edges between each pair in various different ways using their similarities [2], [3], [4]. Using graphs, real-time performance measurement tools can quickly classify new unlabeled samples into their respective neighborhoods based on similarity, thus reducing the need for extensive data collection for accurate clustering. During extensive evaluations, our method yielded up to 61% more accurate performance models compared to the state-of-the-art methods and was more accurate in performance anomaly classification in 60% the times.

**Impact:** Once successful, this method will make real-time decision-making more time- and data-efficient, which can impact cross-cutting domains with real-time decision-making environments, such as disaster response, manufacturing, and autonomous vehicles.

## Area 2: Design Effective and Efficient Downstream Modeling Methodologies; *PI: $750K+; Collaborators: LLNL, BNL, LBNL, AMD.)* 
The objective of my research in this area is to develop novel AI/ML methodologies that are effective in using the data representation from Area 1 to achieve high accuracy predictions and synthesis. The output of these models can guide users, software and facilities in making decisions about how to execute them to achieve multiple objectives such as shorter makespan and higher throughput, power efficiency and system utilization. To accomplish this objective, it is crucial to first understand the resource usage behaviors of applications and pinpoint their bottlenecks, which often involves comparing high-dimensional performance profiles. In this context, my research has developed scalable methodologies and metrics to quantify the impact of parameters and their relationships with underlying systems on various performance objectives such as time to solution, system utilization, and performance variability, leveraging and, where appropriate, extending ML techniques, advancing both HPC and ML [5], [6], [7], [8], [9].

**Impact:** This research has led to several successful technology transfers with national labs and the industry, e.g., integration into Lawrence Berkeley lab's auto-tuner called GPTune and AMD's OmniPerf.

## Area 3: Develop Methods to Efficiently Generalize Models *PI: 450K; Collaborators: AMD, LLNL*
Performance data collection can take days and weeks, and with the advent of a new architecture every six months, the process of data collection, analysis, and optimization need to be repeated, which can delay

✉ tanzima@txstate.edu 📞 +1-765-413-3165 🌐 http://www.tanzimaislam.com in Γ

the scientific code porting tasks in the HPC co-design process. However, with ML, my research team has demonstrated that one can build performance models that can help predict performance in new architectures given an abstract description of the new hardware with high fidelity. However, cross-platform performance prediction is challenging due to the diverse application and platform interactions across platfrom that are not directly comparable. My research develops novel performance modeling methodologies that combine the power of few-shot learning with generative AI, particularly Large Language Models (LLMs) to achieve high prediction accuracy in time- and data-efficient manner compared to the traditional methods [10]. When the number of samples is not sufficient, our method of synthesizing samples using fine-tuned LLMs has been shown to achieve comparable accuracy but with a lot less data ($< 1\%$). This research advances the state-of-the-art in transfer learning for both HPC and ML domains, offering a data-efficient solution for HPC performance modeling.

**Impact:** This research will significantly improve the turn around time of the scientific method by enabling the procurement team to quickly evaluate a future supercomputer from various vendors.

**Area 4: Make Performance Data Reproducible & FAIR** *PI: 300K; Collaborators: BNL, Sandia, Argonne*
In this direction, my research tackles the challenge of reproducible performance in HPC, where various factors can influence the run-to-run execution time. Our novel graph-based ML method [1] has identified crucial but non-obvious metadata for explaining the trade-off between performance optimality and reproducibility, revealing gaps in current metadata collection practices. To address these gaps, we are developing a FAIR (Findable, Accessible, Interoperable, Reusable) framework for performance data collection, annotation, and analysis incorporating new metadata such as I/O device types and network characteristics to better explain and reproduce performance results in HPC environments.

**Impact**: FAIR performance datasets and metadata will enable future research and contributes to making science more reproducible.

**Future Research Directions: Facilitate Interactive Performance-Aware Decision-Making using Heterogeneous Data Sources** To harness the benefits of heterogeneous HPC environments, the next challenge is to develop novel modeling approaches that integrate and model data from multiple, unaligned sources, addressing scalability and heterogeneity to enhance performance analysis.

Opportunity 1: Tackle Heterogeneous, Multi-Modal, and Unaligned Data Sources The objective of this research is to develop modeling techniques that can learn from diverse, distributed sources of performance data collected from various layers of the computing ecosystem, from user-level to system-level. Currently, no methodology in HPC or data science can effectively harness heterogeneous, multi-modal, and unaligned data. This gap presents a significant opportunity for my research to advance the state of the art in ML and HPC.

Opportunity 2: Scaling Performance Modeling with Current and Future Data Sources The objective of this research is to scale performance modeling techniques to handle large volumes of disparate data sources. Current methodologies are inadequate for handling the vast and diverse datasets, leading to inefficiencies and inaccuracies. To bridge this gap, I will develop a modular framework that trains component-specific models on focused datasets and integrate them using parameter correlations into a global model. This approach will capture complex interactions accurately while reducing data and computational requirements. Without novel modeling approaches, harnessing the benefit of data deluge will remain constrained.

**Societal Impact:** My research aims to improve the scientific process by enabling efficient performance analysis and optimization in heterogeneous HPC environments, accelerating advancements in areas such as climate modeling, clean energy, and drug discovery. Additionally, my commitment to integrating research with education is demonstrated by developing AI-powered educational tools that provide personalized, data-driven learning experiences, enhancing HPC education and preparing the next generation of researchers.

# References

[1] Tapasya Patki, Jayaraman J Thiagarajan, Alexis Ayala, and Tanzima Z Islam. Performance optimality or reproducibility: that is the question. In *International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–30, 2019.

[2] Tarek Ramadan, Ankur Lahiry, and Tanzima Z Islam. Novel representation learning technique using graphs for performance analytics. In *2023 International Conference on Machine Learning and Applications (ICMLA)*, pages 1311–1318. IEEE, 2023.

[3] Chase Phelps, Ankur Lahiry, Tanzima Z. Islam, and Line Pouchard. Reimagine application performance as a graph: Novel graph-based method for performance anomaly classification in high-performance computing. In *48th Annual Computers, Software, and Applications Conference (COMPSAC)*. IEEE, 2024. (Accepted). Acceptance rate: 24%.

[4] Tarek Ramadan, Tanzima Z Islam, Chase Phelps, Nathan Pinnow, and Jayaraman J Thiagarajan. Comparative code structure analysis using deep learning for performance prediction. In *2021 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, pages 151–161. IEEE, 2021.

[5] John Tramm, Andrew Siegel, Tanzima Islam, and Martin Schulz. Xsbench-the development and verification of a performance abstraction for monte carlo reactor analysis. *The Role of Reactor Physics toward a Sustainable Future (PHYSOR)*, 2014.

[6] Tanzima Z Islam, Jayaraman J Thiagarajan, Abhinav Bhatele, Martin Schulz, and Todd Gamblin. A machine learning framework for performance coverage analysis of proxy applications. In *SC'16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 538–549. IEEE, 2016.

[7] Tania Banerjee, Jason Hackl, Mrugesh Shringarpure, Tanzima Islam, S Balachandar, Thomas Jackson, and Sanjay Ranka. Cmt-bone — a proxy application for compressible multiphase turbulent flows. *IEEE 23rd International Conference on High Performance Computing (HiPC)*, pages 173–182, Dec 2016. doi: 10.1109/HiPC.2016.029. URL https://ieeexplore.ieee.org/abstract/document/7839682. Acceptance rate: 23%.

[8] J. J. Thiagarajan, R. Anirudh, B. Kailkhura, N. Jain, T. Islam, A. Bhatele, J. Yeom, and T. Gamblin. Paddle: Performance analysis using a data-driven learning environment. In *2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 784–793, May . doi: 10.1109/IPDPS.2018.00088.

[9] Tanzima Islam, Alexis Ayala, Quentin Jensen, and Khaled Ibrahim. Toward a programmable analysis and visualization framework for interactive performance analytics. In *2019 IEEE/ACM International Workshop on Programming and Performance Visualization Tools (ProTools)*, pages 70–77. IEEE, 2019.

[10] Banooqa H. Banday, Tanzima Z. Islam, and Aniruddha Marathe. Perfgen: A synthesis and evaluation framework for performance data using generative ai. In *48th Annual Computers, Software, and Applications Conference (COMPSAC)*. IEEE, 2024. (Accepted). Acceptance rate: 24%.