

Mohammad Tawhidul Hasan Bhuiyan

 mb5332@columbia.edu |  [Tawhid Bhuiyan](#) |  mthbhuiyan.github.io

RESEARCH INTEREST

Compiler, Profile-Guided Optimizations, Systems for Machine Learning, Hardware/Software Co-Design

EDUCATION

Columbia University

PhD in Electrical Engineering

New York, USA

Sept. 2024 – May 2030 (expected)

Palashi, Dhaka, Bangladesh

Bangladesh University of Engineering and Technology (BUET)

Masters of Science in Computer Science and Engineering

July 2021 – Aug. 2024

Bangladesh University of Engineering and Technology (BUET)

Bachelor of Science in Computer Science and Engineering

Palashi, Dhaka, Bangladesh

Mar. 2016 – Feb. 2021

 Class rank: 1st

EMPLOYMENT

Columbia University

Graduate Research and Teaching Assistant

New York, USA

Sept. 2024 – Present

Bangladesh University of Engineering and Technology (BUET)

Lecturer, Department of CSE

Dhaka, Bangladesh

Sept. 2021 – Aug. 2024

PUBLICATIONS

- **Tawhid Bhuiyan**, Sumya Hoque, Angelica Moreira, Tanvir Ahmed Khan, **Stale Profile-Guided Optimizations**, Submitted to ASPLOS-2026
- Ryan Piersma*, **Tawhid Bhuiyan***, Tanvir Ahmed Khan, and Simha Sethumadhavan, **Highly Scalable Power Management Protocol**, Submitted to ISCA-2026, (*Co-first Authors)
- Amir Zarandi, **Tawhid Bhuiyan**, Laurent Schares, and Tanvir Ahmed Khan, **Rebasing GPU Micro-Architectural Power Modeling for Machine Learning Workloads**, Submitted to CAL-2026
- Ryan Piersma, **Tawhid Bhuiyan**, Tanvir Ahmed Khan, and Simha Sethumadhavan, **Reverse Engineering DVFS Mechanisms**, 2025 IEEE International Symposium on Hardware Oriented Security and Trust (HOST), San Jose, CA, USA, 2025
- **Mohammad Tawhidul Hasan Bhuiyan**, Muhammad Rashed Alam and M. Sohel Rahman, **Computing the Largest Common Almost-Increasing Subsequence**, *Theoretical Computer Science*, Volume 930, 2022, Pages 157-178, ISSN 0304-3975
- **Mohammad Tawhidul Hasan Bhuiyan**, Irtesam Mahmud Khan, Sheikh Saifur Rahman Jony, Renee Robinson, Uyen-Sa D. T. Nguyen, David Keellings, M. Sohel Rahman, and Ubydul Haque, **The Disproportionate Impact of COVID-19 among Undocumented Immigrants and Racial Minorities in the US**, *International Journal of Environmental Research and Public Health*, Volume 18(23): 12708, 2021, PubMedID 34886437, ISSN 1660-4601

PROJECTS

Optimizing Datacenter Applications with Stale Profile C++, Python, LLVM	Sept. 2024 – Present
<ul style="list-style-type: none">Developed algorithms to map hardware profiles from deployed datacenter applications to new application binaries, enabling effective optimization despite "stale" profile data.Integrated mapping heuristics with Meta's BOLT and Google's Propeller, demonstrating applicability across industry-standard optimization toolchains.	
Highly Scalable Token-Based Power Management Protocol Python, Numba, Architecture	Aug. 2025 – Present
<ul style="list-style-type: none">Built a simulator for a token-based power management protocol capable of running cycle-accurate power distribution on top of both cycle-accurate and event-driven system components.Designed a novel protocol for system-wide power distribution. Achieved 2.4x more responsive power distribution compared to State-of-the-Art solutions for systems comprising thousands of interconnected components.	
Optimizing GPU energy efficiency for LLM CUDA, NVIDIA Nsight, Architecture	Feb. 2025 – Present
<ul style="list-style-type: none"><i>Thermal-Aware GPU Power Modeling:</i><ul style="list-style-type: none">Developed a high-fidelity power model for LLMs by correlating NVML/CUPTI energy traces with PyTorch layer-level timing.Discovered that standard power models fail for compute-heavy kernels due to thermal throttling; implemented corrections to account for dynamic frequency scaling during long-sequence inference, enabling more accurate energy-per-token predictions.<i>Simulator Scalability & Validation (Accel-sim):</i><ul style="list-style-type: none">Validated Accel-sim traces against real-world LLM workloads on NVIDIA Datacenter GPUs.Identified critical scaling bottlenecks in Accel-sim's handling of transformer architectures and implemented patches to achieve accurate power modeling for large-scale models.<i>Energy Modeling for Multi-GPU Training (Astra-sim):</i><ul style="list-style-type: none">Integrated custom energy models into Astra-sim to quantify energy consumption across distributed multi-GPU training setups.Analyzed the performance-energy trade-offs for different models and network technologies.	

Parallelization Suggestion for Python Programs Python, Compiler	Jan. 2022 – Aug. 2024
<ul style="list-style-type: none">Developed a compiler-based prototype to profile Python programs and suggest optimal parallelization strategies for distributed environments.Achieved performance parity with manually parallelized programs for a set of input programs from Kaggle.	

TECHNICAL SKILLS

Architectural Simulators: Accel-sim, Astra-sim, DRAMsim3, Ramulator2, Calculon

Compiler & Systems: LLVM, BOLT, Propeller, AutoFDO

Performance Tools: perf, NVIDIA Nsight, NVML, CUPTI, PyTorch Profiler

Languages: C/C++, Python, Java, SQL, Bash, Rust, Haskell, JavaScript

Libraries: Polars, Numba, PyTorch, Scikit-Learn, Tensorflow, Pandas, NumPy, Matplotlib

AWARDS

- BUET Alumni Association Award for securing top position in the department
- Dean's Honor List, BUET in all 8 undergraduate semesters
- University Merit Scholarship, BUET

TEACHING & MENTORSHIP

- | | |
|--|-----------------|
| Columbia University | <i>Mentor</i> |
| <ul style="list-style-type: none">• Mentored Amir Zarandi (Undergraduate student, Columbia University) in conducting original research, resulting in a co-authored paper which is under submission.• Mentored 15 students across 4 technical projects. Guided teams through the design, implementation, and debugging complex systems. Some key projects include:<ul style="list-style-type: none">◦ <i>Fine-grained GPU DVFS</i>: Leading the students to update Accel-sim to allow simulating and evaluating new DVFS techniques for GPU cores.◦ <i>HBM DVFS & Co-Optimization</i>: Leading the investigation into HBM idle power; utilizing DRAMsim3 and Ramulator2 to propose co-optimized DVFS mechanisms for GPU compute and memory.◦ <i>Power-Aware Collective Communication</i>: Guiding the integration of different ML models into Astra-sim to easily simulate different model architecture. | |
| Bangladesh University of Engineering and Technology (BUET) | <i>Lecturer</i> |
| <ul style="list-style-type: none">• Developed curriculum and delivered lectures for core undergraduate courses:<ul style="list-style-type: none">◦ <i>Operating Systems</i>: Led hands-on sessions on modifying open-source kernels; taught memory virtualization techniques and file systems.◦ <i>Database Systems</i>: Instructed on schema design, normalization theory, and the internal architecture of DB storage systems. | |

REFERENCES

Tanvir Ahmed Khan
Assistant Professor
Department of Electrical Engineering
Columbia University
✉ tk3070@columbia.edu