

# Rahul Chand

+91-7286063972 | chandrahul0320@gmail.com

Website: [rahulschand.github.io](https://rahulschand.github.io) | LinkedIn: [/RahulSChand](https://www.linkedin.com/in/RahulSChand) | GitHub: [/RahulSChand](https://github.com/RahulSChand)

## EDUCATION

**Birla Institute of Technology, Pilani | 2015-19**

**Bachelor of Engineering (Honors) in Computer Science and Engineering,**

**GPA - 9.62/10, Class Rank - 6/180+**

**Selected Coursework:** Machine Learning, Pattern Recognition, Logic in Computer Science, Neural Networks and Fuzzy Logic, Data Mining, Information Retrieval, Graphs and Networks, Computer Architecture, Operating System

## INDUSTRY AND RESEARCH EXPERIENCE

**Microsoft Research, Bangalore, India | Research Fellow**

**July 2021-July 2023**

- Worked with Microsoft Turing Team on Transformer compression using sparse factorization.
- Worked in the XC team under Manik Verma on Extreme Multi Label Learning (XML). My work primarily involved studying & improving tail performance of extreme classifiers via regularization (Paper under review). Additionally, I also worked on the problem of compressing extreme classifiers.

**Arcesium, Hyderabad, India | Software Engineer & Intern**

**Aug 2019-May 2021**

- Worked in the Performance and Accounting team as a full stack developer.
- Worked as part of the team responsible for developing microservices & frontend using Java, Kotlin, Python, ReactJS & T-SQL for handling large volumes (>100k) of trades daily.

**Indian Institute of Science (IISC), India | Undergraduate Thesis student**

**Jan 2019-July 2019**

- Worked under Prof. Venkatesh Babu at Video Analytics Lab (VAL) for my undergrad thesis on the problem of optical flow estimation using matrix capsule networks.
- Paper: <https://arxiv.org/abs/2304.00306>

**Indian Institute of Remote Sensing (IIRS), Dehradun, India | Research Intern**

**May 2017-July 2017**

- Worked with the Geo-informatics Department on their road-asset management project.
- My work involved developing a deep learning solution on Keras using Faster-RCNN & FCN for road-asset mapping of Indian roads. The model was trained on VOC2012 & images of Dehradun roads obtained from IIRS.
- Report: [github.com/RahulSChand/IIRS-Vehicle-Detection](https://github.com/RahulSChand/IIRS-Vehicle-Detection)

## TEACHING EXPERIENCE

**Teaching Assistant for below courses. Graded assignments, prepared course projects & supervised lab sessions.**

- **Fall 2018:** Data Mining, Principles of Programming Languages, Computer Programming
- **Spring 2018:** Data Structures and Algorithms, Database Systems
- **Fall 2017:** Logic in Computer Science

## PROJECTS

**Open source libraries & contributions**

- **vRAM for LLMs** ([Github](https://github.com/RahulSChand/vRAM))
  - Tool to check GPU vRAM requirements for training & inference of any LLM. Supports frameworks like HuggingFace, vLLM, exLlama, llama.cpp and quantization (bitsandbytes, GGML).
- **llama2.c for dummies** ([Github](https://github.com/RahulSChand/llama2.c))
  - Step by step walkthrough of the inference code of [llama2.c](https://github.com/facebookresearch/llama2.c) written as a starter reference for LLM inference.
- **Fast & tiny datasets for optical flow** ([Github](https://github.com/RahulSChand/optical_flow_datasets))
  - Library to generate tiny optical flow datasets on the fly for sanity testing optical flow estimation models. Written as part of undergraduate thesis at IISC & used in the paper [link](https://arxiv.org/abs/2304.00306)

- **Efficient Batched Torch KSVD** ([Github](#))
  - Library to run sparse dictionary completion algorithm KSVD on batched matrices on GPU. Written using pytorch as part of transformer compression work at Microsoft Research.
- **Language model compression with weighted low-rank factorization** ([Github](#))
  - Pytorch implementation of the ICLR 2022 paper “Language model compression with weighted low-rank factorization”. Code written as part of work at MSR India.
- **Attention network for reading comprehension and question answering** ([Github](#))
  - Tensorflow implementation of the paper “Multi-Granularity Hierarchical Attention Fusion Networks for Reading Comprehension and Question Answering paper”.

---

## ACHIEVEMENTS

- One of 30 students selected from Maharashtra to attend training camp for INMO (Indian National Maths Olympiad) 2015.
- Merit-cum-Need scholarship at BITS Pilani for all 8 semesters.

---

## SKILLS

**Programming Language:** Python, C++, Java, JavaScript, Kotlin

**Libraries and Frameworks:** PyTorch, Numpy, TensorFlow, Keras, HuggingFace, ReactJS, NextJS