

Viresh Duvvuri

Seattle, WA | +1-509-964-5469 | vireshduvvuri@gmail.com | linkedin.com/in/viresh-duvvuri

Senior Software Engineer with 7+ years building AI-first features, distributed systems, and scalable applications from the ground up. Expert in modern AI tooling (LangChain, RAG, vector search) integrated with full-stack development (Python, TypeScript, React) and cloud-native architectures. Proven track record in fast-paced startup environments, delivering production systems serving 200+ daily users with 50-80% efficiency gains through end-to-end ownership from concept through deployment on AWS.

Skills

Programming: Python, TypeScript, JavaScript, C++, SQL, FastAPI, Flask, React, Node.js, NumPy, Pandas, OOP

AI/ML Frameworks: Agentic AI, LangChain, LangGraph, Multi-Agent Systems, MCP (Model Context Protocol), RAG, Context Engineering, Prompt Engineering, Model Evaluation, MLOps, GenAI, FAISS, Pinecone, PyTorch, TensorFlow, Scikit-learn, Feature Engineering, Human-in-the-Loop (HIL), Model Deployment, Responsible AI, Vector Search

Cloud & Infrastructure: AWS, Azure, API Design, Deployment, DevOps, Docker, Kubernetes, Monitoring, Performance Tuning, Scalability, Workflows, Microservices, CI/CD

Data & Analytics: Data Integration, Data Processing, Data Science, Enterprise Integrations, Enterprise Systems, Knowledge Graph, Operational Efficiency

Work Experience

Grid CoOperator

Seattle, WA

AI Engineer

Mar 2025 - Present

- Built multi-agent AI platform from scratch in early-stage startup environment, architecting distributed system using LangChain and MCP frameworks where specialized agents coordinate via APIs to process smart grid analytics, deployed on AWS with microservices architecture and comprehensive observability, reducing analyst workflows by 70% within 2 months through rapid iteration and continuous deployment
- Designed and implemented scalable AI orchestration system with RESTful APIs and event-driven architecture, deployed on cloud infrastructure with monitoring dashboards tracking quality metrics, latency, and cost across 50-100 daily queries, achieving 99%+ uptime through robust error handling and automated recovery mechanisms
- Developed end-to-end production AI system with CI/CD pipelines, Docker containerization, and performance optimization, collaborating cross-functionally with business stakeholders to translate requirements into technical solutions and accelerating deliverables by 60% through agile development practices and continuous feedback loops

Freefly Systems

Woodinville, WA

Senior Software Engineer

Nov 2021 - Oct 2025

- Independently designed, built, and deployed full-stack AI diagnostic system from concept to production serving 200+ daily queries, architecting React frontend with TypeScript and Python Flask REST APIs, integrating foundation model APIs (Ollama, Llama 3.2) with RAG architecture and vector search, containerized with Docker and deployed to production infrastructure reducing expert analysis time from hours to minutes
- Built automated systems to process complex technical data and identify system failures, developing knowledge base enhancements and support tools that streamlined operations
- Contributed to drone platform codebases implementing new features and optimizations for flight control systems and payload integration across multiple product lines, managed software integration projects from planning through release
- Led release management for drone platforms overseeing testing phases from alpha through production deployment, coordinating firmware updates and executing comprehensive testing protocols with cross-functional teams

Lumenier

Sarasota, FL

Drone Software Developer

Jul 2020 - Oct 2021

- Wrote embedded code in C++ to integrate LiDAR and optical flow sensors for obstacle avoidance and position holding with/without GPS under various lighting conditions
- Collaborated with open-source flight control software maintainers for integration, testing, and deployment of autonomous flight algorithms, prototyped innovative features like toss-to-launch for product roadmap development

York Exponential

York, PA

Software Engineer - R&D

Aug 2018 - May 2020

- Developed prototype software for in-house autonomous surveillance mobile robots using ROS2, SLAM, and computer vision technologies
- Built Human Machine Interface for Universal Robot welding applications using Python and Kivy framework, implemented multi-robot control systems with platform independence

Education

Washington State University

Master of Science Computer Science

Pullman, WA

Jan 2015 - Jan 2017

GITAM University

Bachelor of Technology Information Technology

Visakhapatnam, India

Jan 2011 - Jan 2015

Projects

GridCOP: Smart Grid Analytics Agent

- Problem: Power grid analysts needed automated database querying and intelligent insights to understand complex data patterns beyond basic visualizations
- Solution: Developed A2A multi-agent system using LangChain orchestration and MCP where specialized agents coordinate tasks through prompt engineering strategies, implemented RAG and vector search (FAISS) for intelligent querying, implemented model evaluation frameworks to monitor quality and cost metrics, deployed on AWS with observability and logging
- Impact: Enhanced analyst productivity by 70% through AI co-pilot that augments domain experts with automated workflows, implemented human-in-the-loop (HIL) evaluation and testing pipelines for production-ready AI systems with robust error handling through rapid iteration

Production System Optimization Tool

- Problem: Manual system analysis taking hours of expert time, creating bottlenecks in product development and customer support resolution
- Solution: Built full-stack application with React frontend, Python Flask backend, integrated foundation model APIs (Ollama and Llama 3.2) for real-time log processing and interactive analysis using prompt engineering and model evaluation
- Impact: Transformed expert analysis from hours to minutes, deployed to production serving 200+ daily queries with significant performance improvements through rapid iteration and continuous optimization

AI Travel Planner Agent

- Problem: Manual travel planning requiring hours of research across multiple sources with inconsistent and outdated information
- Solution: Built AI agent using Claude 3.5 Sonnet, LangChain, Streamlit, and DuckDuckGo Search API for personalized itinerary generation using prompt engineering techniques
- Impact: Demonstrated end-to-end AI application development, learned conversational AI patterns and real-time data integration techniques through iterative development