

Sudhanshu Agrawal

☎ +1 (310) 500-8571 🏠 San Diego, California ✉ sudhanshuagr27@g.ucla.edu 🌐 sudhanshuagrawal.com

Education

University of California, Los Angeles

Sept 2019 – June 2023

B.S. Computer Science (*Magna cum laude*)

B.S. Mathematics (*Cum laude*)

GPA: 3.92

Undergraduate Research and Teaching

- UCLA Computer Science, Machine Intelligence Group: undergraduate research, advised by Professor Aditya Grover [🔗](#).
- UCLA Mathematics Department: undergraduate research, advised by Professors Levon Nurbekyan [🔗](#) and Samy Wu Fung [🔗](#).
- ACM, UCLA: created and taught ML and deep learning workshops to over 200 students regularly over 2 years. [🔗](#)

Work Experience

Qualcomm : *ML Research Engineer*

Aug 2023 – Present (*San Diego*)

- **LLM efficiency research:** Conducted research on new techniques to accelerate the generation speed of LLMs via speculative decoding. 2 conference papers published and 7 inventions with multiple US/global patent applications pending.
- **Implementation and commercialization:** carried out end-to-end implementation and testing of these algorithms on Qualcomm edge devices, subsequently used for public demos and included in commercial Snapdragon chipset releases.

Qualcomm : *ML Engineering Intern*

June – Sept 2022 (*San Diego*)

- **Novel *Python* profiler:** created a profiling tool for deep learning applications, capable of timing code function-by-function, line-by-line, and profiling multiprocessing applications.

SonicJobs : *ML Engineering Intern*

June – Sept 2021 (*Remote*)

- **Synthetic data-set generation:** created a pipeline to scalably generate over 100,000 *HTML/CSS* webpages with varying contents and styles, solving the challenge of obtaining a diverse, labelled training data-set for a computer vision model.

Reliance Jio : *ML & Data Science Intern*

July – Oct 2020 (*Remote*)

- **Predicting the properties of hydrocarbons:** developed predictive models leveraging lasso regression, ridge regression, support vector regression, recursive feature elimination, linear discriminant analysis, and 1-D convolutions.

Publications

Spiffy: Speculative Decoding for Diffusion LLMs

2025

- Agrawal, Sudhanshu, et al., “Spiffy: Multiplying Diffusion LLM Acceleration via Lossless Speculative Decoding.” *arXiv:2509.18085*. [🔗](#)

VOCABTRIM: Improving Speculative Decoding via LMHead Dimensionality Reduction

2025

- Goel, Raghavv, Agrawal, Sudhanshu, et al., “VOCABTRIM: Vocabulary Pruning for Efficient Speculative Decoding in LLMs.” *ICML 2025 Workshop on Efficient Systems for Foundation Models*. [🔗](#)

AdaEDL: Early Draft Stopping for Speculative Decoding of LLMs

2024

- Agrawal, Sudhanshu, Jeon, Wonseok, and Lee, Mingu, “AdaEDL: Early Draft Stopping for Speculative Decoding of Large Language Models via an Entropy-based Lower Bound on Token Acceptance Probability.” *NeurIPS 2024 Efficient Natural Language and Speech Processing Workshop*. [🔗](#)

ExPT: Synthetic Pretraining for Few-Shot Experimental Design

2023

- Nguyen, Tung, Agrawal, Sudhanshu, and Grover, Aditya, “Expt: Synthetic pretraining for few-shot experimental design.” *Advances in Neural Information Processing Systems 36 (2023): 45856-45869*. [🔗](#)

Efficiently Solving High Dimensional Non-local Mean Field Game Problems

2022

- Agrawal, Sudhanshu, Lee, Wonjun, Fung, Samy W., Nurbekyan Levon, “Random features for high-dimensional nonlocal mean-field games.” *Journal of Computational Physics 459 (2022): 111136*. [🔗](#)

Invited Talks, Judgeships, Reviewing

- **Reviewer:** 2026 AAAI Conference
- **Judge:** 2025 UCSD Graduate Student Research Exposition
- **Judge:** 2025 San Diego State University Student Research Symposium
- **Reviewer:** 2025 ICML Efficient Systems for Foundation Models Workshop
- **Reviewer:** 2025 NeurIPS Efficient Natural Language and Speech Processing Workshop
- **Speaker:** 2024 UCSD and Qualcomm Graduate Students Tech Talk and Recruitment Event
- **Speaker:** 2024 UCSD, IEEE, Qualcomm Careers Panel
- **Panelist:** 2024 UCLA Mathematics Department Alumni Panel