

# Data Governance in the Age of Generative AI: Leveraging Knowledge Graph RAG for Advanced Reasoning



Anthony Alcaraz · Follow

Published in CodeX · 9 min read · Jun 10, 2024



54



As generative AI continues to advance, delivering its transformative benefits will hinge on a solid foundation of data governance.

However, many organizations are finding themselves unprepared for the novel data management and oversight challenges posed by large language models (LLMs) and retrieval-augmented generation (RAG) pipelines.

To truly harness the power of generative AI for advanced reasoning and insight discovery, especially over complex information architectures, a new synthesis is needed — integrating knowledge graphs with RAG models.

This article will explore the data governance imperatives of generative AI, how knowledge graphs and RAG models can complement each other to enable powerful multi-hop reasoning, and a framework for deciding when to invest in building knowledge graph RAG systems.

By marrying the formal knowledge representation of graphs with the flexibility of neural language models, organizations can unlock unparalleled capabilities for querying and making sense of their swelling oceans of data.

# Data Governance Challenges in Generative AI

Generative AI models like LLMs represent a step change in data interaction complexity. Traditional database lookups and business intelligence tools query relatively well-governed structured data through a limited number of explicit interfaces. Generative AI systems, in contrast, are being woven into user-facing applications that must rapidly ingest, retrieve, rank and generate outputs based on huge and dynamic corpuses of unstructured data. This introduces several critical data governance challenges:

**Security, Privacy and Access Control:** Generative AI systems that are exposed to users need fine-grained access policies on the data itself (ensuring private information isn't inadvertently revealed) as well as the model parameters (avoiding extraction of sensitive data the model was trained on). These access policies must be enforced across each step of user request-response loops, from prompt validation to retrieval to final response redaction [1].

**Unstructured Data Management:** LLMs need to be trained on domain-specific sources like documents, knowledge bases, and conversation logs to generate relevant outputs grounded in enterprise knowledge. However, compared to structured databases, unstructured data is often scattered across silos and not governed with consistent metadata, taxonomies and quality controls [1]. Tracking provenance and implementing retention policies is much harder.

**Data Quality:** Ensuring the quality of data ingested by generative AI models is critical, as inaccurate or biased data leads to wrong and even harmful generated content. Model outputs then create more bad data in a vicious cycle. But assessing quality is more complex, requiring natural language processing techniques to parse semantics and domain-specific rulesets to validate claims [2].

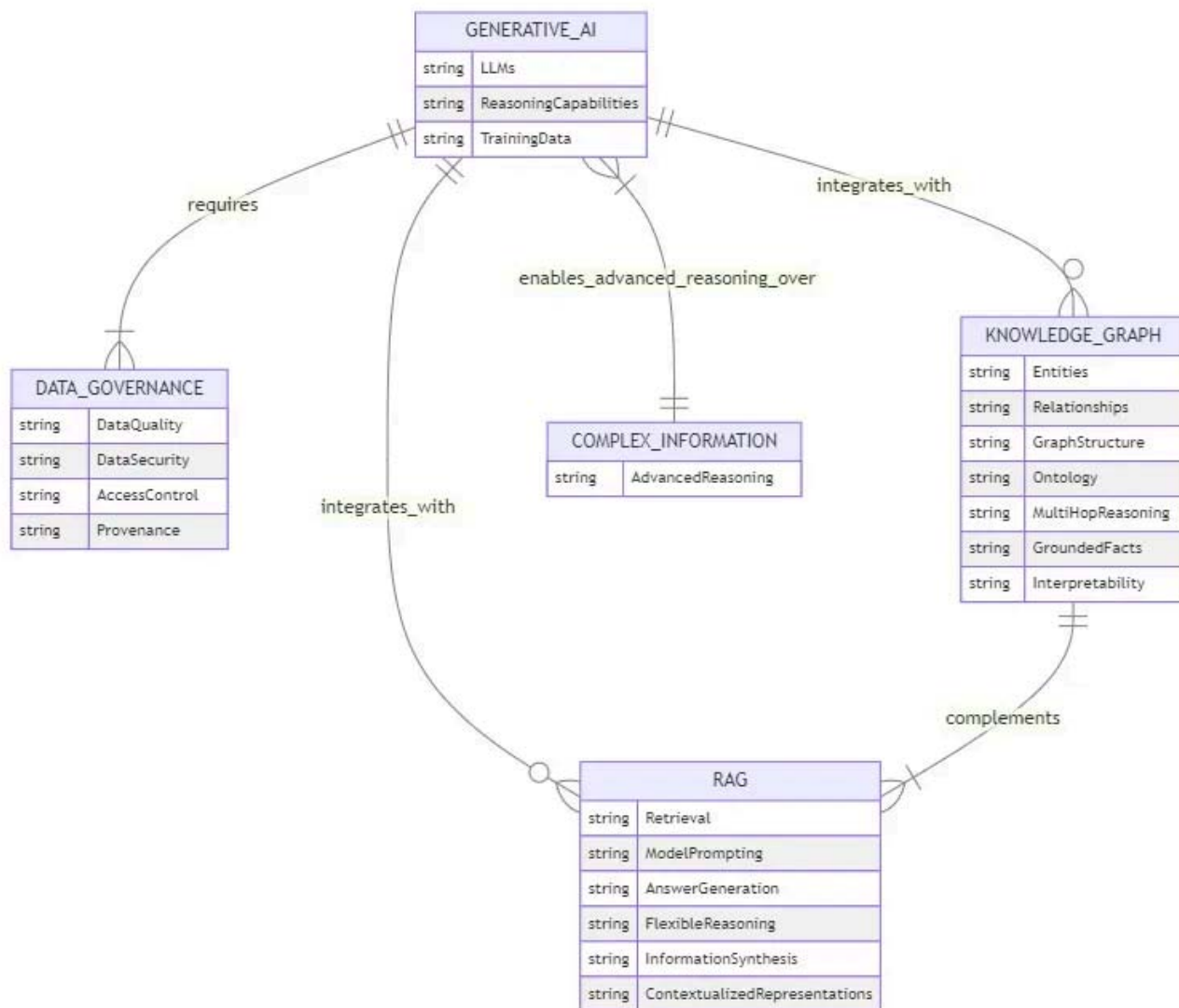
**Complete Lifecycle Governance:** RAG pipelines span data ingestion, wrangling, vectorization, training, finetuning, querying, retrieval, and generation. Each step needs tailored governance — from tracking ingested data through feature stores to auditing model versioning to imposing generation guardrails. Managing this end-to-end lifecycle requires tools for creating observability at scale.

**Mitigating AI Risks:** Immature data governance makes it harder to root out AI risks like data and algorithmic biases that damage model fairness, factual inconsistencies that erode trust, and toxic outputs that create legal liabilities. Tracing risks back to underlying data root causes and having a holistic understanding of data flows and usage is crucial.

Gartner predicts that organizations that operationalize AI transparency, trust and security will see their AI models achieve a 50% result improvement in terms of adoption, business goals and user acceptance.

But getting data governance right for generative AI requires a convergence of technologies, processes and oversight mechanisms that encompasses the unique characteristics of unstructured data and machine learning pipelines.

A key step forward is integrating knowledge graphs with RAG architectures.



Made by the author

## Complementary Strengths: Knowledge Graphs and RAG Models

Knowledge graphs and retrieval-augmented generation have emerged as two different paradigms for question answering — one providing formal guarantees grounded in curated taxonomies and rules, the other flexibility in synthesizing knowledge scattered across unstructured text. By combining them, the factual structure of knowledge graphs can complement the open-ended generation of RAG models:

**Explicit vs Latent Knowledge:** Knowledge graphs represent facts and relationships between entities as an explicit, traversable schema, while RAG

models encode knowledge latently across neural network parameters.

Knowledge graph RAG allows grounding generated text explicitly to facts in the graph, reducing hallucination.

**Logical vs Fuzzy Reasoning:** Knowledge graphs support deterministic logical inferencing based on ontologies and rules, while RAG models reason fuzzily based on patterns and analogies in training data. Using ontologies to guide RAG retrieval and validate generated claims against the graph allows the best of both modes.

**Precise vs Contextual Semantics:** Knowledge graphs precisely define entities and their semantic links, while RAG models learn contextualized representations of fuzzy concepts. Having a symbolic scaffolding from the graph allows RAG to more precisely bind concepts to definitions across contexts.

**Multi-Hop vs Single-Hop Reasoning:** By recursively traversing relationships, knowledge graphs can connect the dots between facts many hops apart. RAG models typically retrieve and attend to single passages, limiting their reasoning horizon. GraphRAG techniques decompose multi-hop queries into chains of local traversals and retrievals.

**Interpretable vs Black Box:** Knowledge graphs make the chain of evidence supporting a conclusion fully inspectable as a path through the graph. RAG models, like most neural networks, are black box systems where the reasoning trace is opaque. Materializing RAG model reasoning as graph traversals promotes interpretability.

Harnessing the complementary strengths of knowledge graphs and RAG models hinges on thoughtful integration. Promising approaches include using knowledge graphs to inform dense passage retrieval for RAG, chaining RAG-retrieved contexts into reasoning trajectories guided by the graph, and using graphs to impose consistency constraints on generated text. A key architectural

decision is when and how to modularize reasoning between knowledge-graph and neural components.

GraphRAG is an example of an architecture that modularizes reasoning at a meso-scale between individual graph nodes and the entire corpus using graph-clustering algorithms. B

y recursively retrieving, ranking and aggregating relevant embeddings of graph clusters (capturing local contexts), the attention horizon of the RAG model is progressively expanded to cover larger swaths of the graph.

This allows dynamically scaling reasoning from single-hop, local connections to arbitrarily long, cross-corpus chains of insights.

## The Case for Knowledge Graph RAG: A Framework

Knowledge graph RAG can be a significant investment, both in terms of the technical infrastructure required and the organizational processes for curating enterprise knowledge in a graph format. So when does it make sense to take on this complexity over using traditional document-oriented RAG or pure knowledge graphs? The key considerations center around data characteristics and business drivers:

**Complex Information Architectures:** Knowledge graph RAG shines when the underlying corpus has a highly interconnected or deeply hierarchical structure. This includes domains like biomedical research, legal contracts, and product documentation where the knowledge is densely linked or categorized. Graphs allow these cross-document connections to be traversed efficiently [1].

**Reasoning-Intensive Use Cases:** Some problem domains inherently require piecing together insights across multiple sources. Key examples include drug discovery (compound/disease relationships), legal investigations (evidence

chains), and supply chain optimization (part/supplier dependencies).

Reasoning over long-range connections tends to break down for vanilla RAG models.

**High-Value, High-Risk Decisions:** The more impactful (and costly if wrong) the decision being assisted by generative AI, the greater need for traceability and verification via grounded facts and evidence trails. Think clinical decision support systems, automated contract analysis, and financial compliance monitoring. Knowledge graphs provide an inspectable substrate to audit conclusions.

**Evolving Knowledge Assets:** In domains where the core knowledge is constantly expanding and changing, like scientific literature and news feeds, keeping generative AI aligned with the latest expert understanding requires frequently retraining models. A knowledge graph provides a single snapshot of versioned truth to flexibly retrain multiple downstream RAG models as needed.

Once the strategic case for knowledge graph RAG is established, the technical implementation requires carefully stitching together multiple complex pipelines:

**Knowledge Graph Construction:** Existing taxonomies, knowledge bases and metadata schemas provide a starting point, but building a high-coverage domain graph often requires information extraction from semi-structured and unstructured sources. This can involve NLP annotation pipelines, entity linking and ontology induction. Human subject matter experts are needed to validate the resulting graph [1].

**Graph Representation Learning:** To make the knowledge graph compatible with neural RAG models, its nodes and edges need to be transformed into embedding vectors that capture semantic similarity. Techniques like graph convolutional networks, node2vec random walks, and

transformers can be applied here. The embeddings should be periodically updated as the graph evolves [6].

**Dense Retrieval and Ranking:** Efficient RAG depends on retrieving only the top-k graph nodes most relevant to a given query using approximate k-nearest neighbor search over the node embeddings. Retrieval quality also depends on robust node and graph-level ranking methods to score relevance using both graph structure and content [11].

**Focused Finetuning and Prompting:** The backbone LLM needs to be trained or finetuned to accept the specific format of retrieved graph contexts and to generate outputs aligned with the desired answer format. Prompts should be carefully engineered to query the knowledge graph effectively. Augmenting graph nodes with unstructured content (e.g. abstracts) the LLM was pretrained on can help.

**Generated Text Grounding:** The model outputs need to be grounded back to nodes and relationships in the knowledge graph, both for answer provenance and for expanding the graph itself with generated knowledge (which can be human-validated). Mention-entity alignment, graph-guided semantic parsing, and entity-aware generation objectives are some of the techniques used.

In addition to the core architecture, surrounding governance mechanisms and tools are critical to a successful knowledge graph RAG implementation:

**Scalable Graph Platforms:** The knowledge graph system should support transactional updates for fast synchronization with upstream sources and handle complex recursive traversal queries. Graph databases like Neo4j, Amazon Neptune and AnzoGraph provide a starting foundation but often need customization for RAG use cases.



**Quality and Testing Infrastructure:** Validating the outputs of generative AI requires a scalable and flexible testing infrastructure. This includes unit tests for knowledge graph upkeep (data completeness, syntactic checks), integration tests for retrieval and generation modules, and end-to-end tests that assess the relevance, factuality and safety of final user-facing outputs.

**Model and Data Monitoring:** The failure modes of knowledge graph RAG systems can be subtle, spanning distributional shifts in input data, retrieval relevancy drifts, and biased generation outputs. Tools for continuously tracking model performance, data quality, and output metrics in production (e.g. Arize, Fiddler, Arthur) and investigating regressions are essential.

**Responsible AI Guardrails:** Generative AI creates novel risks around truthfulness, trustworthiness and bias. Governance processes should assess these risks upfront, implement appropriate architectural guardrails (e.g. safety classifiers, fact-checking APIs, human oversight), and monitor systems for policy conformance and ethical alignment.

## Conclusion

The age of generative AI is upon us, but delivering on its transformative potential requires rethinking data governance from the ground up. The convergence of knowledge graphs and retrieval-augmented generation models offers a powerful new paradigm for grounded, traceable and flexible reasoning over complex information architectures. By representing enterprise knowledge as richly connected graphs and training language models to traverse and synthesize across these graphs, organizations can tap into truly advanced question answering and insight discovery capabilities.

However, implementing knowledge graph RAG is not without its challenges. It requires weaving together multiple complex machine learning and data management pipelines and instilling new governance processes around responsible AI development. Organizations will need to judiciously assess where

the benefits of multi-hop graph reasoning warrant the added technical and organizational overhead based on their specific data characteristics and decision support needs.

As knowledge graph RAG systems mature, they will increasingly blur the lines between neural and symbolic methods.

Future directions include ingesting multimodal data (images, video, equations) into the graphs, training graph representations using augmented language modeling objectives, and building neuro-symbolic architectures that more deeply integrate logical reasoning and natural language generation. The moonshot is to create artificial intelligence that can reason over knowledge with the flexibility of the human mind and the rigor of formal logic.

For business and technology leaders, the strategic question is: what previously intractable queries and unimaginable insights are waiting to be uncovered in your organization's latent knowledge graph?

Unlocking this hidden value will require not just connecting the dots between generative AI and knowledge management, but rewiring data governance for a new era of knowledge-infused, responsible AI systems.

Mapping this uncharted territory will be the defining competitive differentiator of the years ahead.