# Siddhant Ray

Chicago, IL, USA
📱 +1-(773)-457-4156
✉ siddhant.r98@gmail.com
🌐 www.linkedin.com/in/siddhant-ray
https://github.com/Siddhant-Ray (GitHub)

## Education

**2023 – 2028**  **The University of Chicago**, *PhD in Computer Science*
　　Advisor - Junchen Jiang and Nick Feamster

**2020 – 2022**  **ETH Zürich**, *MSc in Electrical Engineering and Information Technology*
　　Advisor - Laurent Vanbever

**2016 – 2020**  **VIT Vellore**, *B.Tech in Electronics and Communication Engineering*

## Experience

**Sep 2023 – Present**  **Graduate Research Assistant**, *Computer Science Department, The University of Chicago*
- Tranformer based model for predicting changes in network latency for use in active queue management and multipath routing.
- Resource allocation and sharing for optimally serving multi-tenant Retrieval Augmented Generation(RAG) LLM systems.

**Sep 2022 – Mar 2023**  **Cloud Networks Researcher**, *Advanced Network Architecures Lab, UPC Barcelona*
- Analysed reinforcement learning based resource sharing, offloading and allocation for cloud-edge systems.
- Developed an approximation for a Mixed-Integer Optimal Matching Algorithm for resource allocation to reduce execution time by 2.5-3x.

**Oct 2021 – Sep 2022**  **Graduate Research Assistant**, *Law, Economics, and Data Science Group, ETH Zurich*
- Research Assistant to Professor Dr. Elliott Ash and worked on improving semantic labelling for text corpora using newer NLP models, sentence simplification and clustering for topic modelling.
- Worked on paraphrase mining to determine clusters of similar narratives in legal corpora and use NLP models to capture underlying narratives in meat policy documents to analyse political discourse.

**May 2019 – July 2019**  **Software Development Intern**, *Capgemini Engineering*
- Developed a K-Shortest Path Searching algorithm for ONOS based Software Defined Layer 2 VPNs.
- Algorithm was subject to dynamic constraints of network resources (e.g.required edges, vertices etc.) to be used for path calculation.

**May 2018 – July 2018**  **Software Development Intern**, *BlueStacks*
- Worked on a machine learning algorithm to predict the App Engine's appropriate display screen based on the customer's past experiences.
- Developed an automation script for generating SVG cards for the App Engine's game front end and an address verification tool using the EasyPost API.

## Publications

**2024**  **Siddhant Ray**, Xi Jiang, Zhuohan Guo, Junchen Jiang, and Nick Feamster. Transformer-based predictions for sudden network changes. In *21st USENIX Symposium on Networked Systems Design and Implementation (Poster Session)*, NSDI '24. USENIX Association, 2024.

**2024**  Yuhan Liu, Hanchen Li, Yihua Cheng, **Siddhant Ray**, Yuyang Huang, Qizheng Zhang, Kuntai Du, Jiayi Yao, Shan Lu, Ganesh Ananthanarayanan, Michael Maire, Henry Hoffmann, Ari Holtzman, and Junchen Jiang. Cachegen: Kv cache compression and streaming for fast language model serving. To appear in SIGCOMM '24, 2024.

**2022**  Alexander Dietmüller, **Siddhant Ray**, Romain Jacob, and Laurent Vanbever. A new hope for network model generalization. In *Proceedings of the 21st ACM Workshop on Hot Topics in Networks*, 2022.

2020  **Siddhant Ray** and Budhaditya Bhattacharyya. Machine learning based cell association for mmtc 5g communication networks. *International Journal of Mobile Network Design and Innovation*, 10(1):10–16, 2020.

## Skills

| | |
|---|---|
| Programming | Python, C++, Java, Bash, Rust, SQL, C, TEX |
| Software | Linux, Git, Docker, P4 switches, ONOS, Google Cloud, AWS, Maven, MATLAB, NetSim, Cadence |
| Frameworks | Mininet, FRRouting, PyTorch, TensorFlow, Sklearn, NLTK, Flask, SciPy, Scapy, BS4, NS-3, Langchain, vLLM |
| Languages | English (C2), Hindi, Bengali, Deutsch (B1) |

## Selected Projects

| | |
|---|---|
| 2022 | Advancing Packet-Level Traffic Predictions with Transformers (Master Thesis) - [code, thesis] |
| 2021 | Towards a New Framework for Integration of Network Planes (Research Project) - [code] |
| 2021 | Attentive Neural Networks for News Classification (Research Project) - [code] |
| 2021 | Investigating Possible Inductive Biases in Local Sparse Attention ViT Architectures Against Traditional CNNs (Course Project) - [code, paper] |
| 2021 | Automatic Certificate Management Environment (Course Project) - [code] |
| 2020 | Maximizing Cross Traffic Flows in a L2/L3 Network with Programmable Switches (Course Project) - [code, poster] |
| 2020 | Machine Learning based Cell Association for 5G Communication Networks (Bachelor Thesis) - [code] |

## Relevant Courses

| | |
|---|---|
| Graduate | Approximation Algorithms, Algorithms, Advanced Computer Networks, System Security, Network Security, Distributed Computing, Discrete Event Systems, Networks Seminar, Introductory Machine Learning, Deep Learning, Learning and Classification Theory, Mathematics of Data Science, Neural Network Theory |
| Undergraduate | Computer Networks, Operating Systems, Wireless Communication, Linear Algebra |

## Honors and Awards

| | |
|---|---|
| 2023 – 2028 | **Liew Family Graduate Fellowship**, University of Chicago |
| 2022 | **Winner at Datathon**, *Microsoft Challenge*, ETH Zurich |
| 2020 | **Best Outgoing Student**, *SENSE department*, VIT Vellore |
| 2019 | **Runner-Up at VIT Hack**, *Education Track*, VIT Vellore |
| 2016 – 2019 | **Merit Scholarship for Academic Excellence**, VIT Vellore |

## Leadership and Volunteering

| | |
|---|---|
| 2019 – 2020 | **Technical Advisor**, IETE VIT |
| 2018 – 2019 | **Organizer**, TEDx VIT Vellore |
| 2017 – 2020 | **President** (2018 – 2019) & **Outreach Worker**, Anokha NGO |