# Vector Store Options for RAG

1] FAISS
   - simple, lightweight ] → similarity search over embeddings

2] Chroma
   - open source Vector DB

3] Qdrant
   - open source rust-based Vector DB
   - docker deployment

→ persistance & metadata filtering

# sentence - transformers

↳ Hugging Face Transformers

↳ embedding models optimized for
   i) semantic similarity
   ii) clustering
   iii) retrieval

# Ollama

↳ framework for running full LLMs locally

↳ generally optimized for generation

# (FAISS) Similarity i.e. Distance Metrics

    ↳ L2 Distance
    ↳ Inner Product
    ↳ Cosine

# Similarity Search Methods (Indexing)

    ↳ Exact Search [uses Euclidean distance]

    ↳ Approx. Nearest Neighbor

        ↳ Quantization based

        ↳ Graph-based

    ↳ GPU accelerated Indexes