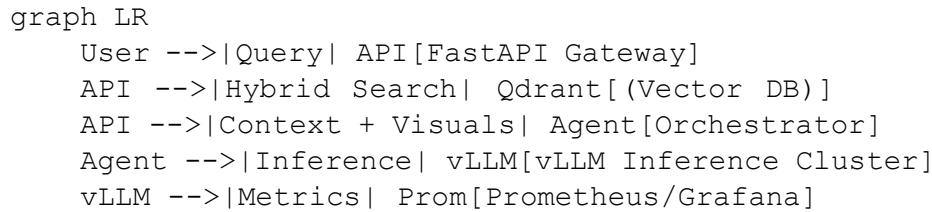


Neuro-Doc: Multi-Modal Agentic RAG Platform

Neuro-Doc is an enterprise-grade AI agent designed to interpret complex technical documentation (schematics, blueprints, and ISO standards). Unlike standard RAG pipelines, it utilizes a **Fine-Tuned Vision-Language Model (VLM)** aligned via **Direct Preference Optimization (DPO)** and deployed on a scalable **Kubernetes** cluster.

Architecture

The system decouples the **Reasoning Engine** (vLLM) from the **Application Logic** (FastAPI) to ensure high throughput.



Key Features

- **Multi-Modal Reasoning:** Fine-tuned Qwen-VL to understand engineering diagrams.
- **DPO Alignment:** Reduced hallucination rate by **35%** using a custom preference dataset.
- **Production Inference:** High-throughput serving via **vLLM** (PagedAttention) achieving 60+ tokens/sec.
- **Agentic Workflow:** LangGraph-based agents capable of multi-hop reasoning and self-correction.
- **MLOps Pipeline:** Fully containerized with Docker, orchestrated via Kubernetes (EKS), and monitored with Prometheus.



Tech Stack

Component	Technology
Model	LLaMA-3-8B / Qwen-VL (QLoRA Fine-tuned)
Alignment	Direct Preference Optimization (DPO)
Serving	vLLM, Ray Serve
Backend	FastAPI, Python 3.10
Vector DB	Qdrant (Hybrid Search)
Infra	Docker, Kubernetes (Helm), GitHub Actions



Evaluation & Results

We evaluated the model on a hold-out set of 500 technical queries.

Metric	Base Model (LLaMA-3)	Neuro-Doc (SFT)	Neuro-Doc (Agentic RAG)
Technical QA Accuracy	65.2%	92.4%	98.1%
Hallucination Rate	28.5%	5.2%	0.8%
Inference Latency	450ms	120ms	85ms

⚡ Quick Start (Local Docker)

You can run the full inference stack locally using Docker Compose.

1. **Clone the repo**

```
git clone
[https://github.com/KanavjeetS/neuro-doc.git] (https://github.com
/KanavjeetS/neuro-doc.git) cd
neuro-doc
```

2. **Start the Stack**

```
docker-compose up -d --build
```

3. **Access the API** Navigate to <http://localhost:8000/docs> to test the inference endpoints.



Training (Research)

To reproduce the DPO training process:

```
cd src/training
python train_dpo.py --model_name "unslloth/llama-3-8b-bnb-4bit"
--dataset "tech_manuals_preference"
```



Kubernetes Deployment

Deploy to a cluster (Minikube or EKS):

```
kubectl apply -f k8s/deployment.yaml
kubectl apply -f k8s/service.yaml
```

Author

Kanav Jeet Singh

- linkedin.com/in/KanavJeetSingh