

Akshay Goindani

asgoindani@gmail.com | 412-696-2588 | [LinkedIn](#) | [Google Scholar](#) | [Personal Website](#)

EDUCATION

Carnegie Mellon University, School of Computer Science	December 2024
Master of Science in Artificial Intelligent Information Systems, GPA: 4.17/4	Pittsburgh, PA
International Institute of Information Technology	July 2022 (Graduated)
B.Tech (Honours) & M.S. (Research) in Computer Science & Engineering, GPA: 8.99/10	Hyderabad, India

WORK EXPERIENCE

Voyage AI	February 2025 - Present
<i>Founding Research Engineer</i>	Palo Alto, CA

- Led the development of instruction-following rerankers, resulting in the release of [voyage-rerank-2.5](#) models, first in the industry that follow instance-level instructions, with improved performance on all benchmarks.
- Driving the development of the next-generation of multimodal embedding models, further enhancing the capabilities of voyage-multimodal-3, a state-of-the-art embedding model.
- Optimized multimodal training infrastructure with improved load-balancing and batch scheduling, boosting throughput by 40–50%.
- Advancing model adaptability through continual learning and model-merging approaches, ensuring long-term scalability and consistent performance across evolving tasks.

<i>Machine Learning Intern</i>	June 2024 - September 2024
--------------------------------	----------------------------

- Developed large-scale data curation, de-duplication, and quality filtering for token-efficient training. Improved data quality leads to the release of state-of-the-art [voyage-3](#) text-embedding models.
- Synthetic multimodal data generation for training embedding and reranking models. Generated high-quality synthetic data improves Vision-Language embedding models, and led to the release of [voyage-multimodal-3](#).

ExaWizards AI platform	July 2022 - June 2023
<i>Machine Learning Engineer</i>	Hyderabad, India

- Deployed a pipeline for Temporal Activity Localization in videos using Natural Language Description.
- Enhanced the efficiency of the object recognition module, improving the inference speed of the platform by 30%.

Amazon International Machine Learning Team	May 2022 – July 2022
<i>Applied Scientist Intern</i>	Bangalore, India

- Built Reinforcement Learning based Attribute Extraction Model to predict missing values in product descriptions.
- Leveraged trajectories of trained online RL agents to design Decision Transformer, improving recall by 8.5 points.

SELECTED PUBLICATIONS

- Liang, Paul Pu*, **Akshay Goindani***, Talha Chafekar, Leena Mathur, Haofei Yu, Ruslan Salakhutdinov, and Louis-Philippe Morency. "Hemm: Holistic evaluation of multimodal foundation models." *Advances in Neural Information Processing Systems* 37 (2024): 42899-42940. [NeurIPS 2024](#) (* Equal Contribution)
- **Goindani, A.** and Shrivastava, M., 2021. A Dynamic Head Importance Computation Mechanism for Neural Machine Translation. [International Conference on Recent Advances in Natural Language Processing](#).
- Sivaprasad, S.*, **Goindani, A.***, Fritz, M. and Gandhi, V., Class-wise Domain Generalization: A Novel Framework for Evaluating Distributional Shift. [NeurIPS 2022 Workshop on Distribution Shifts](#). (* Equal Contribution)

RESEARCH EXPERIENCE

Language Technologies Institute, Carnegie Mellon University	August 2023 – May 2024
<i>Research Student Advisor: Prof. Ruslan Salakhutdinov, Prof. LP Morency, Prof. Paul Pu Liang</i>	Pittsburgh, PA

- Co-led the development of HEMM: A comprehensive benchmark to evaluate generative multimodal models on their ability to perform tasks requiring reasoning, fine-grained alignment and cross-modal understanding. [NeurIPS 2024](#)

Language Technologies Research Centre	March 2019 – July 2022
<i>Research Assistant Advisor: Professor Manish Shrivastava</i>	Hyderabad, India

- Developed a dynamic head importance computation mechanism for LLM based Neural Machine Translation.

- Importance-score based pruning of Attention Heads to minimize redundancy. Proposed approach achieves better performance in low resource conditions, with fewer parameters, and learns better word alignment. [\[Paper\]](#).

Centre for Visual Information Technology

June 2020 – Dec 2022

Research Assistant | Advisor: [Professor Vineet Gandhi](#)

Hyderabad, India

- Proposed a novel class-wise domain generalization framework for evaluating distributional shift. Developed an Iterative Domain Feature Masking method that achieves SOTA performance [\[NeurIPS 2022, Paper\]](#)

Sungkyunkwan University

Dec 2020 – May 2021

Research Intern | Advisor: [Professor Hogun Park](#)

Seoul, South Korea

- Proposed an approach for Augmenting Knowledge Graphs to Question-Answering Systems, using Graph Neural Networks (GNN), to impart commonsense knowledge to QA models, for Open-Domain Question-Answering task

KEY PROJECTS

Needle: A PyTorch-like training framework | LLM Training

- Built a deep learning framework in C++ and Python, supporting tensor operations (convolution, broadcasting, transpose) for training neural networks.
- Optimized training on NVIDIA GPUs with custom CUDA kernels for reduction, matrix multiplication, and other critical operations.
- Designed computation graph execution with backpropagation via topological sort, implementing efficient forward and backward passes for numerous operations.
- Extended capabilities with gradient checkpointing and FP16 mixed-precision training, significantly improving memory efficiency.

Synthetic Data Generation for Off-Policy Preference Optimization | Post-training, Reasoning

- Generated a synthetic preference dataset using a model pool for post-training vision language models for better abductive reasoning. [\[Dataset\]](#)
- Used CLIP for scoring the model generations, the scores are then used to determine the preference order.
- Fine-tuned [PaliGemma-3B](#) on the preference data using Direct Preference Optimization [\[Code\]](#).

MinBERT Classifier | Large Language Models (LLMs)

- Implemented the BERT architecture in PyTorch from scratch with essential components such as Positional Embeddings, Multi-Head attention, etc., and trained it for sentiment classification. Also implemented the AdamW optimizer for training/fine-tuning the LLM [\[GitHub\]](#)

Learning Bilingual Word Embeddings with Minimal Bilingual Data | Natural Language Processing

- Implemented unsupervised method to learn bilingual word embeddings for English & Italian, using a common embedding space via parameterized linear transformation
- Incorporated supervised learning with few known translation pairs to bilingual dictionary iteratively [\[GitHub\]](#)

HONORS & AWARDS

- **Dean's List Award** for Academic Excellence (Top 5% of all students), 2018 - 2021
- [ACM ICPC](#) 2019 (Challenging Programming Contest, Worldwide) - Secured rank in Top 100 at Regional Level
- Recipient of competitive [MITACS Globalink Graduate Fellowship](#) for \$15,000

SKILLS

ML Training/Inference Frameworks: PyTorch, Accelerate, Deepspeed, Megatron, VeRL, vLLM

ML/AI Techniques: Transformers, Diffusion Models, Mixture of Experts

Languages: Python, C/C++, Matlab, Bash, HTML/CSS, JavaScript

COURSES

Artificial Intelligence	Deep Learning Systems, Advanced Natural Language Processing, Multimodal Machine Learning, Machine Learning, Computer Vision, Vision Learning and Reasoning, Artificial Intelligence
Computer Systems	Database Systems, Operating Systems, Software Engineering, Digital Signal Analysis and Applications, Computer Graphics
Math	Discrete Maths, Linear Algebra, Probability & Statistics, Complex Analysis, Multivariate Analysis, Formal Methods
Security & Networks	Advanced Computer Networks, Principles of Information Security
Algorithm & Programming	Data Structures and Algorithms, Computer Programming