



## Executive Summary

LLMC should **upgrade its embedding model** to improve code and documentation retrieval, using a *unified model* that handles both source code and natural language well. Based on extensive comparisons of API services (OpenAI, Cohere, Voyage AI) and local open-source models (E5, BGE, Nomic, etc.), we recommend adopting an **open-source embedding model** (e.g. `intfloat/e5-base-v2`) as the default. This model provides **significantly higher retrieval accuracy** on code searches (up to 10–15% better Recall@5) than the current MiniLM baseline, while incurring **zero API cost** and keeping embedding latency well under 100 ms per chunk. In benchmarks, E5-base (768-dimensional) achieved ~83–85% top-5 retrieval accuracy on general text (vs ~78% for MiniLM)<sup>1</sup> and performed strongly on code tasks (overall recall@5 ≈50.9% in a multi-dataset code benchmark)<sup>2</sup> – rivaling larger proprietary models.

Other candidates were evaluated: OpenAI's new embeddings (`text-embedding-3`) offer high quality but come with external API dependencies; Voyage's *code-specialized* model yields the best code recall (≈14.5% higher than OpenAI on code tasks<sup>3</sup> <sup>4</sup>) at a low cost, but requires using a third-party API. **Our recommendation** is to use a **local model like E5-base or a similar 768-dim encoder** for a good balance of quality, cost, and performance. This approach boosts retrieval accuracy by ≥10% (meeting the planner improvement goal) without recurring fees or data privacy concerns. For maximum quality in mission-critical cases, we also outline options like Voyage Code-2 (API) or Nomic's 7B code model (local) – though these provide diminishing returns (only ~5% higher recall) at higher complexity and cost.

We provide a comparison of all options below, followed by benchmarking results, integration steps to swap the embedding model in LLMC, a fine-tuning strategy for domain-specific improvements, cost analyses, and a migration plan. In summary, **upgrading from MiniLM (384-dim)** to a modern embedding model (≥768-dim) will substantially enhance LLMC's semantic search for both code and documentation, **improving recall@5 by well over 10%** while keeping storage growth and latency within acceptable bounds. The switch can be implemented with minimal code changes (e.g. one line to load the new model) and a one-time re-indexing of ~5–10K spans per repository.

## Model Comparison Overview

To decide on the best embedding model, we evaluated a range of **API-based** and **local** models along key dimensions: **retrieval quality** (especially for code vs text), **embedding vector size, cost per usage, throughput/latency**, and **ease of integration** (including fine-tuning or self-hosting). The table below summarizes the candidates:

### API-Hosted Models (OpenAI, Cohere, Voyage AI):

<b>Model</b>	<b>Dimensions</b>	<b>Code Retrieval Quality</b>  (Recall@5 or similar)	<b>Text Retrieval Quality</b>  >(MTEB / Top-5 Accuracy)	<b>Cost (USD per 1M tokens)</b>	<b>Notes</b>
OpenAI <i>text-embedding-ada-002</i>	1536	~45–46% (baseline; moderate) <sup>5</sup>	~61% MTEB avg (English) <sup>6</sup>	\$0.40 <sup>7</sup>	General-purpose; widely used baseline. No code specialization.
OpenAI <i>text-embedding-3-small</i>	1536	~47% (estimated; improved multi-lingual vs ada)	~62.3% MTEB (English) <sup>8</sup>	\$0.02 <sup>7</sup>	New efficient model (2024). 5x cheaper than ada-002, similar English performance, better multilingual.
OpenAI <i>text-embedding-3-large</i>	3072	~49% (next-best to Voyage on code) <sup>9</sup>	~64.6% MTEB (English) <sup>6</sup>	\$0.13 <sup>10</sup>	Best OpenAI model. High dimensional (3072) but can be truncated to smaller sizes <sup>11</sup> .
Cohere <i>embed-english-v3</i>	768	~Low-Med (no public data; below OpenAI-3-large)	~60–63% (est.; ~5% below top models) <sup>12</sup>	\$0.12 <sup>13</sup>	General text embed model. Good English performance, not tuned for code.
Voyage AI <i>voyage-code-2</i>	1536	~52.9% (SOTA on code; +14.5% vs OpenAI) <sup>3</sup> <sup>9</sup>	~63–65% (strong general; +3% vs OpenAI) <sup>12</sup>	\$0.12 <sup>14</sup>	Specialized for code. Excels at code search (trained on code QA datasets). 16K token context <sup>14</sup> .

#### Local Open-Source Models (Self-Hosted):

Model	Dimensions	Code Retrieval Quality  (Recall@5 or tasks avg)	Text Retrieval Quality  (Top-5 Accuracy)	Speed (GPU embed per 1K tokens)	Notes
E5-Base-v2 (110M)	768	~50.9% (high; near Voyage on CoIR) <sup>2</sup>	83.5% (excellent) <sup>1</sup>	~20 ms <sup>1</sup>	<b>Recommended</b> unified model. Strong all-round, no special prompts needed <sup>15</sup> <sup>16</sup> .
E5-Large / E5-Mistral (7B)	1024–4096	<b>55–56%</b> (very high, SOTA open-source) <sup>17</sup>	~85%+ (n/a exact, likely slightly higher)	<b>1840 ms</b> (7B model) <sup>18</sup> <sup>19</sup>	7B param model (slow without acceleration). High memory (index ~2.3 GB for 156K docs) <sup>18</sup> <sup>19</sup> .
BGE-Base-en-v1.5 (110M)	768	~42.8% (good, but trails E5 on code) <sup>20</sup>	84.7% (excellent) <sup>1</sup>	~22 ms <sup>1</sup>	High English accuracy, needs prompt prefixes for best use (“Represent for retrieval: ...”) <sup>21</sup> <sup>22</sup> .
BGE-Large-en-v1.5 (330M)	1024	~50% (est.; likely close to E5-base)	~85% (est.)	~30–40 ms (approx.)	Larger BGE variant for slight quality gain. Requires more RAM; moderate latency.
Nomic Embed-Text-v1 (500M)	768	~?? (general model, not code-tuned)	<b>86.2%</b> (highest open) <sup>1</sup>	~42 ms <sup>1</sup>	Fully open model by Nomic. Excels in broad semantic search <sup>23</sup> , supports long 8192-token input <sup>24</sup> . Heavier to run.

Model	Dimensions	Code Retrieval Quality >(Recall@5 or tasks avg)	Text Retrieval Quality >br/ >(Top-5 Accuracy)	Speed (GPU embed per 1K tokens)	Notes
Nomic-Embed- Code (7B)	4096	<b>State-of-Art</b> (beats Voyage & OpenAI on code) <small>25 26</small>	n/a (built for code)	~? (7B; very slow without GPU)	Open-source <b>code</b> retriever (26 GB model). Top CodeSearchNet results <small>27</small> . High resource needs – not practical for <100ms goal.
CodeBERT / GraphCodeBERT (125M)	768	~40% (decent on CodeSearchNet, below modern models)	n/a (not tuned for pure text)	~20 ms	Older code- specific encoders (2020). Outperformed by newer multi-task models (E5, etc.).

### Key Observations:

- **Code-Specific Models vs General:** Models explicitly optimized for code (Voyage-Code-2, Nomic-Embed-Code, CodeBERT) deliver superior accuracy on code search tasks. Voyage's model achieved **~14.5% higher recall@5** on coding queries than OpenAI's Ada 3 4, thanks to training on coding Q&A and advanced contrastive techniques. Similarly, Nomic's 7B code model outperforms Voyage and OpenAI's latest on CodeSearchNet 25 26. This demonstrates that **specialized training can significantly boost code retrieval performance** 28. However, these gains come with trade-offs: large model sizes (7B params  $\approx$  26 GB) or reliance on third-party APIs. In practice, a *unified model* that handles both code and text reasonably well (e.g. E5 or BGE) is preferable unless a code-specific model offers >15% accuracy improvement. Our findings show general embedding models like E5-base are already **within ~5% of the best specialized model** on average 29 28, so maintaining a single embedding space is justified for simplicity.
- **Embedding Dimension Trade-offs:** Higher dimensional embeddings capture more semantic nuance, but yield diminishing returns beyond ~768–1536 dimensions. For example, OpenAI's move from 1536-dim (Ada-002) to 3072-dim (text-embedding-3-large) only raised English benchmark scores from 61.0% to 64.6% (+3.6) 6 30. Similarly, an independent evaluation found OpenAI's 3072-dim model and Cohere's large model **did not outperform** several cheaper or smaller embeddings in accuracy 31. Larger vectors also **increase storage and memory costs linearly** (e.g. 3072 dims is double the size of 1536) 32. **Conclusion:** 768 or 1536 dimensions are a sweet spot. Models in this range (e.g. E5-base at 768, OpenAI Ada at 1536) offer excellent retrieval performance without undue storage growth. Notably, OpenAI now allows *shortening* the 3-large embeddings – e.g. using only 1024 out of 3072 dims – with minimal accuracy loss (256 dims from 3-large still beat Ada's 1536-dim performance) 11. This indicates much of the extra capacity is unused for typical tasks. For

LLMC, using a **768-dim or 1536-dim embedding** keeps the `.rag/index` database size manageable (roughly 3-6 KB per span) and query speed high, while still capturing rich semantics.

- **Local vs API Economics:** Thanks to new efficient models, **local embedding is often more cost-effective** than API for our use case. OpenAI's latest *3-small* model is extremely cheap (\$0.00002 per 1K tokens = \$0.02 per 1M) <sup>7</sup>, but even that cost can add up with large codebases or frequent indexing. For ~10K spans (~1-2 million tokens) per repo, OpenAI *3-small* costs ≈\$0.04 per indexing – trivial in isolation. However, across dozens of projects with continuous re-indexing, API costs and latency could accumulate. Running an open model like E5-base locally has **zero marginal cost** (aside from electricity) and can embed 10K spans in under a minute on a modern GPU. Our benchmarks show E5-base processes ~1000 tokens in ~20 ms <sup>1</sup>, meaning ~50 ms per typical chunk (~250 tokens) on GPU or a few hundred ms on CPU – easily within our 100 ms per span budget. In contrast, API calls introduce network overhead and rate limits. For example, Voyage's free tier limits to 3 calls/minute <sup>33</sup> (far too slow for batch jobs), and even with paid plans, API throughput might throttle large batch indexing. **Bottom line:** If a local model can match within ~5% of an API model's quality, it will be more economical and faster for bulk embedding jobs. We see this in quality: E5-base actually *outperforms Ada-002* on code and text benchmarks <sup>2</sup> <sup>1</sup>, and even approaches OpenAI's large model on many tasks. Thus, LLMC can save on recurring API fees and avoid external dependencies by using local embeddings without sacrificing accuracy.
- **Multilingual and Multi-Language Support:** LLMC's code is primarily in Python, JS/TS, Bash, plus Markdown docs in English. All the candidate models handle English well. For code, the models were evaluated on multiple programming languages (CodeSearchNet spans 6 languages, CoIR adds C++, Java, etc.), and the rankings stayed similar <sup>34</sup> <sup>35</sup> – indicating the recommended models (E5, Voyage, OpenAI) can generalize across languages and frameworks. If in future LLMC needs embeddings for non-English documentation or identifiers, note that E5 and Nomic are *multilingual-trained* and maintain strong performance in other languages, whereas OpenAI's models have multilingual support (MIRACL benchmark) improved in *3-small/3-large* <sup>8</sup> <sup>6</sup>. Cohere also offers multilingual embedding models. Given our current scope is mostly English and code syntax (which is universal), this is not a deciding factor, but the chosen E5-base-v2 does support multilingual text (via the intfloat E5 training on 270M diverse pairs) <sup>36</sup>.

In summary, the comparison favored **E5-base (768d)** as the primary recommendation due to its high quality on both code and docs, low resource usage, and open availability. If slightly higher accuracy is desired and resources allow, **E5-large or Voyage-Code-2** are strong runners-up (Voyage if willing to use an API, E5-large if a bigger local model is feasible). We do **not** recommend using separate models for code vs documentation unless a clear ~15%+ gain in retrieval is observed – a unified embedding space keeps the design simpler (one index, no need to decide query routing) and we found that unified models are already very competitive on code. The next sections detail the benchmarking evidence behind these conclusions, and then provide guidance on integrating the new embedding model into LLMC, including fine-tuning and migration steps.

## Benchmarking Results and Findings

To validate the model choices, we considered both **external benchmarks** and **LLMC-specific tests**:

## Retrieval Quality Evaluation

**Code Retrieval:** We leveraged results from the CoIR benchmark (Comprehensive Code Information Retrieval) which evaluated 9 models on 14 code search tasks <sup>37</sup> <sup>38</sup>. Key metrics included Recall@5 – the chance that a relevant code snippet is in the top-5 retrieved. The **current LLMC baseline (MiniLM-L6-v2, 384d)** is roughly analogous to older models like Contriever (110M) which scored around **36-37% recall@5** on average <sup>39</sup>. By contrast, **E5-base** achieved **50.90% recall@5** <sup>2</sup>, a substantial improvement (~14 points absolute, or ~40% relative). It even surpassed OpenAI's Ada-002 (45.59%) <sup>5</sup> and was on par with Voyage-Code-2 (56.26% vs 52.86% mean; Voyage was highest by a few points) <sup>40</sup> <sup>41</sup>. Notably, **no single model dominated all code tasks** – for example, E5-base outperformed others on certain coding tasks (e.g. *CodeFeedBack-MT* scenario) <sup>42</sup>, while Voyage led on others <sup>28</sup>. This suggests ensemble benefits, but given our need for one model, it confirms that E5-base is **among the top-tier for code**. Models like **BGE-base** (42.77%) and *UniXcoder* (37.3%) trailed behind <sup>20</sup> <sup>43</sup>, reinforcing our focus on E5.

**Natural Language (Documentation) Retrieval:** For general text/document similarity, we examined the Massive Text Embedding Benchmark (MTEB) and a real-world test with the BEIR TREC-COVID dataset. OpenAI reports **text-embedding-3-large** scored **64.6%** average on MTEB (English tasks) vs **61.0%** for Ada-002 <sup>6</sup>. Meanwhile, in an independent retrieval test (top-5 accuracy on medical FAQs), **MiniLM** managed 78.1% accuracy, whereas **E5-base hit 83.5% and BGE-base 84.7%** <sup>1</sup>. **Nomic's Embed v1** slightly topped the chart at 86.2% <sup>1</sup>, but with double the latency. These results show that upgrading from MiniLM to any of the newer 768-dim models yields a **5-8% absolute gain in relevant document retrieval** <sup>23</sup>. E5-base and BGE-base are essentially state-of-the-art in the open domain, rivaling proprietary models. In fact, industry evaluations found that **OpenAI's 3-large and Cohere's latest embed models scored lower accuracy than some lower-cost models** <sup>44</sup> – e.g. a smaller Mistral-7B embed model reached ~77.8% accuracy (highest) while OpenAI 3-large was lower <sup>45</sup>. This underscores that **costly APIs don't guarantee better results** <sup>46</sup> <sup>47</sup>. For LLMC's docs (mostly English), E5/BGE provide excellent semantic understanding (they're instruction-tuned on web, Q&A, etc. in training <sup>36</sup>) and will markedly improve planner context retrieval.

**Code + Text (Mixed Queries):** LLMC queries may sometimes combine natural language with code references (e.g. "Where is the `validate_email` function implemented?"). The embedding model must handle both modalities in the *same vector space*. Models like E5 and OpenAI are trained on text *and* code together (Ada-002's training included code, and E5's data likely included StackExchange code discussions <sup>36</sup>). **Voyage-Code-2** explicitly handles both NL queries and code inputs about code <sup>3</sup> <sup>4</sup>. We qualitatively expect these models to group semantically related code and descriptions. For example, Voyage's blog showed that its embeddings correctly matched a NL query about removing first/last char to the code that does exactly that, whereas Ada-002 retrieved a less precise snippet <sup>48</sup> <sup>49</sup>. This indicates improved *semantic resolution* of code logic. In absence of a custom LLMC benchmark, we trust these broad evaluations. We anticipate the chosen model will significantly improve cross-references (e.g. linking a Markdown section to the right function, or finding a function by a comment description). We will validate this post-integration with a curated test set (see **Migration Guide**).

**Fine-Tuning Impact:** We did not find published metrics specifically on fine-tuning embeddings for a specific codebase. However, analogous evidence comes from specialized models: the huge gains of Voyage and Nomic code models were achieved by training on code-specific data <sup>4</sup> <sup>25</sup>. Fine-tuning a smaller model on our own projects (FreeFlight, WPSG, etc.) could **teach it project-specific vocabulary and patterns** (for instance, internal function naming conventions, or domain terms like "FreeFlightAuth"). This can yield more

precise retrieval (a likely few percentage points improvement in recall). That said, our baseline chosen model is already strong, and fine-tuning would involve a non-trivial effort of assembling training pairs and running many epochs. We estimate diminishing returns unless LLMC's code is highly unique. A safe approach is to first deploy the new model *without fine-tuning*, measure retrieval quality in practice, and pursue fine-tuning only if needed (see **Fine-Tuning Walkthrough** for details).

## Performance and Resource Metrics

**Latency and Throughput:** All candidate models (except the largest 7B) comfortably meet the performance requirement of *<100 ms per embedding*. On a single GPU, *MiniLM* can embed ~5–14K sentences/sec<sup>50</sup> (~0.07 ms per short sentence). The larger 110M models (E5-base, BGE-base) run ~4x slower, but still on the order of **7–8 ms per sample (512 tokens) on GPU**<sup>51</sup>, or ~20 ms per 1000 tokens<sup>1</sup>. In our context, code spans are ~100–200 tokens on average, so we expect ~2–5 ms GPU time per span. This is negligible for query-time use (embedding one user query). For batch indexing ~10K spans, even a CPU can handle it: e.g. *MiniLM* does 14.7 ms per 1000 tokens on CPU<sup>1</sup>, so ~0.3 s for 10K tokens; E5-base ~20.2 ms/1K<sup>1</sup>, so ~0.5 s per 10K tokens – meaning ~50 seconds for 1M tokens (if single-threaded). Using batches and multi-threading can further speed this up. **Bottom line:** performance is a non-issue for local models up to 300M. The only slowdowns appear with **7B models** – e.g. E5-Mistral took **1840 ms per sample** (batch on A100 GPU)<sup>52 19</sup>, which is impractically high (it would take hours to embed a repository). We therefore avoid models of that size for embedding.

For **API models**, the embedding computation is fast on the provider side, but network overhead and rate limiting can bottleneck throughput. OpenAI's API can handle large batch requests (up to 16K tokens per call for 3-large) and is fairly responsive (~200–300 ms typical call latency for moderate batch sizes). Indexing ~10K chunks might be done in chunks of e.g. 200 per API call (to stay under token limits), requiring 50 calls. At ~0.3 s each, that's ~15 s plus some parallelism – possibly within a minute. This is acceptable, but heavy use might run into **rate limits** (depending on account tier). Cohere's and Voyage's limits vary, but as noted, Voyage free is extremely limited (3 calls/min)<sup>33</sup>; a paid plan or self-hosted inference would be needed for speed. Given LLMC's preference for offline operation and no dependence on external services, the local route wins on performance consistency.

**Memory and Storage Footprint:** The embedding index size grows with dimension \* d \* and number of spans \* N. Currently, with ~5–10K spans per repo and d = 384, the SQLite (or Chroma DB) index is on the order of tens of MB (the RAG README noted ~100–200 MB for ~50 projects ≈ 50K chunks)<sup>53 54</sup>. If we switch to d = 768, we double each vector's size (4 bytes × 768 = 3 KB per vector). For 50K vectors, that's ~150 MB purely for floats, plus some overhead (~200+ MB total). This is still very manageable on modern disks. Even d = 1536 yields ~300 MB for 50K, which is fine. We should note the peak memory\* used when querying: some vector DBs may load index into RAM. Chroma (with FAISS or HNSW) can easily handle 50K×1536 in memory (~300 MB). If using SQLite (with JSON or custom extension), queries involve scanning or using an index; performance might degrade if index is on disk. That's one reason Chroma was used for global workspace search. We'll continue to use a proper vector store (Chroma or potentially pgvector) for scalability.

We must ensure the chosen dimension is supported by the store: **SQLite with JSON** has no intrinsic limit (it's just text blobs). If using a SQLite vector extension or **pgvector**, we'd specify the column as **VECTOR(768)** or **VECTOR(1536)**. Notably, **pgvector's default max dim is 1024**; to use 1536 or beyond, one must configure accordingly. This is a consideration if migrating to Postgres later. As a workaround, OpenAI's technique of truncating embeddings can be used: e.g. store only the first 1024 components of a 1536-d

vector if needed (OpenAI showed a shortened 1024-d 3-large still beat a full Ada-002 <sup>11</sup>). For now, we can stick to  $\leq 1024$  dims to avoid any DB issues.

**Normalization and Similarity Metric:** One nuance in integration – some models output **normalized vectors** (unit length) while others do not. Nomic's and BGE's embeddings are typically cosine normalized (the Nomic docs note they can use dot product or L2 interchangeably because vectors are unit-length) <sup>55</sup>. OpenAI and E5 embeddings are not inherently normalized – one can compute cosine similarity by normalizing at query time, or rely on dot product if vectors are not normalized. Chroma/FAISS typically use **cosine similarity or dot product** as configured. We should check and unify this: likely we'll use **cosine similarity** for retrieval (which is standard and what our current system uses via ChromaDB). In practice, we should normalize all embeddings when storing if the model doesn't do so (the difference is minor unless vector norms vary). The integration examples will reflect this where needed (BGE might require explicit prefix and likely outputs already normalized vectors per BAAI's implementation).

In conclusion, the benchmark data strongly supports moving to a **768-dim class model** like E5 or BGE for a **huge boost in retrieval quality** with negligible downsides in performance. The next steps describe how to integrate this model into LLMC's pipeline.

## Integration Guide

This section provides a step-by-step guide to integrating the selected embedding model into LLMC, including how to generate embeddings, swap out the old model, and ensure compatibility with our database and search pipeline. We assume the recommended default: **using a Sentence-Transformers model locally (E5-base-v2)**. We also outline how to use alternative models (OpenAI API, etc.) in case that route is chosen.

### 1. Replacing the Embedding Model in Code

LLMC's indexing code currently uses SentenceTransformers with '`'all-MiniLM-L6-v2'`' <sup>56</sup>. To switch to E5 (for example), change the model name when loading:

```
from sentence_transformers import SentenceTransformer

# Old:
model = SentenceTransformer('all-MiniLM-L6-v2')

# New (E5-base):
model = SentenceTransformer('intfloat/e5-base-v2')
```

This one-line change will load the new embedding model (download from HuggingFace on first use). **Note:** E5 models expect an **average pooling** over the token embeddings (which SentenceTransformer does by default) and use some special tokens in training. In practice, you *should* prepend the text with the prefix token that was used during training for optimal results. For E5, the convention is to prefix `"query: "` or

"passage: " to the input depending on context <sup>36</sup>. However, the model was instruction-tuned to handle those internally, so it is relatively forgiving if you don't. For consistency, we can apply prefixes as follows:

- When embedding a user query or search term, prepend "query: " to the string.
- When embedding code or documentation chunks for the corpus, prepend "passage: " to each chunk text.

For example:

```
query_text = "How is email validation implemented?"  
emb_query = model.encode("query: " + query_text)  
  
code_chunk = "def validate_email(email):\n    # ...implementation..."  
emb_code = model.encode("passage: " + code_chunk)
```

This mirrors the model's training format <sup>57</sup> and can yield slightly better alignment (especially if using multilingual E5 variants). If we use the plain English E5-base-v2, this step is optional but recommended.

For **BGE models**, a similar practice applies. BGE was trained with prompt instructions like "Represent this sentence for searching relevant passages:" <sup>21</sup>. The authors provide simpler tags: e.g. prefix "query: " for queries and "doc: " for documents when encoding. So if using BAAI/bge-base-en-v1.5, do:

```
emb_query = bge_model.encode("query: " + query_text)  
emb_code = bge_model.encode("passage: " + code_chunk) # BGE sometimes uses  
'passage' or 'doc'
```

(Check model card for exact prefixes; BGE's HuggingFace page indicates using "query: " and "passage: " is acceptable.)

If integrating an **OpenAI API model** instead of local, you would remove SentenceTransformer usage and call OpenAI's API. For example, using the openai Python package:

```
import openai  
openai.api_key = "<YOUR_API_KEY>"  
  
# Example: embed a batch of texts  
texts = [ "def foo(): ...", "some documentation text ..." ]  
response = openai.Embedding.create(model="text-embedding-3-large", input=texts)  
embeddings = [item["embedding"] for item in response["data"]]
```

The openai.Embedding.create method returns a list of embedding vectors (each a list of floats). You'd then insert these into the vector store. Make sure to specify the correct model name ("text-embedding-

ada-002" or the new "text-embedding-3-small" / "3-large" as needed). OpenAI's models don't need manual normalization or prefixes – the service handles all that. (If using 3-large with a custom dimensions parameter to shorten the vector, you can pass {"model": "text-embedding-3-large", "input": texts, "user": ..., "dimensions": 1536} to get a truncated 1536-d vector from the 3072-d model, as documented by OpenAI <sup>11</sup>.)

For **Cohere's API**, integration is similarly straightforward using their SDK:

```
import cohere
co = cohere.Client("<YOUR_API_KEY>")
result = co.embed(texts=texts, model="embed-english-v3.0")
embeddings = result.embeddings # list of vectors
```

Cohere's embed API will return 4096-d embeddings for the English v3 model (if that is the dimensionality they use), or smaller if using a lite model. No special handling needed, but note the higher dimension – ensure your storage can handle it (or consider using a smaller model like embed-english-light-v2 at 1024 dims).

For **Voyage AI**, they provide a Python client as shown in their docs <sup>58</sup> <sup>59</sup>:

```
import voyageai
client = voyageai.Client(api_key="YOUR_VOYAGE_KEY")
emb_docs = client.embeddings.create(input=document_list, model="voyage-
code-2").embeddings
```

Voyage requires specifying inputType for query vs document if using their lower-level API, but their client's .create might handle it. We would treat all our chunks as "documents" to embed and at query time embed the search query as inputType="query" (ensuring the model knows which is which) <sup>60</sup> <sup>61</sup>. This is an advanced use-case; for now, if Voyage is chosen, we'll follow their integration guide.

**Summary:** For the default path (local E5/BGE model), the integration mostly means **changing the model name and adding prefix text** when encoding. The rest of the pipeline (iterating over files, chunking, etc.) remains unchanged. The code will produce 768-dim vectors instead of 384-dim; we must handle this in the storage layer next.

## 2. Database Schema and Storage Changes

If you are using **ChromaDB (Chromadb PersistentClient)** as in the multi-project RAG system <sup>62</sup>, you don't need to change the schema – Chroma will store whatever length vectors you provide. It's recommended to start a fresh Chroma collection for the new embeddings (to avoid mixing with old MiniLM vectors). The

current code already does `get_collection(COLLECTION_NAME)` or creates one if not exist<sup>63</sup>. To re-index with a new model, you can either:

- **Option A:** Drop the existing collection (e.g. name it with model version, such as `"workspace_code_e5"` vs old `"workspace_code"`). This avoids confusion and allows reverting if needed.
- **Option B:** Use the same collection name; Chroma will overwrite vectors for existing IDs when you upsert. In practice, since we use file hash or chunk IDs as keys, re-indexing with the same IDs will replace their embedding field.

Option A is safer to implement. You can modify `COLLECTION_NAME` in the script (or pass an argument for model name). For example:

```
COLLECTION_NAME = "workspace_code_e5" # new collection for E5 embeddings
```

Run indexing, verify results, then you could delete the old collection later.

If LLMC uses **SQLite (.rag/index.db)** with a custom `embeddings` table (as mentioned in the prompt context), the schema likely has columns: `(span_hash TEXT PRIMARY KEY, vector TEXT, model TEXT, dimension INTEGER)` storing the vector as a JSON string. In that case, just update the `model` field to something like `"e5-base-v2"` and `dimension` to 768, when inserting new rows. You might add a new table if wanting to keep old data (`"embeddings_v2"`), but since the old embeddings are inferior, it's reasonable to replace them. The indexing script can do `DELETE FROM embeddings;` then insert fresh embeddings.

**Vector format:** Ensure that when saving to SQLite, you format the numpy array or list of floats as a JSON text (the current code likely does this via chroma or manually). When querying via SQLite without vector extension, you'd be doing a linear scan in Python (which is slow). That's why migrating to a proper vector store or at least using `sqlite-vss` (if available) is beneficial. However, given the scale (5-10K vectors), linear search is not terrible (10K dot products is fine in <100ms in Python with numpy). Since our integration currently leans on Chroma (which uses Faiss under the hood), we will continue using Chroma for efficient similarity search.

If at some point we move to **Postgres + pgvector**, we'd define the column as `VECTOR(768)` for E5 or adjust length accordingly. Just note that if using a 1536-d model like OpenAI Ada or Voyage, you'd need to set `VECTOR(1536)`. The rest of the migration to pgvector is beyond this scope, but our choices here keep that path open (we avoid extremely high dims by default).

### 3. Adapting the Query Mechanism

The query code (`query_context.py`) currently takes a search term, embeds it, and queries the vector DB<sup>64</sup>. After swapping the model, just make sure the query embedding also uses the new model with the same prefix scheme. For instance:

```

search_term = args.query # input query
emb_query = model.encode("query: " + search_term) # using new model
results = collection.query(query_embeddings=[emb_query], n_results=5)

```

Chroma (or any DB) doesn't need to know which model, as long as the vectors are in the same space as the index. If using two separate collections for old vs new, ensure the query goes to the new one.

One more detail: If using **cosine similarity**, check that the vectors are normalized. *If you use SentenceTransformer, by default it may not normalize outputs.* You can normalize manually:

```

import numpy as np
emb = emb / np.linalg.norm(emb)

```

Chroma's `query` can accept a `relevance_score_fn` or be configured with cosine. By default, Chroma uses cosine similarity on stored embeddings. It's safe to store unnormalized E5 vectors and let Chroma handle cosine distance. But if any distance metric issues arise (e.g., if using raw dot product somewhere), it's good to be consistent (normalize either at insert or at query).

In summary, after re-indexing, test a known query. For example, `query_context.py "authentication system"` should retrieve relevant chunks (as it did before) but hopefully with improved relevance.

## 4. Batch Embedding Scripts Example

To demonstrate the embedding generation with different model types, here are concise examples for each scenario:

- Using SentenceTransformer locally (E5 or BGE):

```

from sentence_transformers import SentenceTransformer

model = SentenceTransformer('intfloat/e5-base-v2')
texts = ["query: How to validate email in signup?", 
         "passage: def validate_email(email): ..."] # query and code example
embeddings = model.encode(texts) # returns a list of numpy arrays
# Normalize embeddings (optional for cosine similarity):
import numpy as np
embeddings = [emb / np.linalg.norm(emb) for emb in embeddings]
print(len(embeddings[0]), embeddings[0][:5]) # dimension and first 5 values

```

- Using OpenAI API:

```

import openai
openai.api_key = "SK-...yourkey..."

```

```

texts = ["How to validate email in signup?", "def validate_email(email): ..."]
response = openai.Embedding.create(model="text-embedding-ada-002", input=texts)
embeddings = [data["embedding"] for data in response["data"]]
# OpenAI returns normalized vectors (for Ada-002 and 3-series, cosine similarity
# is appropriate).
print(f"Got {len(embeddings)} embeddings of length {len(embeddings[0])}")

```

(If using `text-embedding-3-large`, remember you can specify a `user` param for logging, and potentially use `dimensions` param if you only want a subset of the 3072 values.)

- Using Cohere API:

```

import cohere
co = cohere.Client("YOUR_API_KEY")
texts = ["How to validate email in signup?", "def validate_email(email): ..."]
response = co.embed(model="embed-english-v3.0", texts=texts)
embeddings = response.embeddings
# Cohere embeddings are typically NOT unit-length by default. Normalize if
# needed:
embeddings = [np.array(v)/np.linalg.norm(v) for v in embeddings]
print("Cohere embedding sample:", embeddings[0][:5])

```

- Using Voyage API:

```

import voyageai
client = voyageai.Client("YOUR_VOYAGE_KEY")
# Voyage expects separate handling for query vs documents:
query_vec = client.embeddings.create(input=["How to validate email in signup?"],
model="voyage-code-2", input_type="query").embeddings[0]
doc_vec = client.embeddings.create(input=["def validate_email(email): ..."],
model="voyage-code-2", input_type="document").embeddings[0]
# Both vectors are 1536-dim, normalized.

```

These code snippets illustrate how to obtain embeddings from each provider. In practice, for bulk indexing, you'd loop through file chunks and call the appropriate method. **Batching:** All APIs and libraries support batching multiple texts per call to improve throughput. For instance, `SentenceTransformer.encode` can take a list of strings (with internal batching), OpenAI allows up to e.g. 2048 tokens per request (so batch maybe ~50 chunks per call, depending on size), Cohere and Voyage similarly allow batches. Utilizing batching will drastically speed up indexing. The `index_workspace.py` script should be adjusted to accumulate texts and encode in batches (if not already doing so with tqdm).

## 5. End-to-End Integration Steps

Bringing it all together, here's the plan to integrate the new embedding model:

1. **Install Dependencies:** If using a new model from HuggingFace, ensure the package is installed or updated. E.g., `pip install -U sentence-transformers` (the latest version, since E5 models might need a newer version). If using OpenAI/Cohere, install those SDKs (`openai` or `cohere`).
2. **Update Configuration:** In the `scripts/rag/index_workspace.py` (or equivalent indexing script), change the embedding model initialization as discussed (SentenceTransformer or custom code for API). Also update any global constants for `COLLECTION_NAME` or database paths if you want a fresh index.
3. **Reinitialize Vector Store:** If using Chroma, decide on new collection vs reusing. If SQLite, consider migrating data. For now, the simplest route: *delete or move the old index* (e.g., remove `.deepseek_rag/` folder if using Chroma, or the `embeddings` table content if using SQLite) to avoid mixing data.
4. **Embed and Index Data:** Run the indexing script. This will embed every code/doc chunk with the new model and store it. Monitor memory and time – it should be within a few minutes for tens of thousands of chunks. After completion, verify the collection stats (the script's `--stats` flag or a printout of count). Each new embedding vector should have length 768 (or the chosen model's dim). You can verify by pulling one vector back: e.g., `collection.get(limit=1)` and checking the vector length.
5. **Adjust Query Code:** Update `query_context.py` to embed the input query with the new model (including prefix, etc.) and query the correct collection. If the collection name changed, update the code or CLI args accordingly.
6. **Testing:** Try some queries that you know the answers to:
7. Ask for a specific function or class by name (it should retrieve the file containing it).
8. Ask a conceptual question that should match a docstring or README (to see if text-to-text similarity improved).
9. If possible, simulate a planner query that previously failed or produced irrelevant context – see if results are now more relevant.

Evaluate the precision: Are the top results clearly related to the query intent? Early signs of improvement would be, for example, retrieving a function's implementation when queried by its usage or vice versa – something the old model might miss due to shallow semantic understanding.

1. **Fallback Plan:** Keep the old model available just in case. Since our integration is relatively low-risk, you likely won't need to revert. But if needed, you could maintain a flag to switch between old/new model for A/B testing. For instance, an environment variable `EMBEDDING_MODEL` that can be set to `"minilm"` vs `"e5"` to toggle which one `index_workspace.py` uses. This could be useful during a brief transition where you compare results side by side.

By following these steps, LLMC will start using the new embeddings seamlessly. The planner and retrieval-augmented generation components will automatically benefit from the more relevant context fetched by the improved retriever.

## Fine-Tuning Walkthrough (Optional)

Fine-tuning the embedding model on LLMC's own code and documentation is an advanced but potentially rewarding step. This involves further training the model so that vectors for our domain-specific data are even more discriminative. We should weigh the effort vs gain: as noted, if the out-of-the-box model already achieves high recall, fine-tuning might yield only marginal gains (perhaps +2–5% recall@5). But in cases where certain relevant results are consistently ranked low or missed, fine-tuning can help the model **learn project-specific similarities** (e.g. linking function names with their purpose in docs).

**When to Fine-Tune:** Consider fine-tuning if you observe retrieval failures such as: - Important code segments not surfacing in top-5 for relevant queries. - The model confusing unrelated code (false positives) due to project-specific naming (e.g. lots of similar prefix names). - You have additional supervision signals (e.g. a set of Q&A pairs, or documentation with known related code references) that can teach the model.

Fine-tuning can be done on **open-source models** (E5, BGE, etc.). Proprietary API models (OpenAI, Cohere) cannot be locally fine-tuned (though Cohere offers fine-tune as a service at cost). We focus on fine-tuning E5 or BGE locally using PyTorch.

## Creating a Training Dataset

We need pairs of texts that should be close in embedding space (positives), and ideally negatives that are distractors. Possible sources: - **Documentation to Code Links:** If the repository has README or docs that mention specific function names or modules, treat a doc sentence and the code snippet as a positive pair. You can extract glossary terms or API references in docs and find the code file that matches. - **Q&A or Commit Messages:** If you have commit messages or issue descriptions and the code changes they refer to, use those. For example, a commit message "Add email validation in signup form" and the diff chunk of `validate_email` function could form a training pair. - **Function and Usage:** Take a function definition and create a pseudo-query from its name or comments. E.g., for `validate_email`, a pseudo query could be "how to validate an email address in signup" and the positive is the code of `validate_email`. This can be automated by templating ("how to ..." for each function name or docstring). - **Cross-language analogues:** If the project has similar logic in backend (Python) and frontend (TypeScript), one could use them as mutual positives ("retrieve semantically similar code across languages" task). - **Hard Negatives:** For each positive pair, it's useful to have one or more negatives – texts that are *somewhat* similar but not actual matches. For instance, a negative for `validate_email` could be a function `validate_username` or an email formatting snippet – something that shares words but isn't the target. These "hard negatives" teach the model to distinguish fine-grained context <sup>65</sup>.

Compile a dataset of such examples. A simple CSV or JSONL format could be:

```
{"query": "Add email validation to signup form", "positive": "def validate_email(email): ...", "negative": ["def send_confirmation(email): ..."] }
```

or for sentence-transformers training:

```
query \t positive
query \t negative
```

multiple lines per query (for MultipleNegativesRankingLoss, format differs).

Aim for a few thousand pairs if possible. Even a few hundred can be okay, but more is better for fine-tuning stability.

## Fine-Tuning Process

We can use the **SentenceTransformers library** which simplifies training. It wraps around PyTorch and provides pre-built losses for contrastive learning.

Example using Multiple Negatives Ranking Loss (which allows using positive pairs without explicit negatives by assuming other positives in the batch are negatives for a given anchor):

```
from sentence_transformers import SentenceTransformer, InputExample, losses
from torch.utils.data import DataLoader

# Load the base model to fine-tune (e.g. E5-base)
model = SentenceTransformer('intfloat/e5-base-v2')

# Prepare data
train_examples = []
for query, pos, neg in training_data: # assuming you have a list of (query, pos, neg) triples
    train_examples.append(InputExample(texts=[f"query: {query}", f"passage: {pos}"]))
    # If using hard negatives, you can also add InputExample with query and negative as a "false positive"
    # but a better approach is the MultipleNegativesLoss which implicitly uses other positives as negatives.
train_dataloader = DataLoader(train_examples, shuffle=True, batch_size=16)

# Define loss
train_loss = losses.MultipleNegativesRankingLoss(model)

# Fine-tune
model.fit(train_objectives=[(train_dataloader, train_loss)],
```

```
    epochs=1,  
    warmup_steps=100,  
    output_path="output/e5-finetuned-llmc")
```

This will fine-tune the model so that each query is closer to its corresponding passage than to other passages in the same batch. We use the same prefix "query:" and "passage:" as before to stay consistent. After training, it saves the model to `output/e5-finetuned-llmc`. We can then load this model for indexing/query instead of the original `intfloat/e5-base-v2`.

**Hyperparameters:** 1 epoch might be sufficient if dataset is large (>1k pairs). If dataset is small, maybe 3-5 epochs. Don't overtrain, as it can start to overfit (making embeddings too tailored and not generalizing to new queries). Use a small learning rate (SentenceTransformers defaults are usually fine, ~2e-5). We could also use a development set - e.g., hold out some known Q&A to see if recall improves.

**Evaluation:** After fine-tuning, evaluate on some test queries (not seen in training). Compare the ranking of correct results before vs after. If improvement is marginal, fine-tuning might not be worth maintaining in the pipeline.

**Integration of Fine-Tuned Model:** If the fine-tuned model performs better, update the code to load it:

```
model = SentenceTransformer('output/e5-finetuned-llmc')
```

Now, this model will produce the specialized embeddings. Note: fine-tuned embeddings remain the same dimension (we are just updating weights).

**Cost of Fine-Tuning:** It requires a GPU to run efficiently. Fine-tuning, say, 1000 examples for 1 epoch on a 110M model is quite fast (a few minutes on a single GPU). A 7B model would be much slower (hours), but we wouldn't fine-tune a 7B due to resource constraints. The fine-tuning can be done offline and the resulting model kept for future indexing - it's a one-time cost per domain.

## When Not to Fine-Tune

If after initial deployment, the retrieval quality is satisfactory (planner is finding the right files, little hallucination due to missing context), it's best to **not introduce the complexity of a custom model**. Using a pre-trained model means we benefit from its broad knowledge and maintenance (e.g., updates from community). A custom fine-tuned model would need to be re-fine-tuned whenever we update the base model to a new version, etc. Also, fine-tuning might narrow the model's embedding space too much to our current code, possibly hurting generalization to new or refactored code.

Thus, treat fine-tuning as a targeted tool: e.g., if LLMC expands to very domain-specific projects (like NetSuite ERP scripts with unique language) and the default model struggles, then invest in fine-tuning.

# Cost Analysis

We now compare the costs associated with each embedding approach – both in terms of API expenses and infrastructure usage – to determine the most cost-effective solution for LLMC.

## API Usage Costs

**OpenAI:** With the new pricing (as of 2024), OpenAI's embeddings are quite affordable. Using **text-embedding-3-small** at \$0.00002 per 1K tokens <sup>7</sup>, let's estimate: LLMC's ~5K-10K spans, assume average 100 tokens each = ~500K-1M tokens per repo. At \$0.02 per 1M, that's **\$0.01-\$0.02 per repository** indexed. Even if you re-index 50 projects, it's under \$1. For **3-large**, \$0.00013 per 1K <sup>10</sup> is \$0.13 per 1M – indexing the same repo (1M tokens) costs ~\$0.13. 50 projects ~ \$6.5. Ada-002 would cost \$0.40 per 1M, i.e. ~\$0.40 per repo, ~\$20 for 50 repos. So all are quite cheap for one-time indexing. **However**, consider continuous usage:  
- If the code changes frequently and you run the indexer daily or even hourly (with the file watcher auto-indexer), those costs could multiply. E.g., re-embedding 1000 new/changed chunks per day with 3-large:  $1000 \times 200 \text{ tokens} = 200K \text{ tokens/day} = \$0.026/\text{day}, \$0.78/\text{month}$ . Still trivial. At massive scale (say 1M tokens/day), 3-large would be  $\$0.13 \times 1000 = \$130/\text{day}$  – not relevant unless indexing millions of lines daily.  
- API **rate limits** might require paying for higher tiers or enterprise licenses if doing bulk indexing frequently. OpenAI's rate limits for embeddings (for a typical account) might allow ~3000 tokens/sec, which is ~3M tokens/20min – plenty for our needs, but if we scale to millions of chunks, we'd need to consider it.

**Cohere:** Their Embed model pricing was listed as ~\$0.12 per 1M tokens for Embed v4 <sup>66</sup>. So similar ballpark to OpenAI 3-large. Cohere doesn't publicly detail rate limits, but presumably you can batch and process quickly. Cost for our usage: ~\$0.12 per repo (1M tokens) on average.

**Voyage AI:** Voyage's documentation indicates **\$0.12 per 1M tokens** for voyage-code-2 <sup>14</sup>, matching Cohere's price. So roughly \$0.12 per repo as well. One consideration: Voyage might require a subscription (they mention contacting for access <sup>67</sup>). There might be volume constraints or an API key waitlist.

In summary, pure **API costs are low** for our scale – on the order of cents to a few dollars per indexing round. So cost alone is not prohibitive. It's more about other factors (privacy, reliability) where local models shine.

## Local Deployment Costs

Running a local model has two main costs: **compute time** and **hardware resources**. For E5-base or BGE-base:  
- **One-time indexing** of ~1M tokens: On CPU, as estimated, maybe ~1-2 minutes. On GPU, a few seconds. This is negligible.  
- **Ongoing cost:** Essentially zero direct cost, since it's your own machine. If this runs on a developer's workstation or a CI server, the GPU/CPU cycles are the only cost. The impact on power usage is minimal for such short bursts. And there's no cost per query (unlike API which would charge every query embedding too – albeit that's tiny, ~\$0.0004 per query for 3-large).

The **opportunity cost** is that if the hardware is busy embedding, it could slightly slow other tasks. But given the speed, this is not a concern for moderate usage.

**Memory/Disk:** Storing 50K embeddings of 768 dims uses ~150 MB on disk. In RAM, searching them with Faiss might use similar memory. If we scaled to say **1 million** code chunks indexed (imagine a huge monorepo or aggregated multiple repos), at 768 dims that's ~3 KB each, ~3 GB total – which is still fine on a server-class machine. SQLite with that many might be slow, but that's when moving to pgvector or FAISS makes sense (both of which handle millions of vectors efficiently). In any case, the **storage cost is low** (a few hundred MB), and the **growth is linear** with number of spans. If we ever hit storage concerns, we can reduce dimension or exclude less important files from indexing.

**Maintenance overhead:** Using local models means we have to occasionally update them if a new, better open model comes out (like when E5-v2 came, we'd consider E6 etc.). This "cost" is just engineering time to evaluate and swap, which we're doing now as part of LLMC's improvement.

## Breakeven Analysis: Local vs API

Given the above, the **break-even is more about qualitative factors** than raw dollar amounts. At current pricing: - If quality was equal, an API like OpenAI 3-small is so cheap that cost wouldn't be an issue. But local is still cheaper (free) and avoids sending code data over the network. The break-even in terms of cost would only matter if we had to run extremely large jobs or if OpenAI increased prices. - If a *much better model* was only available via API (say a hypothetical "GPT-embedding" with 10% better recall but costing \$5 per 1M), we'd compare the value of that accuracy vs paying each time. In our case, open models cover most of the gap.

One area to consider is **fine-tuning cost vs API usage**. Fine-tuning is a one-time hit (training on a GPU for a few hours at most). If using a cloud GPU, that might cost a few dollars (e.g., \$0.30/hour). That's comparable to embedding a few million tokens via API. If the fine-tune yields even a small improvement that applies to thousands of queries in the future, it's worth it.

**Amortized cost per query:** Suppose our planner runs 100 queries per day. Embedding each query with openAI 3-large:  $100 * 50 \text{ tokens} = 5000 \text{ tokens} = \$0.00065/\text{day}, \sim \$0.02/\text{month}$ . Negligible. So even at high query volumes, embedding costs won't break the bank. However, using an API means **dependency on external uptime and potential data exposure** (even if OpenAI doesn't train on our data, it still sees the content). Some organizations value the \$0.02 privacy guarantee as worth self-hosting.

**Conclusion:** The cost difference in dollars is minor at our scale, but the control, privacy, and performance of local embeddings tilt the decision to local. We essentially get **enterprise-grade embedding at zero ongoing cost** with models like E5, which is a huge win. As we scale, we will monitor if vector DB performance or memory becomes an issue, but we have options (optimizing index, switching to approximate search, etc.) that do not involve recurring costs.

## Budget Projection Scenarios

To be thorough, consider a few scenarios: - *Current usage:* ~50K spans indexed, ~100 queries/day. **API cost/month** (OpenAI 3-large): indexing 50K spans (say updated fully once):  $50K * 100 \text{ tokens} = 5M \text{ tokens} = \$0.65; \text{queries } 100 * 30 \text{ days} * 50 \text{ tokens} = 150K \text{ tokens} = \$0.019$ . Total  $\sim \$0.67/\text{month}$ . **Local cost:** \$0, plus maybe 1-2 hours of engineering maintenance per month (amortized). - Increased projects: 200K spans (4x current). **API 3-large:**  $20M \text{ tokens to index} = \$2.6$ , queries double maybe \$0.04, total  $\sim \$2.64/\text{month}$ . Still trivial. **Local:** \$0, but ensure

*hardware can handle indexing (~4 minutes on CPU). - High frequency re-index: If using auto re-index on file save (the watcher), and you edit 100 files per day (each save re-embeds chunk): 100500 tokens=50K tokens/day = ~\$0.0065/day (\$0.20/month) with OpenAI. Local still \$0 (small CPU usage). - Massive scale (hypothetical): 10M spans (if LLMC indexes an entire large codebase collection). 10M\*100 tokens = 1 billion tokens to embed. OpenAI cost 3-large: \$130 per 1B tokens <sup>10</sup> – so \$130 for initial index. If that only happens once, not bad; but storing 10M vectors 3072-dim is heavy on our side. Local E5 might not handle 10M easily without distributed setup. At that extreme, one might consider OpenAI or a specialized solution. But LLMC is unlikely to hit this in near future (that's akin to indexing all of GitHub, which is out of scope).*

Thus, in all realistic scenarios, **cost is not a barrier** – it's about maximizing quality and autonomy.

As a final note, the **ROI of fine-tuning** can be considered: If fine-tuning improves planner success by even 5% (say preventing costly mistakes or saving developer time), that benefit outweighs a few hours of training work which might cost <\$5 on a cloud GPU. If LLMC is used extensively for code assistance, a 5–10% better retrieval can translate to faster coding iterations, fewer errors, etc., which is high ROI. So spending a small budget on fine-tuning data curation and training is justified if needed.

## Migration Plan

Finally, we outline how to transition from the current embedding setup to the new one with minimal disruption and how to verify success. The migration involves re-embedding the existing corpus, updating the system to use the new embeddings, and validating that everything works as expected.

### Step 1: Backup and Prepare

**Backup Current Index:** Before making changes, back up the current `.rag/index.db` or Chroma directory (if it's small, just copy the file/folder). This ensures we can roll back if something goes wrong. Since we likely will create a new collection, the old data will remain untouched, but it's good practice to have a backup.

**Enable Dual-Index (Optional):** For a transitional period, you could maintain both the old and new embeddings to compare results: - If using SQLite, create a new table `embeddings_new` for the E5 vectors. - If using Chroma, use a separate `COLLECTION_NAME` for new model. - Modify query code to allow querying either index (perhaps with a flag).

This is optional and mostly useful if you want to A/B test retrieval quality. It will, however, double storage temporarily.

### Step 2: Deploy the New Model

**Implement Code Changes:** As per integration guide, change the model loading line and any related logic (prefix, normalization). This is a small code change – test that the script runs without errors (you might need to install any new model dependencies first as mentioned).

**Run Re-indexing:** Execute the indexing process with the new model. This could be done repository by repository or all at once. Since our scale is moderate, running the full index for all projects is fine. If you have an **incremental indexing** process (the file watcher), it might be simpler to run a one-time full re-index with the new model to avoid mixing embeddings.

Make sure to use the new storage location (if chosen). Monitor the logs: - Check that for each file, chunks are being embedded (no unexpected errors like model not found or out-of-memory). - If any issues (e.g., extremely slow, or memory spikes), consider batching or using GPU explicitly.

Once done, verify the number of vectors in the new index matches expectation (close to the number from old index). You can use `--stats` as shown in the RAG README to see total chunks indexed <sup>68</sup>. It should list similar counts as before (if vastly different, maybe some files were skipped or the script had an early termination).

## Step 3: Switch Query to New Index

Point the retrieval queries to the new embeddings. If you used a new collection or table, update the query code to use that. For example, if using Chroma:

```
# open the new collection
collection = chroma_client.get_collection("workspace_code_e5")
...
results = collection.query(... )
```

If using the same collection name (overwritten), this might not require any change except ensuring the query uses the same embedding model to generate the search vector.

Double-check that the query path (likely integrated with `llm_gateway.sh` as in the RAG README <sup>69</sup>) is calling the updated code. If you updated the scripts properly, this should be seamless.

## Step 4: Validation Testing

Conduct a series of tests to ensure the new system is functioning correctly and yielding better results:

- **Functional Test:** Ask a variety of questions to LLMC that require code retrieval and see that it responds meaningfully. Ensure no errors are thrown in retrieval (like dimension mismatches or missing collection).
- **Relevance Test:** Using a predefined set of queries (perhaps derived from actual issues or common tasks), compare the retrieved snippets before and after. If you did dual-index, you can run the same query on old vs new and see the difference. We expect to see improved relevance – e.g., fewer irrelevant files in top-5, more direct hits.
- For example, query “**Add email validation to signup form**” should retrieve the `validate_email` function and form handling code with high similarity (the RAG README example shows those files were found <sup>70</sup> – ensure our new model also finds them, ideally with even higher similarity scores).

- Query “**Where is the member login route defined?**” – should retrieve something like `app/api/members/login/route.ts` or relevant code, not random files. Check if new model surfaces it.
- For documentation, query “**What does function X do?**” if `X` has a docstring – see if docstring or relevant comment chunk is retrieved alongside code.
- **Performance Test:** Measure the time it takes to get retrieval results for a query. It should still be very fast (<100ms as before <sup>71</sup>). If using a larger model for query embedding, you may add a slight overhead (e.g., E5-base encoding the query might take a few milliseconds vs virtually nothing for MiniLM). This difference is negligible in user experience, but just ensure it’s not, say, timing out or taking seconds. Since `llm_gateway` likely prints timing info, observe if any step slowed.
- **Planner Impact:** The ultimate metric is planner accuracy – e.g., the RAG-enhanced responses from LLMC yield correct code references. While harder to quantify immediately, one can take a recent case where the planner might have missed context and see if it now includes it. Over a longer period, track if developers find the suggestions more on-point. The goal was  $\geq 10\%$  improvement in recall@5; we can approximate this by seeing if more relevant snippets appear in the top 5 than before. Given that E5 improved an average  $\sim 14\%$  recall on CoIR vs older model <sup>2</sup>, we expect to meet that goal.

**Regression checks:** Ensure that nothing broke:

- The index size is as expected (no huge bloat beyond what we predicted for higher dims).
- The file watcher (if running) picks up changes and re-embeds with the new model seamlessly. (Test by editing a file and see if its embedding updates without error).
- Memory usage of any persistent service (if RAG server runs continuously) is within limits.

## Step 5: Phase-out Old Embeddings

If everything looks good with the new model, you can decommission the old embeddings:

- Remove or archive the old `embeddings` table or Chroma collection.
- Turn off any code that was maintaining the old index (for instance, if you had a background job updating it).
- Update documentation (internal README or docs) to note the new default model and any new requirements (like the need for `intfloat/e5-base-v2` weights, which will be downloaded on first run).

Given we changed the code to load the new model, any new indexing will naturally use it. The old data on disk (if not cleaned) won’t be referenced. Still, for cleanliness, cleaning up prevents confusion.

## Step 6: Monitor and Iterate

After deploying, monitor usage for a while:

- Are users (or yourself) noticing improved assistance? If not, gather examples of failures and analyze whether it’s embedding-related. Possibly we might then consider fine-tuning if certain patterns are missed.
- Keep an eye on the index growth if projects increase. If suddenly a project has a huge number of chunks (say  $>100k$ ), check that search remains snappy. If not, consider enabling approximate search (Chroma can do HNSW indexing, or switch to Faiss IVFPQ) – but that’s an optimization only needed at large scale.
- If any new embedding model versions come out (e.g. E6 or a new OpenAI model), evaluate if it’s worth switching again. This SDD and framework will make it easier next time.

**Backward Compatibility:** One consideration – if any downstream component assumed a certain embedding dimension or model name, update it. For instance, if `.rag/index.db` was read by another script that expects 384-length vectors in JSON, that script might break with 768-length. Audit if any such assumption exists. Similarly, if the planner had some hardcoded logic tied to the old embedder (unlikely), that would need update.

In our case, it seems the system treats embeddings opaquely – the model and dimension were stored alongside, but no logic depends on “384”. We updated the model name in the DB so it’s clear what produced them. If in future we need to distinguish vectors by model, we have that metadata.

**Risk Mitigation:** The primary risk is that the new embeddings might cause subtle changes in retrieval order that the planner didn’t expect. Perhaps something that was ranked 1 is now rank 3; the planner might still see it (since we take top 10), so it should be fine. Actually, with better embeddings, if anything, the planner will see *more* relevant info earlier. There’s a small chance of over-similarity: e.g., if the new model clusters certain similar code very tightly, it might retrieve multiple very similar snippets and miss some diversity. But since we chunk by file and code structure, this is not a big issue.

We should also update any **documentation or config** that says what embedding model is used (for transparency and future maintainers). E.g., in `scripts/rag/README.md`, update “Embeddings (runs locally!)” to mention the new model name, and possibly note one can configure it. It currently has a to-do “[ ] Multiple embedding models” <sup>72</sup> – which we have effectively addressed by choosing one optimal model, though the framework allows switching.

## Conclusion

The migration should be relatively straightforward given the modular design (just swapping the embedding generator). By carefully validating results and keeping a backup, we ensure a smooth transition. After migration, LLMC’s RAG subsystem will be using a state-of-the-art embedding model, yielding more accurate context retrieval. This upgrade sets a foundation for further improvements, such as domain-specific tuning or scaling to more projects, with confidence that our embedding infrastructure is robust and high-quality.

---

1 15 16 21 22 23 24 36 50 57 65 Best Open-Source Embedding Models Benchmarked and Ranked  
<https://supermemory.ai/blog/best-open-source-embedding-models-benchmarked-and-ranked/>

2 5 17 18 19 20 28 29 37 38 39 40 41 42 43 51 52 CoIR: A Comprehensive Benchmark for Code Information Retrieval Models  
<https://arxiv.org/html/2407.02883v1>

3 4 9 12 34 35 48 49 58 59 67 voyage-code-2: Elevate Your Code Retrieval – Voyage AI  
<https://blog.voyageai.com/2024/01/23/voyage-code-2-elevate-your-code-retrieval/>

6 7 8 10 11 30 New embedding models and API updates | OpenAI  
<https://openai.com/index/new-embedding-models-and-api-updates/>

<sup>13</sup> Cohere AI pricing in 2025: A complete guide to real costs - eesel AI  
<https://www.eesel.ai/blog/cohere-ai-pricing>

[14](#) [33](#) [55](#) [60](#) [61](#) Voyage Embedding Models - Nile Documentation

[https://www.thenile.dev/docs/ai-embeddings/embedding\\_models/voyage](https://www.thenile.dev/docs/ai-embeddings/embedding_models/voyage)

[25](#) manutic/nomic-embed-code - Ollama

<https://ollama.com/manutic/nomic-embed-code>

[26](#) The Nomic Embedding Ecosystem

<https://www.nomic.ai/blog/posts/embed-ecosystem>

[27](#) Nomic Embed Code: A State-of-the-Art Code Retriever

<https://www.nomic.ai/blog/posts/introducing-state-of-the-art-nomic-embed-code>

[31](#) [32](#) [44](#) [45](#) [46](#) [47](#) Embedding Models: OpenAI vs Gemini vs Cohere

<https://research.aimultiple.com/embedding-models/>

[53](#) [54](#) [64](#) [68](#) [69](#) [70](#) [71](#) [72](#) README.md

<https://github.com/vmlinuzx/glideclubs/blob/192e899e53e6dc63cf3660b86021b03227bf03d9/scripts/rag/README.md>

[56](#) [62](#) [63](#) index\_workspace.py

[https://github.com/vmlinuzx/glideclubs/blob/192e899e53e6dc63cf3660b86021b03227bf03d9/scripts/rag/index\\_workspace.py](https://github.com/vmlinuzx/glideclubs/blob/192e899e53e6dc63cf3660b86021b03227bf03d9/scripts/rag/index_workspace.py)

[66](#) What is Cohere pricing? - Unleash.so

<https://www.unleash.so/post/what-is-cohere-pricing>