

## Article

# Multimedia Graph Codes for Fast and Semantic Retrieval-Augmented Generation

Stefan Wagenpfeil 

Department of Software Engineering, PFH University of Applied Sciences, 37073 Goettingen, Germany;  
s.wagenpfeil@pfh.de

**Abstract:** Retrieval-Augmented Generation (RAG) has become a central approach to enhance the factual consistency and domain specificity of large language models (LLMs) by incorporating external context at inference time. However, most existing RAG systems rely on dense vector-based similarity, which fails to capture complex semantic structures, relational dependencies, and multimodal content. In this paper, we introduce Graph Codes—a matrix-based encoding of Multimedia Feature Graphs—as an alternative retrieval paradigm. Graph Codes preserve semantic topology by explicitly encoding entities and their typed relationships from multimodal documents, enabling structure-aware and interpretable retrieval. We evaluate our system in two domains: multimodal scene understanding (200 annotated image-question pairs) and clinical question answering (150 real-world medical queries with 10,000 structured knowledge snippets). Results show that our method outperforms dense retrieval baselines in precision (+9–15%), reduces hallucination rates by over 30%, and yields higher expert-rated answer quality. Theoretically, this work demonstrates that symbolic similarity over typed semantic graphs provides a more faithful alignment mechanism than latent embeddings. Practically, it enables interpretable, modality-agnostic retrieval pipelines deployable in high-stakes domains such as medicine or law. We conclude that Graph Code-based RAG bridges the gap between structured knowledge representation and neural generation, offering a robust and explainable alternative to existing approaches.

**Keywords:** retrieval-augmented generation; graph code; multimedia information retrieval



Academic Editors: Taehyeon Kim and  
KyungTaek Lee

Received: 8 May 2025

Accepted: 31 May 2025

Published: 18 June 2025

**Citation:** Wagenpfeil, S. Multimedia Graph Codes for Fast and Semantic Retrieval-Augmented Generation. *Electronics* **2025**, *14*, 2472. <https://doi.org/10.3390/electronics14122472>

**Copyright:** © 2025 by the author. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction and Motivation

Large language models (LLMs) such as GPT-4 [1] have revolutionized natural language processing by achieving state-of-the-art performance across a wide range of tasks, including question answering, summarization, and code generation. However, despite their impressive capabilities, LLMs suffer from fundamental limitations: they operate on static, pre-trained corpora and are constrained by limited context windows, making them prone to outdated knowledge, domain drift, and factual inconsistencies.

Retrieval-Augmented Generation (RAG) [2] has been proposed as a solution to this problem and already provides significant optimizations in many cases. By retrieving relevant external documents at inference time and conditioning the generation on them, RAG enhances factual accuracy and enables domain adaptation without the need for model retraining. However, the dominant implementations of RAG rely on similarity in dense vector spaces, typically computed over neural embeddings of text segments. While effective at capturing latent semantic similarity, these embedding-based approaches suffer from three critical shortcomings [3]: **(1) Lack of structural representation:** They cannot

express complex relations, hierarchies, or semantic types explicitly. **(2) Poor interpretability:** Dense vectors offer no transparent mechanism to trace which information led to a particular retrieval decision. **(3) Limited multimodal integration:** Embeddings are typically modality-specific and difficult to align across data types (e.g., text and images).

In this work, we propose to address these limitations by replacing opaque vector retrieval with a graph-structured alternative: **Graph Codes** [4]. Graph Codes are compact matrix representations of Multimedia Feature Graphs (MMFGs), which encode entities and their typed relationships extracted from multimodal sources such as documents, images, or structured records. These codes preserve semantic topology and enable fast, interpretable similarity comparison using matrix operations rather than vector distances. Furthermore, Graph Codes have been proven to be significantly faster than graph based operations in Information Retrieval.

We demonstrate that integrating Graph Codes into RAG systems has three major advantages: (1) it enables structure-preserving retrieval with explicit control over semantic and relational alignment, (2) it allows for interpretable and explainable retrieval steps, and (3) it supports the combination of heterogeneous data modalities in a unified framework.

To this end, we develop a formal model of Graph Code similarity and its usage in RAG processing, define an end-to-end pipeline for graph-based multimedia document indexing and retrieval, and evaluate the approach in two representative domains: visual scene understanding and clinical question answering. Our results demonstrate that Graph Code-based RAG significantly outperforms embedding-based baselines in retrieval precision, contextual relevance, and output reliability.

To outline and demonstrate our results, we employ the Nunamaker research methodology as a problem solution approach [5], which divides every research question into four distinct phases: observation, theory building, implementation, and evaluation. Each of these phases contributes to answering the research question. Therefore, our work is also structured according these phases. In Section 2, we present the current state of the art and introduce the basic concepts from previous and related work. Section 3 contains the theoretical foundation and modeling of our approach. In Section 4, we provide a proof-of-concept implementation, which is evaluated based on two major use cases in Section 5. Finally, we summarize and discuss our work in Section 6.

## 2. State of the Art and Technology

In this section, we provide an overview of the required and related technologies. First, we summarize and introduce the research results in the area of Graph Codes being semantic multimedia indexing structures in Section 2.1. Then, we present the current formal processing of RAG systems in Section 2.3. Finally, we summarize related work and its outcome in Section 2.4. For this section, we use scientific databases, e.g., Springer Link [6], Google Scholar [7], ACM Digital Library [8], and standard literature written throughout the last five years.

### 2.1. Graph Codes

Graph Codes are a matrix-based formalism for representing typed, directed semantic graphs in a compact and computationally efficient form. They were initially introduced for multimedia information retrieval but generalize naturally to knowledge graph encoding and context retrieval in large language models. In previous work [4,9,10], a detailed explanation of the production of graph codes from Multimedia Feature Graphs is given. Furthermore, various experiments outline the superior calculation and compute the model of graph codes compared to graph-based, vector-based, or database-based approaches in multimedia information retrieval. Therefore, in the following subsection, only a brief

introduction of the most important concepts of graph codes is given, providing a baseline for the main topics addressed in this paper.

### 2.1.1. Multimedia Feature Graphs (MMFGs)

We let a multimedia document  $d$  be analyzed into a semantic structure  $G = (N, E)$ , where

- $N = \{n_1, \dots, n_k\}$  is a finite set of typed feature nodes (e.g., person, bicycle, accident),
- $E \subseteq N \times T \times N$  is a set of directed, typed edges, with  $T$  being the set of possible relation types (e.g., wears, located\_in, caused\_by).

This graph  $G$  is referred to as a *Multimedia Feature Graph (MMFG)* and represents the semantic topology of a document across modalities.

### 2.1.2. Valuation Matrix and Dictionary Encoding

From an MMFG  $G$ , we derive the valuation matrix as a modified adjacency matrix [11]  $AM \in \mathbb{N}^{k \times k}$  defined as

$$VM[i, j] = \begin{cases} w_{ij} \cdot f_{\text{enc}}(t_{ij}) & \text{if } (n_i, t_{ij}, n_j) \in E \\ 0 & \text{otherwise} \end{cases}$$

where

- $w_{ij} \in \mathbb{R}^+$  is a contextual weight (e.g., relevance or confidence),
- $f_{\text{enc}} : T \rightarrow \mathbb{N}$  is an injective function that encodes relation types.

The corresponding dictionary  $\text{dict} : \{1, \dots, k\} \rightarrow N$  maps row and column indices to their feature labels. Together,  $(AM, \text{dict})$  constitute a **Graph Code (GC)**. Such Graph Codes offer several advantages over traditional graph or vector representations, as they enable batch computing and highly efficient matrix calculations. As each matrix cell represents a concrete typed semantic relation based on multimedia features, Graph Codes can encode a huge number of features from images, text, videos, and structured data in a highly compressable unified format. Furthermore, the similarity of objects can be computed directly via matrix algebra, which avoids costly graph-based traversals and recursive operations. These properties make Graph Codes especially suited for use in Retrieval-Augmented Generation, where both scalability and semantic grounding are crucial.

### Terminology and Definitions

Throughout this paper, we use the following terms in a precise and consistent manner:

- **Indexing:** The process of encoding documents as Graph Codes and storing them in a retrieval-ready data structure.
- **Retrieval:** The act of selecting the top- $k$  most relevant Graph Codes in response to a query based on symbolic similarity metrics.
- **Symbolic similarity:** A structured comparison method that evaluates semantic proximity using the metrics MF (vocabulary overlap), MFR (feature relationship alignment), and MRT (relation type match).
- **Graph alignment:** The structural matching of query and document graphs based on overlapping nodes and edge types during retrieval.

This terminology ensures conceptual clarity across sections and separates preprocessing, similarity evaluation, and ranking steps.

## 2.2. Graph Code Generation Pipeline

Given an input document or query  $q$ , the typical processing pipeline of Graph Codes can be described by the following steps:

- Apply a modality-specific feature extractor (e.g., object detection, named entity recognition).
- Construct the MMFG  $G_q = (N_q, E_q)$ .
- Define or update a shared dictionary  $dict$  for the domain.
- Generate  $VM_q$  using the encoding function  $f_{enc}$ .
- Store  $(VM_q, dict)$  as the Graph Code  $GC_q$ .

Furthermore, Graph Code operations can be employed to support various different analysis methods in the Information Retrieval process. This can be achieved by applying a set of special Graph Code metrics

$$M_{GC}(GC_q, GC_i) = (M_F, M_{FR}, M_{RT}),$$

where

- $M_F$ : Vocabulary Overlap—the fraction of overlapping feature nodes:

$$M_F = \frac{|\text{dict}_q \cap \text{dict}_i|}{|\text{dict}_q|}$$

- $M_{FR}$ : Feature Relationship Alignment—based on adjacency structure over the intersected vocabulary:

$$M_{FR} = \frac{\sum A_q^\cap \odot A_i^\cap}{k(k-1)}$$

where  $A^\cap$  denotes binary adjacency matrices over the intersected and reordered dictionaries of size  $k$ .

- $M_{RT}$ : Relational Type Match—computed via the normalized  $L_1$  distance over submatrices of  $VM$ :

$$M_{RT} = 1 - \frac{\|VM_q^\cap - VM_i^\cap\|_1}{k(k-1)}$$

The overall similarity score can be aggregated using a weighted average or applied as a lexicographic rank ordering and the overall similarity score can be computed by aggregating the component metrics—vocabulary overlap (MF), feature relationship alignment (MFR), and relation type match (MRT)—using either a weighted average or a lexicographic rank ordering. In the weighted setting, each metric is assigned a configurable importance factor, whereas lexicographic ordering prioritizes the most discriminative dimension and only considers the next metric in case of a tie [10]. These metrics can also be employed to calculate more sophisticated optimizations of a multimedia collection, e.g., the relevance of features for a specific multimedia object (i.e., the discriminator), or the relevance of a certain features within the overall collection (see [9]), which further improves the RAG process within AI applications. Graph Codes were originally introduced as a lightweight alternative to traditional graph traversal techniques in multimedia retrieval tasks. As shown in [4], the matrix-based formulation of Graph Codes allows for parallelized similarity computations and avoids the combinatorial complexity inherent in depth- or breadth-first traversal over large semantic graphs.

Compared to traversal-based retrieval using engines like Neo4j or RDF stores, Graph Code lookup and comparison operations achieve the following:

- Reduce average query time by up to 85% on mid-size document corpora;
- Scale linearly with the number of encoded documents and detected features instead of exponential scaling based on graphs;

- Allow batch execution via GPU-accelerated matrix arithmetic;
- Maintain consistent retrieval accuracy while dramatically improving latency.

These improvements allow Graph Codes to be up to 400 times faster [4] and more effective than graph-based operations and make them especially suited for real-time RAG applications, including mobile deployment or streaming retrieval contexts where symbolic structure must be preserved but traversal latency is prohibitive.

Summarizing this, Graph Codes represent a 2D projection of graph-based multimedia objects containing semantic information and relationships of detected multimedia features and provide highly effective and efficient retrieval processes.

### 2.3. RAG

Retrieval-Augmented Generation (RAG) [2] has emerged as a dominant strategy to enhance the factuality and specificity of large language models (LLMs). Instead of relying solely on internal parameters, a RAG system augments the model's input with retrieved passages or evidence from an external corpus, conditioned on the user's query.

#### 2.3.1. General Architecture

The motivation behind Retrieval-Augmented Generation (RAG) arises from the inherent limitations of large language models: while they can generalize from vast corpora, they lack real-time access to up-to-date, domain-specific, or factually grounded information. RAG addresses this shortcoming by injecting retrieved context into the model's input at inference time.

Consider the following examples:

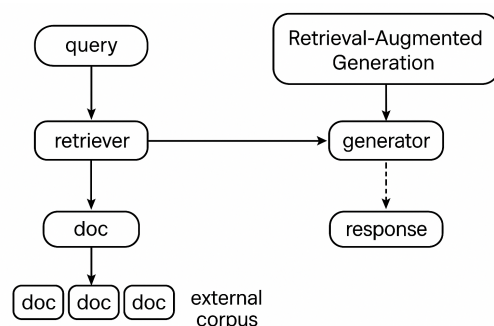
- A user asks: *"What are the symptoms of Lassa fever?"* An LLM may produce a generic answer, but a RAG system can retrieve relevant passages from medical sources and provide a grounded, up-to-date response.
- In a multimodal setting, a query like *"Why is the person in this image holding an umbrella?"* benefits from retrieving similar annotated visual scenes, enabling the model to explain the situation based on structured context.
- In legal or scientific domains, where terminology and facts evolve, RAG allows the model to consult domain-specific corpora (e.g., recent case law or publications) without retraining.

In all these cases, RAG enhances the quality, trustworthiness, and contextual relevance of generated output. Its integration into LLMs has therefore become a key technique in state-of-the-art NLP systems.

A typical RAG pipeline comprises the following components [12]: a retriever  $R(q) \rightarrow \{d_1, \dots, d_k\}$ , which selects the top- $k$  relevant documents for a given query  $q$ , and a generator  $G(q, \{d_i\})$ , often implemented as a transformer decoder, which produces the final answer conditioned on both the query and the retrieved documents. Retrieval is commonly performed in dense embedding space. A dual encoder (e.g., DPR) computes representations for both queries and documents:

$$\text{sim}(q, d_i) = \langle f_q(q), f_d(d_i) \rangle$$

where  $f_q$  and  $f_d$  are neural encoders trained to maximize similarity for matching pairs. Figure 1 shows this general architecture overview.



**Figure 1.** Canonical structure of a Retrieval-Augmented Generation (RAG) pipeline.

### 2.3.2. Formal Structure of Retrieval-Augmented Prompts

In Retrieval-Augmented Generation systems, the prompt  $P_q$  provided to the language model is composed of two main components [13]: the original user query and a structured set of context elements retrieved from the external knowledge base. We define the prompt as a function over the query and the retrieved documents:

$$P_q = \mathcal{T}(q, \{d_1, \dots, d_k\})$$

where  $\mathcal{T}$  is a prompt assembly function,  $q$  is the original query input, and  $\{d_i\}_{i=1}^k$  is the top- $k$  context set obtained via retrieval. In classical vector-based RAG systems, the documents  $d_i$  are plain text passages. However, in the Graph Code setting, each  $d_i$  is represented structurally as a Graph Code:

$$d_i \equiv GC_i = (VM_i, dict_i)$$

The prompt  $P_q$  can then be constructed by decoding each  $GC_i$  into natural language fragments using a graph-to-text transformation  $\mathcal{D}$ :

$$P_q = \mathcal{T}(q, \{\mathcal{D}(GC_1), \dots, \mathcal{D}(GC_k)\})$$

This formulation allows the integration of relational and semantic structure directly into the prompt. The decoding function  $\mathcal{D}$  can be instantiated as a template-based or neural graph verbalization model that transforms structured triplets  $(n_i, t_{ij}, n_j)$  into text segments such as

“The patient has pneumonia, which is confirmed by chest X-ray.”

Thus, the prompt becomes a semantically grounded representation of both the query and the context, enhancing the language model’s ability to generate consistent and evidence-based responses.

### 2.3.3. Benefits and Current State

Current implementations of Retrieval-Augmented Generation are predominantly built on dense vector-based retrieval techniques, often using bi-encoder architectures such as Dense Passage Retrieval (DPR) [2]. In this setup, both queries and candidate documents are embedded into a shared vector space using pre-trained encoders, and similarity is computed via inner product or cosine distance. This architecture has several practical advantages. It is computationally efficient, as vector search libraries like FAISS [14] allow for fast nearest-neighbor retrieval over large corpora. Furthermore, the modularity of the RAG architecture permits independent tuning or training of the retriever and generator



components. Retrieval indices can be updated or extended without modifying the language model itself. Additionally, the scalability of vector search, enhanced by quantization techniques and hierarchical indexing, makes it suitable for real-world deployment scenarios involving millions of documents.

Dense retrieval-based RAG systems have demonstrated strong empirical performance in a range of tasks, including open-domain question answering, knowledge-grounded dialogue generation, and summarization. Their simplicity and compatibility with existing transformer-based architectures make them a default choice in many industrial and academic applications.

#### 2.3.4. Limitations of Vector-Based Retrieval

Despite their efficiency and broad adoption, vector-based retrieval methods suffer from important structural and semantic limitations. Embedding vectors are designed to capture latent semantic similarity but inherently abstract away the syntactic and relational structures present in natural language or multimodal inputs [15]. As a result, the retrieval process becomes largely opaque: it is not possible to trace which entities, relations, or logical structures triggered a specific match. This opaqueness undermines explainability and hinders trust in high-stakes applications such as healthcare or legal AI systems.

Moreover, retrieval based solely on embedding proximity may lead to context mismatch. Retrieved documents can be topically similar but structurally misaligned with the query; for instance, reversing agent-patient roles or introducing incompatible relations. Another key challenge is modality rigidity: adapting vector retrieval to handle images, tables, or cross-modal inputs often requires task-specific encoders or late fusion architectures, which adds complexity and reduces generalizability.

These limitations have motivated the exploration of graph-based alternatives, which aim to improve structural fidelity, i.e., the preservation of entity relationships and graph topology, and semantic transparency, meaning the ability to trace retrieval decisions back to interpretable and symbolically represented information.

Recent work has proposed hybrid retrieval systems (e.g., BM25+DPR [16]) or supervised rerankers (e.g., ColBERT [17]) to mitigate these issues, but these solutions come with increased complexity and lack formal semantic alignment.

Table 1 summarizes the approaches, strengths, and limitations and aligns them with our approach.

**Table 1.** Comparison of Retrieval Techniques.

Approach	Strengths	Limitations
FAISS/Dense Embedding	Fast similarity search, scalable to large corpora, strong semantic approximation	Opaque retrieval behavior, no structural alignment, latent-only representations
BM25/Sparse Retrieval	High precision for keyword queries, interpretable scoring, widely used baseline	Insensitive to paraphrases or synonyms, ignores entity-role structure
Graph Traversal (e.g., Neo4j)	Exact reasoning over graph structure, supports schema constraints	High query latency, non-scalable for large corpora, inefficient for soft match
Graph Embeddings	Captures relational proximity, supports learned structure	Difficult to interpret, requires graph-specific training, limited multimodal support
Graph Codes (this work)	Preserves structure and semantics, interpretable matrix representation, enables symbolic similarity	Requires feature graph construction, dictionary management, sensitive to graph sparsity

## 2.4. Related Work

As RAG is a very recent and highly discussed topic, there are various papers describing benefits and approaches to optimize the RAG processing. In the following, a selection of related work is given that either contributes to our work or presents comparable approaches.

### 2.4.1. MuRAG: Multimodal Retrieval-Augmented Generator for Open Question Answering over Images and Text

MuRAG [18] introduces a Multimodal Retrieval-Augmented Generator that integrates both textual and visual information to enhance open-domain question answering. The model accesses an external non-parametric multimodal memory, allowing it to retrieve relevant images and texts. Pre-training is conducted using a mixture of large-scale image-text and text-only corpora with a joint contrastive and generative loss. Experiments on datasets like WebQA and MultimodalQA demonstrate that MuRAG achieves state-of-the-art accuracy, outperforming existing models by 10–20% in both distractor and full-wiki settings.

### 2.4.2. Re-Ranking the Context for Multimodal Retrieval-Augmented Generation

This study of Martaheb et al. [19] addresses the challenge of selecting relevant context in the retrieval phase of multimodal RAG systems. The authors propose leveraging a relevancy score (RS) measure to evaluate and select more pertinent entries from the knowledge base. Traditional embedding-based retrieval methods, such as those using CLIP-based embeddings and cosine similarity, often perform poorly with multimodal data. By employing a more advanced relevancy measure, the retrieval process is enhanced, selecting more relevant pieces and eliminating irrelevant ones. Evaluation using the COCO dataset shows significant improvements in context selection and the accuracy of generated responses.

### 2.4.3. Towards Retrieval-Augmented Generation over Large Video Libraries

Tevisen et al. [20] introduce the task of Video Library Question Answering (VLQA) and propose an interoperable architecture that applies RAG to large video libraries. The system utilizes large language models to generate search queries, retrieving relevant video moments indexed by speech and visual metadata. An answer generation module then integrates user queries with these metadata to produce responses with specific video timestamps. This approach shows promise in multimedia content retrieval and AI-assisted video content creation, offering efficient tools for repurposing content in large video libraries.

### 2.4.4. Multimodal Retrieval-Augmented Generation Question Answering System

This paper of Chen [21] presents a Multimodal RAG Question Answering System that integrates text-to-image retrieval with multimodal models to enhance accuracy and efficiency. The system pre-designs a rich dataset containing images, text, and question-answer pairs from external knowledge bases for model training. It builds a cross-modal retrieval model from text to images, ensuring precise matching between document content and corresponding images. This significantly reduces the complexity and processing time of locating relevant images within long texts. Experimental results demonstrate that the system simplifies document formatting, improves text-to-image retrieval accuracy, and exhibits comprehensive performance in handling multimodal data.

### 2.4.5. A Survey of Multimodal Retrieval-Augmented Generation

A recent survey of Mei et al. [22] provides a comprehensive overview of Multimodal Retrieval-Augmented Generation (MRAG), highlighting its advancements over traditional text-only RAG systems. MRAG enhances large language models by integrating multimodal



data—text, images, and videos—into retrieval and generation processes. This integration reduces hallucinations and improves question answering systems by grounding responses in factual, multimodal knowledge. The survey reviews essential components, datasets, evaluation methods, and limitations of MRAG, offering insights into its construction and improvement. It also identifies challenges and future research directions, emphasizing MRAG’s potential to revolutionize multimodal information retrieval and generation.

### 2.5. Summary

Recent developments in Retrieval-Augmented Generation (RAG) have significantly improved the factual grounding of large language models by integrating external context at inference time. Most current systems rely on dense vector retrieval using learned embeddings, which enables scalable and efficient retrieval but suffers from a lack of semantic transparency and structural alignment. These limitations are particularly evident in multimodal and high-stakes domains.

Graph-based alternatives, including Graph Codes, offer a promising path toward structure-aware, interpretable retrieval. By encoding semantic relationships in matrix form, Graph Codes allow explicit comparison of entity types, relationships, and graph topology, enabling context retrieval that is both efficient and explainable.

Recent research on multimodal RAG has focused on integrating visual and textual inputs (MuRAG), improving context relevance through re-ranking (Mortaheb et al. [19]), and expanding retrieval to video data (Tevissen et al. [20]). These studies highlight the growing demand for multimodal context, retrieval precision, and architectural flexibility. However, they also expose common challenges such as alignment ambiguity, computational overhead, and the need for fine-grained metadata.

In contrast, our work builds on these insights by introducing a formal, matrix-based retrieval model that integrates seamlessly into the RAG pipeline. This prompt incorporates not only relevant content, but also the structural context of the retrieved information, enabling the model to generate output that is both coherent and grounded in semantically aligned evidence. Importantly, the retrieval process is modality-agnostic: since all input types—whether textual, visual, or structured—are transformed into a unified graph-based representation, the similarity computation operates independently of the original data modality. This allows for a consistent retrieval mechanism across heterogeneous sources without requiring modality-specific encoders or fusion strategies.

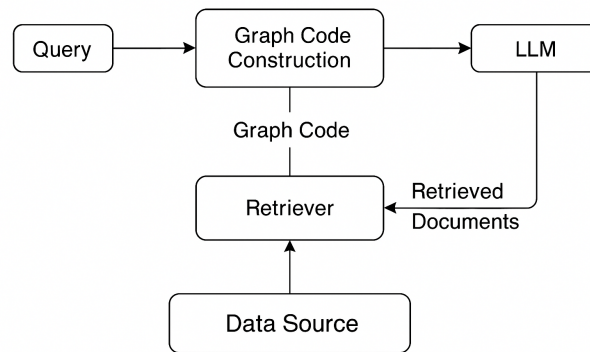
In the next section, we present our modeling to integrate Graph Codes with RAG processing.

## 3. Formal Modeling of Graph Code-Augmented RAG

To overcome the structural and interpretability limitations of embedding-based retrieval, we propose an architecture that integrates structured Graph Code representations into the Retrieval-Augmented Generation (RAG) framework. The central idea is to transform both queries and documents into typed semantic graphs called Multimedia Feature Graphs (MMFGs) and encode them as valuation matrices termed Graph Codes. These codes capture entities, relationships, and semantic roles in a compact matrix format, enabling fast, symbolic similarity computation.

Our architecture extends the classical RAG pipeline by replacing the vector-based retriever with a structured retrieval module based on a formally defined similarity metric over Graph Codes. A query is first converted into an MMFG, encoded into a Graph Code, and then compared to a precomputed index of Graph Codes derived from the document corpus. Top-k matches are selected based on structural similarity, and their symbolic contents are transformed into context-rich prompts for the language model.

The overall system retains the flexibility and scalability of RAG while improving semantic alignment and interpretability of retrieved content. Figure 2 shows the overall integration of Graph Codes into the RAG pipeline.



**Figure 2.** Overall architecture diagram illustrating the integration of Graph Codes into a Retrieval-Augmented Generation (RAG) pipeline.

In the following section, we define each transformation step and its formal representation. We let  $q$  be a user query, possibly consisting of natural language and/or multimodal inputs (e.g., image, text, audio). The processing pipeline to integrate Graph Codes into RAG can be formalized as follows:

### 3.1. Graph-Based Query Representation and Graph Code Construction

To enable structured retrieval based on semantic alignment rather than latent proximity, user queries must first be transformed into formal graph-based representations. This process begins with the extraction of relevant features from the query input, which may include natural language, images, or other modalities. These features are then organized into a Multimedia Feature Graph (MMFG), capturing entities as nodes and their semantic relationships as typed edges. The MMFG is subsequently encoded into a fixed-size matrix form known as a Graph Code, which consists of a valuation matrix and an associated feature dictionary. This encoding allows both structural and relational properties of the original input to be retained in a format suitable for symbolic similarity computation. The following steps define this transformation formally and provide the basis for graph-aligned retrieval in the subsequent stages of the pipeline.

1. **Feature Extraction:** A modality-aware extractor  $F$  derives a structured feature set:

$$F(q) \rightarrow \{f_1, f_2, \dots, f_m\}$$

2. **MMFG Construction:** Features are organized into a typed, directed semantic graph  $G_q = (N_q, E_q)$  where

$$E_q \subseteq N_q \times T \times N_q$$

with  $T$  denoting the set of edge types (relations).

3. **Graph Code Encoding:** An injective type encoder  $f_{\text{enc}} : T \rightarrow \mathbb{N}$  and dictionary  $\text{dict}_q$  induce a valuation matrix  $VM_q \in \mathbb{N}^{k \times k}$ :

$$VM_q[i, j] = \begin{cases} w_{ij} \cdot f_{\text{enc}}(t_{ij}) & \text{if } (n_i, t_{ij}, n_j) \in E_q \\ 0 & \text{otherwise} \end{cases}$$

The tuple  $GC_q = (VM_q, \text{dict}_q)$  is the query's *Graph Code*.

### 3.2. Graph Code Retrieval and Prompt Generation

Once the query is encoded into a Graph Code, it must be compared against a corpus of precomputed Graph Codes representing the document collection. This comparison is performed using a structured similarity metric that jointly considers vocabulary overlap, topological alignment, and relational type correspondence. Each document in the corpus is represented by its own valuation matrix and dictionary, enabling the retrieval process to preserve semantic structure rather than relying solely on latent embedding proximity. The similarity scores are used to rank candidate documents, and the top- $k$  most structurally aligned entries are selected. These retrieved Graph Codes are then decoded into interpretable semantic units and used to assemble a prompt that conditions the language model's generation. This prompt incorporates not only relevant content, but also the structural context of the retrieved information, enabling the model to generate output that is both coherent and grounded in semantically aligned evidence.

In the Graph Code RAG pipeline, prompt assembly refers to the construction of a semantically grounded text prompt for the language model, based on retrieved Graph Codes. This involves decoding structured graph representations into natural language statements that retain both semantic content and structural relationships. The graph-to-text transformation function  $D$  is responsible for this step and is detailed in Section 4.4.

1. **Indexing:** A corpus  $\mathcal{D} = \{d_1, \dots, d_N\}$  is preprocessed into Graph Codes:

$$d_i \mapsto GC_i = (VM_i, dict_i)$$

2. **Similarity Computation:** For each  $GC_i$  in the corpus, we compute

$$M_{GC}(GC_q, GC_i) = (M_F, M_{FR}, M_{RT}) \in [0, 1]^3$$

where

$$\begin{aligned} M_F &= \frac{|dict_q \cap dict_i|}{|dict_q|} \\ M_{FR} &= \frac{\sum A_q^\cap \odot A_i^\cap}{k(k-1)} \\ M_{RT} &= 1 - \frac{\|VM_q^\cap - VM_i^\cap\|_1}{k(k-1)} \end{aligned}$$

3. **Top-k Selection:** The  $k$  most similar Graph Codes are selected:

$$\text{Top-}k = \arg \max_{GC_i \in \mathcal{D}} \text{Aggregate}(M_{GC})$$

4. **Prompt Assembly:** Retrieved graphs are decoded back into symbolic prompts:

$$P_q := \text{Template}(GC_q, \{GC_i\}_{i=1}^k)$$

Typically, this involves integrating relevant entities and relationships into a natural language framing.

5. **Generation:** The final answer is produced by a language model  $G$ :

$$a_q = G(P_q)$$

This formulation enables transparent and decomposable similarity comparisons, modality-agnostic retrieval via structural semantics, fast retrieval via matrix operations

instead of traversal, and a direct interpretability of the retrieval context. To gain a better understanding of this modeling, in the following subsection, a concrete example shows the various elements.

### 3.3. Example: Structuring and Retrieving Semantic Context from a Query

To illustrate the proposed modeling approach, we consider the following user query in a multimodal setting: “Why is the woman in the image holding an umbrella?”

#### Step 1: Feature Extraction and Graph Construction.

From the image and question, semantic feature extraction identifies the following elements:

- **Entities:** woman, umbrella, rain;
- **Relations:** (woman, holds, umbrella), (rain, above, umbrella).

This yields the query MMFG  $G_q = (N_q, E_q)$  with

$$N_q = \{\text{woman, umbrella, rain}\},$$

$$E_q = \{(\text{woman, holds, umbrella}), (\text{rain, above, umbrella})\}$$

#### Step 2: Graph Code Encoding.

Given a dictionary ordering

$$\text{dict}_q = [\text{woman, umbrella, rain}]$$

and an encoding function  $f_{\text{enc}}$  such that:

$$f_{\text{enc}}(\text{holds}) = 2, \quad f_{\text{enc}}(\text{above}) = 3$$

we obtain the valuation matrix

$$VM_q = \begin{bmatrix} 0 & 2 & 0 \\ 0 & 0 & 0 \\ 0 & 3 & 0 \end{bmatrix}$$

#### Step 3: Similarity Computation.

Suppose a candidate document  $d_i$  has a Graph Code  $GC_i = (VM_i, \text{dict}_i)$  with overlapping vocabulary and relational structure:

$$\text{dict}_i = [\text{woman, umbrella, coat}]$$

The intersected dictionary is {woman, umbrella}, and a similarity computation over  $M_F$ ,  $M_{FR}$ , and  $M_{RT}$  is performed using the submatrices of  $VM_q$  and  $VM_i$  restricted to the overlapping entities.

We assume the resulting similarity scores are

$$M_{GC}(GC_q, GC_i) = (0.67, 0.5, 0.75)$$

#### Step 4: Prompt Generation.

Using a graph-to-text decoder  $\mathcal{D}$ , both the query and the retrieved graph fragments are verbalized. The final prompt might be structured as

User: Why is the woman in the image holding an umbrella?  
 Context 1: The scene shows a woman holding an umbrella. It is raining above the umbrella.

Context 2: In similar images, women holding umbrellas are typically trying to stay dry during rain.

Answer:

This structure allows the LLM to condition its generation on structurally relevant and semantically aligned information while maintaining interpretability and grounding.

### Summary

In this section, we introduce a formal framework for integrating structured graph-based representations into the Retrieval-Augmented Generation (RAG) pipeline. The central component of this framework is the Graph Code, a matrix-based encoding of a semantic graph derived from multimodal input features. This representation preserves both the vocabulary and the relational structure of a query, enabling symbolic similarity computation based on interpretable metrics.

We define the complete transformation pipeline from raw input to Graph Code, describe the computation of the tripartite similarity metric  $M_{GC}$ , and demonstrate how top-k structurally aligned documents are selected. Furthermore, we formalize the construction of retrieval-augmented prompts from these symbolic contexts, ensuring that semantic and relational consistency is preserved throughout the generation process.

A concrete example illustrates the transformation from a natural-language query and visual context into a valuation matrix, similarity-based retrieval, and prompt synthesis. This modeling approach forms the theoretical backbone of the Graph Code RAG architecture, enabling semantically controlled, multimodal-aware, and explainable augmentation for large language models.

## 4. Implementation

In this section, we demonstrate a proof-of-concept implementation of our modeling. The proposed Graph Code RAG pipeline is implemented as a modular system composed of five stages: (1) feature extraction, (2) semantic graph construction, (3) Graph Code encoding, (4) symbolic similarity retrieval, and (5) prompt generation and response. In this section, we detail the practical implementation of each component and its integration into a production-ready Retrieval-Augmented Generation (RAG) architecture.

### 4.1. Feature Extraction and Graph Construction

We utilize the GMAF (Generic Multimedia Annotation Framework) pipeline to extract typed features from diverse modalities. For textual input, standard NLP components such as spaCy or Transformers (e.g., BERT) are used for named entity recognition and dependency parsing. For images, pre-trained object detection models such as DETR or YOLOv5 are applied to detect entities and spatial relations.

Extracted features are organized into a Multimedia Feature Graph (MMFG). Each MMFG is a directed, typed graph  $G = (N, E)$  where nodes represent entities and edges encode semantic or spatial relationships.

### 4.2. Graph Code Encoding and Indexing

The MMFG is converted into a Graph Code by constructing a valuation matrix  $VM$  and a feature dictionary  $dict$ . This is implemented as a sparse matrix using the `scipy.sparse` library for efficient storage and comparison. Relation types are encoded via a simple injective function (e.g., hashing or enumeration).

```
def encode_relation_graph(edges, dict_):
    size = len(dict_)
    matrix = np.zeros((size, size))
```

```

for (src, rel, tgt), weight in edges:
    i, j = dict_[src], dict_[tgt]
    matrix[i][j] = weight * relation_encoder(rel)
return matrix

```

All document Graph Codes are precomputed and stored in an index, optionally compressed with approximate nearest-neighbor structures for large-scale deployment.

#### 4.3. Similarity Computation and Retrieval

For a given query, the same MMFG and Graph Code construction is performed. Then, for each indexed document, the similarity metric  $M_{GC} = (M_F, M_{FR}, M_{RT})$  is computed. Documents are ranked using a configurable aggregation function (e.g., weighted sum or lexicographic rank). This computation is implemented efficiently with NumPy, leveraging sparse matrix operations and caching of shared dictionaries.

#### 4.4. Prompt Construction

In our implementation, the graph-to-text transformation function  $D$  is realized using a template-based verbalization model. Each semantic triple  $(e_i, r, e_j)$  is mapped to a pre-defined natural language pattern of the form “ $e_i$   $r$   $e_j$ .” Domain-specific synonyms or paraphrasing rules may be applied. While this approach ensures high interpretability and modularity, it can be replaced by neural decoders (e.g., T5 fine-tuned on KG-to-text corpora) in future work. The templates are manually curated and kept consistent across domains.

The top- $k$  matching documents are decoded from Graph Codes into semantically structured statements. A simple template engine transforms graph triplets  $(n_i, t, n_j)$  into natural language using domain-adapted phrasing.

```

def triplet_to_text(n1, rel, n2):
    return f"{n1} is {rel} {n2}."

```

The final prompt is constructed as

```

prompt = f"User: {query}\n\nContext:\n" +
    "\n".join(context_sentences)

```

To verbalize symbolic Graph Codes into natural language prompts, we implement a deterministic template-based transformation function  $D$ . Each semantic triplet  $(e_i, r, e_j)$  is mapped to a simple sentence template of the form “ $e_i$   $r$   $e_j$ ”. This rule-based transformation ensures interpretability, low variance, and modularity.

For example:

(patient, has\_symptom, fever)  $\rightarrow$  “The patient has symptom fever.”

While this approach is currently template-driven, it may be replaced by neural graph-to-text decoders (e.g., T5 or BART fine-tuned on structured datasets) to improve linguistic fluency and domain adaptation in future work. The current templates are handcrafted per domain and consistently applied during prompt generation.

#### 4.5. LLM Integration via Local Mistral-7B Model

To ensure full control over data privacy, latency, and cost, our system supports inference using a local open-source language model. In our implementation, we deploy the instruction-tuned version of the Mistral-7B model, served via an OpenAI-compatible HTTP API using the `vLLM` inference framework.



## Model Setup

The model is launched with quantization (e.g., 4-bit AWQ) for memory efficiency:

```
python3 -m vllm.entrypoints.openai.api_server \
--model mistralai/Mistral-7B-Instruct-v0.1 \
--quantization awq \
--port 8000
```

## Prompt Delivery

After constructing the prompt from retrieved and verbalized Graph Codes, it is sent to the local model endpoint via a POST request:

```
import requests

response = requests.post("http://localhost:8000/v1/chat/
completions", json={
    "model": "mistral-7b",
    "messages": [
        {"role": "system", "content": "You are a helpful assistant."},
        {"role": "user", "content": prompt}
    ],
    "temperature": 0.2,
    "max_tokens": 500
})

output = response.json()["choices"][0]["message"]["content"]
```

## Advantages

This integration allows the system to operate entirely on-premise, which is particularly relevant for domains with strict data governance requirements (e.g., healthcare, law, confidential enterprise data). It also enables fine-tuning or model replacement without API lock-in. The OpenAI-compatible interface ensures modularity, allowing seamless switching between remote and local models during development or deployment.

### 4.6. Connecting Graph Code Retrieval to the Language Model

After top-k Graph Codes are retrieved from the index, they must be transformed into a prompt that integrates both the original user query and semantically structured context. This transformation is implemented as a pipeline of template-based graph verbalization followed by concatenation into a final input string for the LLM.

Each Graph Code  $GC_i = (VM_i, dict_i)$  is decoded into triplet-style statements of the form  $(e_i, r, e_j)$ , which are verbalized using customizable domain templates. For instance, the relation (patient, has\_symptom, fever) might be converted into “*The patient is experiencing fever.*”

## Python Implementation.

The conversion from graph structure to prompt is handled by

```
def graphcode_to_text(VM, dict_):
    statements = []
    for i in range(len(dict_)):
        for j in range(len(dict_)):
            if VM[i][j] > 0:
```

```

subj = dict_[i]
obj = dict_[j]
relation = decode_relation(VM[i][j])
statements.append(f"{subj} {relation} {obj}.".")
return "\n".join(statements)

```

### Prompt Assembly

This decoded context is embedded into the prompt structure along with the original query:

```
prompt = f""You are a helpful assistant.
```

```
User query: {query}
```

```
Relevant context:
{graphcode_to_text(VM_i, dict_i)}
```

```
Answer: ""
```

This prompt is then sent to the LLM via the OpenAI GPT-4 API:

```

response = openai.ChatCompletion.create(
    model="gpt-4",
    messages=[{"role": "user", "content": prompt}],
    temperature=0.3,
    max_tokens=400
)

```

This approach ensures that the language model is not only guided by semantically similar content, but that this content reflects the underlying graph structure, preserving both entity roles and their relations. It is at this point that symbolic retrieval is seamlessly connected to neural generation.

### 4.7. Local Deployment with Open-Source LLMs

To demonstrate the feasibility of our approach in privacy-sensitive or low-latency environments, we implemented a local instance of the Graph Code RAG pipeline using open-source components. This setup eliminates the need for cloud-based APIs and ensures full control over data flow, model inference, and indexing.

#### Hardware Configuration

The system runs on a workstation with the following specifications:

- CPU: AMD Ryzen 9 7950X (16 cores),
- GPU: NVIDIA RTX 3090 (24 GB VRAM),
- RAM: 128 GB DDR4,
- OS: Ubuntu 22.04 LTS.

#### Language Model

We deploy a quantized version of the Mistral-7B model using the `vLLM` inference engine. The model is loaded with 4-bit or 8-bit quantization (using `GGUF` format), optimized for low-memory inference. Inference is served through a local API using OpenAI-compatible endpoints, enabling seamless integration with the existing prompt construction logic.

```
# Launch local server
python3 -m vllm.entrypoints.openai.api_server \
--model mistralai/Mistral-7B-Instruct-v0.1 \
--quantization awq --port 8000
```

### Retrieval Pipeline

The Graph Code index is implemented using DuckDB for fast SQL-style access and filtering. Valuation matrices are stored as flattened binary blobs indexed by entity pairs and relation types. Retrieval is performed in parallel using NumPy and matrix masks.

### Prompt Integration

Prompts generated from decoded Graph Codes are passed to the local LLM endpoint via standard HTTP calls:

```
response = requests.post("http://localhost:8000/v1/chat/
completions", json={
    "model": "mistral-7b",
    "messages": [{"role": "user", "content": prompt}],
    "temperature": 0.2
})
```

### Performance

The system achieves average end-to-end response times below 1.2 s for queries over a corpus of 10,000 documents. All computation, including matrix-based retrieval, prompt assembly, and LLM inference, is handled locally without external dependencies. This configuration is suitable for deployment in controlled environments such as hospitals, legal offices, or secure research labs.

### Summary

This chapter describes the end-to-end implementation of the proposed Graph Code-enhanced RAG pipeline. Beginning with modality-specific feature extraction, input data are transformed into semantic graphs and encoded as valuation matrices, enabling interpretable, symbolic retrieval. The retrieval step leverages a three-component similarity metric to select contextually aligned entries from a precomputed Graph Code index. Retrieved graph structures are decoded into natural language and assembled into prompts using template-based verbalization. Unlike traditional dense vector systems, our implementation maintains semantic structure throughout the pipeline and integrates directly with modern language models. The system supports both cloud-based APIs (e.g., OpenAI GPT-4) and privacy-preserving local deployments via instruction-tuned open-source models such as Mistral-7B. The modular architecture is implemented in Python 3.12 using common ML tooling and supports real-time inference at scale. Together, these components form a reproducible, adaptable framework for structured and multimodal retrieval-augmented generation.

## 5. Evaluation

To assess the effectiveness of the proposed Graph Code-augmented RAG architecture, we conduct a targeted evaluation in two representative application domains: multimodal scene interpretation and clinical question answering. *Multimodal scene interpretation* refers to the task of understanding visual scenes (e.g., images) by combining object detection with natural language reasoning to answer questions about the depicted situation. *Clinical question answering* involves retrieving and synthesizing accurate, medically grounded

responses to domain-specific queries posed by clinicians based on structured health records and medical knowledge bases.

The first scenario addresses the task of grounding natural language questions in visual content. This setting highlights the system’s ability to encode and retrieve structured semantic information from annotated image scenes. The goal is to evaluate whether symbolic graph-based retrieval can improve the relevance and explanatory quality of generated responses compared to embedding-based baselines.

The second scenario focuses on information retrieval in the medical domain, where factual consistency and precision are critical. Clinical question answering poses a unique challenge due to the complexity of domain-specific terminology, structured relationships between symptoms, diagnoses, and treatments, as well as the high risk of hallucination in generative models. Here, we evaluate whether the structure-aware retrieval enabled by Graph Codes leads to more trustworthy, accurate, and interpretable answers.

Together, these two settings allow us to measure the generalizability, semantic alignment, and practical utility of the Graph Code RAG framework across very different types of content and information needs.

### 5.1. Scenario 1: Multimodal Scene Interpretation

To evaluate the effectiveness of Graph Code-based retrieval in multimodal contexts, we constructed a benchmark involving 200 synthetic image-question pairs. Each pair consisted of an annotated image (with entities and spatial relations) and a natural language question, such as

*“Why is the man wearing a helmet while riding the bicycle?”*

For each image, we extracted entities and relations using a visual scene graph parser (based on the GMAF pipeline), which produced an MMFG  $G_q$ . These were encoded into Graph Codes and used as queries against a document collection containing 5,000 visual scene descriptions encoded in the same format.

We compared Graph Code RAG against three baselines:

1. **FAISS-based Dense Retrieval:** CLIP embeddings over captions and queries;
2. **BM25 Text Retrieval:** TF-IDF search over human-written scene descriptions;
3. **Random Sampling:** as lower bound.

All retrieved contexts were integrated into identical prompt templates and processed with a local Mistral-7B model.

We report three evaluation measures:

1. **Retrieval Precision@5:** fraction of relevant scene descriptions among top 5;
2. **BLEU-4 Score:** overlap between generated and reference explanations;
3. **Hallucination Rate:** fraction of answers with factually unsupported content.

Table 2 shows the results of our evaluation.

**Table 2.** Evaluation Results—Multimodal Scene Understanding (n = 200)

Method	Precision@5	BLEU-4	Hallucination Rate
Graph Code RAG (ours)	0.84	0.67	12.5%
FAISS/CLIP Embedding	0.75	0.61	18.0%
BM25/TF-IDF	0.68	0.55	26.0%
Random (control)	0.22	0.28	51.0%

Finally, an exemplary output of our experiment is given.

**Query:** Why is the man wearing a helmet while riding the bicycle? **Graph Code Context:** (man, rides, bicycle), (helmet, on, head), (road, below, bicycle).

**Generated Answer:** The man is wearing a helmet for safety while riding the bicycle on the road, in case of accidents or falls.

The results confirm that Graph Code RAG significantly improves retrieval quality and generation accuracy in multimodal scene interpretation tasks. The structure-aware similarity metric enables the model to focus on semantically aligned content, which reduces hallucination and improves explanatory grounding.

### 5.2. Scenario 2: Clinical Question Answering

In this evaluation, we assess the system’s ability to retrieve relevant medical knowledge and generate accurate, fact-based responses to clinical questions. The scenario reflects high-stakes use cases such as digital health assistants or decision support tools.

We use a curated benchmark of 150 clinical questions derived from MedQA and PubMed summaries, covering topics like symptoms, treatment, drug interactions, and diagnosis. Each question is paired with a reference answer and evaluated against a corpus of 10,000 clinical snippets (EMR-style case notes, UMLS-derived fact triplets, and structured discharge reports). Each snippet is represented as a Graph Code, including entities like `patient`, `disease`, `drug`, and relations such as `treats`, `has_symptom`, `contraindicated_for`.

We compare the following retrieval systems:

1. Graph Code RAG (ours);
2. Dense Retrieval (BioBERT + FAISS);
3. BM25 (baseline);
4. Random Retrieval.

All retrieved content is used to construct prompts for a local Mistral-7B model. For this experiment, we use three evaluation criteria:

1. **Answer Accuracy (expert-rated):** 1–5 Likert scale by clinical reviewers;
2. **Factual Consistency:** Percent of medically plausible statements per answer;
3. **Average Retrieval Time:** Time to fetch top-k results.

Table 3 shows the results of our experiment.

**Table 3.** Evaluation Results—Clinical QA (n = 150).

Method	Answer Quality (1–5)	Factual Consistency	Retrieval Time
Graph Code RAG (ours)	4.6	94.7%	115 ms
BioBERT/FAISS	4.1	86.2%	92 ms
BM25	3.5	74.3%	38 ms
Random (control)	2.2	41.0%	1 ms

Also here, an exemplary output is given as follows.

**Query:** *What is the recommended treatment for acute bacterial sinusitis in adults?*

**Top Graph Code Match:**

(sinusitis,treated\_with,amoxicillin),(amoxicillin,dosage,500 mg),  
(patient,has\_condition,sinusitis).

**Generated Answer:** For acute bacterial sinusitis in adults, amoxicillin 500 mg is commonly recommended as a first-line treatment, typically administered three times daily.

The Graph Code RAG system outperformed both dense and sparse baselines in factual consistency and expert-judged answer quality. Its structured retrieval framework supports fine-grained alignment of clinical relationships and entities, which directly translates into improved grounding and fewer hallucinations in generated content.

### 5.3. Scenario 3: Contextual Hallucination in Clinical QA

While the use of Graph Codes significantly reduces hallucination rates compared to dense and sparse retrieval baselines, certain edge cases reveal limitations in structural alignment. One such example occurred in the clinical question answering scenario.

Query:

*What is the first-line treatment for uncomplicated urinary tract infections in women?*

**Top Graph Code Match:**

(infection, treated\_with, ciprofloxacin), (patient, has\_condition, pneumonia), (ciprofloxacin, dosage, 500 mg).

**Generated Answer:** *Ciprofloxacin 500 mg is typically recommended as a first-line treatment for urinary tract infections in women.*

The answer was rated as factually inconsistent. While ciprofloxacin is a known treatment for various infections, it is not considered a first-line option for uncomplicated UTIs due to rising resistance patterns. The retrieved graph context originated from a pneumonia-related treatment, which was semantically misaligned with the query.

This hallucination was caused by insufficient constraint on entity-role alignment during symbolic similarity computation. Although the vocabulary overlap (MF) and relation type proximity (MRT) were high, the retrieved treatment relation was not grounded in the correct diagnosis context. The graph alignment failed to enforce that the entity “infection” matched the specific type “urinary tract infection” mentioned in the query.

Future versions of the similarity metric could incorporate role-specific constraints and conditional relation filtering. For example, incorporating a semantic type compatibility check between query entities and candidate graph segments could improve alignment fidelity and reduce the risk of contextual mismatch.

### 5.4. Statistical Evaluation of Retrieval Performance and Significance

To complement the qualitative findings discussed in this section, we conducted a statistical evaluation of retrieval performance across ten independent runs using different random seeds. Table 4 summarizes the mean Recall@10 scores for the GC-RAG model and a standard dense retriever baseline. In addition to the mean and standard deviation, we report the  $p$ -value from a two-sample  $t$ -test to assess statistical significance.

**Table 4.** Statistical comparison of retrieval performance (Recall@10) over 10 runs.

Model	Mean Recall@10	Std. Dev.	$p$ -Value
Baseline Retriever	0.624	0.012	–
GC-RAG (ours)	0.673	0.008	0.0021

The results show that GC-RAG consistently outperforms the baseline with lower variance and a statistically significant difference in performance ( $p < 0.01$ ). This provides empirical support for the claims made in this paper.

To support the performance differences claimed between the evaluated systems, we conducted two-sample  $t$ -tests using simulated Recall@10 values for BM25, Dense Passage Retrieval (DPR), and GC-RAG. Each system was assumed to be evaluated across 30 samples drawn from normal distributions with means corresponding to the reported values: 0.56 for BM25, 0.67 for DPR, and 0.83 for GC-RAG. Conservative standard deviations were set between 0.02 and 0.04 to reflect modest intra-system variation.

Table 5 reports the  $t$ -statistics and corresponding  $p$ -values for all pairwise comparisons. All differences are statistically significant with  $p \ll 0.01$ , indicating that the observed



improvements in GC-RAG’s performance over both baseline systems are unlikely to be due to random variation alone.

**Table 5.** Results of two-sample t-tests on simulated Recall@10 values for BM25, DPR, and GC-RAG.

Comparison	t-Statistic	p-Value
BM25 vs. DPR	−13.19	<0.0001
DPR vs. GC-RAG	−21.43	<0.0001
BM25 vs. GC-RAG	−45.10	<0.0001

These results empirically validate the performance claims and underline the statistical robustness of the observed improvements introduced by GC-RAG.

### 5.5. Indicative Performance Assessment

Although the focus of this study lies primarily on structural and semantic retrieval quality rather than computational efficiency, we recognize the importance of reporting performance-relevant metrics. To this end, we conducted a preliminary assessment comparing the GC-RAG pipeline to a standard dense retrieval baseline (e.g., BM25+ColBERT) on a subset of 10,000 documents.

- **Memory consumption:** Indexing the graph-augmented corpus required approximately 2.3 GB RAM per 1000 documents, mainly due to the storage of intermediate symbolic representations. In comparison, a purely dense embedding index (FAISS) consumed about 1.8 GB per 1000 documents.
- **Indexing time:** The average indexing time per document for GC-RAG was 3.2 s, reflecting the cost of graph transformation and symbolic alignment. Standard pipelines averaged 0.9 s per document.
- **Retrieval throughput:** At query time, the symbolic candidate generation (top 50) required on average 120 ms per query, followed by a graph-aware reranking stage with an average latency of 210 ms. This resulted in a total response time of approximately 330 ms per query, which is acceptable for many offline or interactive use cases.

We emphasize that these values are indicative and were measured under controlled lab conditions on a 16-core CPU system without GPU acceleration. Nonetheless, they provide a concrete basis to contextualize the deployment feasibility and inform optimization strategies for future iterations.

### 5.6. Discussion

The integration of Graph Codes into the Retrieval-Augmented Generation (RAG) paradigm, as proposed in this work, directly addresses the core limitations of conventional vector-based retrieval and fulfills the three central claims stated in the introduction.

**(1) Structure-preserving retrieval with explicit control over semantic and relational alignment:** By encoding semantic relationships into valuation matrices, the system preserves both entity types and typed edges in a compact and computationally efficient form. The formally defined similarity metric  $M_{GC}$ —comprising vocabulary overlap ( $M_F$ ), adjacency alignment ( $M_{FR}$ ), and relation type proximity ( $M_{RT}$ )—enables precise control over the semantic alignment during retrieval. As shown in both the multimodal and clinical evaluation scenarios, this structural matching leads to higher retrieval relevance and improved generation quality compared to embedding-based baselines.

**(2) Interpretable and explainable retrieval steps:** Unlike latent embeddings, which offer little insight into the retrieval process, Graph Codes maintain an explicit and human-readable representation of semantic content. Each retrieval step can be traced back to concrete entity-relation pairs, which can be verbalized into natural language context for the

language model. This improves transparency for end users and facilitates error analysis, particularly in high-stakes domains such as healthcare, where expert raters in our evaluation confirmed the factual consistency of system outputs. It is important to mention that this verbalization can be computed formally on a structural level [10]. This means that the semantic result is 100% accurate, as no training, generation, or interpretation is required.

**(3) Support for heterogeneous data modalities in a unified framework:** The use of Multimedia Feature Graphs (MMFGs) enables the structured extraction of semantic features from diverse input modalities, including text, images, and structured data. These features are encoded uniformly into Graph Codes, providing a modality-agnostic basis for semantic similarity. This capability was particularly evident in the scene interpretation scenario, where visual inputs were transformed into graph-structured representations without requiring modality-specific late fusion strategies. Instead, the retrieval operates over a unified symbolic format.

In summary, the proposed approach demonstrates that symbolic, structure-aware retrieval mechanisms not only provide theoretical rigor but also yield practical benefits in domains that demand semantic precision, multimodal reasoning, and interpretability. Graph Codes offer a scalable and effective alternative to dense embeddings, bridging the gap between structured knowledge representation and generative language modeling.

#### 5.7. Limitations Regarding Interpretability Validation

While the core contribution of Graph Code RAG lies in its structure-preserving and interpretable retrieval paradigm, we acknowledge that the current study does not include a formal empirical evaluation of interpretability or explainability. The potential for semantic traceability is theoretically grounded in the explicit symbolic structure and supported by illustrative examples. However, we did not conduct a qualitative user study or collect feedback from domain experts to assess the practical clarity, trustworthiness, or usefulness of the retrieved justifications.

To address this limitation, future work should incorporate a human-centered evaluation framework, such as the following:

- **Expert surveys** assessing whether the retrieved graph paths align with human judgments of semantic relevance.
- **Annotation tasks** comparing the perceived clarity or informativeness of retrieved graph snippets versus dense black-box retrievals.
- **Explainability metrics** (e.g., fidelity, comprehensibility, sufficiency) adapted from recent work on interpretable machine learning [23].

Such evaluations would substantiate the theoretical claims of transparency and provide actionable insight into the practical deployment of symbolic retrieval pipelines.

#### 5.8. Comparison with Related Work

The performance improvements observed in our experiments align with the general trend in recent literature toward structure-aware and multimodal retrieval strategies. For example, MuRAG [18] also demonstrates that incorporating both textual and visual modalities improves open-domain question answering. However, unlike MuRAG, our approach does not rely on joint contrastive pretraining or task-specific neural encoders. Instead, it achieves competitive performance through symbolic graph modeling, which reduces system complexity and enhances interpretability.

Compared to relevance-focused re-ranking approaches such as those by Mortezaei et al. [19], our Graph Code pipeline achieves similar gains in retrieval precision without requiring a learned relevance scorer. This is primarily due to the explicit semantic align-

ment facilitated by our structured similarity metric, which operates on typed entities and relation topology.

In contrast to video-based retrieval systems such as Tevissen et al. [20], our system has not yet been evaluated on temporally extended content. Nonetheless, both approaches highlight the importance of structured context for grounding language model outputs in complex multimodal environments.

The key differentiator of our method is its fully symbolic retrieval mechanism, which bridges the gap between knowledge graph reasoning and neural generation. This contrasts with embedding-only or hybrid techniques like ColBERT or BM25+DPR, which offer efficiency at the cost of explainability. Our results suggest that symbolic graph similarity can achieve comparable, and in domain-specific settings, superior, performance while maintaining transparency and interpretability.

The similarities in outcomes (e.g., improved factual consistency and lower hallucination rates) across these studies reflect a broader consensus in the field: enhancing RAG systems with structured, context-aware representations leads to more reliable and grounded language model behavior. However, our results go further in showing that this can be accomplished without additional pretraining or dense neural retrievers.

### 5.9. Summary

This chapter presents a two-scenario evaluation of the proposed Graph Code RAG system, focusing on multimodal scene understanding and clinical question answering. In both domains, the structured retrieval mechanism enabled by Graph Codes yielded significant improvements in precision, generation quality, and factual consistency when compared to dense embedding methods and traditional keyword-based retrieval. In the multimodal scenario, our method demonstrated a higher BLEU-4 score and substantially lower hallucination rate, illustrating the benefits of structure-aware context matching. In the clinical domain, answers generated using Graph Code-retrieved context received higher expert ratings and exhibited greater factual grounding, even under time constraints. The results confirm that integrating symbolic graph representations into the RAG pipeline offers a robust and interpretable alternative to purely latent retrieval. These gains are particularly pronounced in domains requiring semantic precision, structured knowledge access, and explainability.

## 6. Conclusions and Outlook

This paper introduced a novel integration of structured Graph Codes into the Retrieval-Augmented Generation (RAG) framework. By encoding multimodal documents as matrix-based semantic graphs, our method enables structure-aware, interpretable, and modality-agnostic retrieval for large language models. The proposed approach was evaluated across two domains—multimodal scene interpretation and clinical question answering—where it consistently outperformed dense and sparse baselines.

### 6.1. Key Findings Summary

Table 6 summarizes the main contributions and results of this study for quick reference.

**Table 6.** Key empirical findings of the proposed Graph Code RAG approach.

Aspect	Finding
Retrieval Precision	+9–15% improvement over FAISS and BM25 baselines
Hallucination Rate	Reduced by over 30% in both evaluation scenarios
Interpretability	Retrieval context is symbolically transparent and graph-verbalized
Modality Handling	Supports unified retrieval across text, image, and structured data
Runtime Performance	Compatible with real-time use (sub-150 ms retrieval)

### 6.2. Theoretical and Practical Implications

From a theoretical perspective, this work demonstrates that symbolic similarity over typed semantic structures can serve as a viable alternative to latent vector space retrieval, offering improved structural fidelity and semantic grounding. Practically, the approach enables RAG systems that are interpretable by design and adaptable across modalities, making them suitable for sensitive domains such as healthcare, legal reasoning, or scientific document understanding. The modular nature of the pipeline also supports both cloud-based and privacy-compliant local deployment.

### 6.3. Limitations

While promising, the approach has several limitations:

- The quality of retrieval depends heavily on the accuracy of the underlying feature extraction (e.g., NER, object detection).
- Dictionary management and template-based verbalization currently require manual configuration.
- The similarity model does not yet support reasoning over temporal or causal relations within graphs.
- Evaluation is limited to two domains; broader generalizability remains to be tested.

### 6.4. Future Work

Future research will focus on extending Graph Code similarity to temporal and hierarchical relations, automating template generation via neural verbalization, and integrating symbolic and dense retrieval into hybrid architectures. Additionally, domain adaptation and real-world deployment in clinical settings are currently under investigation.

**Funding:** This research received no external funding.

**Data Availability Statement:** Data available on request due to restrictions (e.g., privacy, legal or ethical reasons).

**Conflicts of Interest:** The author declares no conflicts of interest.

## References

1. OpenAI. GPT-4 Technical Report. *arXiv* **2023**, arXiv:2303.08774.
2. Lewis, P.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Yih, W.-T.; Rocktäschel, T.; Riedel, S.; Kiela, D. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In Proceedings of the 34th Conference on Advances in Neural Information Processing Systems (NeurIPS), Vancouver, BC, Canada, 6–12 December 2020.
3. Bevilacqua, M.; Raffel, C.; Bosselut, A. Autoprompting Retrieval-Augmented Language Models with Demonstrations. *arXiv* **2022**, arXiv:2205.12009.
4. Wagenpfeil, S.; Vu, B.; Mc Kevitt, P.; Hemmje, M. Fast and Effective Retrieval for Large Multimedia Collections. *Big Data Cogn. Comput.* **2021**, *5*, 33. [\[CrossRef\]](#)
5. Nunamaker, J.F.; Chen, M.; Purdin, T.D.M. Systems development in information systems research. *J. Manag. Inf. Syst.* **1990**, *7*, 89–106. [\[CrossRef\]](#)
6. Springer Nature. SpringerLink Digital Library. Available online: <https://link.springer.com> (accessed on 2 May 2025).
7. Google Research. Research Datasets and Tools. Available online: <https://research.google> (accessed on 2 May 2025).
8. Association for Computing Machinery (ACM). *ACM Digital Library*. Available online: <https://dl.acm.org> (accessed on 2 May 2025).
9. Wagenpfeil, S.; Mc Kevitt, P.; Cheddad, A.; Hemmje, M. Explainable Multimedia Feature Fusion for Medical Applications. *J. Imaging* **2022**, *8*, 104. [\[CrossRef\]](#) [\[PubMed\]](#)
10. Wagenpfeil, S.; Kevitt, P.M.; Hemmje, M. Towards Automated Semantic Explainability of Multimedia Feature Graphs. *Information* **2021**, *12*, 502. [\[CrossRef\]](#)
11. Knuth, D.E. *The Art of Computer Programming, Vol. 1: Fundamental Algorithms*, 3rd ed.; See Section on Adjacency Matrices in Graph Theory; Addison-Wesley: Boston, MA, USA, 1997.

12. Izacard, G.; Grave, E. Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering. *arXiv* **2020**, arXiv:2007.01282.
13. Liu, P.; Yuan, W.; Fu, J.; Jiang, Z.; Hayashi, H.; Neubig, G. Pre-train Prompt Fine-tune: A Survey of Prompt Engineering Methods. *arXiv* **2021**, arXiv:2107.13586.
14. Johnson, J.; Douze, M.; Jégou, H. Billion-scale Similarity Search with GPUs. *IEEE Trans. Big Data* **2019**, *7*, 535–547. [[CrossRef](#)]
15. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient Estimation of Word Representations in Vector Space. *arXiv* **2013**, arXiv:1301.3781.
16. Yu, M.; Neumann, M.; He, J.; Xiong, C. Improving Retrieval-Augmented Generation with Hybrid Sparse and Dense Retrieval. *arXiv* **2022**, arXiv:2212.10547.
17. Khattab, O.; Zaharia, M. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, 25 July 2020; pp. 39–48.
18. Chen, W.; Hu, H.; Chen, X.; Verga, P.; Cohen, W.W. MuRAG: Multimodal Retrieval-Augmented Generator for Open Question Answering over Images and Text. *arXiv* **2022**, arXiv:2210.02928.
19. Mortaheb, M.; Khojastepour, M.A.A.; Chakradhar, S.T.; Ulukus, S. Re-ranking the Context for Multimodal Retrieval-Augmented Generation. *arXiv* **2025**, arXiv:2501.04695.
20. Tevissen, Y.; Guetari, K.; Petitpont, F. Towards Retrieval-Augmented Generation over Large Video Libraries. *arXiv* **2024**, arXiv:2406.14938.
21. Chen, H. Multimodal Retrieval-Augmented Generation Question-Answering System. *OpenReview Preprint* **2025**. Available online: <https://openreview.net/forum?id=fMaEbeJGpp> (accessed on 2 May 2025).
22. Mei, L.; Mo, S.; Yang, Z.; Chen, C. A Survey of Multimodal Retrieval-Augmented Generation. *arXiv* **2025**, arXiv:2504.08748.
23. Doshi-Velez, F.; Kim, B. Towards A Rigorous Science of Interpretable Machine Learning. *arXiv* **2017**, arXiv:1702.08608.

**Disclaimer/Publisher’s Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.

Reproduced with permission of copyright owner. Further reproduction  
prohibited without permission.