




Dav Vrat Chadha

✉ davvrat.chadha@mail.utoronto.ca  in/davvratchadha  github.com/davvratchadha  davvratchadha.com

System Design Engineer and AI enthusiast driven by a passion for solving real-world challenges. Committed to leveraging cutting-edge technology to create innovative solutions that make a meaningful impact on society.

Education

B.A.Sc. in Engineering Science, University of Toronto

Major: Machine Intelligence [GPA: 3.93]

Toronto, ON, Canada

Sept 2020 – Apr 2025

Work Experience

System Design Engineer Intern, Memory Team

AMD - Datacenter GPU

Markham, ON, Canada

May 2023 - Aug 2024

- Designed Char_wizard, a memory characterization tool that reduced characterization time by 70% (650 hours) using innovative grid search algorithms. Improved runtime complexity from $O(n^2)$ to $O(n)$ and incorporated concurrency, resulting in over 2,000 hours of usage.
- Engineered a 4pt HBM-PHY margin testing tool for filtering **MI325** parts during manufacturing, to reduce customer RMAs.
- Executed electrical validation and memory characterization for HBM3 and HBM3E across all major memory vendors, tuning memory settings to optimize performance and reliability for MI300 and MI325 platforms.
- Performed F_{max} correlation analysis to validate reliable, high-speed data transfers in HBM3 stacks, ensuring alignment across different MI300 platforms.
- Contributed to HBM3E bring-up for the MI325 product with memory validation, automation, and tool support, and assisted in HBM3 debugging during the El Capitan supercomputer's development phase.
- Conducted Bit Error Rate (BER) testing for the HBM-PHY interface on all MI300 platforms, optimizing memory voltage settings.
- Developed SuperScript, a **Python** library for memory validation and debugging, integrated into AMD tools to boost efficiency.
- Designed and developed a **Python** and **C++** tool, amd-address-mapper, to translate physical addresses from ROCm workloads to failing locations in HBM, enhancing debugging efficiency and currently used internally and at customer sites.

Research Student - Computational Linguistics Lab

University of Toronto

Toronto, ON, Canada

Sept 2024 - Present

- Conducting thesis research in Prof. Gerald Penn's Computational Linguistics lab on mechanistic interpretability of **large language models** through circuit discovery with DiscoGP.
- Exploring model merging and editing techniques for circuits to optimize LLMs for multi-task performance in resource-constrained environments, enhancing efficiency for on-device AI in wearables and smartphones.

ML Engineer - FINCH Satellite Mission

University of Toronto Aerospace Team

Toronto, ON, Canada

Sept 2023 - Aug 2024

- Implemented a novel **diffusion model** conditioned on neighboring spectral frames for destriping hyperspectral images resulting in PSNR = 39.2274, LPIPS = 0.2214, SSIM = 0.8817, and SAM = 0.0423, for the ICVL-HSI dataset.
- Worked on the hyperspectral data augmentation pipeline in **PyTorch**, utilizing a modified image patch extraction technique to create new images out of existing ones to increase the size of the training and validation dataset.
- Operated in an **agile** environment, contributing to continuous improvement and innovation within the team.

Publications

- Dav Vrat Chadha et al., *Optimizing Shmooing for AI HPC HBM Memory Characterization Using Decision Trees*, AMD Internal Publication, July 2024. Reviewed by AMD Fellows Committee.
- Ian Vyse et al., *Beyond the Visible: Jointly Attending to Spectral and Spatial Dimensions with HSI-Diffusion for the FINCH Spacecraft*, Presented in 38th Annual Small Satellite Conference, 2024. DOI: 10.48550/arXiv.2406.10724

Projects

TalentRank (Engineering Capstone)

PyTorch, ChromaDB, Python, IR, Recommendation systems, LLMs

Toronto, ON, Canada

Sept 2024 - Present

- Leading a team to design and develop a **PyTorch** and **ChromaDB**-based information retrieval and **recommendation system**, successfully identifying top candidates from a large application pool.
- Implementing dense retrieval techniques with **LLM** integration to rank candidates based on job relevance, enhancing selection accuracy and reducing processing time.
- Coordinating with stakeholders to align system capabilities with hiring goals, ensuring the model effectively meets real-world recruitment needs.

Image Rendering Engine

Toronto, ON, Canada

C, Performance analysis and Optimization, Linux

Sept 2024 - Oct 2024

- Optimized a high-performance image rendering engine for real-time object manipulation based on preprocessed sensor data, achieving smooth 60 FPS rendering from 1500Hz oversampled input.
- Enhanced C-based performance on **Linux** by implementing efficient data structures and profiling with perf, gprof, and GCov, achieving a 700x speedup over the baseline implementation.
- Developed and optimized frame-buffer transformations (e.g., shifting, rotating, mirroring) to enable real-time, accurate image manipulation on constrained hardware.

NoPunIntended

Toronto, ON, Canada

Repository; PyTorch, Python, LLMs, NLP

Jan 2023 - Apr 2023

- Fine-tuned an ensemble of transformer-based **LLMs** (DeBERTa and RoBERTa) using **PyTorch** to accurately detect and locate puns within text by applying contextual masking techniques and clustering with K-means for enhanced interpretability and localization.
- Advanced Amazon's research on pun detection by implementing and optimizing model architecture and training strategies, resulting in an improved method that reached a 75.58% test accuracy—a performance approaching that of GPT-4 (82.77%) on similar tasks.

Skills

Python, PyTorch, NumPy, Tensorflow, C, C++, NLP, AI, ML, IR, Sklearn, Keras, CUDA, Jax, Objax, spaCy, Cython, CI/CD, Git, Linux, Java, MySQL

Honors & Awards

- Two-time AMD Executive Spotlight Award winner - Aug 2024, Dec 2023
- AMD Intern Innovation Showcase Award - Mar 2024
- Two-time AMD Spotlight Award winner - Dec 2023, Aug 2023