

Avi Singhal

avisinghal6@gmail.com|9452491966 |[linkedin.com/in/avi-singhal99/](https://www.linkedin.com/in/avi-singhal99/)Fremont, California|github.com/avisinghal6 |avisinghal6@gmail.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.: Machine Learning Model Development Engineer

Feb 2024 - Present

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

California, USA

- Leading development of edge inference framework of LLMs, devised integer inference & quantization methodology for non-linear layers with **<1%** quantization loss. Developing adaptive quantizer with float fallback to maintain high accuracy when integer computations lower quantized model accuracy.
- Prototyping & finetuning compressed versions of vision transformers, compressing state space-based MAMBA model. Developing architectures and search spaces for audio applications, object, and face detection.
- Model development and training for gauge reading applications.

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 and post training quantization framework for edge devices.
- 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.
- 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%**.

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 in python which **reduced** test hardware redesign **time, cost** by **30%**.

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.

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 with pretrained models. Executed parallel dense video captioning for event proposals, enhanced frames with YOLO, used BLIP for image captioning. Employed ChatGPT for caption summarization.
- Evaluated performance with BERTScore, achieved impressive results without model training with latency **~5s**.