

# Ranvir Virk

ranvirsv@gmail.com | <https://www.linkedin.com/in/ranvir-singh-virk> | <https://www.github.com/ranvirsv>

## Education

---

**University of Technology Sydney**  
Master of Data Science and Innovation

2026 – Present

**Indiana University Bloomington**  
Master of Science in Computer Science  
Bachelor of Science in Computer Science

08/2021 – 05/2025

GPA: 3.6

GPA: 3.8

## Skills

---

**Languages:** Python, R, SQL

**Machine Learning:** PyTorch, TensorFlow, scikit-learn, YOLO, OpenCV, spaCy, NLTK, Transformers

**Data Analysis:** Pandas, NumPy, Matplotlib, Time Series Forecasting, Predictive Modeling, SciPy

**AI & LLM Tooling:** CrewAI, LangChain, ChromaDB, Vector Databases, Pydantic

## Experience

---

**Software Engineer - AI Intern** | Postman | New York, NY 06/2025 - 09/2025

- Architected and delivered an end-to-end agentic AI workflow, defining task decomposition, agent provisioning, execution loops, and feedback reconciliation using CrewAI multi-agent orchestration and LangChain pipelines.
- Integrated ChromaDB (vector database) with OpenAI's text-embedding-3-large, and jina-embeddings-v2-base-code models to power semantic context retrieval, reducing agent query latency by 40% and improving relevance of retrieved code snippets.
- Implemented a codebase indexing pipeline and chunking mechanism using Tree-sitter for AST-style parsing and metadata extraction; generated and stored 10k+ fine-grained embeddings to support sub-file granularity in agent prompts.
- Automated CI/CD with GitHub Actions and Postman Flows, enabling zero-downtime multi-region rollouts and reducing manual release effort by 2 hours per cycle.

**Research Assistant - Data Science** | Indiana University | Bloomington, Indiana 02/2023 - 05/2025

- Developed and optimized dual-channel particle tracking pipelines by integrating MATLAB's Deep Learning Toolbox, improving tracking accuracy by 20%.
- Preprocessed and cleaned 20M+ row geological and imaging datasets using Python (Dask) and SQL; built reliable ETL workflows to ensure data integrity and scalable storage in PostgreSQL.
- Built and validated machine learning models (CNNs for particle classification, LSTM/ARIMA for time-series forecasting) using scikit-learn and TensorFlow; applied hyperparameter tuning, cross-validation, and statistical diagnostics to improve model reliability.

**Machine Learning Engineer Intern** | Nod.ai - AMD | Santa Clara, California 05/2023 - 08/2023

- Implemented production-grade C++ extensions for PyTorch linear algebra operators, integrated into the CI/CD pipeline, improving computational throughput and reducing GPU kernel latency by 20%.
- Debugged and refactored 10+ IREE compiler passes, developed custom MLIR dialects for INT8 quantization, and implemented LLVM IR optimization passes to fuse matrix multiplications, achieving performance uplift in inference workloads while lowering model memory footprint.
- Contributed upstream to the IREE compiler by fixing a CPU-side 16-bit fmaxf bug, authoring an ExpandF16MaxFToF32 pass in MLIR and adding LIT tests, ensuring reliable 16-bit float support across CPU backends.

## Projects

---

### Deep Dive into Vision Transformers

- Utilized YOLOv5 to generate facial and body crop annotations across 15k Kaggle images, enabling evaluation of ViT sensitivity to data quality and yielding an 18% precision boost.
- Performed systematic ablation on augmentation techniques and dataset scaling, quantifying a 5% accuracy improvement per additional 1k high-fidelity samples and stabilizing ViT training dynamics.
- Fine-tuned Vision Transformer models, applying NMF-based deep feature factorization and Grad-CAM to visualize attention shifts under varied data conditions.

### Safe-Feed Recommender: Multi-Task Learning Evaluation

- Addressed competing recommender objectives by deploying a Shared Bottom Network to quantify negative transfer, establishing a baseline that drove a 5% performance gain.
- Resolved cross-task gradient conflicts by architecting a Multi-gate Mixture-of-Experts (MMoE) model, utilizing task-specific routing to effectively balance competing signals and boost multi-task prediction accuracy by 4%.