

Professional Experience

AI Software Solutions Engineer Intel Corporation January'24 – Ongoing

- Developing **high-performance** deep learning **kernels** with **dynamic shape support** for Intel's next-generation GPU using SYCL, optimizing latency, memory bandwidth, I/O access, and compute utilization.
- Implemented an efficient **cumsum kernel**, achieving **2x perf improvement** over **IPEX**.
- Designed and **implemented** complex operations like **TopK** and media operators such as **Brightness and Contrast** as graphs in C++ using MLIR types and attributes, enabling efficient GPU execution.
- Developed a novel machine learning algorithm combining **VAEs** and **Diffusion Models** for NLP and CV.
- **Co-authored two papers**; one **submitted to CVR 2025** and actively seeking conferences for the second.

Senior ML Engineer Qualcomm Corporate R&D December'22 – January'24

- Spearheaded **ONNX optimizations** on Qualcomm's **AI100 accelerator**, achieving an **8.5% performance boost** for large language models (LLMs) like **ChatGLM2-6B** through **node-fusions, graph simplifications**.
- Enhanced **GPT model** efficiency by **2x** through caching Key-Value matrices and minimizing DDR reads/writes.
- Developed a Graph Neural Network algorithm to enhance compiler efficiency, resulting in a filed patent.
- **Led a three-member team** in optimizing and deploying the top 120 models from Hugging Face library.

ML Engineer Qualcomm Corporate R&D November'20 – November'22

- Engineered software modules in C++ & Python for AI/Deep Neural Network frameworks.
- Implemented auto-detection of post-processing in CV models, replacing them with **optimized kernels**.
- Achieved a **28.2% perf improvement** for (**BERT and variants**) through **Graphcore's packing strategy**.
- Enhanced operator support in the **GLOW compiler** for the **Cloud AI100 SDK**.

Patent and Publication

U.S. Patent application 18/330,253 and 18/500,014 (*Pending*)

- "Pre-Processing For Deep Neural Network Compilation Using Graph Neural Networks," filed on June 06, 2023.

Machine-Style Handwriting Generation with Diffusion *CVR 2025 Conference (Submitted)*

- Developed an algorithm leveraging VAEs and Diffusion models for handwriting generation.

Home Automation Using Panoramic Image Using IoT (*Published in 2018*)

Technical Blogging & Content Creation

Technical blogger

- Write in-depth technical articles on **Understanding GPUs** and **Parallel Programming with CUDA**.

YouTube Channel Co-Owner & Python Instructor

Tech Stack for Software Development and Machine Learning

- **Programming:** Python, C++, SYCL (DPC++), CUDA
- **Machine Learning Frameworks:** PyTorch, ONNX, ONNX Runtime
- **ML Domain & Techniques:** NLP, CV, Quantization, Pruning, Node Fusion, Graph Optimization, GNN
- **Others:** GPU Optimization, Git, Docker, GLOW (Machine Learning Compiler)

Education

Mumbai, IN IIT-Bombay July'18 - July'20

- M.Tech in Computer Science and Engineering, July 2020. CPI: **8.43** (on scale of 10).

Mumbai, IN University of Mumbai June'14 – June'18

- B.E. in Computer Engineering, June 2018. CPI: **9.07** (on scale of 10).

Master Thesis

- Computational Model to Understand and Predict Emotions. (2020)