# MUDIT DHAWAN

Pittsburgh, PA | 412-721-8074

mdhawan@cs.cmu.edu | linkedin.com/mudit-dhawan | mudit-dhawan.github.io/ | Google Scholar

## EDUCATION

**Carnegie Mellon University** | Master of Science in **Machine Learning**                    December 2025
Courses: Probabilistic Graphical Models, Convex Optimization, Advanced Introduction to Machine Learning (PhD)
Independent Study with **Prof. Michael J Tarr**: Analysis of Category Representation in Pre-trained Models and Brain

**IIIT Delhi, India** | Advisor: *Prof. Ponnurangam Kumaraguru* [thesis]                    May 2022
BTech. in Electronics and Communication Engineering                                    **Dept Rank: 7/81**

## EXPERIENCE

**Adobe** | Machine Learning Engineer Intern                                    May 2025 - Aug 2025
- Built a scalable training pipeline using **Supervised Finetuning** followed by **Group Relative Policy Optimization** (GRPO) to train **SLM-based conversational agents** for **tool orchestration**, with modular support for LoRA, quantization, and vLLM, achieving over **95%+ reduction in time-to-market** for production-ready agents.
- Developed a role-based **synthetic data generation pipeline** using LLMs (GPT, LLaMA 70B), incorporating an **MCP**-like framework to streamline tool integration and simulate diverse user interactions for grounded data.

**Microsoft Research** | Research Fellow                                    July 2022 - July 2024
*Extreme Classification (XC) Group, Advisor: Dr. Manik Varma*
- Deployed a novel reformulation of Query AutoComplete as an XC task on Bing AI Chat and Search platforms. Led to **1% gain in Click Through Rate (CTR)** in production. On Bing AI Chat Platform, it led to an **8% increase in Keystrokes Saved Per Query**, and **6% increase in acceptance rate** for long queries [publication: ICLR'24].
- Highlighted in-efficiency of current pair-wise rankers in recommendation systems and devised an algorithm with **30x lower latency** than SOTA theoretical baseline, **10% more accurate** than production system on an offline set, and led to **0.5% increase** in **Impression Yield Click Yield** in en-markets in online A-B tests.
- Executed a novel scoring method to use GPT-4 as a large scale oracle to de-noise hard-negatives for retriever training. Led to **30x increase in throughput**, reduced API costs and trained retriever led to **0.6% absolute gain in clicks** during **online A-B tests** in English markets.

**IIIT Hyderabad** | Research Assistant                                    August 2019 - June 2022
*Precog Research Group, Advisor: Prof. Ponnurangam Kumaraguru*
- Proposed novel algorithms for multimodal fake news detection [publications: MMAsia'21, SNAM'24] [code].
- Introduced **Multi-Task Learning Framework** for Bail prediction (**2% more accurate** than SOTA) with an auxiliary extractive summarization task for grounding predictions in Hindi documents [publication: ACL Findings'22].

## SELECTED PUBLICATIONS                                    CITATIONS: 145 | H-INDEX: 4

**Accurate and Efficient Cross-encoders for Ranking** *Neural Information Processing Systems (NeurIPS) ENLSP 2024 Workshop, Empirical Methods in Natural Language Processing (EMNLP) WiNLP 2024 Workshop*

**Enhancing Tail Performance in Extreme Classifiers by Label Variance Reduction** *International Conference on Learning Representations (ICLR) 2024*

**HLDC: Hindi Legal Documents Corpus** *Association for Computational Linguistics (ACL) Findings 2022*

## SELCTED HONORS AND AWARDS

Granted a **US Patent** for optimized ranking model: *CROSS-JEM: Cross-encoder Joint Efficient Modeling for ranking in large-scale search and recommendation systems*.

## SELECTED PROJECTS

**Multi-Intent Session-Based Recommendation Systems by LLMs**: Performed task-specific distillation leveraging Low-Rank Adaptation of GPT-4 to increase **inference speedup by 100x** by using smaller LLMs with minimal loss in relevance and diversity of predictions.

## SKILLS

**Languages**: Python, Java, C/C++, BASH; **Software**: PyTorch, Tensorflow, HuggingFace, Keras, MATLAB
**ML/DL Skills**: Retrieval Augmented Generation (RAG), Multimodal ML, Large Language Models (LLMs), Natural Language Processing (NLP), Diffusion, Computer Vision, Audio Modelling, Distillation, Retrieval, Search, Recommendation and Ranking, Extreme Classification (XC), Generative Models, Data Structures and Algorithms.