

# Mohammad Taufeeque

🌐 taufeeque9.github.io · ✉ 9taufeeque9@gmail.com · 🐙 github/taufeeque9 · in mtaufeeque

## EDUCATION

### Indian Institute of Technology Bombay

2018–2022

B.Tech (with Honors) in Computer Science and Engineering

Mumbai, India

- GPA - **9.24**/10.0
- Bachelor's Thesis - Fianchetto: Speed, Belief, Guile, Caution to Win at Reconnaissance Blind Chess

## PUBLICATIONS

- [1] K. Pelrine\*, **M. Taufeeque\***, M. Zając, E. McLean, and A. Gleave, “Exploiting Novel GPT-4 APIs”, 2023. arXiv: 2312.14302 [cs.CR].
- [2] A. Tamkin, **M. Taufeeque**, and N. D. Goodman, “Codebook Features: Sparse and Discrete Interpretability for Neural Networks”, 2023. arXiv: 2310.17230 [cs.LG].
- [3] A. Gleave\*, **M. Taufeeque\***, J. Rocamonde\*, E. Jenner, S. H. Wang, S. Toyer, M. Ernestus, N. Belrose, S. Emmons, and S. Russell, “imitation: Clean Imitation Learning Implementations”, 2022. arXiv: 2211.11972 [cs.LG].
- [4] **M. Taufeeque\***, N. Tongia\*, and S. Kalyanakrishnan, “Fianchetto: Speed, Belief, Guile, Caution to Win at Reconnaissance Blind Chess”, 2022.
- [5] G. Perrotta, R. W. Gardner, C. Lowman, **M. Taufeeque**, N. Tongia, S. Kalyanakrishnan, G. Clark, K. Wang, E. Rothberg, B. P. Garrison, P. Dasgupta, C. Canavan, and L. McCabe, “The Second NeurIPS Tournament of Reconnaissance Blind Chess”, in *Proceedings of the NeurIPS 2021 Competitions and Demonstrations Track*, ser. Proceedings of Machine Learning Research, vol. 176, 2022, pp. 53–65.
- [6] **M. Taufeeque\***, S. Koita\*, N. Spicher, and T. M. Deserno, “Multi-camera, multi-person, and real-time fall detection using long short term memory”, in *Medical Imaging 2021: Imaging Informatics for Healthcare, Research, and Applications*, vol. 11601, SPIE, 2021, pp. 35–42.

## INDUSTRY & RESEARCH EXPERIENCE

### FAR AI - Research Engineer

Berkeley, USA

Supervisor: Dr. Adam Gleave

Aug '22 - Present

- Benchmarked and improved the performance and documentation of the open-source library *imitation*
- **Codebook Features**: Developed and published a novel neural network interpretability tool
- Red-teaming: Found novel ways to exploit the **GPT-4** Fine-tuning and Assistants API

### Microsoft Research - Research Intern

Bangalore, India

Guide: Prof. Sunita Sarawagi (IIT Bombay) & Dr. Sriram Rajmani (Microsoft Research)

Dec '21 - Aug '22

- Developed an online algorithm against OOD data to integrate noisy feedback rules to a trained ML text classifier
- Deployed the algorithm to **Microsoft's Ads** system & improved the compliance of the system by over **55%**

### Goldman Sachs - Summer Analyst

Bangalore, India

Guide: Kesavan Mukunthan

Summer 2021

- Created a webapp module to compute and display performance exposures to different factors for each stock in a portfolio of mutual funds that helps portfolio managers to analyse the drivers of performance of every fund

- Developed an application that **detects falls in real-time** using human pose keypoints from multiple cameras
- Maintainer of the open-source project on GitHub with **200+ stars & 50+ forks**

**AI Agent for Reconnaissance Blind Chess**

IIT Bombay

Bachelor's Thesis | Guide: Prof. Shivaram Kalyanakrishnan | Won NeurIPS competition (RBC) Aug '21 - Dec '21

- Developed an AI Agent for RBC, a variant of Chess where only a 3x3 region can be sensed before making a move
- Won the **NeurIPS 2021** competition on RBC with **91.3%** win rate & **100 Elo pts** margin from the runner-up

**Randomized Planning Algorithms for POMDPs**

IIT Bombay

Guide: Prof. Shivaram Kalyanakrishnan

Spring 2021

- Designed planning algorithms for **POMDPs** that achieved **20%** higher rewards than the SoTA algorithms
- Combined random subsets of nodes from Finite State Controllers of weak policies to obtain a strong policy

**SAFE App Vulnerabilities**

IIT Bombay

Guide: Prof. Bhaskaran Raman

Aug '20 - Apr '21

- Found **severe vulnerabilities** in the SAFE App IITB, used by many institutions to conduct remote exams
- Reported **data-leak**, **APK signature verification** and **timing-based** vulnerabilities in the Android app

**SCHOLASTIC ACHIEVEMENTS**

<b>Attended Google Research Week</b>	Selected among 50 undergraduates nationwide	2022
<b>Best AI Agent out of 18 Bots</b>	RBC chess competition in NeurIPS 2021	2021
<b>AP (Advanced Performer) Grade</b>	Best performance in Machine Learning course (GNR 638)	2020
<b>All India Rank 303</b>	JEE Advanced (230,000 aspirants)	2018
<b>All India Rank 330</b>	JEE Mains (1.2 Million aspirants)	2018
<b>Merit-cum-Means (MCM) Scholarship</b>	IIT Bombay	2018
<b>National top 1%</b>	Indian National Physics Olympiad (INPhO)	2017
<b>National top 1%</b>	Indian National Chemistry Olympiad (INChO)	2017
<b>KVPY Science Fellowship</b>	Government of India	2016

**TECHNICAL SKILLS**

<b>Programming</b>	C++, C, Python, Java, Bash, Racket, Prolog, MIPS, PostgreSQL
<b>ML Libraries</b>	PyTorch, Tensorflow, Transformers (by Huggingface), Keras, Scikit-learn, Pandas

**TEACHING & MENTORSHIP****Teaching Assistantships**

- Intro to ML Safety (Summer 2022 & Spring 2023)
- Artificial Intelligence & Machine Learning (Autumn 2021)
- Medical Image Computing (Spring 2022)
- Calculus II (Autumn 2021)

**AISCF IIT Delhi**

Autumn 2023

- Conducted an in-person semester-long reading group on **AI Safety** with 30+ students across 3 groups

**English Language Tutor at ELIT IIT Bombay**

Autumn 2020

- Designed an English language curriculum and held classes for 50+ students registered for the program

**Member in Developer's Community (DevCom) IIT Bombay**

Jan '19 - Aug '20

- Maintained, developed & updated features of **InstiApp**, the institute app with over **10,000+** downloads

**Mentor of SoC project - Intrusion Detection System | Github: taufeeque9/IDS**

Summer 2020

- Mentored a team of **9 developers** in building a real-time system to monitor networks for malicious activity