

JOHNNY HUANG

San Mateo, CA, 94403 | (650)-278-6570 | h.johnny@wustl.edu | [Portfolio](#) | [Linkedin](#) | [GitHub](#)

EDUCATION

St. Louis, MO

Washington University in St. Louis

- **B.S/M.S: Mathematics + Computer Science** *Aug. 2022 – May 2026*
- **Honors and Activities:** Chancellor's Career Fellow, Taylor Scholar (full tuition), 3x WashU Hackathon (2023 Co. Organizer)
- **Relevant Coursework:** Machine Learning, Object Oriented Design, Advanced DSA, System Design, Bayesian Statistics, Optimization

PROFESSIONAL EXPERIENCES & INVOLVEMENT

Software Engineering Intern

June – Sep. 2025

City and County of San Francisco

San Francisco, CA

- Scaled **Flask**-based webserver used by all city administrators from ~40 to ~800 daily users; coordinated workers & threading using **Gunicorn** on **Linux** w/ load balancing; resulting server is non-blocking & stable for **500+** concurrent RPM; designed careful server locks for STA processes & offloading CPU tasks to **Celery** ensuring **100%** server up time.
- Deployed a light-weight **Redis** cache in backend for frequent & heavy DB queries; configured HTTP headers for caching static assets client side; drastically reducing Largest Content Paint time by **84%** to **<900ms** for page loads on average; introduced asynchronous models in **JS** frontend to eliminate UI freezes.
- Optimized **SQL** database queries & managed DB sessions to alleviate existing pool exhaustion issues; established a new read-only DB replica for higher throughput reads & improved DB availability.

Machine Learning Intern

Jan. – Apr. 2025

SF Office Of The Chief Medical Examiner

San Francisco, CA

- Designed a LLM for generating concise impressions for complex forensic cases; applied 8-bit quantization to Qwen-235B & performed parameter-efficient tuning w/ **HF's** LoRA library; pipeline fully configured for **CUDA**; achieved training speedup of **400%** & memory overhead reduction of **50%**; resulting impressions sped up SF's forensic case turnaround time by **15%** MoM.
- Built a plug-and-play training repo w/ **HF** & **PyTorch** for validating key lab reports & court documents sent to permanent storage; image-based validation using OpenAI's CLIP; resulting workflow is **100%** unsupervised & virtually eliminates **100%** of false negatives.

Large Language Models Intern

May – Aug. 2024

Rad AI

San Francisco, CA

- Led developments of a full-stack search feature using a RAG framework on **300k+** internal patient records with 2 other interns; stored embedded documents in a **Pinecone** vector database w/ **Langchain** pipeline for optimized retrieval; **pydantic** for API validation
- Established lightning-speed asynchronous API's for inference using **FastAPI**; deployed through **Docker** containers on **AWS** ec2 **Linux** instances; achieved inference speed of **<2000ms** per request + model query.

PERSONAL PROJECTS (See [More!](#))

- **Diary App:** [Website](#) for users to take notes/write diaries; built w/ **ASP.NET** MVC framework in **C#**; setup REST API's in backend; stores user information persistently using a document-based JSON **NoSQL** database.
- **URL Shortner:** **Golang** [webserver](#) for shortening any URL designed to be fast; setup go routines & channels for fully non-blocking & concurrent-safe model; **JS** for real-time DB changes; MRU cache & **NoSQL** DB loaded in RAM for instant access.
- **Petrichor:** [Mental health application](#) aiming to match users with their perfect therapist; implementing user login, modules, calendars; backend logic built with **NodeJS** & **ExpressJS** app; data stored w/ **MySQL**; containerized app and deployed on **AWS**.

SKILLS

- **ML:** PyTorch, TensorFlow, HuggingFace, Pydantic, Langchain, PySpark, Pinecone, CUDA, Pandas
- **Backend:** Java, C++, Python, NodeJS, ExpressJS, FastAPI, SQL, NoSQL, C#, .NET, Golang, Flask, Celery
- **Frameworks:** Jupyter, AWS, MongoDB, Docker, Kubernetes, Git, Bash, Linux, Gunicorn & Uvicorn, Redis