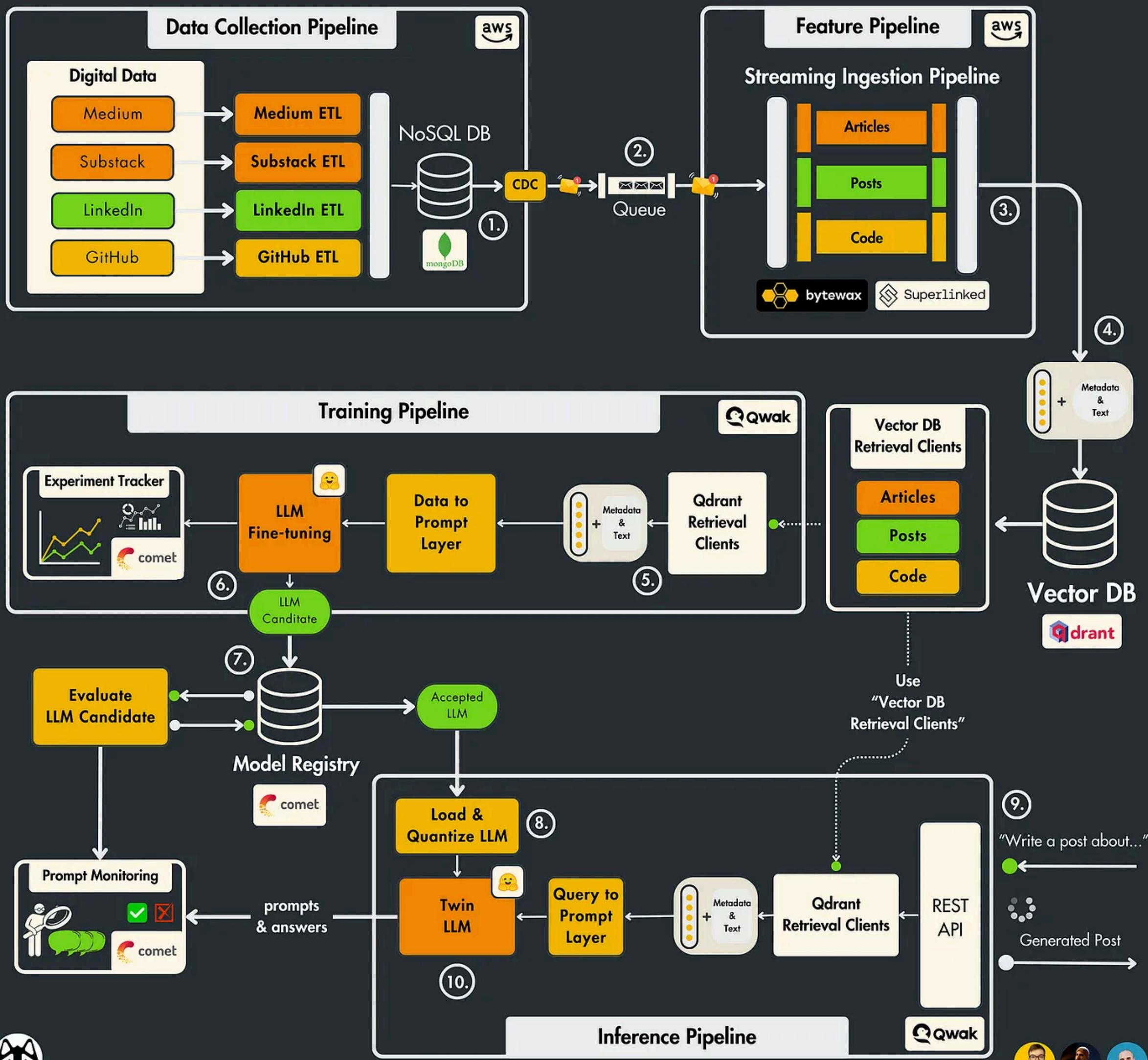
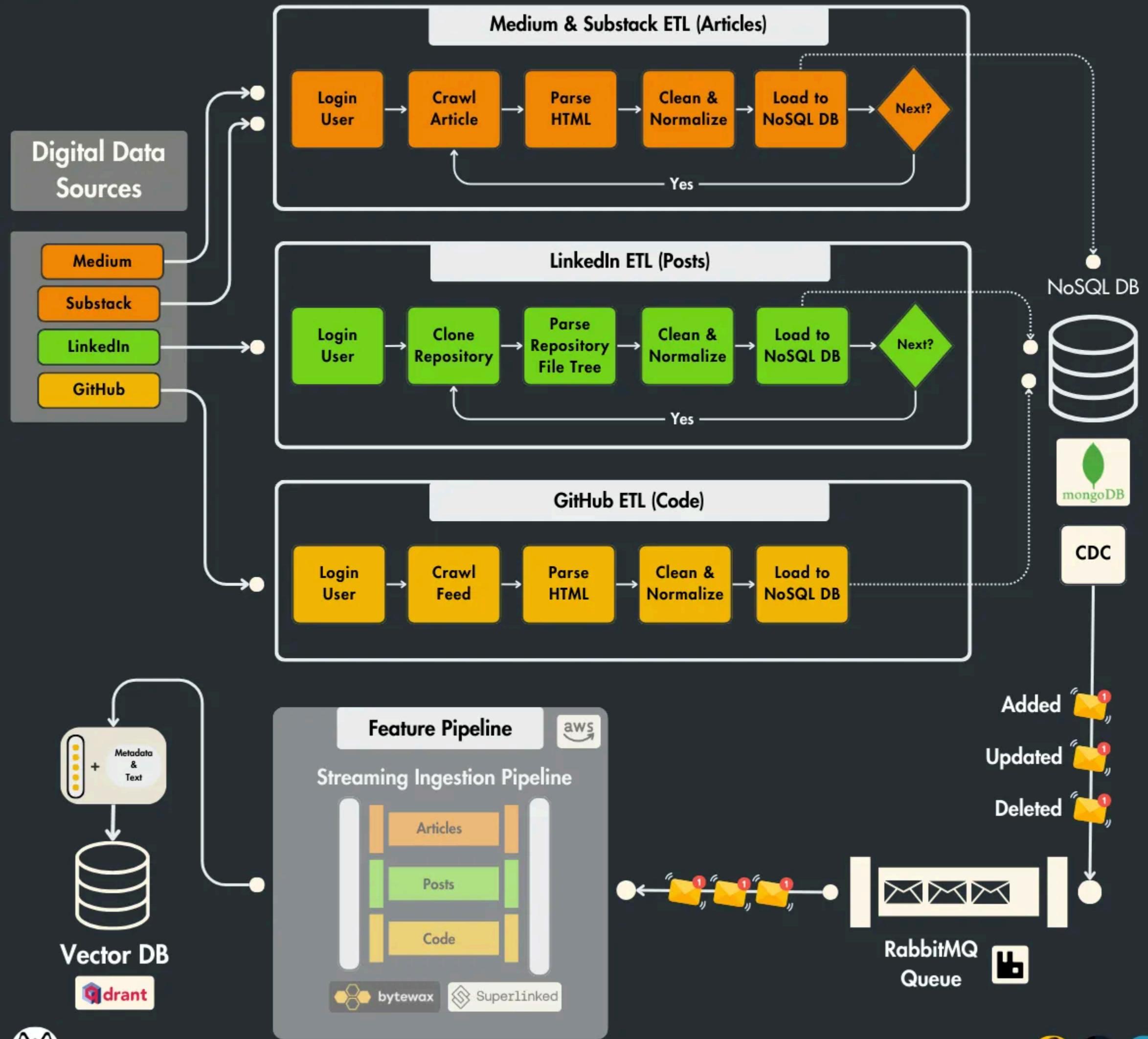


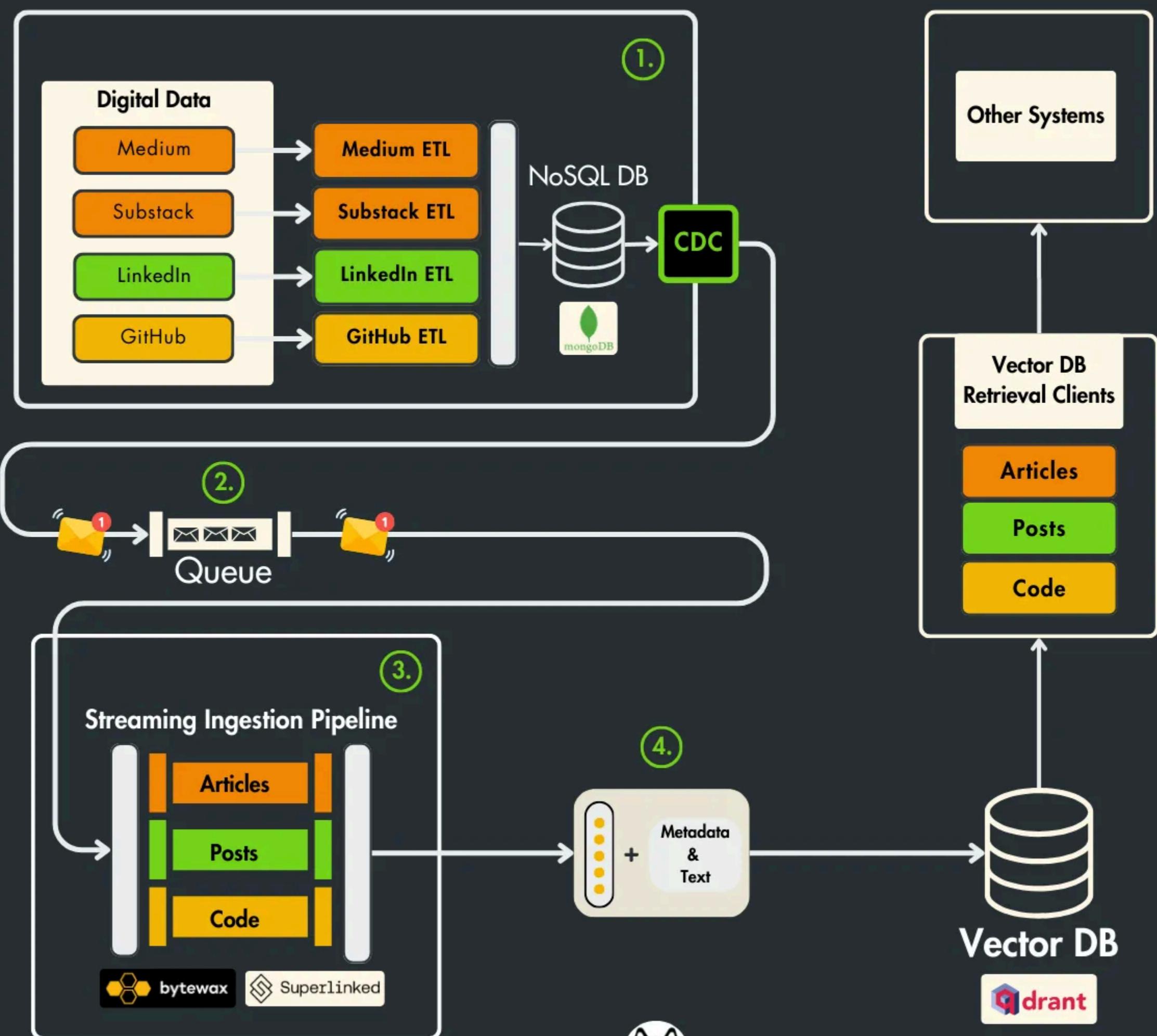
An end-to-end framework for production-ready LLM systems by building your LLM twin



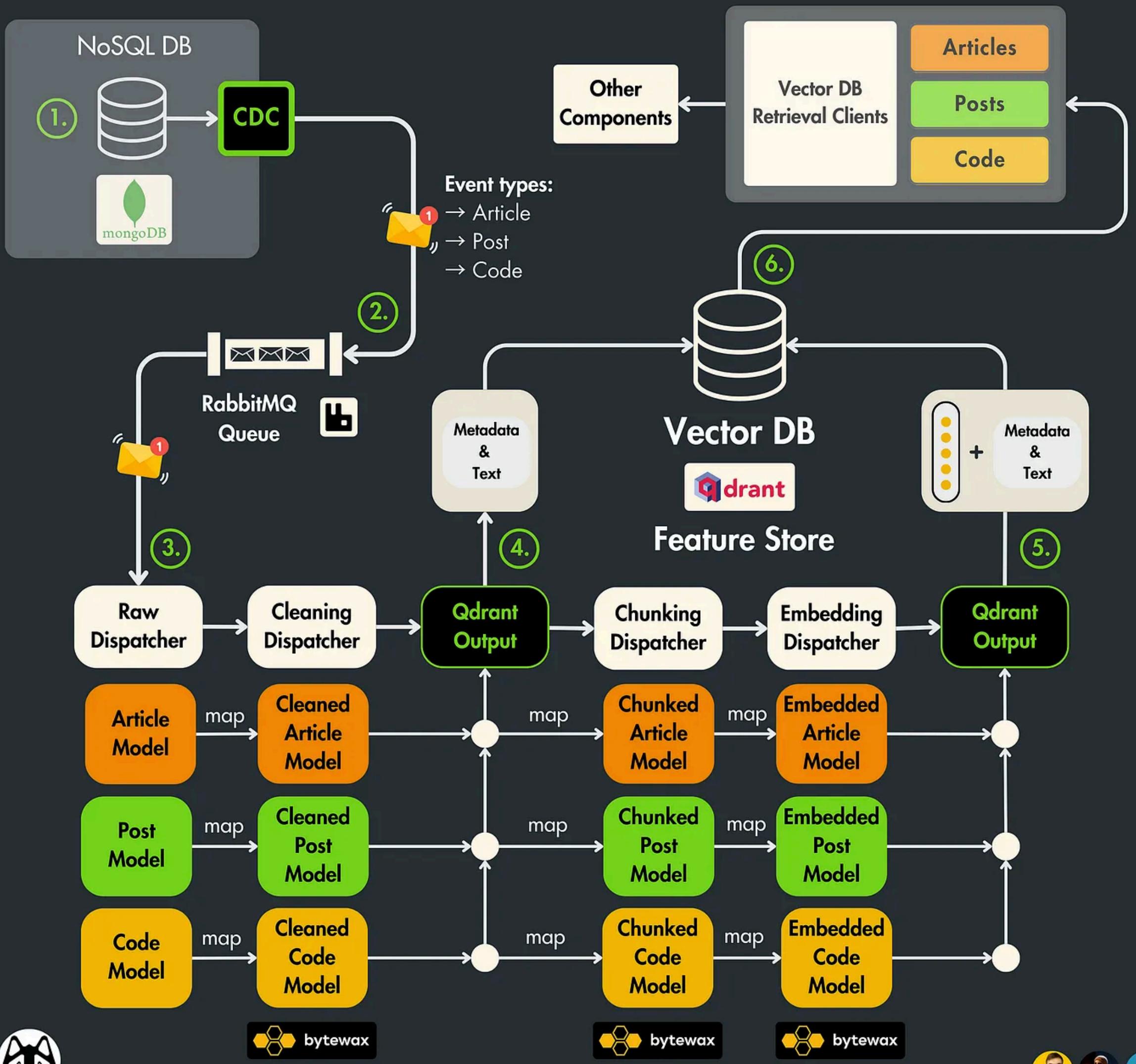
Data Collection Pipeline



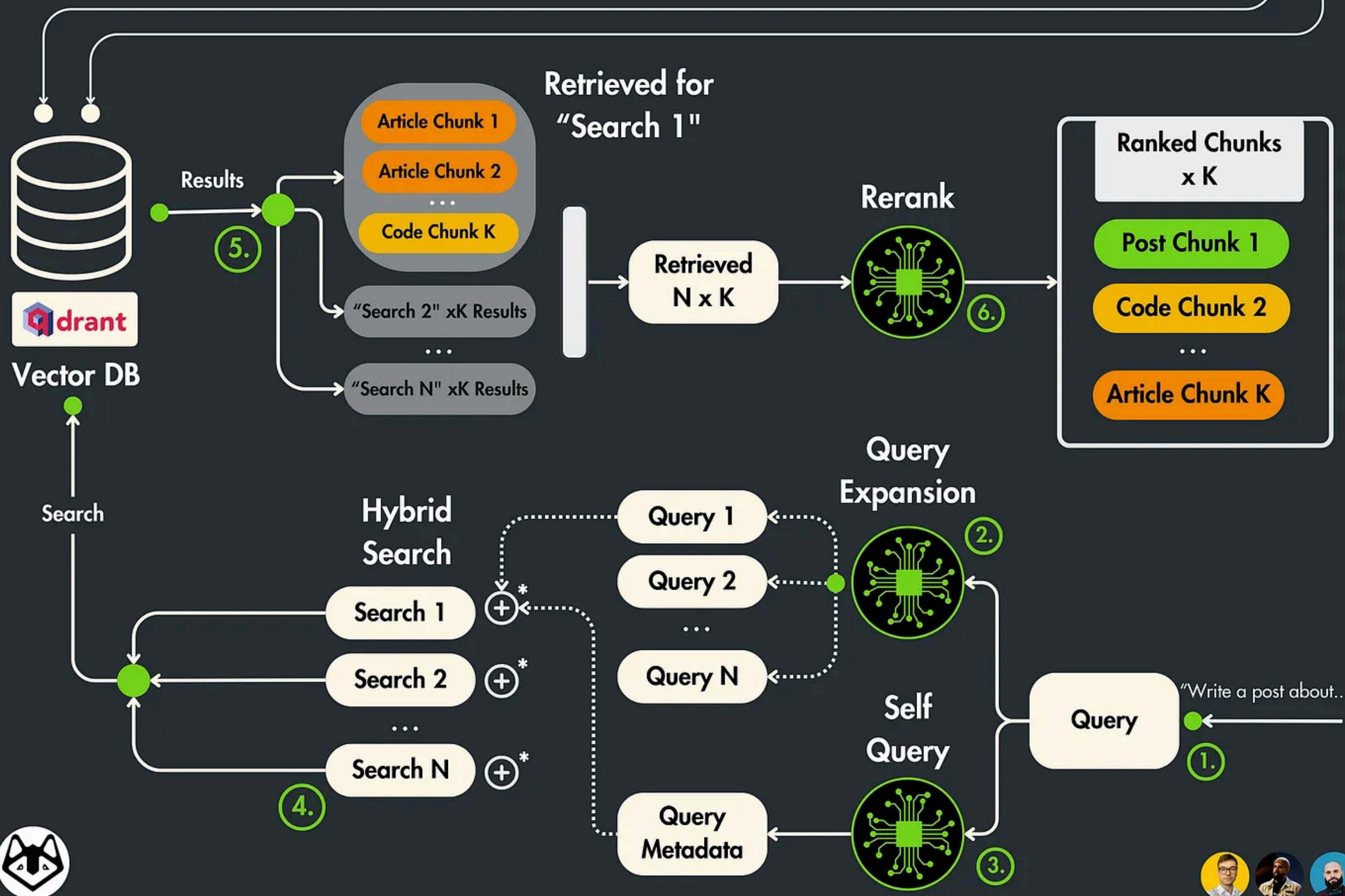
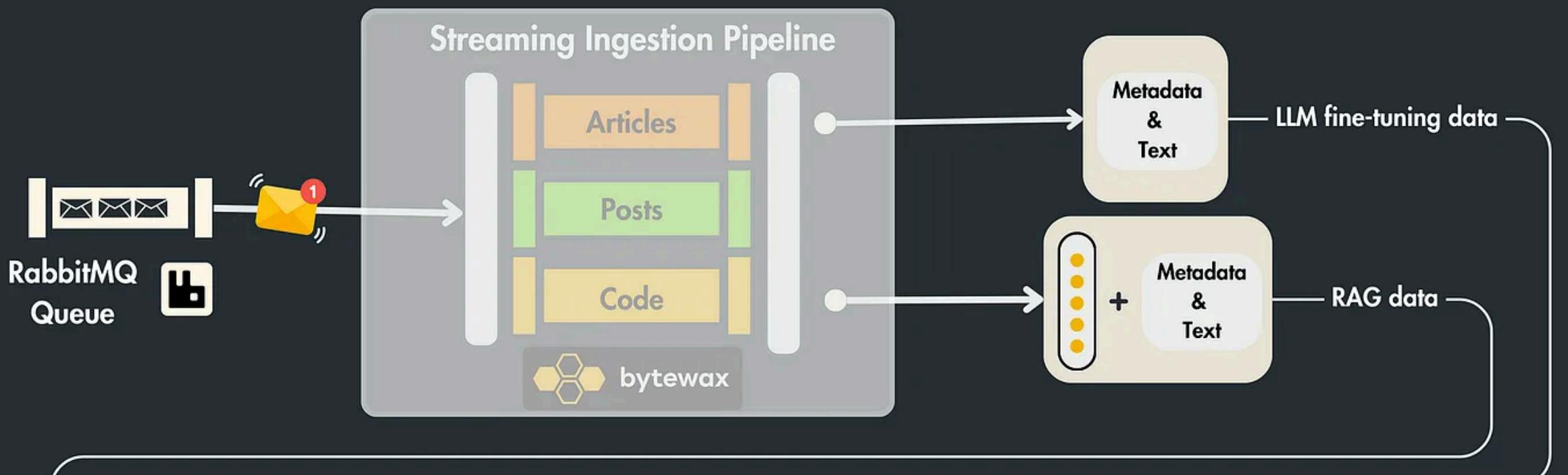
An end-to-end framework for production-ready LLM systems by building your LLM twin



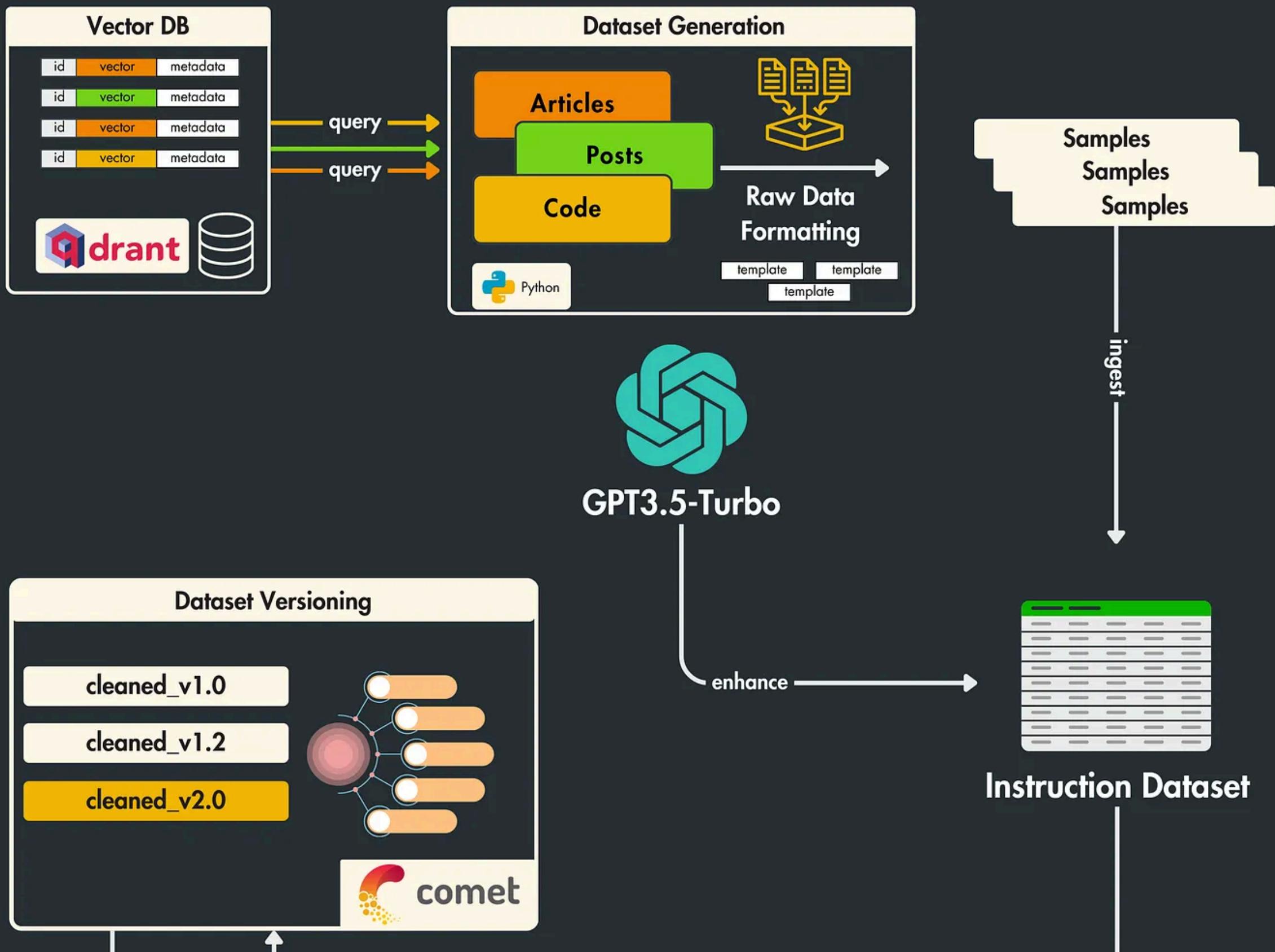
SOTA Python Streaming Pipelines for Fine-tuning LLMs and RAG – in Real-Time!



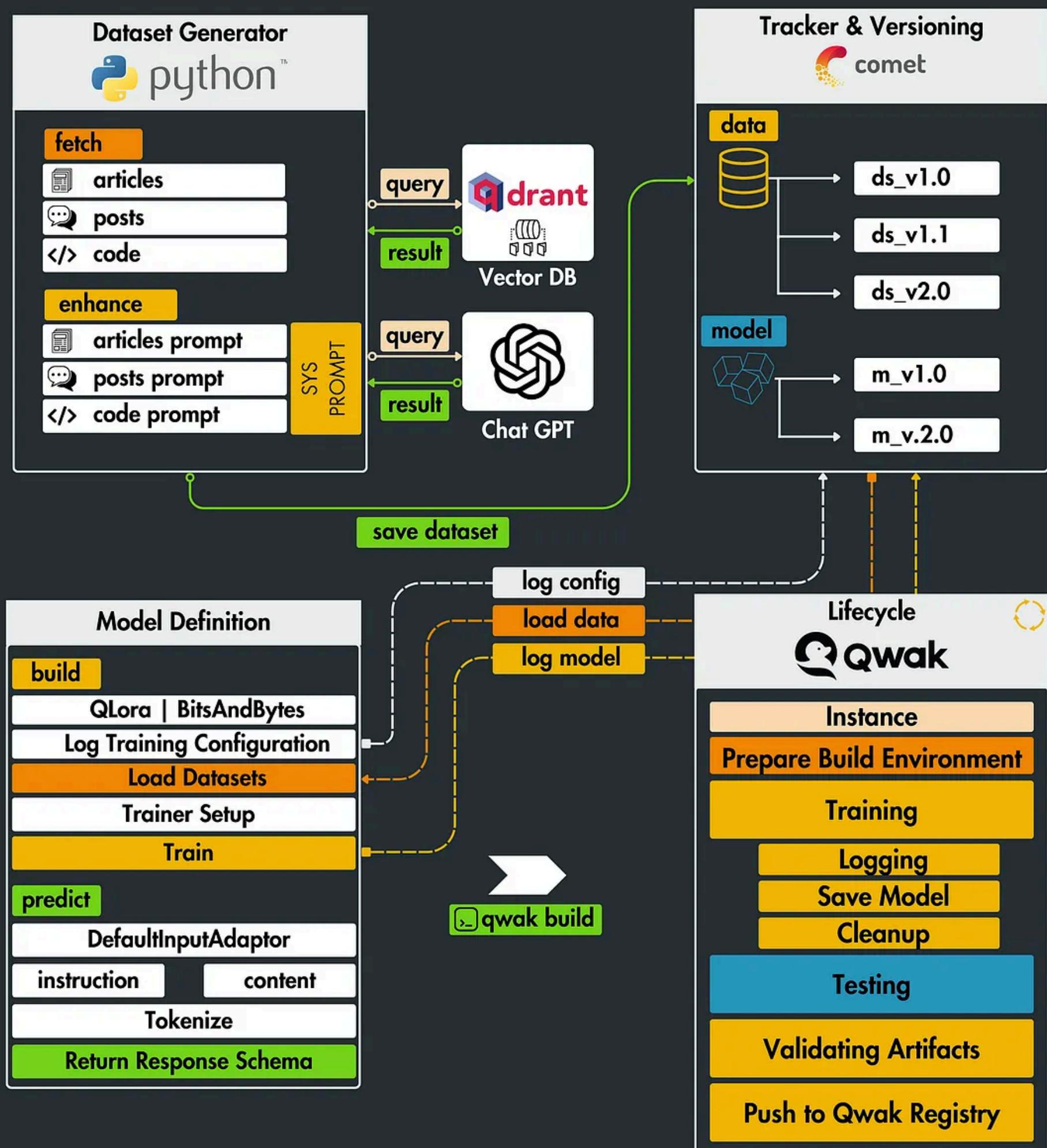
Advanced Retrieval-Augmented Generation (RAG): Retrieval Python Module Architecture



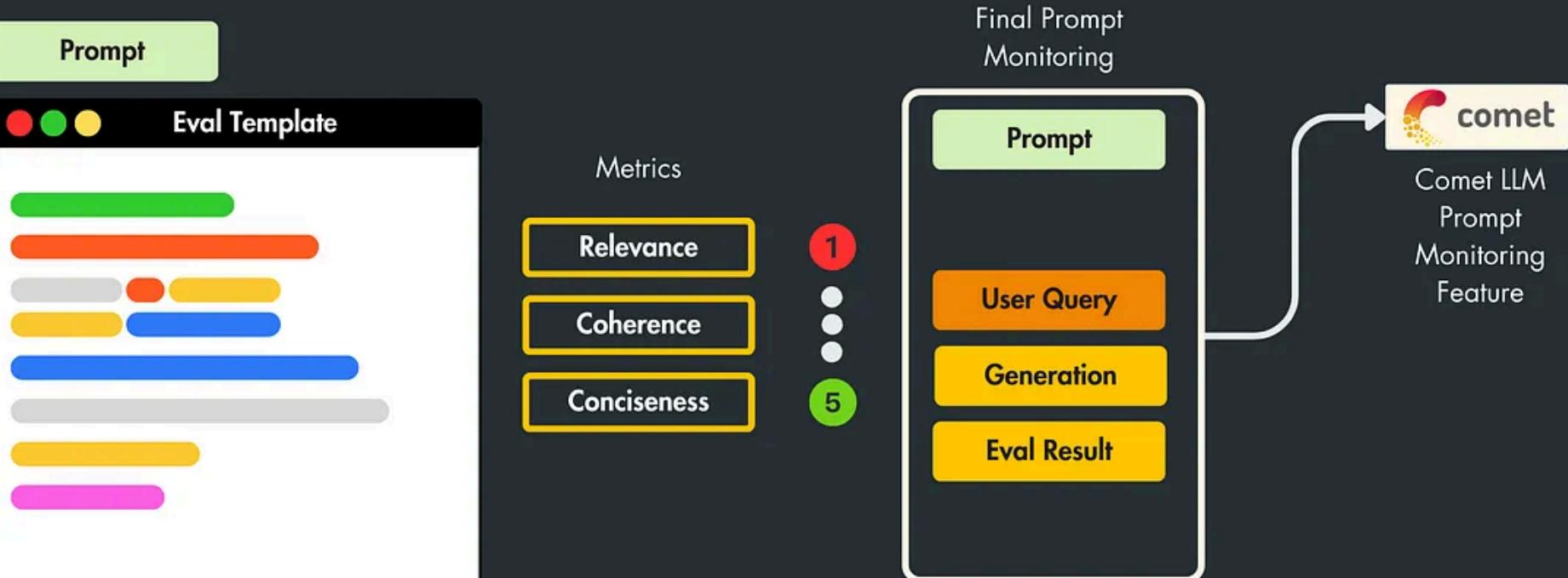
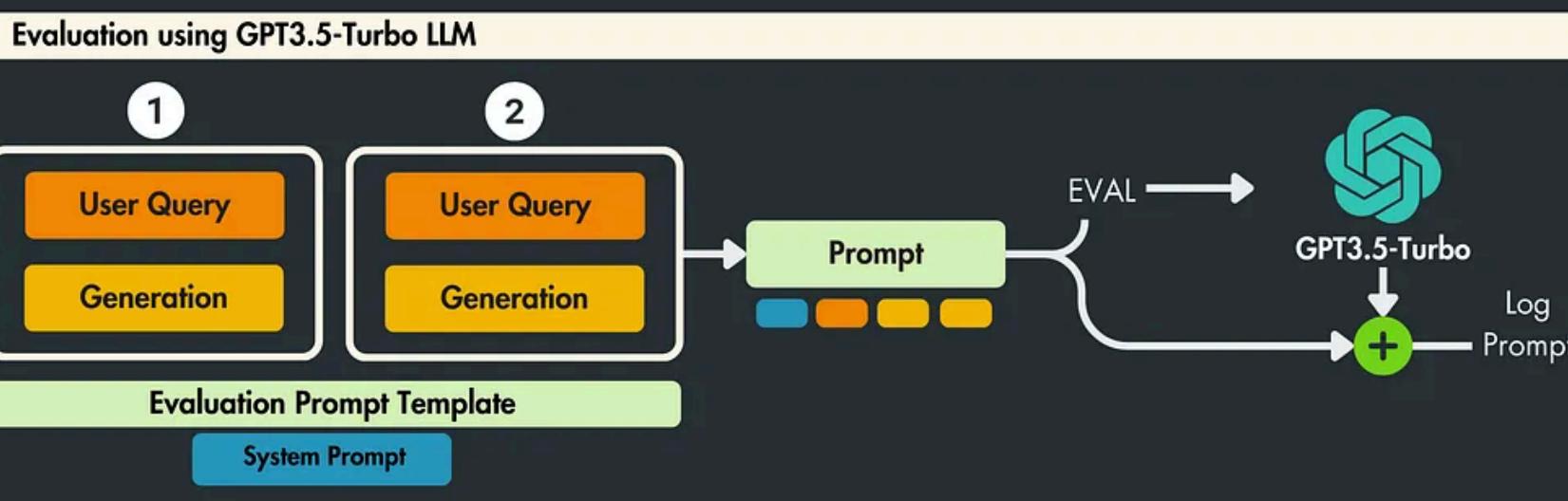
Custom LLM Finetuning Dataset Generation using Knowledge Distillation



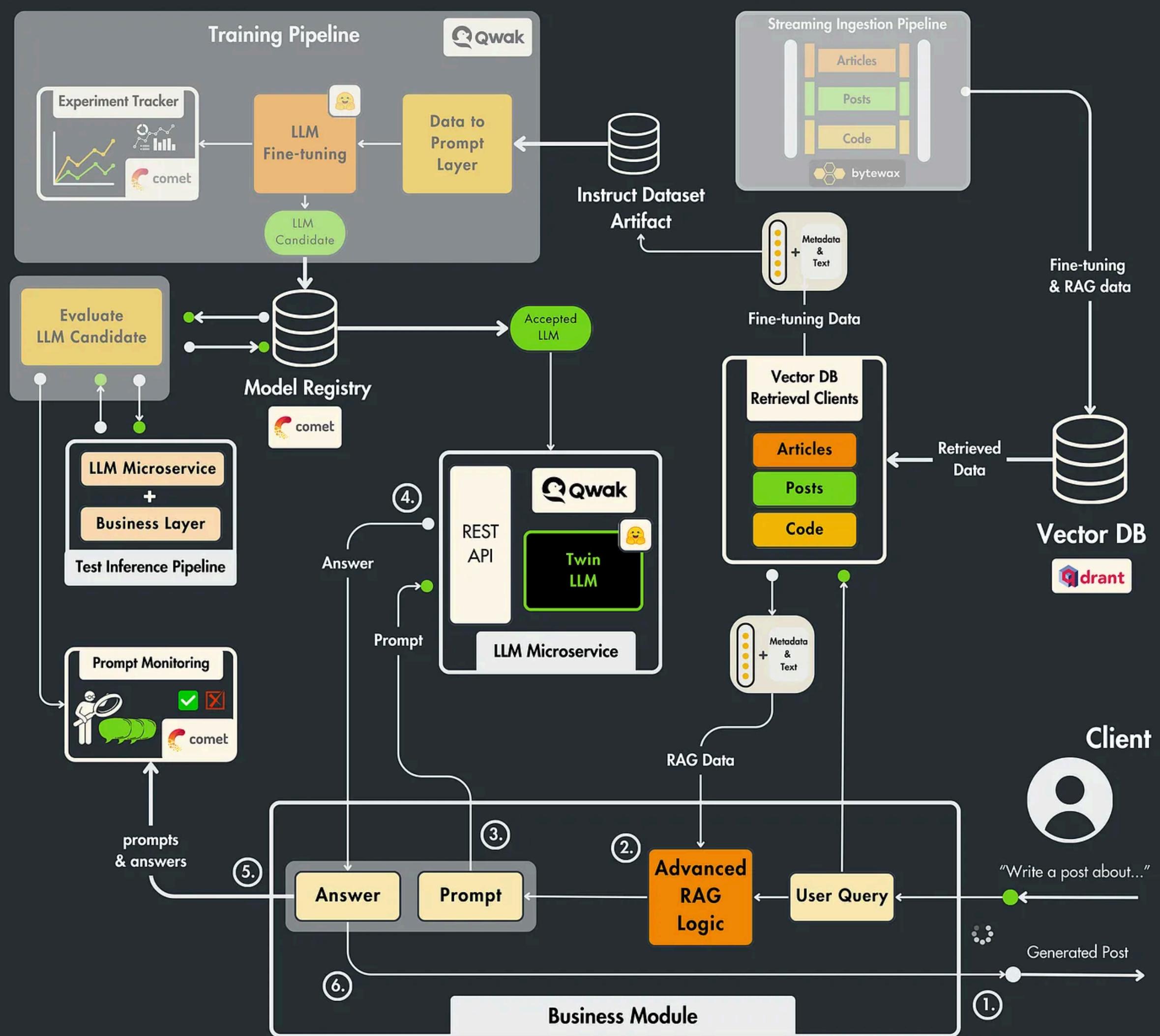
LLM Twim finetuning workflow on Qwak



Evaluating the fine-tuned LLM Twin model

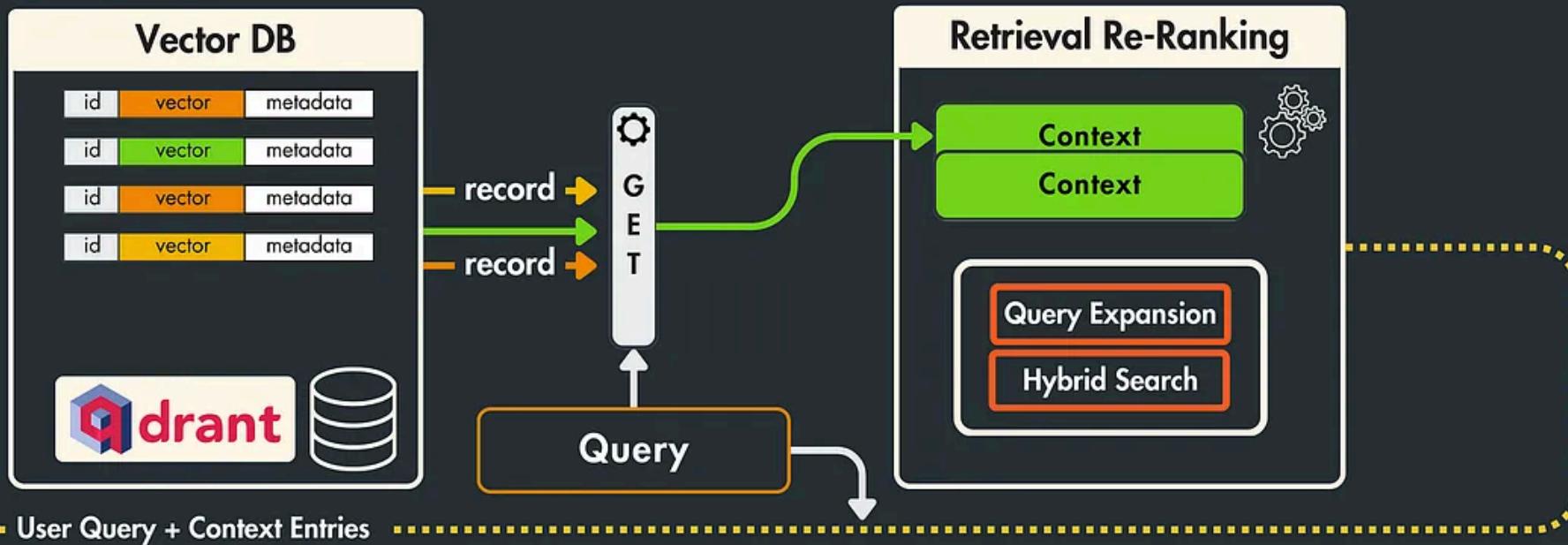


Architect scalable and cost-effective LLM & RAG inference pipelines



Evaluating the LLM-Twin RAG

Context Retrieval



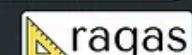
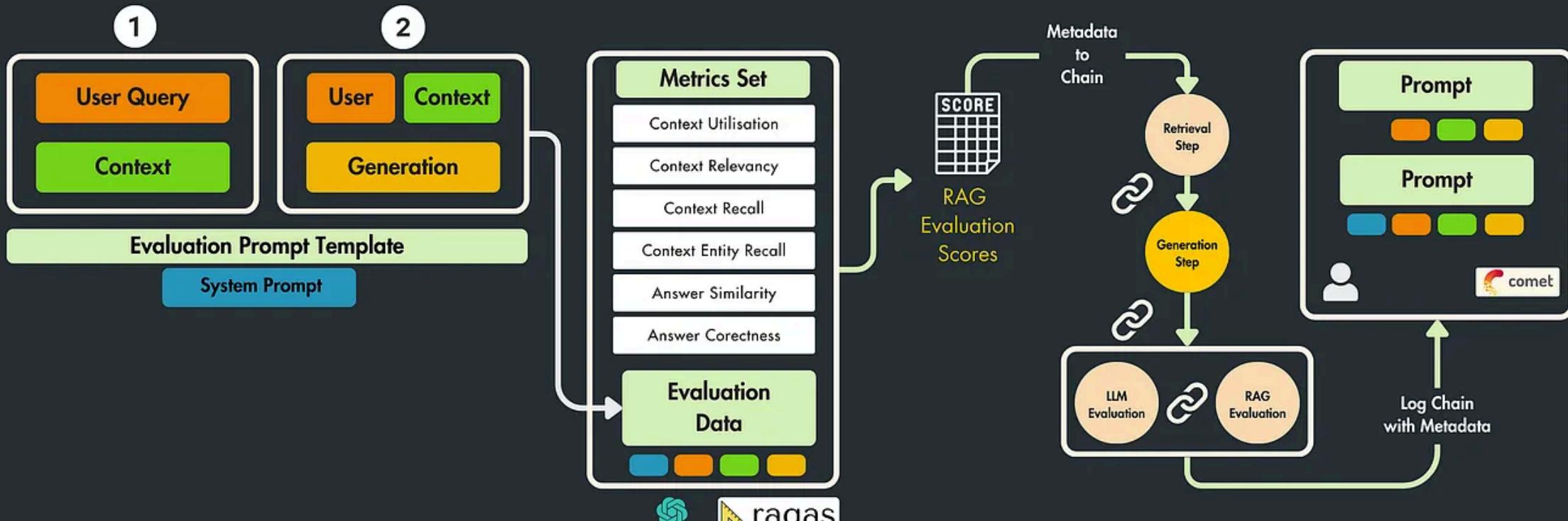
User Query + Context Entries

Inference with fine tuned LLM



LLM Generation Response

RAG Evaluation



Build a scalable RAG ingestion pipeline using 74.3% less code

