

AYUSH KUMAR MALIK

📞 +1(930) 333-4460 | 📩 ayushkumarmalik10@gmail.com | 💬 linkedin.com/in/ayush67

Education

Indiana University Bloomington
MS in Computer Science, GPA: 3.9/4.0

Bloomington, IN
Graduated: 2025

Shiv Nadar University, Delhi NCR
BSc in Computer Science, GPA: 3.8/4.0

Greater Noida, India
Graduated: May 2020

Work Experience

Brckt (Peristyle Labs)

AI/ML Engineer

Dec 2024 – Present

Indianapolis, IN (Remote)

- Built **real-time tennis match analysis system** using **Llama 3.3-70B** via Venice.ai API, generating professional head-to-head predictions with streaming responses.
- Developed **web scraping infrastructure** using **Playwright** headless browser with anti-detection measures (user-agent rotation, heading-based navigation), extracting H2H stats from matchstat.com.
- Implemented **TTL caching layer** with thread-safe operations, order-independent keys, and automatic eviction, reducing redundant scraping by caching H2H data for 2 hours.
- Deployed **FastAPI** backend with **Server-Sent Events (SSE)** for real-time streaming, **Docker** containerization, and **Caddy** reverse proxy handling HTTPS termination.

Riverside Global LLC

AI Engineer

Jun 2025 – Dec 2025

Hampton, IL (Remote)

- Architected production **RAG system** with 5-stage pipeline: query routing, reformulation, **hybrid retrieval** (BM25 + semantic), cross-encoder reranking, and GPT-4 generation with hallucination checking, reducing document research time by 60%.
- Built **hybrid search engine** combining Sentence-BERT embeddings with BM25 keyword matching using **Reciprocal Rank Fusion (RRF)**, achieving 94% retrieval relevance on 10,000+ environmental documents.
- Developed **LLM-powered data extraction** pipeline using GPT-4 function calling with custom JSON schemas, achieving 95% accuracy and reducing manual extraction from 3 hours to 15 minutes per document.
- Implemented document classification system for permits, water quality reports, and EPA notices with automated validation checks and human-in-the-loop review for compliance workflows.

Projects

CodePilot: Multi-Agent AI Coding System

[GitHub] | [Live Demo]

- Architected **multi-agent orchestration system** with 4 specialized agents (Planner, Coder, Reviewer, Explorer) using **Claude Sonnet 4.5** and function calling for autonomous code generation.
- Designed **hybrid retrieval engine** combining BM25 keyword search with semantic embeddings using **Reciprocal Rank Fusion (RRF)**, improving code retrieval precision by 25% over single-method approaches.
- Implemented **token-efficient context tools** (file outlines, code chunk extraction) achieving **40x reduction** in token consumption vs. naive full-file reads.
- Built iterative **feedback loop** between Coder and Reviewer agents with state machine orchestration; deployed on **HuggingFace Spaces** with real-time agent visualization.

VerbaQuery: Semantic Search & RAG System

[GitHub]

- Built production **RAG pipeline** using **Sentence-BERT (all-MiniLM-L6-v2)** embeddings to encode 1,000+ document pages into 384-dimensional dense vectors for semantic retrieval.
- Implemented **hybrid retrieval** combining BM25 sparse vectors with dense embeddings, achieving **92% recall@10** compared to 78% with keyword search alone.
- Integrated **cross-encoder reranking** (ms-marco-MiniLM) to refine top-k candidates, reducing hallucination rate by 40% in generated responses.
- Deployed vector storage using **FAISS** with IVF indexing for sub-100ms query latency at scale.

LLM Evaluation Framework & Fine-Tuning Experiments

[GitHub]

- Developed comprehensive **evaluation framework** measuring LLM performance across accuracy, latency, and cost metrics for function calling and structured output tasks.
- Benchmarked GPT-4, GPT-3.5-Turbo, and Llama models on **JSON schema compliance** and instruction following, achieving **95% cost reduction** through optimal model selection.
- Built modular **fine-tuning pipeline** with dataset preparation, training loops, and automated metric collection for model specialization experiments.
- Conducted statistical analysis of model outputs, generating reports on accuracy distributions, failure modes, and confidence calibration.

Technical Skills

Languages: Python, SQL

ML Frameworks: PyTorch, HuggingFace Transformers, Sentence-Transformers

LLM & RAG: GPT-4 API, Claude API, Function Calling, Hybrid Search (BM25 + Semantic), Cross-Encoder Reranking (ms-marco-MiniLM), Prompt Engineering

Vector Databases: ChromaDB, FAISS

Infrastructure: AWS (EC2, S3), Docker, FastAPI, PostgreSQL, pdfplumber