

Soumen Kumar Mondal Male
Center for Machine Intelligence and Data Science (CMIInDS) M.S. by Research
Indian Institute of Technology Bombay, India DOB: 12/08/1996
Specialization: Data Science and Artificial Intelligence Mobile: +91 8759240557
Email-ID : mondalsoumen00@gmail.com or soumenkm@iitb.ac.in Website: <https://soumenkm.github.io>

Qualification	Specialization	Institute (all located in India)	Year	CPI/% (R)
MS by Research	Data Science & Artificial Intelligence	Indian Institute of Technology, Bombay	2023 - 2026	9.83 (1 st)
MTech	Structural Engineering	Indian Institute of Science, Bangalore	2018 - 2020	9.30 (4 th)
BTech	Civil Engineering	Jadavpur University, Kolkata	2014 - 2018	8.74 (6 th)
12 th Standard	Science (PCM), Languages	Haldia High School (West Bengal Board)	2012 - 2014	89.60 (1 st)
10 th Standard	Science, Arts, Languages	Haldia High School (West Bengal Board)	2011 - 2012	86.85 (2 nd)

Technical Skills

- **Programming & Scripting Languages:** Python, C, C++
- **Tools and Technologies:** PyTorch, HuggingFace, LangChain, TensorFlow, Scikit-Learn, LaTeX, Git, Linux

List of Publications, [Google Scholar ID: 154CUKcAAAAJ](#)

1. **Mondal, Soumen K.**, Sen, S., Singhania, A., & Jyothi, P. (2025). Language-Specific Neurons Do Not Facilitate Cross-Lingual Transfer. In Proceedings of the InsightsNLP in NAACL 2025 (oral). [ArXiv: 2503.17456](#).
2. Sona, SE., **Mondal, Soumen K.**, Sen, S., Singhania, A., & Jyothi, P. (2025). LoFTI: Localization and Factuality Transfer to Indian Locales. In Findings of the ACL 2025. [ArXiv: 2407.11833](#).
3. **Mondal, Soumen K.**, Varmora, A., Chanda, P., & Ramakrishnan, G. (2025). FairPO: Robust Preference Optimization for Fair Multi-Label Learning. Under review in NeurIPS 2025. [ArXiv: 2505.02433](#).

Work Experience

- **Fujitsu Research, Bengaluru, India**

(AI Research Intern, Project Title: Proactive RCA – A Hypergraph-Abstracted and Mamba-based RAG Framework for Causal Error Prediction in Warrior System Logs) (May 2025 - July 2025)

- Engineered a novel **ProactiveRCA** framework to accurately forecast future system errors by analyzing partial log data (20-30% of log). Designed and implemented a multi-stage pipeline using **Hypergraphs** and **PageRank** algorithm to compress massive logs (**500M+** tokens) into a format compatible with **Mamba**.
- Implemented domain-adaptive pre-training using **LoRA** to adapt the **Mamba** model to the non-natural language syntax of logs, improving the BLEU score for next-token prediction by 42%.
- Modified the **InstructRAG** pipeline by aligning **hyper-GNN** and text embeddings using **Cross Attention** module, achieving a semantic correctness score of 4.4/5.0 when evaluated by **LLM-as-a-Judge**.

- **General Electric (GE Vernova), Bengaluru, India**

(System Value Optimization Engineer) (Aug 2020 - July 2023)

Engineered end-to-end wind turbine load optimization strategies by implementing **CatBoost** model for predictive load assessments, reducing extreme load by 70%, leading to a **GE Spotlight Impact Award**.

M.S. by Research in Data Science and Artificial Intelligence at IIT Bombay

- **Controlling Large Language Models for Low-Resource Languages using RAG and Preference Optimization**

(M.S. Thesis, Advisor: Prof. Preethi Jyothi, CSE, IIT Bombay) (Spring 2025 - Present)

- Designed and implemented a data-efficient framework to improve low-resource Neural Machine Translation by fine-tuning **Llama 3.1** with **DPO**. Exploring the explainability of DPO for 70% BLEU score degradation on **Flores-200** dataset. [GitHub: soumenkm/ModelEditing](#)
- Proposed a unified **agentic RAG** and **DPO** architecture for multilingual QA on the **MCPQA** dataset to ensure factually grounded responses in a cross-market context. [GitHub: soumenkm/MarketQA](#)

- **Improving Downstream Task Performance in Multi-lingual LLMs by Intervening Language Specific Neurons**

(M.S. Thesis, Advisor: Prof. Preethi Jyothi, CSE, IIT Bombay) (Autumn 2024)

Investigated the mechanistic interpretability of language-specific neurons in **LLMs** to assess their utility for cross-lingual transfer. Demonstrated through neuron specific **LoRA** fine-tuning and test-time interventions that neurons are *polysemantic*, yielding 1% improvement [paper 1]. [GitHub: soumenkm/LangSpecificNeurons](#)

- **IIT Bombay - Amazon Collaboration: Localizing Text Across Domains Using RARR Attribution Technique**
(M.S. R&D Project, Advisor: Prof. Preethi Jyothi, CSE, IIT Bombay) (Grade: 10, Spring 2024)
Developed LoFTI, a novel benchmark of **1,100** factual entity pairs across **99** categories, to evaluate the factual and localization capabilities of LLMs in a cross-geographical context. Utilized RARR attribution to enhance the performance, boosting factual correctness by **13%** [paper 2]. [GitHub: soumenkm/RnD_Project](#)
- **Cross-lingual Factual Knowledge Transfer in Multi-lingual Language Models**
(M.S. Seminar, Advisor: Prof. Preethi Jyothi, CSE, IIT Bombay) (Grade: 10, Spring 2024)
Analyzed factual knowledge representation and transfer in **multilingual BERT** and **XLM-R** models across **53** languages using **Probeless** method to inspect internal neuron activity. [GitHub: soumenkm/TracingRootFacts](#)

Machine Learning Course Projects at IIT Bombay

- **FairPO: Robust Preference Optimization for Fair Multi-Label Learning**
(Course Project, Optimization for ML, Prof. Ganesh Ramakrishnan, CSE, IIT Bombay) (Spring 2025)
Proposed FairPO, a novel framework to address label fairness in Multi Label Classification using **DPO**, **CPO**, **SimPO**, and **GRPO**, achieving a **3.44%** mAP gain on **MS-COCO** dataset [paper 3]. [GitHub: soumenkm/FairPO](#)
- **Vision Transformer (ViT) Model Fine-Tuning with MillionAID Dataset using LoRA**
(Course Project, Advanced Deep Learning for CV, Prof. Biplab Banerjee, CSRE, IIT Bombay) (Autumn 2024)
Implemented **LoRA** from scratch for efficient fine-tuning of a **DINO-ViT** on **Million AID** dataset, achieving **97%** classification accuracy. [GitHub: soumenkm/IITB-GNR650-ADLCV/CodingProject](#)
- **Learning to Classify Images under Noisy Labels using Turtle**
(Course Project, Advanced Deep Learning for CV, Prof. Biplab Banerjee, CSRE, IIT Bombay) (Autumn 2024)
Achieved **88%** accuracy on **CIFAR-100** with **40%** label noise by implementing a **CLIP & DINO-ViT** framework for label denoising and subsequent supervised fine-tuning. [GitHub: soumenkm/IITB-GNR650-ADLCV/Project1](#)
- **Zero Shot Learning (ZSL) for Image Classification on AwA2 Dataset**
(Course Project, Advanced Deep Learning for CV, Prof. Biplab Banerjee, CSRE, IIT Bombay) (Autumn 2024)
Designed a ZSL pipeline with ViT, FastText and **Class Normalization**. Attained **40%** test accuracy on AwA2 with a challenging **50:50** train-test split. [GitHub: soumenkm/IITB-GNR650-ADLCV/Project2](#)
- **Deep Learning based System to Estimate the Calorie Content in Food from Images**
(Course Project, Foundations of Machine Learning, Prof. Sunita Sarawagi, CSE, IIT Bombay) (Autumn 2023)
Developed an automated calorie estimation system using **YOLOv8** (detection) and **GrabCut** (segmentation), achieving **7.6%** mean absolute error across **19** food classes. [GitHub: soumenkm/CS725-FML-Project](#)

Machine Learning "From Scratch" Self Projects

- **Build GPT2 and BERT from Scratch:** Developed all the core components of **GPT-2** and **BERT** (multi-head self-attention, MLPs, trainer class, pre-training and instruction fine-tuning) from scratch without **HuggingFace**. [GitHub: soumenkm/Build-LLM-from-scratch](#), [GitHub: soumenkm/Build-BERT-from-scratch](#) (Autumn 2024)
- **Build Diffusion Model from Scratch:** Developed a **DDPM** from scratch, implementing diffusion/ reverse processes and sampling. [GitHub: soumenkm/Diffusion-Model-from-Scratch](#) (Autumn 2024)

Courses at IIT Bombay

- | | |
|---|---|
| CS 725: Foundations of Machine Learning (Grade: 10) | GNR 638: Deep Learning for CV (Grade: 10) |
| SC 607: Convex Optimization (Grade: 10) | GNR 650: Advanced Deep Learning for CV (Grade: 10) |
| CS 601: Algorithms and Complexity (Grade: 10) | CS 769: Optimization for Machine Learning (Grade: 10) |
| EE 635: Applied Linear Algebra (Grade: 9) | GNR 602: Advanced Satellite Image Processing (Grade: 10) |
| IE 621: Probability and Stochastic Process (Grade: 9) | BB 610: Biomedical Microsystems (Inst. Elec.) (Grade: 10) |

Teaching Assistant Positions at IIT Bombay

- | | |
|---|--|
| IITB e-PG Diploma in AI: Statistics (Class Size: 150) | CS 6106: Statistical Learning Theory (Class Size: 50) |
| CS 725: Foundations of Machine Learning (Class Size: 180) | DS 303: Introduction to Machine Learning (Class Size: 180) |

Achievements

- Maintained a perfect **10/10 CGPA** across three consecutive academic semesters at IIT Bombay. (2025)
- Awarded the **Institute Academic Prize** for outstanding academic performance at IIT Bombay. (2024)
- Received the **GE Spotlight Impact Award** for reducing operational costs at GE Vernova. (2022)
- Won the **Innovate 2021 AI/ML Challenge** hosted by GE Vernova. (2021)