

Aru Sharma

+91-7452029206 | arusharmazxx000@gmail.com | [personal-page](#)

EDUCATION

Panjab University

Bachelor of Engineering in Information Technology

Chandigarh, India

Oct. 2022 – Present

Little Scholars

Mathematics and Computer Science; (Central Board of Secondary Education); 95%

Kashipur, India

May 2019 – Jun 2021

EXPERIENCE

ML Engineering Intern

Nannie.ai

Sep 2025 – Present

London, UK

- Working on testing and deploying **SOTA Vision algorithms** for classification, segmentation and pose detection for animals.
- Deployed **OSS text to video generation** models for in-house testing and benchmarking against **Veo3**.

AI Engineer (contract)

Deskree Inc.

Nov 2025 – Jan 2026

Toronto, Canada

- Worked on **Tetrix** and building AI agents for your infrastructure including cloud services like AWS .
- Developed **Tetrix CLI** - a tool to review architecture, and security issues and enforce code quality.

OSS Developer at Google Summer Of Code

Mifos Initiative

Jun 2025 – Sep 2025

Seattle, WA

- Developed a **multi-agent bot** let users know the status of **Jira tickets**, questions related to **Slack discussions**.
- Developed a **full-stack web** application using **FastAPI, NextJs and Firestore** as database and Auth client.

OSS Developer at Summer Of Bitcoin

Bitcoin-dev-project

May 2025– Aug 2025

Manhattan , NY

- Designed and prototyped **AI-assisted coding tools for Bitcoin** using small language models and domain-specific Retrieval-Augmented Generation (RAG).
- Improved **data pipelines** to ingest bitcoin related knowledge from Bitcoin Conference talks, correct and summarize them using LLMs

Software Engineering Intern

CNCF WasmEdge

Sep 2024 – Dec 2024

Austin, TX

- Developed a **RAG** based chatbot for code assistance using **opensource LLMs** with **WasmEdge runtime**.
- Created a pipeline to ingest data from **Github repository**, augmented it using QnA pairs, summary and then embed this into a **Qdrant vector database**.

PUBLICATIONS

Robust Speech Emotion Recognition Across Diverse Datasets: A Comparative Study of Deep Learning and Transformer-Based Approaches for VoIP Applications,16th International Conference on Computing Communication and Networking Technologies, 2025 Accepted for publication.

PROJECTS

Hidden-state-extractor | *vLLM, Pytorch*

- Created a custom plugin in vLLM for hidden state extraction from LLMs for Mechanistic Interpretability.
- Uses Pytorch Forward hooks for extraction and sub-processes for consumption via CUDA IPCs.

Memory-Augmented-Agents | *Long Horizon Reasoning, Memory Routers*

- Built a chat agent designed with long-term memory capabilities to make them personalised and adaptive.
- Leverages a semantic, keyword and reranking mechanism for intelligent retrieval and a sophisticated mechanism for continuous memory consolidation.

ACHIEVEMENTS

- Ranked 15 globally on the **NTIRE Image Dehazing and Denoising challenge** at **CVPR 2024**.