

# HARSH MANKODIYA

• +1 (602)-517-6750 • hmankodi@asu.edu • linkedin.com/in/harshmankodiya • github.com/hmankodiya

## Education

<b>Arizona State University</b> <i>Master of Science, Computer Science: GPA 3.70</i> <i>Courses: NLP, Statistical Learning, Artificial Intelligence, Data Mining</i>	<b>August 2023 - May 2025</b> <i>Tempe, USA</i>
--	--

<b>Institute of Technology, Nirma University</b> <i>Bachelor of Technology, Computer Science Engineering</i> <i>Courses: Machine Learning, Deep Learning, Data Structures, Linear Algebra, Calculus, Probability and Statistics</i>	<b>August 2019 - May 2023</b> <i>Ahmedabad, India</i>
---	--

## Professional Experience

<b>Nomic AI</b> <i>Machine Learning Engineer</i>	<b>July 2025 - Present</b> <i>NYC, USA</i>
---	---

- Training Transformer-based document parsing models on complex, high-density datasets and building scalable pipelines for evaluation, benchmarking, and deployment, driving agentic document AI capabilities and multimodal understanding.

<b>Cellino Biotech</b> <i>Machine Learning Intern</i>	<b>May 2024 - August 2024</b> <i>Cambridge, USA</i>
--	--

- Developed a proof of concept for a central embedding model for patch selection, anomaly detection, cell segmentation and cell classification.
- Fine-tuned **DinoV2** using **Vision Transformer** based heads for downstream 2D-segmentation tasks, achieving average F1-Score of **82%**. Utilized **Weights & Biases** for experiment tracking, artifact logging and hyperparameter sweep.
- Streamlined dataset preparation for ML pipelines by integrating **Google Cloud APIs** with **PyTorch Dataset** utilities to convert **Zarr arrays** to **Tensors** and implemented **local caching** to boost data loading throughput.
- Decomposed embeddings using **t-SNE** and **PCA**, leveraged **GMM clustering** for zero-shot cell artifact detection. Optimized with **JAX** for on-device GPU inference, achieving a **10x** speedup.
- Containerized the inference pipeline with **Docker**, enabling real-time data processing and developed automated testing using **PyTest** and **BASH**.

<b>Lens Lab, Arizona State University</b> <i>Research Assistant</i>	<b>August 2023 - May 2024</b> <i>Tempe, USA</i>
--	--

- Integrated eXplainable AI techniques with RL agents in **Gymnasium** environments to enhance decision explainability.
- Trained **Proximal Policy Optimization (PPO)** using **StableBaselines3** with **VAE**-based feature extraction for image stream processing.
- Leveraged **Multi-Modal CLIP** models for **zero-shot segmentation** and concept sampling for policy rollouts.
- Published findings at **NeurIPS 2024 SATA Workshop**, focusing on explainability in robotic decision-making.

<b>Bosch</b> <i>Research Intern</i>	<b>January 2023 - May 2023</b> <i>Bangalore, India</i>
--	---

- Designed a **GradCAM**-based **Knowledge Distillation** pipeline to train **SegNet** and **U-Net** models for image segmentation, achieving **IoU** scores exceeding **85%** on multiple datasets using **NVIDIA DGX A100** systems.

<b>Samyak Infotech Pvt Ltd</b> <i>Machine Learning Intern</i>	<b>May 2022 - July 2022</b> <i>Ahmedabad, India</i>
--	--

- Fine-tuned **BERT**-based **LayoutLM**, for structured information extraction from scanned business invoices, achieving an **F1-score of 81%**.
- Summarized long invoices by using **T5 Transformer** and obtained a strong **BERTScore of 0.95**
- Performed **KMeans clustering** on layout-aware embeddings to organize invoices into structurally similar groups reducing manual annotation effort by **10%**.

## Relevant Projects

<b>Continual Knowledge Expansion for Book Retrieval Systems</b>   <i>Python, FAISS, HuggingFace</i>	<b>Dec 2024</b>
• Developed a pipeline to continuously ingest and embed new book content, enabling a <b>RAG</b> system to stay updated.	
• Utilized <b>DPR</b> and <b>FAISS</b> to extract new knowledge chunks from book sections, improving retrieval precision.	
• Implemented a <b>Curriculum Learning</b> -based pretraining strategy, reducing perplexity of <b>Pythia-2.8B</b> and <b>LLaMA2-7B</b> by <b>20%</b> compared to vanilla pretraining.	

## Technical Skills

<b>Languages</b>	- Python, C++, Shell, Docker, Git
<b>ML Frameworks</b>	- PyTorch, HuggingFace, Jax, TensorFlow, scikit-learn, XGBoost, Stable-Baseline3, Gym, LangChain, LangGraph, Ollama
<b>Python Libraries</b>	- NumPy, SciPy, Pandas, OpenCV, Pillow, Zarr, Dask, Seaborn, Matplotlib, Plotly, W&B, MLFlow, PySpark
<b>ML Techniques</b>	- LLMs, RAG, Knowledge Distillation, Reinforcement Learning, SSL, CLIP, Image Captioning, Image Classification, Image Segmentation, VAE, GANs, Style Transfer