

# Evaluating RAG Retriever Layer Options for Qdrant and Neo4j

## Context and Criteria

The goal is to choose a **Retrieval-Augmented Generation (RAG)** retriever component that fits a tech stack with multiple data stores and services, including Neo4j (graph DB), Qdrant (vector DB), PostgreSQL, Redis, RabbitMQ/Dramatiq, MinIO, and agent-driven LLM workflows. Key requirements include:

- **Direct integration with Qdrant (vector embeddings) and Neo4j (knowledge graph)** for hybrid retrieval (combining unstructured text similarity with graph-based context).
- **Support for agent-driven workflows**, including multi-hop or iterative retrieval (where an LLM agent can perform *search* -> *reason* -> *search* cycles) and streaming responses.
- **Summarization and incremental context building**, i.e. ability to condense information or build up answers from multiple pieces of context.
- **Hybrid search pipelines** that utilize both vector similarity and knowledge graph relationships.
- **Python compatibility** to integrate with the existing stack and orchestration (likely Python-based).

Below we compare several frameworks and tools – both popular and lesser-known – against these criteria, and then recommend the most suitable option.

## Frameworks Considered

### LlamaIndex (GPT Index)

LlamaIndex is a flexible framework purpose-built for connecting LLMs with external data. It **natively supports Qdrant** as a vector store (via a plugin package) <sup>1</sup> <sup>2</sup>, and also provides **graph store integration** for Neo4j. In fact, LlamaIndex *“excels in RAG applications, providing robust capabilities for knowledge graphs, document indexing, and structured data access”* <sup>3</sup>. Key points for LlamaIndex:

- **Data Integration:** It can ingest data from various sources (documents, databases, APIs) and build composite indices. For example, you can use a `QdrantVectorStore` to handle embeddings <sup>2</sup> and simultaneously use a `Neo4jGraphStore` to persist or query a knowledge graph extracted from text. LlamaIndex's KnowledgeGraphIndex can even extract triplets from text and store/query them in Neo4j <sup>4</sup> <sup>5</sup>. This means it's well-suited to hybrid pipelines combining vector search with graph queries.
- **Retrieval and Summarization:** LlamaIndex offers various query modes (retrieval, synthesis, **summarization**) and index structures (vector indices, tree indices, keyword tables, knowledge graphs). It supports **incremental summarization and hierarchical context**: e.g. you can retrieve many documents and then use an LLM to summarize them (it has built-in `response_mode="tree_summarize"` for hierarchical summarization <sup>5</sup>). This is useful if a direct answer requires condensing information from multiple nodes or documents.

- **Agent/Workflow Support:** While LlamaIndex is not primarily an agent framework, it can be composed in multi-step workflows. You can route queries to different indices or do sequential retrieval (for instance, first query the graph for a specific entity, then use that result to query the vector index). It also integrates with LangChain, so you could wrap LlamaIndex queries as tools for an LLM agent if needed. Out-of-the-box streaming token support is limited (it typically returns completed answers), but one can integrate it with streaming LLM APIs by customizing the LLM callback in LlamaIndex.
- **Integration Effort:** LlamaIndex provides high-level abstractions, which means less manual coding to combine Qdrant and Neo4j. You'd use its API to connect to Qdrant (as a vector index) <sup>2</sup> and connect to Neo4j via a GraphStore. If a direct Neo4j query is needed (e.g. a Cypher query), that might require a custom tool or LangChain agent, since LlamaIndex's knowledge graph interface is geared toward graphs it builds. Overall, it's a powerful option if you need to seamlessly integrate *embeddings + knowledge graph* with support for complex data types.

**Trade-offs:** LlamaIndex is very flexible but can be complex to configure optimally. Its learning curve is moderate. It is a rapidly evolving project, so ensure version compatibility. It may not have built-in "agent loops" like LangChain, but it excels at index construction and has *first-class support for knowledge graphs and hybrid data* <sup>3</sup>. This makes it ideal when you need sophisticated data integration and summarization.

## Haystack (deepset)

Haystack is a mature open-source framework for search and QA pipelines. It now supports Qdrant via a `QdrantDocumentStore` integration <sup>6</sup> <sup>7</sup>, allowing you to use Qdrant for dense vector retrieval. Haystack's strength lies in building **production-ready pipelines for search-oriented applications** – as noted, *"Haystack stands out in search-oriented applications, offering modular pipeline construction for question-answering systems and document processing."* <sup>8</sup> Key points for Haystack:

- **Vector Search:** Using `qdrant-haystack`, you can store and query embeddings in Qdrant just like other Haystack stores (it's maintained by Qdrant team) <sup>6</sup>. Haystack also supports hybrid retrieval (dense + sparse) by combining retrievers, and can re-rank results using a separate model.
- **Knowledge Graph Integration:** Haystack does not have native Neo4j integration in its pipeline *by default*. However, it provides a concept of a **knowledge graph** and `GraphRetriever` in older versions (for example, using a small in-memory graph or a Neo4j integration library). In fact, an official Neo4j-Haystack connector exists (`Neo4jDocumentStore`) which leverages Neo4j's indexes for storing documents/embeddings <sup>9</sup>. One approach is to use Neo4j's latest capability as a vector index (Neo4j 5+ supports vector similarity on nodes) as a document store. Alternatively, you could use a **custom Node** in a Haystack pipeline: e.g., one branch queries Qdrant, another branch queries Neo4j via Cypher (perhaps retrieving related facts by entity), and then a **JoinNode** merges the results before passing to the generator. This is possible, but not as seamless as in specialized graph-aware frameworks. It likely requires writing a custom Python node that runs a Neo4j query given the input.
- **Agent/Multi-hop:** Haystack's pipeline is typically a directed acyclic graph (DAG) of components rather than an *LLM agent loop*. It's great for one-pass retrieval and answer generation. Doing multi-hop reasoning (where an LLM decides follow-up queries) is not its primary design, though you can craft a pipeline with multiple retrievers or an iterative prompt. Haystack does support **streaming responses** from generative models (the `HaystackChatGPT` component, for instance, can stream token by token), which is useful for real-time agent responses.
- **Summarization and Context:** Haystack provides components like Summarizers and Transformers-based Readers which can summarize documents or do generative QA. You can incorporate a

summarizer to compress retrieved docs if needed. However, these are separate nodes – Haystack will not automatically build hierarchical summaries; you have to design that into the pipeline (e.g., retrieve many docs, summarize them, then feed to final answer step).

- **Integration Effort:** If you strictly want to combine Qdrant and Neo4j, Haystack may require more glue code. For example, you might maintain your knowledge graph in Neo4j but extract textual triples or context to feed into Haystack's documents. Alternatively, use Neo4j indirectly: store relationships or metadata in the Qdrant payload and let the retriever filter by those (not fully utilizing graph traversal though). An advanced approach is to use an LLM in the loop: e.g., retrieve documents via Qdrant, extract an entity, then have a second step where a custom node queries Neo4j for that entity's info. This is doable but not provided out-of-box.

**Trade-offs:** Haystack is **robust and scalable** (used in many production QA systems) and is friendly to deploying pipelines. It excels at pure text **search and QA pipelines** <sup>8</sup>, and it now cleanly integrates with Qdrant. However, for *deep integration of a knowledge graph*, it is not as straightforward as LlamaIndex or Neo4j's own tools. If your use-case is very graph-centric or requires iterative reasoning, Haystack might feel a bit rigid without custom extensions. On the plus side, it has a large community and is proven in enterprise settings.

## LangChain (LLM Orchestration & Agents)

LangChain is a popular framework for composing LLM "chains" and agents. It provides a wide range of integrations and could serve as the glue to incorporate both Qdrant and Neo4j in an **agent-driven workflow**. While LangChain itself is more of an orchestration layer than a standalone retriever, it offers **standardized interfaces for vector stores and tools** that can fulfill our needs:

- **Qdrant Integration:** LangChain has a Qdrant vector store wrapper (`langchain.vectorstores.Qdrant`), which can be used to create a retriever object for semantic search. This makes it trivial to plug Qdrant into a chain or agent as a tool.
- **Neo4j Integration:** LangChain does not have a built-in "Neo4j retriever," but it *does* support querying databases through natural language. One common pattern is using an LLM agent with a **Cypher search tool**. For example, you can give the agent the ability to call a tool that executes Cypher queries on Neo4j. The agent can then decide, given a question, to formulate a Cypher query (often by having the LLM generate it) to fetch structured data from the graph. There are community examples of a *GraphCypherQChain* or an agent that translates questions to Cypher for Neo4j. This allows leveraging the knowledge graph when appropriate.
- **Agents and Multi-hop Reasoning:** LangChain's biggest strength is facilitating complex **multi-step reasoning** with LLMs. It provides an agent loop where the LLM can observe intermediate results and decide the next action (e.g., "search the vector DB", then "query the graph for related info", then "summarize and answer"). This directly addresses the requirement for live agent workflows and multi-hop search. Using LangChain, you can define **multiple tools**: one tool might do a semantic similarity search in Qdrant, another tool might do a structured query in Neo4j, another might do an SQL lookup in PostgreSQL, etc. The agent (powered by an LLM like GPT-4 or Claude) can dynamically choose which tools to call in order to assemble the answer. This is very powerful for complex queries that need both unstructured and structured data. LangChain's agent system also supports **streaming** token outputs (if the LLM API supports it), so you can stream the final answer to the user.
- **Summarization & Memory:** LangChain provides chain templates for tasks like summarization and also a concept of *memory* (storing conversational context). For incremental context building, one could use a summary chain to distill long retrieved texts, or use the agent to chunk and summarize

on the fly. However, these would be custom-designed chains; LangChain gives you the building blocks.

- **Integration Effort:** Using LangChain will likely require writing some code to set up the chains/agent. For instance, configuring a Neo4j tool requires coding the interface to run a query against the Neo4j Python driver and return results. The upside is that LangChain has many examples and a vibrant ecosystem. Its “*integration with popular tools and databases streamlines development*”, letting developers focus on logic rather than low-level connections <sup>10</sup>. Indeed, LangChain already has integrations for Qdrant, and Neo4j integration can be achieved either via existing community chains or a few lines of custom tool code.

**Trade-offs:** LangChain is very **flexible and feature-rich** (with a “first-mover advantage” in the LLM framework space) <sup>11</sup> <sup>12</sup>. It’s likely the best choice if you need an **LLM agent orchestrating multiple data sources**. The trade-off is complexity: you are assembling the pipeline yourself, which can introduce more points of failure if not done carefully. Also, LangChain’s abstractions can sometimes add runtime overhead. In summary, LangChain would shine in scenarios where an intelligent agent must decide how to query Neo4j vs Qdrant vs other sources dynamically. If you prefer a more static or straightforward pipeline, a purpose-built retriever (like GraphRAG or LlamaIndex) might be simpler. Many teams use **LangChain in combination with other tools** (for example, LangChain agent calling LlamaIndex indices as tools) to get the best of both.

*(Note: An extension of LangChain called LangGraph specifically targets complex agent workflows with stateful control flows. LangGraph provides a way to define cyclic workflows and persistent state for agents, which can be useful for reliable multi-turn interactions <sup>13</sup> <sup>14</sup>. It works with LangChain’s components (including Qdrant retrievers) and could be considered if you need very fine-grained control over an agent’s planning loop. This is an emerging tool, so it adds complexity but addresses cases where LangChain’s standard agent might falter in long reasoning chains.)*

## Neo4j GraphRAG (Graph Retrieval-Augmented Generation)

Neo4j’s **GraphRAG** is a specialized framework developed by Neo4j to combine graph data with vector search. It is a **first-party Python library** from Neo4j designed for exactly this use case <sup>15</sup>. GraphRAG offers a “*robust, feature-rich, and high-performance solution, with long-term support from Neo4j, and it natively provides a Qdrant retriever*” <sup>15</sup>. Key features:

- **Seamless Qdrant + Neo4j Retrieval:** GraphRAG’s `QdrantNeo4jRetriever` bridges the two systems. In a single call, you can encode a query, perform a **vector search in Qdrant**, and then automatically **fetch the corresponding Neo4j nodes** for the top hits <sup>16</sup> <sup>17</sup>. Each vector record in Qdrant is expected to carry a reference (like a Neo4j node ID in its payload), so the retriever links them: “*Each Qdrant match references a Neo4j node ID... The retriever uses that ID to fetch the corresponding node in Neo4j.*” <sup>16</sup> <sup>17</sup>. This yields a list of results enriched with graph data (properties/relationships from Neo4j).
- **Graph-Based Context:** Because results are actual Neo4j nodes, you can easily leverage the graph structure for context. For instance, once you have a node, you could use Neo4j queries to pull its neighbors (related entities) as additional context for the LLM. GraphRAG is built to “*use Neo4j’s graph model to add context and relationships to your results*” <sup>18</sup>. It effectively lets Qdrant handle the heavy vector similarity search, while Neo4j provides the knowledge connections – “*the best of both worlds for RAG*” <sup>19</sup>.

- **Integrated Pipeline:** GraphRAG includes not just retrieval but also generation support. It has an `GraphRAG` class that takes a retriever and an LLM interface (it's compatible with LangChain's LLM wrappers, or you can use OpenAI API, etc.) <sup>20</sup> <sup>21</sup> . This means you can call `rag.search(query)` and it will internally retrieve context and then call the LLM to generate an answer. The library allows configuration of the prompt used for the LLM and can return the answer directly. You can also get the retrieved context back (by setting `return_context=True`) for inspection <sup>22</sup> <sup>23</sup> .
- **Support for Multiple Vector Stores:** While Qdrant is a primary target, the GraphRAG package also supports other vector DBs (Weaviate, Pinecone, etc.) <sup>18</sup> . Alternatively, if you prefer, it can use Neo4j's own vector indexing (Neo4j 5.x) via a `VectorRetriever` that queries an index on the graph nodes <sup>24</sup> . But using Qdrant externally is advantageous for scaling and was a key design goal <sup>25</sup> <sup>26</sup> .
- **Agent Workflow:** Out-of-the-box, GraphRAG executes a **single-step retrieval + answer**. It doesn't provide an agent loop for multi-hop reasoning (it assumes one query leads to one set of retrieved context and one answer). However, you can still incorporate GraphRAG into an agent manually. For example, you could use a LangChain agent where one tool calls GraphRAG's retrieval+LLM in one go, or have the agent use GraphRAG's retriever as a sub-step. If more complex reasoning is needed (e.g. ask one question, get nodes, then follow-up question on those nodes' neighbors), you might need to orchestrate that on your own. GraphRAG's strength is in simplifying the typical RAG query, rather than enabling arbitrary tool use by an agent.
- **Summarization/Context Building:** The GraphRAG library doesn't explicitly provide hierarchical summarization or chunking (its philosophy is to let the LLM handle the returned context). But you could configure the prompt or do a follow-up step to summarize if needed. Since it's a Neo4j-supported tool, it's optimized for performance and integration rather than fancy prompting techniques. If you have extremely large context (many graph nodes), you may need to perform an extra summarization step outside GraphRAG.

**Trade-offs:** For a stack that *already includes Neo4j and Qdrant*, Neo4j GraphRAG is a very **natural fit**. It was literally built to *"make combining graph data and vector similarity easier"* <sup>18</sup> , and to connect Neo4j and Qdrant with just a few lines of code <sup>18</sup> . The big advantage is **low integration effort** and official support. You'll get updates and maintenance from Neo4j's team, which is valuable for long-term reliability <sup>15</sup> <sup>27</sup> . The main caution is that GraphRAG is a newer library (ensure it meets your specific needs and consider the community size – though backed by Neo4j, it might be less community-driven than LangChain or Haystack). Also, if your use case requires steps beyond the straightforward retrieve-and-answer (like complex multi-turn dialogues with reasoning), you might end up combining GraphRAG with an agent framework. In summary, GraphRAG offers an **excellent hybrid retriever** that directly satisfies the Qdrant+Neo4j integration requirement.

## Other Notable Frameworks (Agent & Workflow Focused)

In addition to the above, a couple of less common frameworks may be worth mentioning for completeness:

- **Stanford DSPy:** Stanford's DSPy (Declarative Sequential Processes for LLMs) is an AI framework for defining complex decision-making and reasoning workflows with LLMs. In the context of RAG, DSPy can orchestrate retrieval and chain-of-thought reasoning. For example, a recent demo combined DSPy with Qdrant to build a medical QA bot that performs multi-step retrieval and reasoning <sup>28</sup> <sup>29</sup> . In that pipeline, *"DSPy searches against the Qdrant vector database to retrieve top documents... the retrieved passages are then reranked... DSPy uses these passages to guide the LLM through a chain-of-thought reasoning to generate the answer."* <sup>29</sup> . This showcases DSPy's ability to handle *multi-hop logic*

(retrieve, re-rank, reason, etc.). While powerful, DSPy is more of a research-centric toolkit and doesn't directly provide Neo4j integration. You would likely use DSPy in tandem with code that queries Neo4j (similar to how one would with LangChain). Its strength lies in structuring the agent's reasoning process and ensuring more **transparent, step-by-step logic** (which can improve correctness for complex queries). If your focus is on **advanced reasoning and you're open to a newer framework**, DSPy could be an option – but expect to write custom code for database interfacing.

- **Semantic Kernel or Custom Orchestration:** Microsoft's Semantic Kernel and other orchestration libraries allow plugin mechanisms to connect to databases. For example, you could create a Semantic Kernel function to query Neo4j and another to query Qdrant, and have a semantic planner decide which to call. These are more custom and lower-level compared to LangChain or LlamaIndex. They give you flexibility in C# or Python, but you'll be implementing more manually.
- **Hugging Face Transformers Tools / AutoGen:** There are emerging "agent" systems like HuggingFace's transformers agent or frameworks like AutoGen (by Microsoft) that enable multi-agent or tool-using LLMs. These can be quite powerful (AutoGen, for instance, focuses on multi-agent conversations for problem solving <sup>30</sup>). However, adopting them would mean custom integration with Qdrant/Neo4j and a relatively higher complexity. They might be overkill unless you have a specific need for multi-agent collaboration.

In summary, beyond the mainstream options, there are frameworks that excel in certain niches (LangGraph for complex workflow state management <sup>31</sup>, DSPy for guided CoT reasoning, etc.). These could be layered on top of the retrieval layer if needed. For choosing the core retriever component, though, the focus should remain on the tools that directly address Qdrant+Neo4j integration and the RAG pattern.

## Feature Comparison of Top Options

Below is a comparison of how the top candidates meet the key criteria:

Framework	Qdrant Integration	Neo4j Integration	Agent/ Multi-hop Support	Summarization & Context	Hybrid (KG + Vector)	Python Friendly
LlamaIndex	Yes (native vector store plugin for Qdrant) <sup>1</sup>	Yes (GraphStore for Neo4j, builds or queries knowledge graphs) <sup>4</sup>	Partial – not an agent loop itself, but can be used in multi-step flows; integrates with LangChain for agents	Strong – built-in hierarchical indexing and summary capabilities (tree summarize, etc.)	Yes – designed for integrating structured data (KG) with text <sup>3</sup>	Yes (pure Python)

Framework	Qdrant Integration	Neo4j Integration	Agent/ Multi-hop Support	Summarization & Context	Hybrid (KG + Vector)	Python Friendly
Haystack	Yes (via <code>qdrant-haystack</code> <code>DocumentStore</code> ) <sup>6</sup>	Partial – no out-of-box Neo4j in pipeline (Neo4j integration exists for storage, or use custom node) <sup>9</sup>	Limited – uses static pipelines; multi-hop would require custom pipeline logic (no LLM agent deciding)	Moderate – can add summarizer nodes; primarily focuses on retrieval + optional reader (generative QA)	Indirect – can combine sources in pipeline but requires custom merging (no native KG traversal support)	Yes (Python)
LangChain	Yes (Qdrant VectorStore class available)	Yes (indirectly via tools/ agents – e.g., LLM-generated Cypher queries)	Excellent – supports agents with tool use for iterative search, multi-hop reasoning, and streaming outputs <sup>10</sup>	Custom – provides utilities to summarize or manage memory, but you construct these chains manually	Yes – through agent tools you can incorporate both vector search and graph queries in one workflow	Yes (Python)
Neo4j GraphRAG	Yes (first-party Qdrant retriever built-in) <sup>15</sup>	Yes (direct Neo4j driver usage; retrieves actual Neo4j nodes) <sup>16</sup>	Partial – easy one-shot queries; for multi-hop, would need external orchestration (can be combined with agents if needed)	Decent – retrieves relevant graph nodes; relies on LLM to use that context. No automatic summarizer, but context size is manageable via top-k.	Yes – <b>native hybrid:</b> purpose-built to combine Neo4j relationships with vector similarity <sup>19</sup>	Yes (Python)

Framework	Qdrant Integration	Neo4j Integration	Agent/ Multi-hop Support	Summarization & Context	Hybrid (KG + Vector)	Python Friendly
Stanford DSPy	Yes (via Python code – e.g., use Qdrant client in DSPy functions)	Partial (no native driver, but can call Neo4j via Python in the chain)	Excellent – focuses on stepwise CoT reasoning, you explicitly script multi-step retrieval and thinking <sup>29</sup>	Moderate – not specific to summarization, but chain-of-thought approach can incorporate summary steps if scripted	Possible – not built-in, but you can design a workflow that queries both as needed	Yes (Python)

*Table: Comparison of RAG Retriever Layer Options* for integration, agent capabilities, and hybrid search support.

## Recommendation and Implementation Strategy

Considering the trade-offs, the **most suitable RAG retriever layer** for this tech stack is **Neo4j’s GraphRAG, augmented with LangChain for agent orchestration**. This combination plays to the strengths of each tool:

- **Neo4j GraphRAG** as the core retriever gives a straightforward, supported way to query Qdrant and Neo4j together. It directly fulfills the need to unify vector search with graph context: *Qdrant handles vector similarity while Neo4j provides relationships — “the best of both worlds for RAG”* <sup>19</sup>. With GraphRAG, you can, for example, store documents or facts as Neo4j nodes (with embeddings in Qdrant) and easily retrieve them by semantic similarity, then leverage Neo4j to pull connected information. This will drastically simplify implementation compared to crafting a custom pipeline from scratch. Since it’s maintained by Neo4j, you gain reliability and updates <sup>15</sup>.
- **LangChain** (or a similar agent framework) can be layered on top to handle the **agent-driven workflow** requirements. While GraphRAG can answer direct queries on its own, LangChain would allow the system to have a dialogue or multi-step interactions. For instance, you can expose a tool to the agent that wraps GraphRAG’s `search` function. The agent (an LLM) could decide to call this tool when it needs information, and you could also have other tools (e.g. a direct Cypher query tool for complex graph queries, or a SQL tool for Postgres if needed). This way, you enable **multi-hop reasoning** – the LLM might ask GraphRAG for an initial piece of info, then follow up by drilling down into the graph via another Cypher query, etc., all in a single conversation. LangChain will also easily allow streaming the LLM’s answer to the user as it’s generated.
- This hybrid approach meets all criteria: direct Qdrant/Neo4j integration (via GraphRAG), agent-based multi-hop search (via LangChain’s agent), summarization if needed (the agent can invoke a summarization chain or you can prompt the LLM to summarize). It is Python-based and fits well with your stack.

**Why not solely LlamaIndex or Haystack?** LlamaIndex is a close contender – it offers similar hybrid capabilities and might even be used in this solution for specific sub-tasks (e.g., building an index of documents or performing a complex graph query via its KG index). If your team is very comfortable with



LlamaIndex, it could replace GraphRAG as the retrieval layer, since it “provides robust capabilities for knowledge graphs... making it ideal for sophisticated data integration” <sup>3</sup>. However, GraphRAG gets the edge due to its purpose-built integration and official support, which reduces development friction. Haystack, on the other hand, would require more customization to incorporate Neo4j, and it shines less in an agent setting (it’s better for static pipelines). Haystack could still be useful for certain components (for example, using its `DocumentStore` abstraction for other data or its reader for extractive QA), but it’s not the optimal choice for the centerpiece of a Neo4j+Qdrant RAG system.

**Implementation Pattern:** The recommended implementation is to use **Neo4j GraphRAG’s QdrantNeo4jRetriever** in conjunction with a **LangChain Agent**:

1. **Set up GraphRAG Retriever:** Use `QdrantNeo4jRetriever` from `neo4j_graphrag`. Provide it the Neo4j driver and Qdrant client, and ensure each vector in Qdrant has a payload field linking to the Neo4j node ID (e.g., `neo4j_id`). This linking is crucial – as described, “make sure your Qdrant payload (e.g. `neo4j_id`) matches the node property in Neo4j so the retriever can join data seamlessly.” <sup>32</sup> <sup>33</sup>. Once configured, test that `retriever.search(query_text="...")` returns the top-k Neo4j nodes relevant to a query.
2. **Incorporate LLM and Prompt:** Decide on the LLM to use (could be OpenAI GPT-4, or local models via HuggingFace). You can use GraphRAG’s `GraphRAG` pipeline class to generate answers directly <sup>21</sup>, or simply use the retriever in a custom chain. For more control, you might manually prompt the LLM with the retrieved context. For example: “Given the following information from our knowledge base: {retrieved\_context}, answer the user’s question...”. GraphRAG’s default pipeline will handle this for you, but in an agent scenario you might handle it in the agent’s prompt.
3. **LangChain Agent with Tools:** Create a LangChain agent and register tools:
4. A “**KnowledgeSearch**” tool that calls the GraphRAG retriever (you could wrap `retriever.search` so that the agent, given a query, gets back a snippet of knowledge or a set of facts). The result from this tool would include the content of the Neo4j nodes (and possibly their relationships if you choose to fetch some neighbors as context).
5. Optionally, a “**GraphQuery**” tool that allows the agent to execute arbitrary Cypher queries on Neo4j. This could be useful for questions that require a specific graph traversal or when the agent deduces it needs structured data (for instance, counting relationships, finding a path between entities, etc.).
6. A “**DatabaseQuery**” tool for PostgreSQL if the use-case might involve structured data lookup from SQL.
7. Any other tools (maybe a “WebSearch” or a “Calculator” etc., depending on your application needs).
8. The agent’s prompt will instruct it on how and when to use these tools.
9. **Streaming Responses:** Ensure the chosen LLM supports streaming. With LangChain, you can enable streaming easily on the LLM wrapper (e.g., `StreamingStdOutCallbackHandler`). This way, as the agent formulates the final answer, you can stream it token-by-token to the user for a responsive experience.
10. **Summarization & Context Management:** If the retrieved context is too large (say the graph returns many nodes or lengthy text), consider adding a step to summarize it before passing to the final answer. This could be a simple call to an LLM to compress the info, or using LlamaIndex’s summarizer on the fly. Also implement conversation memory if this is a chatbot – LangChain’s `ConversationBufferMemory` or a summary memory can help keep track of dialogue history without exceeding token limits.
11. **Testing and Tuning:** Finally, rigorously test the system. Evaluate how the agent decides between using Qdrant vs direct Cypher queries. You might need to fine-tune the agent instructions so that it

properly uses the GraphRAG tool for most semantic lookups and uses direct Cypher for specific structured queries (if at all). Monitor for errors (like the agent hallucinating a Cypher query that doesn't run – you'll need to catch exceptions and possibly feed that back to the agent).

By following this approach, you leverage GraphRAG's optimized retrieval to get relevant, graph-enriched context, and LangChain's agent to integrate that retrieval in a flexible, intelligent workflow. This addresses all the requirements: **direct Qdrant/Neo4j integration**, **multi-hop reasoning**, **streaming answers**, and **hybrid data sources**.

In conclusion, **Neo4j GraphRAG** provides a purpose-built RAG retrieval backbone for your Neo4j + Qdrant data, while **LangChain (or a similar agent framework)** provides the brains to use that backbone effectively in an interactive LLM application. This combination is well-suited to your stack and should minimize custom glue code, as most integration points are supported out-of-the-box. It offers a powerful yet maintainable solution for building a RAG system that can draw on both **graph knowledge** and **vector semantics** to assist your LLM in generating accurate, context-rich answers.

#### Sources:

- Neo4j GraphRAG official documentation and Qdrant integration example <sup>15</sup> <sup>16</sup>
- Neo4j Developer Blog on combining Neo4j and Qdrant for RAG <sup>19</sup> <sup>18</sup>
- Analysis of LlamaIndex, LangChain, and Haystack capabilities in RAG contexts <sup>3</sup> <sup>8</sup> <sup>10</sup>
- Qdrant and Neo4j GraphRAG usage guide <sup>34</sup> <sup>32</sup>
- Stanford DSPy example for chain-of-thought retrieval reasoning <sup>29</sup>

---

#### <sup>1</sup> <sup>2</sup> LlamaIndex - Qdrant

<https://qdrant.tech/documentation/frameworks/llama-index/>

#### <sup>3</sup> <sup>8</sup> <sup>10</sup> <sup>11</sup> <sup>12</sup> <sup>30</sup> <sup>31</sup> AI Frameworks in Action: Building RAG Systems with LangChain, LlamaIndex, and Haystack!

<https://www.linkedin.com/pulse/ai-frameworks-action-building-rag-systems-langchain-pavan-belagatti-wmg3c>

#### <sup>4</sup> <sup>5</sup> Neo4j Graph Store - LlamaIndex

[https://docs.llamaindex.ai/en/stable/examples/index\\_structs/knowledge\\_graph/Neo4jKGIndexDemo/](https://docs.llamaindex.ai/en/stable/examples/index_structs/knowledge_graph/Neo4jKGIndexDemo/)

#### <sup>6</sup> <sup>7</sup> Haystack - Qdrant

<https://qdrant.tech/documentation/frameworks/haystack/>

#### <sup>9</sup> Haystack Neo4j Integration - Neo4j Labs

<https://neo4j.com/labs/genai-ecosystem/haystack/>

#### <sup>13</sup> <sup>14</sup> LangGraph - Qdrant

<https://qdrant.tech/documentation/frameworks/langgraph/>

#### <sup>15</sup> <sup>27</sup> Neo4j GraphRAG - Qdrant

<https://qdrant.tech/documentation/frameworks/neo4j-graphrag/>

#### <sup>16</sup> <sup>17</sup> <sup>18</sup> <sup>19</sup> <sup>25</sup> <sup>26</sup> <sup>32</sup> <sup>33</sup> <sup>34</sup> Integrate Qdrant and Neo4j to Enhance Your RAG Pipeline - Graph Database & Analytics

<https://neo4j.com/blog/developer/qdrant-to-enhance-rag-pipeline/>

20 21 22 23 24 **User Guide: RAG — neo4j-graphrag-python documentation**

[https://neo4j.com/docs/neo4j-graphrag-python/current/user\\_guide\\_rag.html](https://neo4j.com/docs/neo4j-graphrag-python/current/user_guide_rag.html)

28 29 **Building a Chain-of-Thought Medical Chatbot with Qdrant and DSPy - Qdrant**

<https://qdrant.tech/documentation/examples/qdrant-dspy-medicalbot/>