# Ayush Luhar

13ayushjluhar@gmail.com | linkedin.com/in/ayushluhar
13aluminium.netlify.app | github.com/13aluminium | (562)-819-4405 | Long Beach, CA

## Education

| | |
|---|---|
| **California State University, Long Beach** – Masters of Science in Computer Science | Aug 2025 - May 2027 |
| **Charotar University of Science and Technology** – B.Tech in Computer Engineering | Oct 2021 - May 2025 |

## Skills

**Languages:** Python, Java, SQL

**ML & Data:** PyTorch, TensorFlow, scikit-learn, Hugging Face, FAISS, Pandas, NumPy

**Systems & MLOps:** PySpark, Ray, MLflow, Docker, Distributed Training (DDP, NCCL)

**Software Engineering:** Git, CI/CD pipelines, testing, debugging, model evaluation, A/B testing

**Concepts:** Ranking & Retrieval, Transformers, Embeddings, Tokenizers, Overfitting/Regularization, ML Pipelines

## Experience

**Research Assistant –** California State University, Long Beach (Long Beach, USA)     Sep 2025 – Present

- Implemented scalable ML pipelines using PySpark to process **100K+ behavioral records**, reducing end-to-end processing time by 60%.
- Worked with stakeholders to develop neural ranking models, achieving **0.89 NDCG@10** in large-scale retrieval tasks.
- Built automated **50+ feature engineering** workflows and evaluated models using offline metrics and controlled experiments.
- Applied responsible AI practices by analyzing bias, data distribution, and model behavior to improve fairness and reliability of ranking systems.

**Research Intern –** CHARUSAT (Changa, India)     Dec 2024 – May 2025

- Bootstrapped a **2.5M bilingual sentence-pair dataset** using RLHF-based data generation and filtering in early low-data regimes.
- Fine-tuned **T5 transformer models** on the curated dataset, achieving **+42% BLEU** and **−46% WER** improvements.
- Scaled training via distributed multi-GPU pipelines (PyTorch DDP + NCCL), reducing runtime from **72h → 18h**.
- Applied mixed-precision optimization and tracked **200+ experiments with MLflow** to ensure reproducible post-training.

**Machine Learning Intern –** HyperVect (Gandhinagar, India)     May 2023 – June 2023

- Built GPU-accelerated pipelines for **image and audio datasets (50K+)** using PyTorch and TensorFlow, improving it by **3.2×**.
- Developed LLM-based multimodal data workflows with inference caching, reducing latency by **25%**.
- Scaled pipelines using Ray-based parallel processing, cutting runtime by **~40%**.

## Projects

**Java 8 to Java 17 Migration Agent |** Python, LLMs, Agentic Systems huggingface.co/datasets/13Aluminium/JAVA_8-JAVA_17

- Built a master–worker agentic system for Java 8 to Java 17 migration using LLMs and AI coding tools (Claude, GitHub Copilot), with validation and review.
- Created a supervised dataset from **100+ GitHub** repositories through manual Java version migration.
- Achieved strong migration quality with **BLEU 72.26**, ROUGE-L 0.858, and **69.54% XML recall**, validating structural accuracy.

**Content Recommendation Engine |** Python, PyTorch, FAISS

- Built collaborative filtering and ranking models on 1M+ user–item interactions, achieving **0.85 Precision@10**.
- Implemented approximate nearest neighbor (ANN) retrieval with FAISS, reducing latency from 2s to 50ms.
- Designed offline evaluation and A/B testing framework, improving click-through rate by 12%.

**Sanskrit-Morpheme-Tokenizer |** Python, NLP, Hugging Face  huggingface.co/spaces/snskrt/Sanskrit-morpheme-Tokenizer

- Designed a large-scale NLP pipeline over a **5GB corpus (600K+ morphemes)** using PySpark and distributed processing.
- Developed a custom tokenizer expanding vocabulary to 500K+, improving downstream transformer performance by 23%.
- Deployed on Hugging Face Spaces, serving 5K+ monthly users with **<200ms p95 latency**.

## Publications

**Luhar Ayush**, et al. (2024). "From Pixels to Text: A Comparative Evaluation of Convolution Neural Network and Recurrent Neural Network Models in Image Processing" *International Conference on Smart Trends in Computing and Communications, Springer, Singapore.*     doi.org/10.1007/978-981-97-1326-4_32

**Vyakaran Dataset**, Hugging Face     doi.org/10.57967/hf/6073