

# Amey Agrawal

CS Ph.D. Student, Georgia Tech

[ameya.info](mailto:ameya.info) [@agrawalamey12@gmail.com](mailto:agrawalamey12@gmail.com) [github.com/agrawalamey](https://github.com/agrawalamey) [Google Scholar](https://scholar.google.com/citations?user=ameya.info)

## Education

<b>Present</b> <b>Aug 2022</b>	<b>Georgia Institute of Technology</b> Ph.D., Computer Science. GPA 4.00/4.00. Advisor: <i>Prof. Alexey Tumanov</i> / Area: Systems for machine learning, LLM inference systems.	<b>Atlanta, USA</b>
<b>Jul 2018</b> <b>Aug 2014</b>	<b>Birla Institute of Technology and Science Pilani</b> B.E. (Hons.), Computer Science	<b>Pilani, India</b>

## Experience

<b>Aug 2024</b> <b>May 2024</b>	<b>Microsoft Azure Systems Research</b> Research Intern / Mentor: <i>Dr. Esha Choukse</i> Building systems to serve large language models with multimillion context length requests.	<b>Redmond, USA</b>
<b>Aug 2023</b> <b>May 2023</b>	<b>Microsoft Research</b> Research Intern / Mentors: <i>Dr. Ramchandran Ramjee, Dr. Bhargav Gulavani</i> Designing high-throughput, low-latency inference systems for large language models.	<b>Bangalore, India</b>
<b>Aug 2022</b> <b>Jan 2021</b>	<b>Microsoft Research</b> Research Software Engineer-II / Mentor: <i>Dr. Muthian Sivathanu</i> Built parts of the elasticity sub-system that leveraged efficient time sharing of GPUs to provide transparent scaling of deep learning training workloads. This work was done as a part of the Singularity project, Microsoft's planet-scale AI infrastructure service.	<b>Bangalore, India</b>
<b>Nov 2020</b> <b>Jul 2018</b>	<b>Qubole Inc.</b> Member of Technical Staff-II / Mentor: <i>Rohit Karlupia</i> Worked on various applied machine learning and software engineering problems to enhance Qubole's data science platform. Published research in several top-tier venues.	<b>Bangalore, India</b>
<b>Dec 2017</b> <b>Jul 2017</b>	<b>Software Engineering Intern</b> / Mentor: <i>Bharath Bhushan</i> Built core data-plane components for Qubole's Deep Learning clusters based on TensorFlow and Apache Spark.	

## Publications

- DynaQuant: Compressing Deep Learning Training Checkpoints via Dynamic Quantization** [pdf]  
Amey Agrawal, Sameer Reddy, Satwik Bhattamishra, Sarath Nookala, Vidushi Vashishth, Kexin Rong, and Alexey Tumanov  
Proceedings of 15th ACM Symposium on Cloud Computing, Redmond, 2024, Redmond [SoCC'24]
- Mnemosyne: Parallelization Strategies for Efficiently Serving Multi-Million Context Length LLM Inference Requests Without Approximations** [pdf]  
Amey Agrawal, Junda Chen, Íñigo Goiri, Ramachandran Ramjee, Chaojie Zhang, Alexey Tumanov, Esha Choukse  
Preprint: arXiv:2409.17264 (2024) [CoRR]
- Metron: Holistic Performance Evaluation Framework for LLM Inference Systems** [pdf][code]  
Amey Agrawal\*, Anmol Agarwal\*, Nitin Kedia, Jayashree Mohan, Souvik Kundu, Nipun Kwatra, Ramachandran Ramjee, Alexey Tumanov  
Preprint: arXiv:2407.07000 (2024) [CoRR]
- Taming Throughput-Latency Tradeoff in LLM Inference with Sarathi-Serve** [pdf][code][video]  
Amey Agrawal, Nitin Kedia, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav S. Gulavani, Alexey Tumanov, Ramachandran Ramjee  
Proceedings of 18th USENIX Symposium on Operating Systems Design and Implementation, 2024, Santa Clara [OSDI'24]
- Vidur: A Large Scale Simulation Framework For LLM Inference** [pdf][code][video]  
Amey Agrawal, Nitin Kedia, Jayashree Mohan, Ashish Panwar, Nipun Kwatra, Bhargav S. Gulavani, Ramachandran Ramjee, Alexey Tumanov  
Proceedings of 7th Annual Conference on Machine Learning Systems, 2024, Santa Clara [MLSys'24]
- Sarathi: Efficient LLM Inference by Piggybacking Decodes with Chunked Prefills** [pdf]  
Amey Agrawal, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav S. Gulavani, Ramachandran Ramjee  
Preprint: arXiv:2308.16369 (2023) [CoRR]

- Sybill: Deep Learning Workload Tuning with Virtual GPUs** [\[poster\]](#)  
 Srihas Yarlagadda\*, Amey Agrawal\*, Sarath Nookala, Pranavi Bajjuri, Shivam Mittal, Alexey Tumanov  
*ACM Symposium on Cloud Computing Poster, 2023* [SoCC'23]
- Singularity: Planet-Scale, Preemptible, Elastic Scheduling of AI Workloads** [\[pdf\]](#)  
 Singularity Team, Microsoft  
*Preprint: arXiv:2202.07848 (2022)* [CoRR]
- Logan: A Distributed Online Log Parser** [\[pdf\]](#)  
 Amey Agrawal, Rajat Gupta, and Rohit Karlupia  
*Proceedings of IEEE International Conference on Data Engineering, 2019, Macau* [ICDE'19]
- Learning Digital Circuits: A Journey Through Weight Invariant Self-Pruning Neural Networks** [\[pdf\]](#)[\[code\]](#)  
 Amey Agrawal, and Rohit Karlupia  
*Sparsity in Neural Networks Workshop 2021; New in ML Workshop, NeurIPS, 2019, Vancouver* [SNN'21]
- Delog: A Privacy Preserving Log Filtering Framework for Online Compute Platforms** [\[pdf\]](#) [\[dataset\]](#)  
 Amey Agrawal, Abhishek Dixit, Namrata Shettar, Darshil Kapadia,  
 Rohit Karlupia, Vikram Agrawal, and Rajat Gupta  
*Proceedings of IEEE International Conference on Big Data, 2019, Los Angeles* [BigData'19]

## Honours and Awards

---

- Center for Research into Novel Compute Hierarchies (CRNCH) Fellowship, 2023** [\[🌐\]](#)  
 › For research of automatic hardware-aware optimization of deep learning training workloads.
- School of Computer Science Fellowship, 2022**  
 › PhD fellowship from Georgia Institute of Technology.

## Teaching Roles

---

- Systems for Machine Learning** *Head Teaching Assistant | Georgia Institute of Technology* Fall'24  
 › Helping revise the curriculum and assignments to reflect the changing landscape in AI systems.
- Systems for Machine Learning** *Guest Lecturer | Georgia Institute of Technology* Fall'23  
 › Conducted a three-part lecture series on large-language model inference systems.
- Neural Networks & Fuzzy Logic** *Lead Teaching Assistant | BITS Pilani* [\[assignments\]](#) Aug'17 - May'18  
 › Introduced Python programming assignments along with a new custom-built evaluation platform. Other responsibilities included coordinating the team of seven teaching assistants to conduct labs, designing assignments and helping students with the term project.
- Machine Learning** *Teaching Assistant | BITS Pilani* Jan'18 - May'18  
 › Conducted introductory sessions on the scientific Python ecosystem, and organized tests and programming assignments for over 100 students in the class.