# MD. SHARIAR KABIR

✉ shariar1405076@gmail.com  🌐 shariar076.github.io  ⌂ Shariar076  𝒢 Google Scholar

## RESEARCH INTERESTS

My research focuses on developing reliable, interpretable, and aligned large language models, with particular emphasis on understanding and improving their behavior in long-context and real-world settings.
Currently I am working on:

- Red-teaming LLMs to elicit unintended behaviors in long-context settings as well as controlling them.
- Improving explainability by identifying features and circuits responsible for model behaviors.
- Integrating interpretability with intervention methods for controlled model editing.

## SELECTED PUBLICATIONS

[5] Md. Shariar Kabir, and Muhammad Abdullah Adnan. AgnoSVD: Dynamic resource allocation for serverless workloads using collaborative filtering. In *Array (2025): 100662*.

[4] Shariar Kabir, Kevin Esterling, and Yue Dong. PReSS: A Black-Box Framework for Evaluating Political Stance Stability in LLMs via Argumentative Pressure. [paper]. Under Review at LREC.

[3] Shariar Kabir, Kevin Esterling, and Yue Dong. Beyond the surface: Probing the ideological depth of large language models. *arXiv preprint arXiv:2508.21448*, 2025. In Progress.

[2] Shariar Kabir, Nazmun Nahar, Shyamasree Saha, and Mamunur Rashid. Automatic speech recognition for biomedical data in Bengali language. *arXiv preprint arXiv:2406.12931*, 2024.

[1] Syed Mostofa Monsur*, Shariar Kabir*, and Sakib Chowdhury*. Synthnid: Synthetic data to improve end-to-end Bangla document key information extraction. In *EMNLP 2023 Workshop on Bangla Language Processing*, 2023.

*\* indicates co-first author.*

## RESEARCH EXPERIENCE

**University of California, Riverside**                                   **January 2025 - December 2025**
*Research Intern – NLP Lab (Prof. Yue Dong, and Prof. Kevin Esterling)*                        *Riverside, CA*

*Understanding LLMs' response instability over longer context.*

  - Finding the correlation between instability and epistemic uncertainty over multiturn conversations.
  - Evaluating stability after model finetuning and activation steering.

*Mechanistic Interpretability of LLM in Socio-Political Reasoning.*

  - Analyzing activation pathways linked to ideological responses using SAE features from Neuronpedia.
  - Evaluating steerability and understanding ideological depth formation mechanisms in LLMs.

*LLMs' Social Epistemology using Bayesian Statistics.*

  - Implementing Multidimensional IRT from scratch in Stan (validated with 98% correlation to DW-NOMINATE scores).
  - Applying psychometric methods to model LLM ideological positioning compared to humans.

## INDUSTRY RESEARCH EXPERIENCE

**Celloscope Ltd.**                                                                     **2020 – Present**
*Lead AI Research Engineer*                                                      *Dhaka, Bangladesh*
  - Led a number of NLP and computer vision-based projects deployed across multiple industrial domains.
  - Engineered *private and self-hosted* Conversational AI solutions using open-sourced LLMs and RAG.

**MedAI Pvt. Ltd.**                                                                         **2021 – 2024**
*Data Scientist (Part-time)*                                                           *Cambridge, UK*
  - Collected and curated Bengali biomedical audio data for fine-tuning Whisper.
  - Built multimodal disease prediction pipelines incorporating structured & unstructured clinical data.
  - Evaluated Conversational AI for mental health symptoms, integrating symptom classification models.

## EDUCATION

**Bangladesh University of Engineering and Technology (BUET)**                          2026
*M.Sc. in Computer Science & Engineering (Part-time)*

**Thesis**: Dynamic Resource Allocation for Workloads in Serverless Architecture. [dissertation]
**Coursework**: Bioinformatics Algorithms, Distributed Computing Systems, Data Mining, Data Management in the Cloud, etc.

**Bangladesh University of Engineering and Technology (BUET)**                          2019
*BSc in Computer Science & Engineering*                                      *CGPA: 3.53/4.00*

**Thesis**: Active Learning on Big Data for scalable classification using distributed infrastructure. [dissertation]
**Selected Coursework**: Machine Learning, Pattern Recognition, Artificial Intelligence, Digital Image Processing, etc.

## SKILLS

**Research**:            Mechanistic Interpretability, Topic Modeling, Item Response Theory, Model Visualization
**Programming**:         Python, Shell, C, C++, STAN, LaTeX, SQL, TypeQL
**Machine Learning**:    PyTorch, ScikitLearn, OpenCV, Pandas, Datasets, Transformers, SpaCy
**Tools**:               LangChain, Neuronpedia, SGLang, Ollama, OpenAI, Spark, PySpark, Docker
**Soft-Skills**:         Communication, Collaboration, Presentation, Technical Writing

## RECENT INDUSTRY PROJECTS

**Medical Classification by Probing LLMs:** Multi-label classification of medical disciplines by training linear probes on activations from LLMs pretrained on Medical data. We extract layer-wise attention head activations from medical-domain LLMs and use Ridge regression classifiers to predict relevant medical disciplines from clinical descriptions. [code]

**Exercise Monitoring System:** Inspired by research works like VidDiff and HuMMan. We created a system for LG Nova's Real-Time AI Fitness Coaching. Our system leverages Vision-Language Models (VLMs) to assist users in performing exercises correctly by comparing their execution against reference videos of expert demonstrations.

**Drawing Checker:** An initiative to automate the design-error detection and verification in engineering drawings. The system used computer vision techniques and generative models to evaluate technical drawings, identify inconsistencies, and flag deviations from design constraints. I directed the model training and dataset curation pipelines, ensuring that the models achieved consistent accuracy across diverse geometric and structural inputs.

**Resume Shortlister:** An NLP-driven retrieval and ranking system designed to automate candidate selection for enterprise recruitment. By designing a hybrid RAG approach combining rule-based filtering with semantic retrieval, we developed a system capable of aligning candidate attributes with organizational requirements.

## AWARDS & ACHIEVEMENTS

**Industry Coding Assessment**                                                          2025
*CodeSignal General Coding Assessment (ICA): 510/600 ($\approx$ 722/850 equivalent GCA, top 15%)*

**Global Health Equity Challenge Award**                                                2024
*MIT Solve*

Recognized for innovative approach to accessible healthcare (Top 6/2200+). [link]

## REFERENCES

**Prof. Yue Dong**                                                        ✉ yue.dong@ucr.edu
*Assistant Professor of Computer Science and Engineering, University of California, Riverside*

**Prof. Kevin Esterling**                                           ✉ kevin.esterling@ucr.edu
*Professor of Public Policy and Political Science, University of California, Riverside*

**Prof. Muhammad Abdullah Adnan**                                     ✉ adnan@cse.buet.ac.bd
*Professor in the Department of CSE, Bangladesh University of Engineering and Technology*