

PRAHALADH CHANDRAHASAN

prahald92@gmail.com | [Github](#) | [LinkedIn](#) | [Website](#)

EDUCATION

CARNEGIE MELLON UNIVERSITY (SCHOOL OF COMPUTER SCIENCE)

Master's in Information Technology (GPA : 3.92/4.33)

Pittsburgh, PA

Aug 2024 - Dec 2025

SKILLS

Programming Languages : C++, Python, SQL, JavaScript, Typescript, Shell Scripting (Git & Bash)

Frameworks and Libraries : Pytorch, FastAPI, nodejs, Express JS, HuggingFace, OpenAI, Tensorflow, vLLM, wandb, MLFlow, LangChain, LlamalIndex, Anthropic Agent SDK, React

Cloud & Devops : AWS(EC2, S3, EKS, ECR, RDS, IAM, Bedrock, SageMaker), Azure, Redis, Kubernetes, Docker, GitHub Actions

WORK EXPERIENCE

Founding Forward Deployed Engineer

Circle AI

San Francisco, CA

Jan 2025-Present

- Built and maintained end-to-end CI/CD infrastructure and staging environments : designed a complete pre-production environment on **Azure Container Apps** with isolated **OAuth configs**, multi-tier deployment strategies, and PR-aware workflow gating that eliminated redundant deployments for release PRs
- **Optimized AI agent container performance** : replaced HTTP polling with persistent **WebSocket** connections for agent-server communication, pre-baked container base layers for faster cold starts, removed unnecessary cloud credential provisioning from ephemeral containers, and migrated file uploads to direct browser-to-storage to bypass server-side payload limits
- **Hardened agent security against prompt injection and privilege escalation** : implemented PreToolUse hook system to intercept destructive bash commands and block environment variable exfiltration attempts, restricted container filesystem access, removed embedded cloud credentials, and enforced multi-tenant data isolation via Row-Level Security with impersonation-aware policies
- Shipped full-stack features and large-scale codebase improvements — built a Google Docs-style task sharing system (RLS, REST APIs, real-time search, org-wide sharing), fixed critical container file sync bugs, migrated Python services to **structured logging**, and removed 20,000+ lines off dead code to reduce maintenance burden and attack surface

Language Technologies Institute

Pittsburgh, PA

Machine Learning Engineer

Jan 2025 - Present

- Engineered and deployed **DeepResearch Comparator**, a large-scale agentic evaluation platform on an **EKS Cluster**, enabling **fine-grained human feedback collection** and benchmarking of closed- and open-source agents
- Productionized **deep research agents** and **multimodal RAG systems** by exposing them as **scalable APIs**, integrating visualization and monitoring of outputs and metrics using **ZenoML**
- Benchmarked DeepResearch agents (e.g., WebThinker, custom agents) on **BrowseComp** and **HealthBench** using **HPC clusters**, optimizing throughput with the **vLLM** library by tuning model serving parameters for various GPU configurations
- Built a live arena-style platform for **NeurIPS MMU-RAG (Text-to-Video track)** as a cloud-native application to host baseline and participant VideoRAG submissions, developed baseline VideoRAG systems, and a fast evaluation pipeline for submissions on **VBench**

Bank of America Continuum India

Chennai, India

Software Engineer

Jul 2022 - Jul 2024

- Accelerated testing efficiency by implementing Tosca-based automation for comprehensive end-to-end payment flows from initiation to clearing, eliminating **1,000+** manual regression testing hours annually and improving release velocity by **65%**
- Architected reusable Tosca UI and API modules deployed across 5 regional payment landscapes, reducing development cycles by **40%**
- Detected and resolved 5+ critical production defects through rigorous testing protocols, preventing potential revenue loss of **\$5 million**
- Spearheaded innovation by developing and filing a **patent** for a payments fraud detection algorithm and presented work across various business verticals

RedHat

Bangalore, India

Software Engineer Intern

Jan 2022 - Jul 2022

- Engineered and delivered two critical features ([ENTESB-18633](#) and [ENTESB-18785](#)) successfully integrated into [Hawtio release 7.11](#) using **AngularJS** and **PatternFly** framework, enhancing real-time monitoring capabilities for RedHat Fuse environments and improving enterprise customer experience

PROJECTS

Evaluation methods for identifying gender bias in LLMs

- Designed a prompt-conditioned activation steering pipeline that routes incoming queries through a lightweight classifier at inference time, applying targeted latent suppression only to malicious prompts — significantly outperforming unconditional clamping baselines on both forgetting effectiveness and knowledge retention
- Built a synthetic adversarial dataset generation pipeline using Llama 3.3-70B to produce diverse labeled prompts spanning harmful and benign intent classes, enabling robust classifier training that generalized across unseen evaluation distributions without relying on real-world hazardous data