

Vashisth Tiwari

☎ (585) 524-8385 | ✉ vashistt@andrew.cmu.edu | in www.linkedin.com/in/vashistht/ | 🌐 <https://vashistht.github.io/>

EDUCATION

Carnegie Mellon University

Master of Science in Artificial Intelligence Engineering (GPA: 4.0)

Pittsburgh, PA

Aug. 2023—Dec. 2024

- Courses: Advanced NLP, Deep Learning, Machine Learning, AI Systems, Stochastic Processes, Deep Learning Systems*
- Teaching: Advanced Natural Language Processing (Fall'24) w/ Prof. Graham Neubig

University of Rochester

Bachelor of Science in Physics, Bachelor of Arts in Mathematics (GPA: 3.97)

Rochester, NY

Aug. 2019—Aug. 2024

- Awards: Phi Beta Kappa, Semi-Finalist Rhodes Scholarship India, Harry W. Fulbright Prize, Undergrad Teaching Award
- Courses: Modern Statistics & Exploration, Data Structures & Algorithms, Probability, Real Analysis

RESEARCH EXPERIENCE

SLAB Lab@CMU w/ Prof. Emma Strubell

Pittsburgh, PA

Research Assistant

Aug. 2024—Present

- Designing and exploring the efficacy of active synthetic data generation for LLMs by incorporating student model feedback

InfiniAI Lab@CMU w/ Prof. Beidi Chen

Pittsburgh, PA

Research Assistant

Feb. 2024—Aug. 2024

- Demonstrated that speculative decoding significantly improves both throughput and latency in LLM inference, a key insight for performance optimization. Achieved up to 2x speedup for LLaMA-3 inference with batch sizes of 128 and above
- Investigated integration of diverse model sparsities (weight, attention, activation) in speculative decoding draft models

Dark Energy Spectroscopic Instrument w/ Prof. Segev Benzvi

Rochester, NY

Research Assistant

Jan. 2020—May 2021

- Designed multi-class CNNs for spectral data with TensorFlow, scikit-learn to find galaxies with supernovae
- Enhanced network performance by applying noise-removal techniques like binning and filtering to preprocess spectral data
- Achieved 95%+ accuracy and high precision for supernovae classification tasks in the DESI data pipeline

WORK EXPERIENCE

Mana Finance Corporation

Hillsborough, CA

Quantitative Research Intern

May 2022—Aug. 2022

- Utilized statistical techniques to analyze stock price distributions and quantify investment risk; developed ML models for assessing expected yields on potential investments
- Prototyped a tool demonstrating direct tracking of Ethereum blockchain data on UniSwap (ingestion, indexing, & visualization)

Los Alamos National Laboratory

Los Alamos, NM

Research Intern

Jun. 2021—Aug. 2021

- Modeled complex quantum system using Python; utilized numerical differential equation solvers in Mathematica and Python
- Discovered optimal laser pulse parameters through high-dimensional optimization and parameter estimation using SciPy, CvxPy
- Improved the system performance by 5% beyond the current state-of-the-art pulse parameters

PROJECTS

Pruning while Preserving Reasoning Capabilities

Mar. 2024—May 2024

- Improved upon Bonsai (forward pass only structured LLM pruning, compatible with consumer hardware) in math-reasoning tasks
- Demonstrated that novel task-aware pruning metric better retains reasoning abilities than the standard perplexity baseline

End-to-End NLP System Building

Feb. 2024—Mar. 2024

- Engineered a Retrieval Augmented Generation (RAG) based chatbot on CMU utilizing webpages and semantic scholar data
- Implemented core RAG components: LangChain embedder, Faiss+ColBERT retriever, and reader using open-source LLMs

LLaMA-2 from Scratch

Feb. 2024

- Built a 42M LLaMA-2 model and trained on TinyStories dataset: implemented ROPE embeddings, AdamW optimizer, attention
- Continued pre trained on CFIMDB and fine-tuned on SST-5 datasets to enable zero-shot movie review sentiment analysis

MLBareBones

Aug. 2023—Dec. 2024

- Implemented ML algorithms from scratch using NumPy: neural networks, SVMs, linear regression, decision trees, AdaBoost, etc.

SKILLS & LEADERSHIP

Languages: Python, Java, C++, Mathematica, Bash

Libraries / Frameworks: PyTorch, Hugging Face, TensorFlow, NumPy, Pandas, Scikit-Learn, Spark, Apache Kafka, Cuda

Leadership: Peer Mentor, Department of ECE (CMU); President, Society of Physics Students (University of Rochester)

PUBLICATIONS

[1] Chen J*, **Tiwari V***, Sadhukhan R*, et al. "MagicDec: Breaking the Latency-Throughput Tradeoff for Long Context Generation with Speculative Decoding." *ECCV 2024 Efficient Foundation Model Workshop* (2024). (*Equal contribution)

[2] Uzun C, Pandey S, **Tiwari V**, et al. "Improved Bragg splitting of Bose-Einstein condensates into high-order momenta wave-packets." *APS Division of Atomic, Molecular and Optical Physics* (2023).

[3] Wasserman A, **Tiwari V**, et al. "Using ML to Develop a Transient Identification Pipeline for DESI." *AAS*.