




# VAISHNAVI PANCHAVATI

Mountain View, CA

✉ vpanchavati10@gmail.com    vaishnavi-panchavati    vaishnavi    portfolio

Work Authorization: Eligible to work in the U.S with valid **H4 EAD** (no sponsorship required)

## Education

### University of Massachusetts Amherst

Feb 2023 – Dec 2024

*Master of Science in Computer Science (Minor: Machine Learning), GPA : 4.0/4.0*

*Amherst, MA*

### National Institute of Technology Karnataka (NITK)

Aug 2014 – May 2018

*Bachelor of Technology in Electrical and Electronics Engineering, CGPA : 9.26/10*

*Surathkal, Karnataka*

**Relevant Coursework:** Advanced Natural Language Processing, Neural Networks, Advanced Algorithms, Reinforcement Learning, Distributed Systems, Data Structures and Algorithms

## Technical Skills

Languages: C, *C++*, *Python*, Go, JavaScript, HTML, CSS, Shell Scripting  
Frameworks & Libraries: *Spark*, *FastAPI*, *Flask*, *PyTorch*, *gRPC*, RESTful APIs, vLLM  
Tools & Platforms: *Docker*, *Kubernetes*, *AWS (EC2)*, Spark, Git, Linux  
Databases: *PostgreSQL*, *SQLite*, *MongoDB*  
ML and Systems: Model Distillation, *Distributed ML Inference*, Model Serving, *Inference Optimization*

## Work Experience

### Texas Instruments

Aug 2019 – May 2022

*Software Engineer (5G Transceiver)*

*Bangalore, Karnataka*

- Designed and implemented *C/C++ based functional test frameworks* for evaluating multiple embedded subsystems, ensuring *high system reliability* and precise feature validation.
- Developed *scalable Python based orchestration* to manage regression test workflows, *reducing manual validation time by 40%*.
- Migrated signal-processing logic from MATLAB to modular Python packages and *integrated it into a versioned local cloud test system, reducing datapath errors by 60%*

### Qualcomm

Sep 2018 – April 2019

*Software Engineer (4G/5G Small Cells)*

*Hyderabad, Telangana*

- Implemented features for LTE/5G NR PHY layer protocols in CATM1 and NB IoT devices in C++ and Python.
- Identified corner cases to *boost performance by up to 90%* in key tests like throughput and SINR for both FDD and TDD systems.

## Projects

### Model Serving System | *Python, FastAPI, Kubernetes*

April 2025 – Present

- Building a *distributed, fault-tolerant LLM* serving system with *gRPC and FastAPI* based control APIs, supporting autoscaling, model caching, fault tolerance, dynamic model loading, and Redis backed persistence.
- Deployed the system on AWS with Docker containers and prometheus based monitoring, achieving *70%* reduction in operational overhead

### LLM Quantization & Inference Optimization | *Python, Spark*

Feb 2025 – April 2025

- Implemented *GPTQ* based 4-bit *quantization pipeline* with on the fly dequantization, reducing model size with minimal degradation in performance.
- Leveraged *Spark* for efficient parallel data preprocessing prior to quantization pipeline and added custom quantized linear layers and packed state compression, cutting memory footprint by *50%* for a 70GB model.

### Stock Bazaar Application | *Flask, Docker, AWS, Python*

Feb 2023 — April 2023

- Built a *fault-tolerant, scalable microservices* based stock trading system using *thread-per-request concurrency and locking* to handle coordinated client requests.
- Implemented *server push cache consistency* and leader based Order replication with container deployment on *AWS EC2*, *reducing* latency by *5ms*.

### Quote-MI | *Pytorch, LLM, NVIDIA A100*

Feb 2024 – April 2024

- Fine-tuned *BERT, RoBERTa, GPT-2, T5* on manually curated quote datasets for multi-label classification obtaining *70% accuracy* and evaluated T5 and Gemma using zero and few-shot prompting.
- Deployed GPU accelerated BERT* with TensorFlow Serving and Docker, *reducing* inference latency by *5s*.

### miniLlama | *Pytorch*

Jan 2025 – Feb 2025

- Built and fine tuned a compact Llama 2 model with custom LoRA for sentence completion and classification, achieving *30%* performance improvement on CFIMDB dataset