

Avi Singhal

avisinghal6@gmail.com|+19452491966 |[linkedin.com/in/avi-singhal99/](https://www.linkedin.com/in/avi-singhal99/)Houston, Texas, USA |github.com/avisinghal6 |avisinghal.com

EDUCATION

Rice University, Houston

Aug 2022 - Dec 2023

Master of Computer Science, GPA- 4.0/4.0

Texas, USA

Coursework: Probabilistic Algorithms and Data Structures, Deep Learning for Vision and Language, Machine Learning, Machine Learning with Graphs, Design and Analysis of Algorithms, Software Engineering Methodology, Computer Architecture

Delhi Technological University

Aug 2017 - Jun 2021

Bachelor of Technology- Electronics and Communication, GPA- 8.75/10

New Delhi, India

Coursework: Microprocessors and Interfacing, Embedded Systems, Pattern Recognition, Web Development

SKILLS

Programming Languages: Python, C++, Java, SQL, JavaScript, Matlab

Technologies and Frameworks: PyTorch, TensorFlow, Hugging Face, AWS, GCP, ONNX, Docker, Linux, Microsoft NNI

WORK EXPERIENCE AND INTERNSHIPS

TetraMem Inc.: Software/ML Intern

May 2023 - Nov 2023

Skills: Python, C++, PyTorch, TensorFlow, Hugging Face, ONNX, AWS, Docker, Linux, Git

California, USA

- Achieved high accuracy (~**85%**), low quantization (**uint8**) loss (~**1%**), low latency (<**5K** MAC operations), compressed model size(<**350KB**) by building neural architecture search (NAS) and post training quantization framework for edge devices. Improved model performance by ~**15%** by injecting noise during training to promote robustness.
- Improved accuracy by ~**3%** to reach ~**88%** accuracy for CIFAR10 on resource constrained devices using **joint optimization** of NAS and Hyperparameter optimization (HPO) inspired by CVPR 23's MA2ML using reinforcement learning.
- Conducted extensive debugging & rigorous tests for quantization, increasing test coverage from **65%** to **87%**.
- Introduced support for **10+** intricate ONNX operators and simulation of noise to ML compiler to enhance model inference on AI accelerator and **improve** accuracy by **at least ~5%**.
- Led development for edge inference framework of LLMs, developed integer inference for non-linear activation functions with <**2%** quantization loss. Finetuned vision transformers to achieve high performance for image classification.

Texas Instruments (TI): Test Engineer

Jan 2021- Jul 2022

Skills: Python, C++, git

Bangalore, India

- Reviewed large C++ code base, crafted scalable and efficient test program for production release. Resolved **50+** bugs, incorporated **20+** features to **enhance debugging** in the verification tool crafted at TI, **recognized** as top contributor.
- Constructed a parasitic extraction tool with user interface with python scripts. The tool **reduced** test hardware redesign **time, cost** by **30%** and better output correlation. Work published at TI conference.

RELEVANT PROJECTS

Adaptive Learning with Dynamic Batch Creation Using Near-Neighbors

Aug 2022 - May 2023

- Created an adaptive batch creation algorithm to create new batches using near neighbors of samples with highest gradients from previous batch like the paper from ICLR 2016. Attained ~**5%** faster convergence compared to random batches.
- Conducted extensive experiments with several thresholds & variations, monitored experiments using weights and biases.

Graph Based Recommender System

Jan 2023 - May 2023

- Performed comparative analysis based on scalability, precision, latency on YouTube dataset using: Pixie Random Walk (PRW), Random walk based embeddings and link prediction without GNN with the latter outperforming all methods.

Mechanistic Interpretability of Transformers

Aug 2022 - Dec 2022

- Trained a decoder only transformer with only attention layers from scratch inspired by ChatGPT. Analyzed attention scores via heatmaps, revealing insights including copying mechanism, skip gram behavior in ~**30%** attention heads.
- Performed extensive debugging and parameter tuning to achieve successful training of transformer model.

Recipe Generation from Videos

Jan 2023 - May 2023

- Engineered a recipe generation pipeline leveraging pretrained models. Executed parallel dense video captioning for event proposals, reduced redundancy with cosine similarity and elevated image analysis of frames using YOLO, harnessed BLIP for image captioning, Resnet for frame-to-feature conversion. Employed ChatGPT for caption summarization.
- Evaluated performance with BERTScore, achieved impressive results without model training with latency ~**5s**.