# Amogh Patankar

linkedin.com/in/apatankar22 | github.com/apatankar22 | apatankar22.github.io

## EDUCATION

**University of California, San Diego** — September 2023 – March 2025
*M.S., Computer Science & Engineering; Concentration:* **AI & Machine Learning** — *La Jolla, CA*
- **Growth Strategist, Perplexity AI**
- **Co-President, CSE Graduate Student Council**

**University of California, San Diego** — September 2020 – March 2023
*B.S., Data Science* — *La Jolla, CA*

## PROFESSIONAL EXPERIENCE

**AI Applications Development Intern** — September 2024 – December 2024
*Advanced Micro Devices (AMD)* — *San Jose, CA*
- Optimized generative AI workload (Llama-2, Llama-3, Stable Diffusion, etc.) execution on RyzenAI neural processing unit (NPU) using strategies such as caching, batching, quantization, as well as model and data parallelism.
- Benchmarked generative AI workloads on Ryzen[AI] and competitor hardware (Snapdragon X Elite, Lunar Lake), specifically time to first token, and tokens/sec, i.e. latency and throughput, as metrics.
- Evaluated NPU latency and utilization during concurrent execution of generative AI workloads and high-fidelity AI effects (Microsoft Studio Effects).

**Data Scientist, Generative AI Intern** — June 2024 – September 2024
*Marvell Technology Inc.* — *Santa Clara, CA*
- Architected an end-to-end ETL pipeline consisting of data engineering, analysis and visualization using numpy, pandas, sklearn, Tableau, etc. Data was ingested using Apache and AWS Kinesis, and stored in SnowflakeDB, and AWS S3.
- Utilized generative pretrained transformer (GPT) models to generate synthetic data, leveraging parameter efficient fine-tuning (PEFT) techniques such as low-rank adaptation (LoRA).
- Implemented reinforcement learning (RL) and deep learning methods for DSP parameter optimization.
- Integrated large language models (LLM) and retrieval augmented generation (RAG) to automate hardware modeling process.

**Researcher** — June 2023 – August 2024
*Stanford University School of Medicine* — *Palo Alto, CA*
- Performed data engineering, analysis, and visualization through various methods in Python and R; specifically, pandas, geopandas, numpy, matplotlib, scikit-learn and ggplot. Developed unique statistical packages composed of chi-squared and Fisher tests.
- Led research teams mentored by Dr. Gross, Dr. Palaniappan, and Jin Long. **Opioid overdose research** published in British Journal of Anaesthesia, and diabetes research in preprint at Journal of Asian Health.

**Software Development Engineer Intern** — June 2022 – September 2022
*Amazon Web Services (AWS)* — *Seattle, WA*
- Led architectural changes in Lex ASR (Automatic Speech Recognition) Services and DataHub, improving latency for conversational AI models. Conducted A/B tests to evaluate architectural changes for customer use.
- Performed cohort analysis to segment critical and non-critical customer data in DataHub, enhacing Lex ASR schemas, and enabled compliant storage of all data using AWS Kinesis, S3, and Lambda.
- Applied time series analysis on AWS CloudWatch metrics to track and optimize the performance of ASR Service, leading to lower response time and accelerating customer request resolution by up to ∼75%.

**Research Intern** — June 2021 – August 2021
*Scripps Research Translational Institute* — *La Jolla, CA*
- Developed an R package to estimate genetic regulatory variation using a confidence interval estimation methods, applying simulation-based techniques like bootstrapping to evaluate the estimates' reliability.
- Implemented binomial distributions and various other statistical concepts, and used sensitivity analysis to quantify the impact of different assumptions on the genetic data from the Genotype Tissue Expression Project (GTEx).

## SKILLS

- **Languages:** Python, Java, R, C++, SQL, JavaScript
- **Frameworks and Tools:** PyTorch, Tensorflow, ONNX, LangChain, Pinecone, Keras, numpy, sklearn, pandas, seaborn, plotly, ggplot, AWS [S3, EC2, SageMaker, Kinesis, Lambda], SnowflakeDB, Tableau, Streamlit, git

## PROJECTS

- **Capstone Project**: Active Learning with Neural Processes for Epidemiology Modeling
- Prediction of Causes of Patient Readmission using Large Language Model(s) (**OpenAI, Llama-7B, BART-Large**)
- Data Science Interview Tool (**GPT-3**, PyTorch, Python)
- Autonomous Vehicle Trajectory using Deep Learning (**Argoverse 2 Dataset**, PyTorch)