

MULTI-HEAD RAG: SOLVING MULTI-ASPECT PROBLEMS WITH LLMs

Anonymous authors

Paper under double-blind review

ABSTRACT

Retrieval Augmented Generation (RAG) enhances the abilities of Large Language Models (LLMs) by enabling the retrieval of documents into the LLM context to provide more accurate and relevant responses. Existing RAG solutions do not focus on queries that may require fetching multiple documents with substantially different contents. Such queries occur frequently, but are challenging because the embeddings of these documents may be distant in the embedding space, making it hard to retrieve them all. This paper introduces Multi-Head RAG (MRAG), a novel scheme designed to address this gap with a simple yet powerful idea: leveraging activations of Transformer’s multi-head attention layer, instead of the decoder layer, as keys for fetching multi-aspect documents. The driving observation is that different attention heads learn to capture different data aspects. Harnessing the corresponding activations results in embeddings that represent various facets of data items and queries, improving the retrieval accuracy for complex queries. We provide an evaluation methodology and metrics, multi-aspect datasets, and real-world use cases to demonstrate MRAG’s effectiveness. We show MRAG’s design advantages over 18 RAG baselines, empirical improvements of up to 20% in retrieval success ratios, and benefits for downstream LLM generation. MRAG can be seamlessly integrated with existing RAG frameworks and benchmarks.

1 INTRODUCTION

Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) emerged as a promising remedy for several key limitations of Large Language Models (LLMs). By decoupling knowledge from model weights, RAG reduces the risk of leaking confidential data (Yan et al., 2025), a critical concern when training on sensitive corpora. It also mitigates hallucinations (Huang et al., 2025b) by grounding LLM outputs in retrieved, verifiable information. The core mechanism involves augmenting a generative LLM with a retrieval module that fetches relevant passages from an external corpus in response to a query. Rather than relying solely on static, parametric knowledge, RAG dynamically incorporates retrieved content into the model’s context, enabling more accurate and up-to-date responses. While early RAG systems required training task-specific retrievers and readers (Humeau et al., 2020), the current trend favors lightweight, in-context learning (ICL) approaches (Gao et al., 2024), which avoid the cost and complexity of retraining and allow for rapid knowledge updates without modifying the underlying LLM parameters.

A RAG pipeline consists of two main stages: data preparation and query execution. In the preparation stage, a vector database (DB) is constructed by embedding a collection of documents and storing these embeddings alongside their associated content. At inference time, the query is similarly embedded, and nearest-neighbor search retrieves the most relevant data items, which are then passed to the LLM for final answer generation. Ongoing developments in RAG have led to variety of different RAG designs (Gao et al., 2024).

Yet, no existing RAG method or benchmark explicitly targets *multi-aspectual* problems, that is, queries requiring the integration of multiple, semantically distinct aspects. For example, answering “What car did Alexander the Great drive?” (assuming no historical pretraining) requires retrieving unrelated documents on Alexander the Great and on car manufacturing, whose embeddings may lie far apart in the vector space. Such multi-aspect queries are common in industrial settings, as confirmed by extensive discussions with our industry collaborators and further supported by our analysis of over 35 industry reports (details are in Appendix A; we considered accident prevention,

healthcare, airport management, and others). For example, in a chemical plant accident, determining the cause might require accessing diverse and confidential documents related to worker psychology (“Was it mismanagement?”), equipment records (“Was a part outdated or rusty?”), weather conditions (“Were there power spikes due to a storm?”), or even microclimate (“Was prolonged humidity a factor?”). As shown in Section 5, such cases have been unaddressed by modern RAG schemes and benchmarks (Chen et al., 2024b; Xiong et al., 2024; Lyu et al., 2025; Es et al., 2024).

In this work, we propose Multi-Head RAG (MRAG): a scheme that addresses the above problem (**contribution 1**). Common practice in many modern RAG designs is to use embeddings derived from *last-layer decoder block activations* of a decoder-based *embedding LLM* (a language model specifically fine-tuned to provide high-quality embeddings). Examples of such embedding models are SFR-Embedding-Mistral and E5-Mistral-7B. Our key idea is to extend this design by incorporating activations from the *multi-head attention (MHA) modules of decoder blocks* as embedding sources. This enables the representation of multiple distinct aspects of the input text. Specifically, a Transformer consists of a stack of blocks (e.g., 96 in GPT-3 (Wang et al., 2025b)), each containing an MHA module with multiple *heads* that are trained with separate parameter sets. Through a survey of literature into Transformer design and interpretability, we find empirical and theoretical support for the conjecture that *different heads specialize in different aspects of the input* (details are in Section 2.4). This enables efficient *multi-aspect embeddings* without increasing space or compute costs compared to standard RAG, and without requiring *any* additional fine-tuning or architectural modifications to the base model.

Considering multi-aspectuality comes with challenges. For example, it is unclear how to assess whether a RAG solution does indeed harness multiple aspects when fetching documents. For this, we develop a multi-aspect RAG pipeline that includes data preparation and query processing with both multi-aspect retrieval and ranking schemes (**contribution 2**). We also establish an evaluation methodology and provide multi-aspect datasets, complementing existing RAG benchmarks (Chen et al., 2024b) (**contribution 3**). We ensure the relevance of our RAG datasets in real use cases by working directly with tech leaders (e.g., a generative AI division head) from 3 corporations, all of which actively use RAG in their own LLM infrastructures. We illustrate the advantages of MRAG over 18 traditional and modern RAG designs in various design criteria and in both time and space complexities (**contribution 4**). In evaluation, MRAG enhances the relevance of retrieved documents by up to 20% over modern RAG baselines, offers comparable performance without degradation for single-aspect queries, and benefits the downstream LLM generation (**contribution 5**). Thanks to its simplicity, MRAG can be seamlessly integrated into any stores while its benchmarking methodology can straightforwardly extend benchmarks such as RAGAs (**contribution 6**).

2 MRAG: DESIGN & IMPLEMENTATION

A typical RAG scheme (see Figure 2) consists of two main parts: **data preprocessing** ④ and **query execution** ⑤; both parts heavily use an **embedding model** ③ and the **data store** ②. During preprocessing, each document in the database is encoded into one or more embeddings using the embedding model; these embeddings are stored with that document (usually as a key-value pair in a vector DB). During query execution, the embedding of the user-provided query is constructed using the same embedding model; then, the *retriever* fetches candidate documents based on embedding similarity, the *reranker* refines their ordering using more precise scoring, and the *reader* uses a downstream generative model to synthesize the final answer from the top-ranked documents.

2.1 CONSTRUCTING MULTI-ASPECT EMBEDDINGS ①

An embedding is constructed for each data item in a pre-existing database (part ④) and for the user query (part ⑤). In standard modern RAG, given an input chunk of n tokens constituting a document

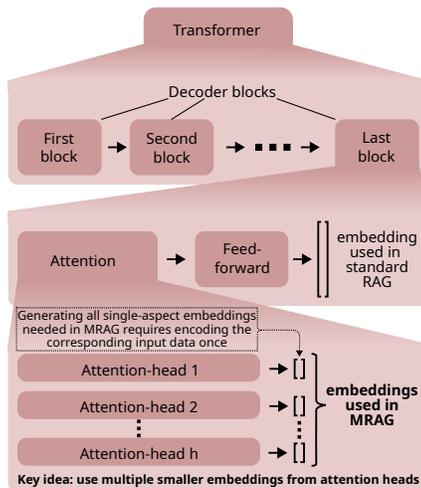


Figure 1: An overview of the decoder architecture, and a comparison of how standard RAG and Multi-Head RAG embeddings are generated.

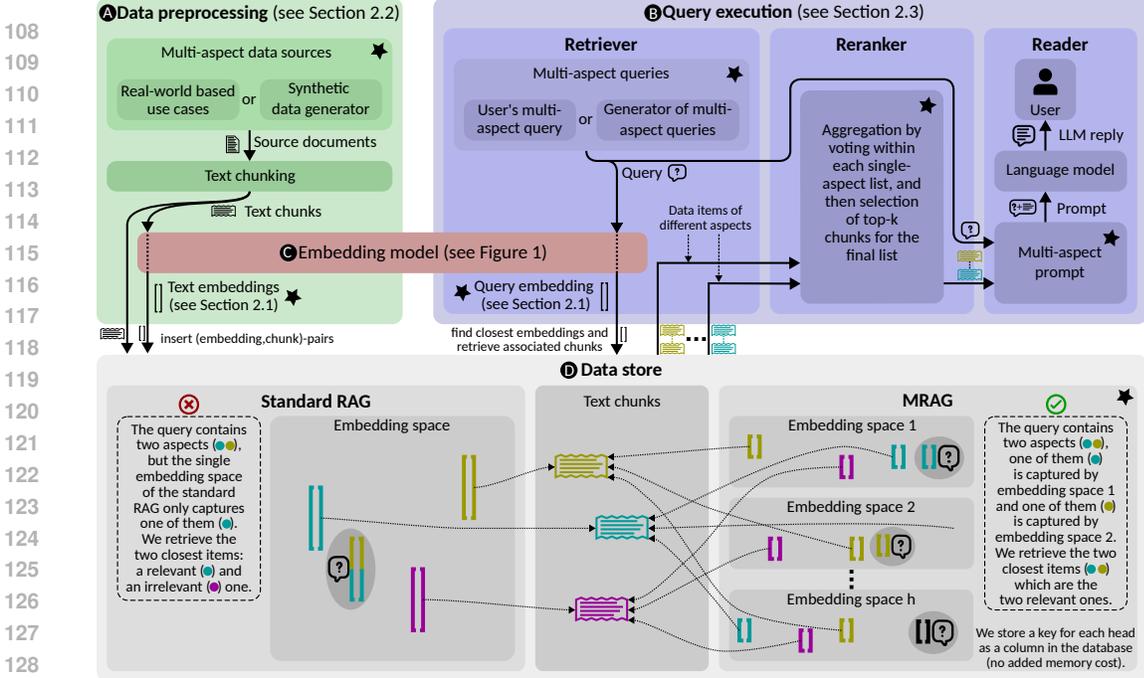


Figure 2: Overview of the MRAG pipeline, consisting of two parts: data preparation **A** and query execution **B**. The embedding model **C** and the data store **D** are used by both parts. The data store **D** contains text embeddings **E** linking to text chunks **F** reflecting three different aspects (cyan, magenta, yellow). Blocks marked by a star \star are a novelty of this work.

or chunk to be embedded, one obtains a corresponding embedding by applying the embedding model to the data item and extracting the activation of the final feed-forward layer for the last token \mathbf{x}_n , i.e., $\text{feed-forward}(\text{multi-head}(\mathbf{x}_n)) \in \mathbb{R}^d$. Appendix B provides additional mathematical details.

In MRAG, instead of relying on this single embedding, the activations of all h attention heads are obtained *before* they are merged by the final projection layer. Specifically, for the last token \mathbf{x}_n , we extract the set of head-specific vectors $\mathcal{S} = \{\mathbf{e}_k\}_{k=1}^h$, where each $\mathbf{e}_k = \text{head}^k(\mathbf{x}_n) \in \mathbb{R}^{d/h}$ is referred to as a “single-aspect embedding” and \mathcal{S} is called a “multi-aspect embedding”. This provides h semantically diverse embeddings per input, reflecting the different perspectives captured by each attention head. Since these vectors are extracted from the same internal activation used by standard RAG, the overall space and time requirements remain unchanged.

2.2 DATA PREPROCESSING **A**

We populate a data store **D** with multi-aspect embeddings **E** for their corresponding documents **F** or text chunks **G**. Unlike in standard RAG, where a single embedding **E** points to a single text chunk **G**, in MRAG, each out of h single-aspect embeddings **E** points to the original text chunk **G** (i.e., the data store **D** contains h embedding spaces). This crucial feature allows MRAG to compare query **H** and text chunks **G** in multiple embedding spaces that capture multiple aspects of the data.

MRAG performs one-time *importance scoring* for each embedding space i , capturing its relevance to the dataset. Each score s_i combines: (i) the average L2 norm a_i of embeddings in space i , reflecting *attention strength*, and (ii) the average cosine distance b_i between sampled embeddings, approximating its *semantic spread*. The final score is $s_i = a_i \cdot b_i$, encouraging both head relevance (through high a_i) and high representational diversity (through high b_i). Full details are provided in Algorithm 1, Appendix C.1.

2.3 QUERY EXECUTION

During query execution, MRAG first generates a multi-aspect embedding \mathcal{S} of the input query (cf. Section 2.1). These embeddings are computed fully in parallel during the same inference pass, so this multi-vector representation introduces no additional latency.

Retriever. Given \mathcal{S} , the retriever conducts parallel retrieval of the top- c nearest chunks within each embedding space, and the aggregation of the resulting candidates across all h spaces.

Reranker. Once hc candidate chunks are retrieved across all heads, the ranker stage consolidates them into a single list of top- k results using a simple but effective voting strategy. For each candidate chunk at position p in the ranked list for space i , we assign a score of $s_i \cdot 2^{-p}$, where s_i is the precomputed importance score. This exponentially discounts lower-ranked candidates and balances influence across heads. The final top- k list is obtained by globally sorting all candidates by these scores. The voting procedure is described in Algorithm 2 in Appendix C.2.

Reader. The top- k retrieved results are inserted into the LLM context using a multi-aspect prompt template (prompts are fully specified in Appendix D). Each result is placed in a separate section of the prompt. Stored metadata can be included alongside each chunk to provide additional context, such as the aspect or the chunk category.

2.4 CAPTURING MULTI-ASPECTUALITY WITHOUT ADDITIONAL TRAINING

A key design decision in MRAG is to use the hidden representations immediately after the attention block in the last decoder layer, *without additional fine-tuning*, to avoid training overhead and make deployment easy. This is motivated by growing evidence that attention heads in Transformer models naturally converge during training to focus on distinct aspects of the input. For example, Wang et al. (2025a) show that heads diverge into clusters specialized for different input patterns, while Olsson et al. (2022) and others (McDougall et al., 2024; Gould et al., 2024) discover that heads focus on repeated sequences or named entities. Similar trends are observed in BERT (Clark et al., 2019; Kovaleva et al., 2019; Htut et al., 2019). We provide more details in a brief literature survey (Appendix E.1). Given these findings, we assume that even without additional fine-tuning, the embeddings output by MHA already encode multi-aspectuality suitable for downstream retrieval. Our own brief attention pattern analysis (Appendix E.2) confirm shifting token focus across heads.

2.5 PARALLELIZATION AND SYSTEMS CONSIDERATIONS

MRAG evenly distributes the embedding dimensionality across attention heads, keeping total storage and compute $O(nd)$ (same as standard RAG; we show a more detailed complexity analysis in the following section). It uses off-the-shelf ANN indexes that support parallel subspace search, and it avoids dynamic or variable-length vectors that complicate indexing and caching.

Modern vector databases already support the parallel multi-vector search required by our approach. Milvus can execute distributed, parallel ANN queries across shards with low latency (Wang et al., 2025c; Clavié et al., 2024; Wang et al., 2021). Pinecone provides cascading retrieval with fixed-size multi-vector encodings (such as ConstBERT) to maintain predictable costs. Furthermore, ESPN tackles multi-vector retrieval at SSD and GPU levels by implementing prefetching pipelines and storage bypass mechanisms, achieving near-memory query latency, even when dealing with large indices (Shrestha et al., 2024). Taken together, these technologies illustrate that multi-vector retrieval, especially when using fixed-size vectors as MRAG does, can be implemented at scale with negligible practical overhead.

3 COMPUTE & STORAGE COMPLEXITY ANALYSIS

We analyze the runtime and space complexity of MRAG and 18 RAG baselines, the results are in Table 1 (the runtime metrics of work/depth and the notation are explained in the table caption). Derivation details and a broader discussion are in Appendix F. Overall, MRAG achieves competitive results in all aspects. At inference time, it extracts h attention-head embeddings in parallel from a single forward pass, resulting in the same latency as standard RAG. Preprocessing is lightweight, requiring only one pass per document and simple statistics for head scoring, avoiding the cost of training or complex structure construction. Storage overhead is minimal, as h single-aspect embeddings per data item have the same dimensionality as a standard embedding. In contrast, prior schemes like Poly-encoder (Humeau et al., 2020) and ColBERT (Khattab & Zaharia, 2020) incur significant cost due to using many token-level embeddings or time-consuming training rounds.

4 BENCHMARKING MULTI-ASPECTUALITY

To assess how well MRAG performs on multi-aspect queries, we need (1) datasets of multi-aspectual documents, (2) queries to the LLM that require retrieving multi-aspect documents, and (3) metrics that assess how well a RAG scheme retrieves such multi-aspect data. We now summarize these three elements; details are in Appendix G.

Table 1: **Time and storage complexity of different RAG schemes; MRAG is the only scheme to match the results of a plain vanilla RAG.** **Work** is the total number of all operations, **Depth** is runtime complexity assuming unlimited parallel processing units; both are established measures of analyzing parallel algorithms. **Retrieval** describes steps during inference, including any postprocessing before returning documents or summaries. **Preprocessing** includes all steps before inference. n : number of data items in a database (i.e., document chunks), d : embedding dimensionality, k : number of retrieved top chunks per single user query, l_q : average token length of the query, l_d : average token length of the document, W_m/D_m : work/depth to run a transformer-based model, W_e/D_e : work/depth to embed a graph, W_i/D_i : work/depth to index a graph, s : polynomial function that models the complexity of various additional scheme-specific design decisions, heuristics, etc., which cannot be straightforwardly modeled with closed-form expressions (e.g., count of graph communities, toy-graph size, self-note length, BM25 matching cost, keyword matching cost); it scales at most linearly with n and is typically considerably larger than $O(1)$.

| Scheme | Retrieval | | Preprocessing | | Storage |
|-------------------------------------|-------------------------|----------------------------|-----------------------------|-------------------|-----------------|
| | Work | Depth | Work | Depth | |
| Vanilla RAG | $O(W_m + nd)$ | $O(D_m + \log d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |
| Poly-encoders (Humeau et al., 2020) | $O(W_m + nd + s(n)d)$ | $O(D_m + \log d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |
| Lewis et al. (2020) | $O(2W_m + nd)$ | $O(2D_m + \log d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |
| ColBERT (Khattab & Zaharia, 2020) | $O(W_m + ndl_q l_d)$ | $O(D_m + \log d l_d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd l_d)$ |
| EMDR (Singh et al., 2021) | $O((k+1)W_m + nd)$ | $O(2D_m + \log d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |
| Self-RAG (Asai et al., 2024) | $O((2k+1)W_m + nd + k)$ | $O(2D_m + \log d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |
| Chain-of-Note (Yu et al., 2024) | $O((ks(n)+1)W_m + nd)$ | $O((ks(n)+1)D_m + \log d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd + s(n)d)$ |
| RAPTOR (Sarithi et al., 2024) | $O(W_m + nd + s(n)d)$ | $O(D_m + \log d)$ | $O(W_m n + W_m s(n))$ | $O(D_m)$ | $O(nd + s(n)d)$ |
| RAGraph (Jiang et al., 2024) | $O(W_m + nd + s(n)d)$ | $O(D_m + \log d)$ | $O(W_m n + W_e)$ | $O(D_m + D_e)$ | $O(nd + s(n)d)$ |
| RQ-RAG (Chan et al., 2024) | $O(s(n)(W_m + nd))$ | $O(s(n)(D_m + \log d))$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |
| ActiveRAG (Xu et al., 2024) | $O(3W_m + nd)$ | $O(2D_m + \log n)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |
| HiQA (Chen et al., 2024c) | $O(W_m + nd + s(n)d)$ | $O(D_m + \log d)$ | $O(W_m n + s(n)d)$ | $O(D_m)$ | $O(nd + s(n)d)$ |
| GraphRAG (Edge et al., 2025) | $O(s(n)(W_m + d))$ | $O(D_m + s(n)d)$ | $O(W_m n + W_m s(n) + W_e)$ | $O(D_m + D_e)$ | $O(nd + s(n)d)$ |
| Fusion RAG (Rackauckas, 2024) | $O(s(n)W_m + knd)$ | $O(2D_m + \log d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |
| Meta-chunking (Zhao et al., 2025b) | $O(W_m + (n + s(n))d)$ | $O(D_m + \log d)$ | $O(l_d W_m + s(n)W_m)$ | $O((l_d + 1)D_m)$ | $O(nd + s(n)d)$ |
| MoC (Zhao et al., 2025a) | $O(W_m + (n + s(n))d)$ | $O(D_m + \log d)$ | $O(W_m n + s(n))$ | $O(D_m)$ | $O(nd + s(n)d)$ |
| Parametric RAG (Su et al., 2025) | $O(W_m + nd + s(n))$ | $O(D_m + \log d + s(n))$ | $O(W_m n s(n))$ | $O(D_m s(n))$ | $O(n s(n))$ |
| SuperRAG (Yang et al., 2025) | $O(W_m + nd + W_i)$ | $O(D_m + \log d + D_i)$ | $O(W_m n + s(n) + W_e)$ | $O(D_m + D_e)$ | $O(nd + s(n))$ |
| HiRAG (Huang et al., 2025a) | $O(W_m + nd + W_i)$ | $O(D_m + \log d + D_i)$ | $O(W_m n + s(n) + W_e)$ | $O(D_m + D_e)$ | $O(nd + s(n))$ |
| MRAG [This Work] | $O(W_m + nd)$ | $O(D_m + \log d)$ | $O(W_m n)$ | $O(D_m)$ | $O(nd)$ |

We construct three **multi-aspect datasets**: (1) a synthetic Wikipedia dataset with documents sampled from clearly distinct categories (e.g., countries, shipwrecks, board games); (2) a real-world-based legal document dataset annotated by a legal areas (energy law, family law, criminal law, etc.) or document language style (aggressive, mild, neutral, etc.); and (3) a real-world-based dataset of industry accident reports, categorized by cause. For each dataset, we generate **multi-aspect queries** using GPT-4o, combining n distinct categories into single queries (with $n \in \{1, 2, 3, 4, 5, 6, 10, 15, 20, 25\}$). For example, a query with 10 aspects must contain a question about 10 different documents from 10 different categories.

To evaluate retrieval performance, we introduce three **metrics for assessing multi-aspectuality**. Let Q be a query, S a Reranker scheme, and Q_{rel} the ideal set of n relevant documents. The *Retrieval Success Ratio* is defined as $\Xi(Q, n) = \frac{|S(Q, n) \cap Q_{rel}|}{|Q_{rel}|}$, measuring the fraction of exactly matched documents. The *Category Retrieval Success Ratio* Ξ_c extends this by also accepting matches from the same categories as those in Q_{rel} even if the exact document has not been matched. Finally, we define a tunable *Weighted Success Ratio* $\Xi_w = \frac{w \cdot \Xi + \Xi_c}{w + 1}$, allowing the user to adjust the trade-off between exact and category-level matches via the weight w .

We motivate category-level retrieval using both previous works (V et al., 2025; Liu et al., 2021; Chen et al., 2024a; Xiao et al., 2024) as well as our hands-on experience in legal and industrial accident analysis. Namely, even when exact matches are missing, retrieving documents from the same semantic category can enhance generation, as supported by the classic *clustering hypothesis* in information retrieval (V et al., 2025; Liu et al., 2021): documents in the same “neighborhood” are likely to be relevant to the same query. Recent studies further validate this across open-domain QA and ontology-guided RAG (Chen et al., 2024a; Xiao et al., 2024).

5 EVALUATION

We now illustrate the advantages of MRAG over the state of the art. Further details on evaluation setup and additional results are in – respectively – Appendix H and I.

Comparison Baselines. We consider three main baselines: **Standard RAG**, **Split RAG**, and **Fusion RAG** (Rackauckas, 2024). The first represents a modern RAG pipeline in which each document uses the activations of the last decoder layer as its embedding. The second is a blend between Standard RAG and MRAG; it splits the activation of the last decoder layer in the same way as MRAG and applies a voting strategy. The purpose of Split RAG is to show that *MRAG’s benefits come from using the multi-head output as embedding and not merely using multiple embedding spaces*. Additionally, we consider **Fusion RAG** (Rackauckas, 2024), an optional mechanism that we harness to *further enhance the benefits of MRAG at the tradeoff of additional tokens* (detailed in Section 5.3).

Embeddings & models. While MRAG allows for extracting multi-aspect embeddings from *any* block, we found that the last MHA works best in our experiments. MRAG can use any embedding

model with MHA; we consider two embedding models from the MTEB leaderboard (Huggingface, 2025), the SFR-Embedding-Model (Meng et al., 2024) and the e5-mistral-7b-instruct (Wang et al., 2024), both based on the Mistral 7B architecture with 32 decoder blocks and 32 attention heads.

5.1 MRAG DELIVERS SUPERIOR PERFORMANCE FOR MULTI-ASPECT QUERIES

We start with the **Wikipedia multi-aspect dataset** and corresponding queries (cf. Section 4). In each query, we mention the documents to be fetched in the text and then assess the success ratio of different RAG strategies in finding these documents and their categories (a full example of such a query is in Figure 10 in Appendix G). We show first the absolute retrieval performance of MRAG over Standard RAG in Figure 3. We fix the number of aspects present in the queries to 10, and vary the total number of retrieved documents from 10 to 30. MRAG consistently outperforms Standard RAG ($> 10\%$ increase in the retrieval success ratio on average for exact document matches). Moreover, the retrieval performance increase is even more significant on category matches ($> 25\%$ increase in the retrieval success ratio on average). The performance increase is further detailed in the histograms on the right side. Here, for a specific number of documents fetched, MRAG’s histogram indicates a better distribution of retrieval success ratios (across all 25 queries). This gain stems from MRAG’s ability to decompose query semantics across multiple attention heads, increasing the likelihood of matching each aspect with a semantically aligned document.

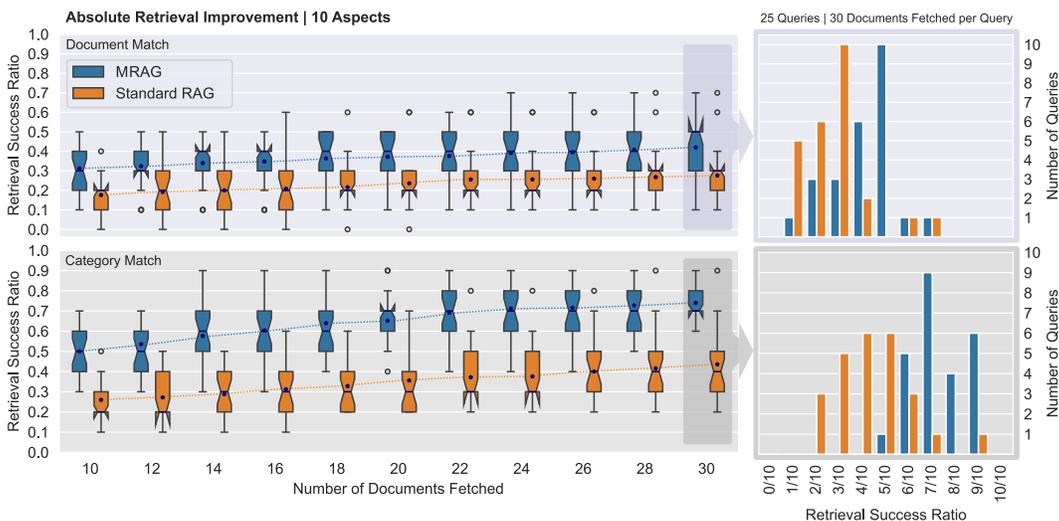


Figure 3: Retrieval success ratio over 25 queries between MRAG and Standard RAG; each query uses 10 different aspects. The top part presents exact document matches, the bottom part presents category only matches (we explain the metrics in Sec. 4). A histogram is presented for a specific sample to showcase the detailed distribution among the 25 queries (there are 30 documents fetched for each query).

Next, Figure 4 shows the relative weighted performance improvement of MRAG with respect to Standard RAG as we vary the number of aspects present in the queries. We show data for two different embedding models (SFR and e5). MRAG consistently outperforms the Standard RAG by 10-20% on average, not only across the number of documents fetched, but also across the increasing counts of aspects present in the replies, and does it for both embedding models. This robustness suggests that MRAG scales better than Standard RAG with query complexity, as its multi-head representation distributes semantic load more evenly than the single-vector baseline.

To further illustrate advantages of MRAG, we also consider two **real-word use cases** from in-house industry data analytics projects, namely, the synthesis of legal documents and the analysis of causes of chemical plant accidents. The results are in Figure 5. In the former (the left side), the task is to create a document based on user requirements that may be related to different *aspects*, for example to the law being considered (e.g., the British or the US one), the subject (e.g., energetic or civil), the style of the document (e.g., aggressive or mild), and others. This task is executed with RAG that can fetch documents from a database. In the latter (the right side), the task is to discover a cause of an accident. Here, one also wants to retrieve documents from a database that should be used in the LLM context to facilitate discovering the cause of the accident. The causes are grouped in categories such as utility impact due to severe weather, lack of preparedness and planning, incorrect installation

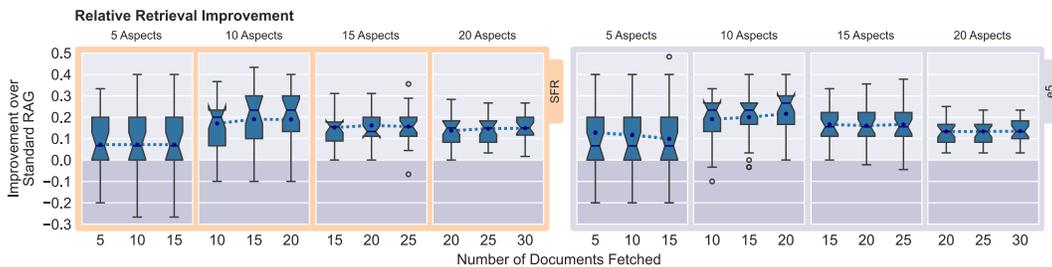


Figure 4: **Relative retrieval improvement of MRAG over Standard RAG** across queries with different numbers of aspects and different embedding models (SFR in the left side, e5 in the right side).

of equipment, lack of maintenance, et cetera. Similarly to the previous analyses, we measure the retrieval success ratio over corresponding databases. MRAG offers advantages over other schemes.

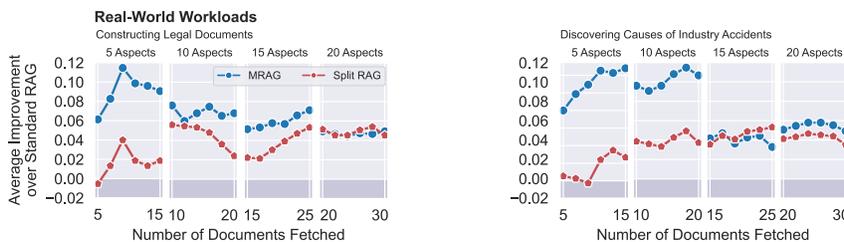


Figure 5: **Average improvement of the retrieval success ratio of MRAG and Split RAG over Standard RAG** for two real-world workloads *constructing legal documents* (left) and *discovering causes of industry accidents* (right).

We also delve deeper into the underlying factors for MRAG’s performance gains. For this, we compare MRAG to the Split RAG baseline in Figure 6. The blue plots show the relative weighted performance of MRAG and Split RAG over Standard RAG. MRAG performs better than Split RAG, illustrating that its *high accuracy is due to the actual multi-head part*, and not merely just partitioning the vector and using multiple embedding spaces.

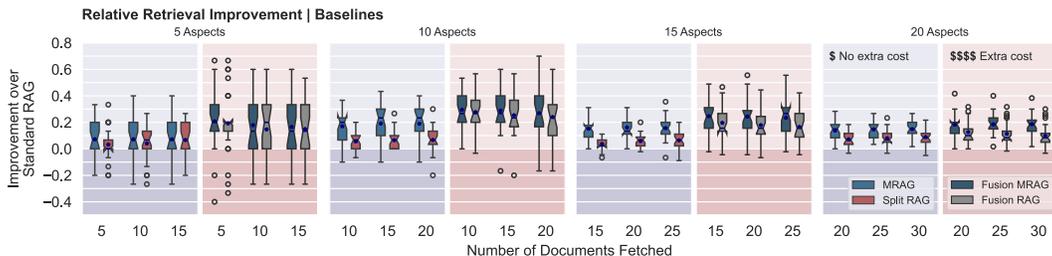


Figure 6: **Relative retrieval improvements of MRAG over Standard RAG** for the SFR embedding model compared with **Split RAG** (the blue plots), and the **relative retrieval improvements of Fusion MRAG over Standard RAG** compared with **Fusion RAG** (the red plots).

5.2 MRAG ENSURES HIGH PERFORMANCE FOR SINGLE-ASPECT QUERIES AS WELL

We additionally show in Table 2 that MRAG performs on-par with Standard RAG on queries where only a single aspect is expected. This confirms that MRAG’s design generalizes well: in low-aspect settings, the aggregation over heads still captures dominant semantics, effectively collapsing into a strong single-vector representation with only a negligible accuracy drop for single-aspect tasks.

5.3 MRAG SEAMLESSLY ENHANCES EXISTING RAG SCHEMES

MRAG’s simplicity ensures that it can be seamlessly integrated with other RAG approaches. As an example, we combine MRAG with *Fusion RAG*, which uses an LLM (additional token cost) for even more accurate retrieval. Fusion RAG uses an LLM to create a fixed number of questions about the RAG query. Each question is separately applied through an embedding model using Standard RAG. This relies on multiple LLM calls and heuristic reranking, inflating latency and cost.

We apply MRAG to each of these questions and denote the combined scheme as *Fusion MRAG*. Red plots of Figure 6 show that Fusion MRAG consistently outperforms pure Fusion RAG, indicating

Table 2: **Retrieval success ratio** (the exact document match) for 25 **single-aspect** queries on different datasets (Multi-Aspect Dataset, Legal Dataset, Accidents Dataset), using different embedding models (SFR, e5). For every query, a specific single document (single-aspect) is expected to be among the fetched documents for the retrieval to be classified as successful.

| Documents Fetched | Wikipedia Dataset | | | | Legal Dataset | | Accidents Dataset | |
|-------------------|-------------------|--------------|-------|--------------|---------------|--------------|-------------------|--------------|
| | SFR | | e5 | | SFR | | SFR | |
| | MRAG | Standard RAG | MRAG | Standard RAG | MRAG | Standard RAG | MRAG | Standard RAG |
| 1 | 24/25 | 25/25 | 24/25 | 25/25 | 24/25 | 24/25 | 25/25 | 25/25 |
| 2 | 25/25 | 25/25 | 25/25 | 25/25 | 25/25 | 25/25 | 25/25 | 25/25 |
| 3 | 25/25 | 25/25 | 25/25 | 25/25 | 25/25 | 25/25 | 25/25 | 25/25 |

that these optimizations can be combined together. However, both Fusion strategies introduce a greater variance than MRAG and additional costs in terms of compute, latency, and tokens.

5.4 MRAG SEAMLESSLY ENHANCES DOWNSTREAM LLM GENERATION

We also illustrate that MRAG enhances the downstream LLM generation with its improved multi-aspect retrieval. To show this, we use a sampled subset of multi-aspect Wikipedia queries, for which we applied both Standard RAG and MRAG to retrieve supporting documents. These documents are then integrated into the prompt templates and are passed to the LLM to answer the original query with the aid of the retrieved data. To quantify the effect from RAG, we count facts in the LLM output (e.g., years or named entities such as locations). On average, MRAG generations contain on average 15.4 pieces of factual information per query, compared to 11.2 for Standard RAG. The average improvement further confirms MRAG’s advantages. By explicitly leveraging multi-aspectual embeddings from MHA, one increases the likelihood that the LLM generation includes diverse and complementary facts, especially for complex queries spanning multiple domains (e.g., combining historical events with technological timelines). Overall, MRAG helps the LLM in delivering richer, more comprehensive responses.

5.5 MRAG ENSURES NO LATENCY AND STORAGE OVERHEADS

MRAG introduces no latency overhead at query embedding time, as all head-level embeddings are extracted in parallel from a single standard forward pass. Retrieval across embedding subspaces is also fully parallelizable with modern vector databases (Pan et al., 2024; Ma et al., 2025), and our use of a modest number of heads (e.g., 16–32 for SFR-Embedding-Model) ensures this parallelism is within easy reach. Moreover, as the total dimensionality of embeddings remains unchanged (e.g., 1024 split across 32 heads), MRAG needs no additional storage compared to standard RAG.

5.6 FURTHER ANALYSES & ABLATION STUDIES

Additional analyses are in Appendix I, they include analyzing the impact of using embeddings from **different decoder blocks** (rather than the last one), **scalability of preprocessing**, and **additional voting strategies for reranking**. These analyses all confirm the previous findings of MRAG broadly outperforming other baselines.

6 RELATED WORK & ADVANTAGES OF MULTI-HEAD RAG

RAG Solutions. We compare MRAG to a large number of both traditional and modern RAG solutions in Table 3 in terms of design advantages (complexity analysis and empirical evaluation are in, respectively, Sections 3 and 5). While prior RAG systems support retrieving multiple documents for a single query (e.g., RAG (Lewis et al., 2020), EMDR (Singh et al., 2021)), none of them generates *multi-aspect* embeddings per document, and therefore do not offer multi-aspectuality. Similarly, although these systems employ Transformer architectures with MHA, they make no use of the MHA internal structure. Other methods such as Poly-encoder (Humeau et al., 2020) and ColBERT (Khattab & Zaharia, 2020) do produce multiple embeddings per document, but for a different reason: they are based on models such as BERT, which inherently yields token-level embeddings. Thus, these models require multiple vectors per document simply to cover its content at the token level, not to represent distinct semantic aspects. In contrast, MRAG leverages the powerful internal structure of modern decoder-based LLMs, where a small number of vectors derived from MHA (or even a single vector from the final decoder block, as in standard RAG) suffices to represent the semantics of an entire document or a chunk, as indicated by recent work (Lee et al., 2025; Besta et al., 2025).

Building on this foundation, MRAG achieves further practical advantages. It introduces *scalable preprocessing* since (1) global embeddings are computed in a single forward pass per document and (2) scoring of heads is straightforwardly parallelizable. At inference time, it incurs *no additional cost* relative to standard RAG – in contrast to token-level approaches, which require computing nu-

Table 3: **Comparison of the advantages of different RAG schemes** (sorted top to bottom chronologically). **No additional training**: does a given scheme require additional training beyond standard pre-training and fine-tuning? **Works with any MHA LLM**: does a given scheme can seamlessly work with any multi-head attention (MHA) LLM? **Extensibility to other modalities**: could a given scheme be relatively easily used beyond LLMs, e.g., with Graph Foundation Models (GFMs), Vision Models, and others? **No overhead at inference**: does a given scheme enable zero additional overhead at the inference? **Scalable preprocessing, low storage overhead**: does a given scheme enable scalable preprocessing and little storage overhead, respectively? (details in Section 3 and Table 1). **Multi-Aspectuality**: does a given scheme enable extracting multiple aspects of indexed data? “”: full support, “”: partial support, “”: no support, “”: unknown.

| Scheme | No additional training | Works with any MHA LLM | Extensibility to other modalities | No overhead at inference | Scalable preprocessing | Low storage overhead | Multi-Aspectuality |
|-------------------------------------|---|--|---|---|---|---|---|
| Poly-encoders (Humeau et al., 2020) |  |  (BERT based) |  |  |  |  |  |
| Lewis et al. (2020) |  |  |  |  |  |  |  |
| CoBERT (Khattab & Zaharia, 2020) |  |  (BERT based) |  |  |  |  |  |
| EMDR (Singh et al., 2021) |  |  |  |  |  |  |  |
| Self-RAG (Asai et al., 2024) |  |  |  |  |  |  |  |
| Chain-of-Note (Yu et al., 2024) |  |  |  |  |  |  |  |
| RAPTOR (Sarthi et al., 2024) |  |  |  |  |  |  |  |
| RAGraph (Jiang et al., 2024) |  |  (only GFMs) |  (only GFMs) |  |  |  |  |
| RQ-RAG (Chan et al., 2024) |  |  |  |  |  |  |  |
| ActiveRAG (Xu et al., 2024) |  |  |  |  |  |  |  |
| HiQA (Chen et al., 2024c) |  |  |  |  |  |  |  |
| GraphRAG (Edge et al., 2025) |  |  |  |  |  |  |  |
| Fusion RAG (Rackaukas, 2024) |  |  |  |  |  |  |  |
| Meta-chunking (Zhao et al., 2025b) |  |  |  |  |  |  |  |
| MoC (Zhao et al., 2025a) |  |  |  |  |  |  |  |
| Parametric RAG (Su et al., 2025) |  |  |  |  |  |  |  |
| SuperRAG (Yang et al., 2025) |  |  |  |  |  |  |  |
| HIRAG (Huang et al., 2025a) |  |  |  |  |  |  |  |
| MRAG [This work] |  |  |  |  |  |  |  |

merous pairwise token similarities during retrieval, our method only compares compact, chunk-level vectors. *Storage overhead is also minimal*: the MHA-derived embeddings match the dimensionality of standard last-layer embeddings (detailed time and storage complexity analyses are in Section 3). Finally, MRAG requires *no training* and is *highly versatile*: it can be applied to any Transformer model with MHA, and is easily extensible to other modalities such as vision and graph data, where Transformer architectures are now common.

Reranking. Retrieval is sometimes enhanced by a **reranking** phase (Rosa et al., 2022; Nogueira & Cho, 2020; Nogueira et al., 2020; Li et al., 2023; Gao et al., 2021; MacAvaney et al., 2019). Here, after retrieving a set of relevant chunks, they are re-ranked using specialized models. *In this work, we provide a heuristic reranker that considers multi-aspectuality, but we design MRAG specifically so that it can be seamlessly used in conjunction with any other existing cross-encoders.*

Multi-Head Embeddings Outside RAG. Several methods propose a new model variant or architectural modification to the Transformer in order to better exploit MHA (Park et al., 2020; Huang et al., 2019; Wang et al., 2020; Xue & Aletras, 2023). For example, MHSAN (Park et al., 2020) introduces a novel visual-semantic embedding network that extracts multiple region- and phrase-sensitive features using MHA. Wang et al. (2020) develop a speaker embedding network that explicitly enforces head-level diversity via contrastive learning across resolutions. Xue & Aletras (2023) propose PIT, a Transformer variant that composes attention heads across layers to reduce redundancy and enable efficient inference. Vashisht et al. (2025) introduce MAGE, a technique that mixes heads across models to improve generalization through stochastic head combinations. *In contrast, MRAG is the first method to harness embeddings from MHA in pre-trained, decoder-based embedding LLMs for the purpose of more effective RAG. Unlike the above works, MRAG requires no architectural modifications, no training, and no custom modules, being fully plug-and-play.*

7 CONCLUSION

We introduced Multi-Head RAG (MRAG), a simple yet powerful extension to RAG that leverages the multi-head attention (MHA) activations of decoder models to capture multiple semantic aspects of a query or document. Motivated by the challenge of retrieving semantically diverse documents for multi-aspect queries, which is a common need in real-world applications, MRAG generates a set of aspect-specific embeddings without requiring extra training, model calls, or increased storage. Through a comprehensive evaluation including synthetic and industrial datasets, and tasks ranging from the accuracy of retrieval to the quality of downstream LLM generations, we demonstrate that MRAG consistently improves retrieval relevance (up to 20%) and yields richer, more factual outputs. Unlike many advanced RAG baselines, MRAG remains plug-and-play, scalable, and efficient—making it a practical and principled solution for high-precision multi-aspect retrieval in LLM-driven systems.

Ethics statement. This work uses datasets and models that are publicly available or licensed for use. We did not conduct research involving human subjects and did not collect or share any personal or sensitive information. All data in our case studies was de-identified and used according to its license terms. Our work focuses on retrieval and evaluation methods for multi-aspect queries and does not create new risks.

Reproducibility statement. We make our results easy to reproduce. We release code, configuration files, and scripts to rebuild indexes, run retrieval/reranking, and execute all evaluations, along with detailed instructions, hyperparameters, prompts, and random seeds. We document software and hardware versions and model/dataset identifiers. Where licensing prevents sharing specific corpora, we provide preprocessing scripts and exact instructions to reconstruct them. All figures and tables in the paper can be regenerated.

REFERENCES

- 486 **Ethics statement.** This work uses datasets and models that are publicly available or licensed for
487 use. We did not conduct research involving human subjects and did not collect or share any personal
488 or sensitive information. All data in our case studies was de-identified and used according to its
489 license terms. Our work focuses on retrieval and evaluation methods for multi-aspect queries and
490 does not create new risks.
- 491 **Reproducibility statement.** We make our results easy to reproduce. We release code, configura-
492 tion files, and scripts to rebuild indexes, run retrieval/reranking, and execute all evaluations, along
493 with detailed instructions, hyperparameters, prompts, and random seeds. We document software and
494 hardware versions and model/dataset identifiers. Where licensing prevents sharing specific corpora,
495 we provide preprocessing scripts and exact instructions to reconstruct them. All figures and tables
496 in the paper can be regenerated.
- 497
- 498 REFERENCES
- 499
- 500 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-RAG: Learning
501 to Retrieve, Generate, and Critique through Self-Reflection. In B. Kim, Y. Yue, S. Chaudhuri,
502 K. Fragkiadaki, M. Khan, and Y. Sun (eds.), *Proceedings of the Twelfth International Conference
503 on Learning Representations*, ICLR '24, pp. 9112–9141, Vienna, Austria, May 2024. Interna-
504 tional Conference on Learning Representations. URL https://proceedings.iclr.cc/paper_files/paper/2024/hash/25f7be9694d7b32d5cc670927b8091e1-Abstract-Conference.html.
- 505
- 506
- 507 Maciej Besta, Lorenzo Paleari, Marcin Copik, Robert Gerstenberger, Ales Kubicek, Piotr Nyczyk,
508 Patrick Iff, Eric Schreiber, Tanja Srin dran, Tomasz Lehmann, Hubert Niewiadomski, and Torsten
509 Hoefler. CheckEmbed: Effective Verification of LLM Solutions to Open-Ended Tasks, July 2025.
510 URL <https://arxiv.org/abs/2406.02524>. arXiv:2406.02524.
- 511 R. S. Bridger. A Guide to Human Factors in Accident Investigation. In *Sensemaking in Safety
512 Critical and Complex Situations*, pp. 13–32. CRC Press, Boca Raton, FL, USA, 2021. doi: 10.1
513 201/9781003003816-2. URL <https://www.taylorfrancis.com/chapters/oa-edit/10.1201/9781003003816-2/guide-human-factors-accident-investigation-bridger>.
- 514
- 515 Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. RQ-
516 RAG: Learning to Refine Queries for Retrieval Augmented Generation. In *Proceedings of the
517 First Conference on Language Modeling*, COLM '24, Philadelphia, PA, USA, October 2024.
518 OpenReview. URL <https://openreview.net/forum?id=tzE7VqsaJ4>.
- 519
- 520 Jianfa Chen, Emily Shen, Trupti Bavalatti, Xiaowen Lin, Yongkai Wang, Shuming Hu, Harihar
521 Subramanyam, Ksheeraj Sai Vepuri, Ming Jiang, Ji Qi, Li Chen, Nan Jiang, and Ankit Jain.
522 Class-RAG: Content Moderation with Retrieval Augmented Generation, December 2024a. URL
523 <https://arxiv.org/abs/2410.14881>. arXiv:2410.14881.
- 524 Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. Benchmarking Large Language Models in
525 Retrieval-Augmented Generation. *Proceedings of the AAAI Conference on Artificial Intelligence*,
526 38(16):17754–17762, March 2024b. doi: 10.1609/aaai.v38i16.29728. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29728>.
- 527
- 528 Xinyue Chen, Pengyu Gao, Jiangjiang Song, and Xiaoyang Tan. HiQA: A Hierarchical Contextual
529 Augmentation RAG for Multi-Documents QA, September 2024c. URL <https://arxiv.org/abs/2402.01767>. arXiv:2402.01767.
- 530
- 531
- 532 Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What Does BERT
533 Look at? An Analysis of BERT’s Attention. In Tal Linzen, Grzegorz Chrupała, Yonatan Be-
534 linkov, and Dieuwke Hupkes (eds.), *Proceedings of the 2019 ACL Workshop BlackboxNLP: An-
535 alyzing and Interpreting Neural Networks for NLP*, BlackboxNLP '19, pp. 276–286, Florence,
536 Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4828.
537 URL <https://aclanthology.org/W19-4828/>.
- 538 Benjamin Clavié, Antoine Chaffin, and Griffin Adams. Reducing the Footprint of Multi-Vector
539 Retrieval with Minimal Performance Impact via Token Pooling, September 2024. URL <https://arxiv.org/abs/2409.14683>. arXiv:2409.14683.

- 540 Columbia Accident Investigation Board. *Columbia Accident Investigation Board Report*, volume 2.
541 October 2003. URL <https://digital.library.unt.edu/ark:/67531/metadc1282015/>.
542
- 543 Bruce G. Coury, Vernon S. Ellingstad, and Joseph M. Kolly. Transportation Accident Investigation:
544 The Development of Human Factors Research and Practice. *Reviews of Human Factors and*
545 *Ergonomics*, 6(1):1–33, November 2010. doi: 10.1518/155723410X12849346788624. URL
546 <https://journals.sagepub.com/doi/10.1518/155723410X12849346788624>.
547
- 548 Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt,
549 Dasha Metropolitanaky, Robert Osazuwa Ness, and Jonathan Larson. From Local to Global: A
550 Graph RAG Approach to Query-Focused Summarization, February 2025. URL <https://arxiv.org/abs/2404.16130>. arXiv:2404.16130.
551
- 552 Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. RAGAs: Automated Eval-
553 uation of Retrieval Augmented Generation. In Nikolaos Aletras and Orphee De Clercq (eds.),
554 *Proceedings of the 18th Conference of the European Chapter of the Association for Computa-*
555 *tional Linguistics: System Demonstrations*, EACL ’24, pp. 150–158, St. Julians, Malta, March
556 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.eacl-demo.16/>.
557
- 558 Luyu Gao, Zhuyun Dai, and Jamie Callan. Rethink Training of BERT Rerankers in Multi-Stage
559 Retrieval Pipeline. In *Advances in Information Retrieval: Proceedings of the 43rd European*
560 *Conference on IR Research, Part II*, ECIR ’21, pp. 280–286, Virtual Event, March 2021. Springer.
561 ISBN 978-3-030-72239-5. doi: 10.1007/978-3-030-72240-1_26. URL [https://doi.org/10.1](https://doi.org/10.1007/978-3-030-72240-1_26)
562 [007/978-3-030-72240-1_26](https://doi.org/10.1007/978-3-030-72240-1_26).
563
- 564 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng
565 Wang, and Haofen Wang. Retrieval-Augmented Generation for Large Language Models: A Sur-
566 vey, March 2024. URL <https://arxiv.org/abs/2312.10997>. arXiv:2312.10997.
- 567 Rachael Gordon, Rhona Flin, and K. Mearns. Designing and Evaluating a Human Factors Inves-
568 tigation Tool (HFIT) for Accident Analysis. *Safety Science*, 43(3):147–171, March 2005. doi:
569 10.1016/j.ssci.2005.02.002. URL [https://www.sciencedirect.com/science/article/abs/](https://www.sciencedirect.com/science/article/abs/pii/S0925753505000068)
570 [pii/S0925753505000068](https://www.sciencedirect.com/science/article/abs/pii/S0925753505000068).
571
- 572 Rhys Gould, Euan Ong, George Ogden, and Arthur Conmy. Successor Heads: Recurring, Inter-
573 pretable Attention Heads in the Wild. In B. Kim, Y. Yue, S. Chaudhuri, K. Fragkiadaki, M. Khan,
574 and Y. Sun (eds.), *Proceedings of the Twelfth International Conference on Learning Representa-*
575 *tions*, ICLR ’24, pp. 9236–9261, Vienna, Austria, May 2024. International Conference on Learn-
576 ing Representations. URL [https://proceedings.iclr.cc/paper_files/paper/2024/hash](https://proceedings.iclr.cc/paper_files/paper/2024/hash/2722a0ccf6acfe3d144fdbb0dedd80b5-Abstract-Conference.html)
577 [/2722a0ccf6acfe3d144fdbb0dedd80b5-Abstract-Conference.html](https://proceedings.iclr.cc/paper_files/paper/2024/hash/2722a0ccf6acfe3d144fdbb0dedd80b5-Abstract-Conference.html).
- 578 Peter G. Harle. Investigation of Human Factors: The Link to Accident Prevention. In *Aviation*
579 *Psychology in Practice*, pp. 127–148. Routledge, London, UK, 1997. doi: 10.4324/9781351218
580 825-7. URL [https://www.taylorfrancis.com/chapters/edit/10.4324/9781351218825-](https://www.taylorfrancis.com/chapters/edit/10.4324/9781351218825-7/investigation-human-factors-link-accident-prevention-peter-harle)
581 [7/investigation-human-factors-link-accident-prevention-peter-harle](https://www.taylorfrancis.com/chapters/edit/10.4324/9781351218825-7/investigation-human-factors-link-accident-prevention-peter-harle).
582
- 583 Phu Mon Htut, Jason Phang, Shikha Bordia, and Samuel R. Bowman. Do Attention Heads in BERT
584 Track Syntactic Dependencies?, November 2019. URL <https://arxiv.org/abs/1911.12246>.
585 arXiv:1911.12246.
- 586 Haoyu Huang, Yongfeng Huang, Junjie Yang, Zhenyu Pan, Yongqiang Chen, Kaili Ma, Hongzhi
587 Chen, and James Cheng. HiRAG: Retrieval-Augmented Generation with Hierarchical Knowl-
588 edge, June 2025a. URL <https://arxiv.org/abs/2503.10150>. arXiv:2503.10150.
589
- 590 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong
591 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A Survey on Hallucination in
592 Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *ACM Trans.*
593 *Inf. Syst.*, 43(2):42:1–42:55, January 2025b. ISSN 1046-8188. doi: 10.1145/3703155. URL
<https://doi.org/10.1145/3703155>.

- 594 Po-Yao Huang, Xiaojun Chang, and Alexander Hauptmann. Multi-Head Attention with Diver-
595 sity for Learning Grounded Multilingual Multimodal Representations. In Kentaro Inui, Jing
596 Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical*
597 *Methods in Natural Language Processing and the 9th International Joint Conference on Natu-*
598 *ral Language Processing*, EMNLP-IJCNLP '19, pp. 1461–1467, Hong Kong, China, Novem-
599 ber 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1154. URL
600 <https://aclanthology.org/D19-1154/>.
- 601 Huggingface. Massive Text Embeddings Benchmark Leaderboard, 2025. URL <https://huggingface.co/spaces/mteb/leaderboard>. Accessed: 2025-09-21.
- 602
603
- 604 Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux, and Jason Weston. Poly-Encoders: Architec-
605 tures and Pre-Training Strategies for Fast and Accurate Multi-Sentence Scoring. In *Proceedings*
606 *of the Eighth International Conference on Learning Representations*, ICLR '20, Virtual Event,
607 April 2020. OpenReview. URL <https://openreview.net/forum?id=SkxgnnNFvH>.
- 608 Xinke Jiang, Rihong Qiu, Yongxin Xu, Wentao Zhang, Yichen Zhu, Ruizhe Zhang, Yuchen Fang,
609 Chu Xu, Junfeng Zhao, and Yasha Wang. RAGraph: A General Retrieval-Augmented Graph
610 Learning Framework. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak,
611 and C. Zhang (eds.), *Proceedings of the Thirty-Eighth Annual Conference on Neural Information*
612 *Processing Systems (NeurIPS '24)*, volume 37 of *Advances in Neural Information Processing*
613 *Systems*, pp. 29948–29985, Vancouver, Canada, December 2024. Curran Associates. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/34d6c7090bc5af0b96aea92fa074899-Abstract-Conference.html.
- 614
615
- 616 Madhu Kalia. Personalized Oncology: Recent Advances and Future Challenges. *Metabolism*, 62:
617 S11–S14, January 2013. doi: 10.1016/j.metabol.2012.08.016. URL <https://www.sciencedirect.com/science/article/abs/pii/S0026049512003204>.
- 618
619
- 620 Omar Khattab and Matei Zaharia. ColBERT: Efficient and Effective Passage Search via Context-
621 tualized Late Interaction over BERT. In *Proceedings of the 43rd International ACM SIGIR*
622 *Conference on Research and Development in Information Retrieval*, SIGIR '20, pp. 39–48, Vir-
623 tual Event, July 2020. Association for Computing Machinery. ISBN 9781450380164. doi:
624 10.1145/3397271.3401075. URL <https://doi.org/10.1145/3397271.3401075>.
- 625
626
- 627 Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. Revealing the Dark Se-
628 crets of BERT. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceed-*
629 *ings of the 2019 Conference on Empirical Methods in Natural Language Processing and the*
630 *9th International Joint Conference on Natural Language Processing*, EMNLP-IJCNLP '19, pp.
631 4365–4374, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:
632 10.18653/v1/D19-1445. URL <https://aclanthology.org/D19-1445/>.
- 633
634
- 635 Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catan-
636 zaro, and Wei Ping. NV-Embed: Improved Techniques for Training LLMs as Generalist Em-
637 bedding Models. In Y. Yue, A. Garg, N. Peng, F. Sha, and R. Yu (eds.), *Proceedings of*
638 *the Thirteenth International Conference on Learning Representations*, ICLR '25, pp. 79310–
639 79333, Singapore, April 2025. International Conference on Learning Representations. URL
640 https://proceedings.iclr.cc/paper_files/paper/2025/hash/c4bf73386022473a652a18941e9ea6f8-Abstract-Conference.html.
- 641
642
- 643 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
644 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe
645 Kiela. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In H. Larochelle,
646 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Proceedings of the Thirty-Fourth Annual*
647 *Conference on Neural Information Processing Systems (NeurIPS '20)*, volume 33 of *Advances in*
Neural Information Processing Systems, pp. 9459–9474, Virtual Event, December 2020. Curran
Associates. URL https://proceedings.neurips.cc/paper_files/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html.
- 648
649
- 650 Canjia Li, Andrew Yates, Sean MacAvaney, Ben He, and Yingfei Sun. PARADE: Passage Represent-
ation Aggregation for Document Reranking. *ACM Trans. Inf. Syst.*, 42(2):36:1–36:26, September
2023. ISSN 1046-8188. doi: 10.1145/3600088. URL <https://doi.org/10.1145/3600088>.

- 648 Ye Liu, Kazuma Hashimoto, Yingbo Zhou, Semih Yavuz, Caiming Xiong, and Philip Yu. Dense Hi-
649 erarchical Retrieval for Open-Domain Question Answering. In Marie-Francine Moens, Xuanjing
650 Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Findings of the Association for Computa-
651 tional Linguistics: EMNLP 2021*, pp. 188–200, Punta Cana, Dominican Republic, November
652 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-emnlp.19.
653 URL <https://aclanthology.org/2021.findings-emnlp.19/>.
- 654
655 Yuanjie Lyu, Zhiyu Li, Simin Niu, Feiyu Xiong, Bo Tang, Wenjin Wang, Hao Wu, Huanyong Liu,
656 Tong Xu, and Enhong Chen. CRUD-RAG: A Comprehensive Chinese Benchmark for Retrieval-
657 Augmented Generation of Large Language Models. *ACM Trans. Inf. Syst.*, 43(2):41:1–41:32,
658 January 2025. ISSN 1046-8188. doi: 10.1145/3701228. URL <https://doi.org/10.1145/3701228>.
- 660 Le Ma, Ran Zhang, Yikun Han, Shirui Yu, Zaitian Wang, Zhiyuan Ning, Jinghan Zhang, Ping Xu,
661 Pengjiang Li, Wei Wei Ju, et al. A Comprehensive Survey on Vector Database: Storage and
662 Retrieval Technique, Challenge, June 2025. URL <https://arxiv.org/abs/2310.11703>.
663 arXiv:2310.11703.
- 664
665 Sean MacAvaney, Andrew Yates, Arman Cohan, and Nazli Goharian. CEDR: Contextualized Em-
666 beddings for Document Ranking. In *Proceedings of the 42nd International ACM SIGIR Con-
667 ference on Research and Development in Information Retrieval, SIGIR '19*, pp. 1101–1104,
668 Paris, France, July 2019. Association for Computing Machinery. ISBN 9781450361729. doi:
669 10.1145/3331184.3331317. URL <http://doi.org/10.1145/3331184.3331317>.
- 670 Callum Stuart McDougall, Arthur Conmy, Cody Rushing, Thomas McGrath, and Neel Nanda. Copy
671 Suppression: Comprehensively Understanding a Motif in Language Model Attention Heads. In
672 Yonatan Belinkov, Najoung Kim, Jaap Jumelet, Hosein Mohebbi, Aaron Mueller, and Hanjie
673 Chen (eds.), *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neu-
674 ral Networks for NLP, BlackboxNLP '22*, pp. 337–363, Miami, FL, USA, November 2024. As-
675 sociation for Computational Linguistics. doi: 10.18653/v1/2024.blackboxnlp-1.22. URL
676 <https://aclanthology.org/2024.blackboxnlp-1.22/>.
- 677 Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. SFR-
678 Embedding-Mistral: Enhance Text Retrieval with Transfer Learning. Salesforce AI Research
679 Blog, October 2024. URL <https://www.salesforce.com/blog/sfr-embedding/>. Accessed:
680 2025-09-21.
- 681
682 Christina Messiou, Richard Lee, and Manuel Salto-Tellez. Multimodal Analysis and the Oncology
683 Patient: Creating a Hospital System for Integrated Diagnostics and Discovery. *Computational
684 and Structural Biotechnology Journal*, 21:4536–4539, 2023. ISSN 2001-0370. doi: 10.1016/j.cs
685 bj.2023.09.014. URL <https://www.sciencedirect.com/science/article/pii/S2001037023003264>.
- 686
687 Paul Michel, Omer Levy, and Graham Neubig. Are Sixteen Heads Really Better than One? In
688 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.),
689 *Proceedings of the Thirty-Third Annual Conference on Neural Information Processing Systems
690 (NeurIPS '19)*, volume 32 of *Advances in Neural Information Processing Systems*, pp. 14014–
691 14024, Vancouver, Canada, December 2019. Curran Associates. URL [https://proceedings.
692 neurips.cc/paper_files/paper/2019/hash/2c601ad9d2ff9bc8b282670cdd54f69f-Abst
693 ract.html](https://proceedings.neurips.cc/paper_files/paper/2019/hash/2c601ad9d2ff9bc8b282670cdd54f69f-Abst.html).
- 694
695 Jordan Multer, Joyce Ranney, Julie Hile, Thomas Raslear, and A. John. Developing an Effective
696 Corrective Action Process: Lessons Learned from Operating a Confidential Close Call Reporting
697 System. In Nastaran Dadashi, Anita Scott, John R. Wilson, and Ann Mills (eds.), *Rail Human
698 Factors: Supporting Reliability, Safety and Cost Reduction*, pp. 659–669. Taylor & Francis, Lon-
699 don, UK, 2013. URL <https://www.govinfo.gov/content/pkg/GOVPUB-TD3-PURL-gpo48075/pdf/GOVPUB-TD3-PURL-gpo48075.pdf>.
- 700
701 Rodrigo Nogueira and Kyunghyun Cho. Passage Re-Ranking with BERT, April 2020. URL <https://arxiv.org/abs/1901.04085>. arXiv:1901.04085.

- 702 Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. Document Ranking with a Pre-
703 trained Sequence-to-Sequence Model. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings*
704 *of the Association for Computational Linguistics: EMNLP 2020*, pp. 708–718, Virtual Event,
705 November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-em
706 nlp.63. URL <https://aclanthology.org/2020.findings-emnlp.63/>.
- 707 Office of Railroad Safety. Collaborative Incident Analysis and Human Performance Handbook.
708 Technical report, Federal Railroad Administration, May 2014. URL https://railroads.dot.gov/sites/fra.dot.gov/files/fra_net/14293/FRA_ORS_HumanPerformanceManual.pdf.
- 709 David O’Hare. The ‘Wheel of Misfortune’: A Taxonomic Approach to Human Factors in Accident
710 Investigation and Analysis in Aviation and Other Complex Systems. *Ergonomics*, 43(12):2001–
711 2019, December 2000. doi: 10.1080/00140130050201445. URL <https://www.tandfonline.com/doi/abs/10.1080/00140130050201445>.
- 712 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
713 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli,
714 Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane
715 Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,
716 and Chris Olah. In-Context Learning and Induction Heads, September 2022. URL <https://arxiv.org/abs/2209.11895>. arXiv:2209.11895.
- 717 OpenText. Cognitive Computing Reshapes Enterprise Decision-Making. Technical report, 2022.
718 URL <https://www.opentext.com/assets/documents/en-US/pdf/opentext-wp-cognitive-computer-reshapes-enterprise-decision-making-en.pdf>.
- 719 Nigel J. Packham. Columbia Crew Survival Investigation Report. Technical report, NASA, November
720 2017. URL <https://ntrs.nasa.gov/api/citations/20170011659/downloads/20170011659.pdf>.
- 721 James Jie Pan, Jianguo Wang, and Guoliang Li. Vector Database Management Techniques and
722 Systems. In *Companion of the 2024 International Conference on Management of Data, SIGMOD*
723 *’24*, pp. 597–604, Santiago AA, Chile, June 2024. Association for Computing Machinery. ISBN
724 9798400704222. doi: 10.1145/3626246.3654691. URL <https://doi.org/10.1145/3626246.3654691>.
- 725 Geondo Park, Chihye Han, Wonjun Yoon, and Daeshik Kim. MHSAN: Multi-Head Self-Attention
726 Network for Visual Semantic Embedding. In *Proceedings of the IEEE Winter Conference on*
727 *Applications of Computer Vision, WACV ’20*, pp. 1518–1526, Snowmass, CO, USA, March 2020.
728 IEEE Press. doi: 10.1109/WACV45572.2020.9093548. URL <https://ieeexplore.ieee.org/document/9093548>.
- 729 Zackary Rackauckas. RAG-Fusion: A New Take on Retrieval-Augmented Generation. *International*
730 *Journal on Natural Language Computing*, 13(1):37–47, February 2024. ISSN 2319-4111. doi:
731 10.5121/ijnlc.2024.13103. URL <https://doi.org/10.5121/ijnlc.2024.13103>.
- 732 J. Reason, Erik Hollnagel, and Jean Paries. Revisiting the Swiss Cheese Model of Accidents. Technical
733 Report 2006-017EEC Note 2006/13, Eurocontrol Experimental Centre, October 2006. URL
734 <https://www.eurocontrol.int/publication/revisiting-swiss-cheese-model-accidents>.
- 735 Guilherme Rosa, Luiz Bonifacio, Vitor Jeronymo, Hugo Abonizio, Marzieh Fadaee, Roberto Lotufo,
736 and Rodrigo Nogueira. In Defense of Cross-Encoders for Zero-Shot Retrieval, December 2022.
737 URL <https://arxiv.org/abs/2212.06121>. arXiv:2212.06121.
- 738 Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning.
739 RAPTOR: Recursive Abstractive Processing for Tree-Organized Retrieval. In B. Kim,
740 Y. Yue, S. Chaudhuri, K. Fragkiadaki, M. Khan, and Y. Sun (eds.), *Proceedings of the Twelfth*
741 *International Conference on Learning Representations, ICLR ’24*, pp. 32628–32649, Vienna,
742 Austria, May 2024. International Conference on Learning Representations. URL https://proceedings.iclr.cc/paper_files/paper/2024/hash/8a2acd174940dbca361a6398a4f9df91-Abstract-Conference.html.

- 756 Susav Shrestha, Narasimha Reddy, and Zongwang Li. ESPN: Memory-Efficient Multi-Vector In-
757 formation Retrieval. In *Proceedings of the 2024 ACM SIGPLAN International Symposium on*
758 *Memory Management*, ISMM 2024, pp. 95–107, Copenhagen, Denmark, June 2024. Association
759 for Computing Machinery. ISBN 9798400706158. doi: 10.1145/3652024.3665515. URL
760 <https://doi.org/10.1145/3652024.3665515>.
- 761 Devendra Singh, Siva Reddy, Will Hamilton, Chris Dyer, and Dani Yogatama. End-to-End
762 Training of Multi-Document Reader and Retriever for Open-Domain Question Answering. In
763 Joakim Nivre, Leon Derczynski, Filip Ginter, Bjørn Lindi, Stephan Oepen, Anders Søgaard,
764 and Jörg Tidemann (eds.), *Proceedings of the Thirty-Fifth Annual Conference on Neural Infor-*
765 *mation Processing Systems (NeurIPS '21)*, volume 34 of *Advances in Neural Information Pro-*
766 *cessing Systems*, pp. 25968–25981, Virtual Event, December 2021. Curran Associates. URL
767 [https://proceedings.neurips.cc/paper_files/paper/2021/hash/da3fde159d754a255](https://proceedings.neurips.cc/paper_files/paper/2021/hash/da3fde159d754a2555ea198d2d105b2-Abstract.html)
768 [5ea198d2d105b2-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2021/hash/da3fde159d754a2555ea198d2d105b2-Abstract.html).
- 769 Weihang Su, Yichen Tang, Qingyao Ai, Junxi Yan, Changyue Wang, Hongning Wang, Ziyi Ye, Yujia
770 Zhou, and Yiqun Liu. Parametric Retrieval Augmented Generation. In *Proceedings of the 48th*
771 *International ACM SIGIR Conference on Research and Development in Information Retrieval*,
772 SIGIR '25, pp. 1240–1250, Padua, Italy, July 2025. Association for Computing Machinery. ISBN
773 9798400715921. doi: 10.1145/3726302.3729957. URL [https://doi.org/10.1145/3726302.](https://doi.org/10.1145/3726302.3729957)
774 [3729957](https://doi.org/10.1145/3726302.3729957).
- 775 Venkatesh V, Mandeep Rathee, and Avishek Anand. SUNAR: Semantic Uncertainty Based Neighbor-
776 hood Aware Retrieval for Complex QA. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.),
777 *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for*
778 *Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, NAACL
779 '25, pp. 5818–5835, Albuquerque, NM, USA, April 2025. Association for Computational Lin-
780 guistics. ISBN 979-8-89176-189-6. URL [https://aclanthology.org/2025.naacl-long.3](https://aclanthology.org/2025.naacl-long.300/)
781 [00/](https://aclanthology.org/2025.naacl-long.300/).
- 782 Varun Vashisht, Samar Singh, Mihir Konduskar, Jaskaran Singh Walia, and Vukosi Marivate.
783 MAGE: Multi-Head Attention Guided Embeddings for Low Resource Sentiment Classification,
784 February 2025. URL <https://arxiv.org/abs/2502.17987>. arXiv:2502.17987.
- 785 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
786 Łukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In I. Guyon, U. Von Luxburg,
787 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Proceedings of the*
788 *Thirty-First Annual Conference on Neural Information Processing Systems (NIPS '17)*, volume 30
789 of *Advances in Neural Information Processing Systems*, pp. 5998–6008, Long Beach, CA, USA,
790 December 2017. Curran Associates. URL [https://proceedings.neurips.cc/paper_files](https://proceedings.neurips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html)
791 [/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html).
- 792 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing Multi-Head
793 Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned. In Anna
794 Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of*
795 *the Association for Computational Linguistics*, ACL '19, pp. 5797–5808, Florence, Italy, July
796 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1580. URL [https:](https://aclanthology.org/P19-1580/)
797 [//aclanthology.org/P19-1580/](https://aclanthology.org/P19-1580/).
- 798 George Wang, Jesse Hoogland, Stan van Wingerden, Zach Furman, and Daniel Murfet. Differentia-
799 tion and Specialization of Attention Heads via the Refined Local Learning Coefficient. In Y. Yue,
800 A. Garg, N. Peng, F. Sha, and R. Yu (eds.), *Proceedings of the Thirteenth International Con-*
801 *ference on Learning Representations*, ICLR '25, pp. 101037–101082, Singapore, April 2025a.
802 International Conference on Learning Representations. URL [https://proceedings.iclr.cc/p](https://proceedings.iclr.cc/paper_files/paper/2025/hash/fad7c708dda11f3e72cc1629bb130379-Abstract-Confere)
803 [aper_files/paper/2025/hash/fad7c708dda11f3e72cc1629bb130379-Abstract-Confere](https://proceedings.iclr.cc/paper_files/paper/2025/hash/fad7c708dda11f3e72cc1629bb130379-Abstract-Confere)
804 [nce.html](https://proceedings.iclr.cc/paper_files/paper/2025/hash/fad7c708dda11f3e72cc1629bb130379-Abstract-Confere).
- 805 Jianguo Wang, Xiaomeng Yi, Rentong Guo, Hai Jin, Peng Xu, Shengjun Li, Xiangyu Wang, Xi-
806 angzhou Guo, Chengming Li, Xiaohai Xu, Kun Yu, Yuxing Yuan, Yinghao Zou, Jiquan Long,
807 Yudong Cai, Zhenxiang Li, Zhifeng Zhang, Yihua Mo, Jun Gu, Ruiyi Jiang, Yi Wei, and
808 Charles Xie. Milvus: A Purpose-Built Vector Data Management System. In *Proceedings of*
809

- 810 *the 2021 International Conference on Management of Data, SIGMOD '21*, pp. 2614–2627, Vir-
811 tual Event, June 2021. Association for Computing Machinery. ISBN 9781450383431. doi:
812 10.1145/3448016.3457550. URL <https://doi.org/10.1145/3448016.3457550>.
- 813
814 Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Ma-
815 jumder, and Furu Wei. Text Embeddings by Weakly-Supervised Contrastive Pre-Training, Febru-
816 ary 2024. URL <https://arxiv.org/abs/2212.03533>. arXiv:2212.03533.
- 817 Zhiming Wang, Kaisheng Yao, Xiaolong Li, and Shuo Fang. Multi-Resolution Multi-Head Attention
818 in Deep Speaker Embedding. In *Proceedings of the IEEE International Conference on Acoustics,*
819 *Speech and Signal Processing, ICASSP '20*, pp. 6464–6468, Barcelona, Spain, May 2020. IEEE
820 Press. doi: 10.1109/ICASSP40776.2020.9053217. URL [https://ieeexplore.ieee.org/do](https://ieeexplore.ieee.org/document/9053217)
821 [cument/9053217](https://ieeexplore.ieee.org/document/9053217).
- 822
823 Zichong Wang, Zhibo Chu, Thang Viet Doan, Shiwen Ni, Min Yang, and Wenbin Zhang. History,
824 Development, and Principles of Large Language Models: An Introductory Survey. *AI Ethics*, 5:
825 1955–1971, June 2025b. ISSN 2730-5961. doi: 10.1007/s43681-024-00583-7. URL [https://](https://link.springer.com/article/10.1007/s43681-024-00583-7)
826 link.springer.com/article/10.1007/s43681-024-00583-7.
- 827 Zikai Wang, Qianxi Zhang, Baotong Lu, Qi Chen, and Cheng Tan. Towards Robustness: A Critique
828 of Current Vector Database Assessments, July 2025c. URL [https://arxiv.org/abs/2507.0](https://arxiv.org/abs/2507.00379)
829 [0379](https://arxiv.org/abs/2507.00379). arXiv:2507.00379.
- 830 Jinfeng Xiao, Linyi Ding, James Barry, Mohab Elkaref, Geeth De Mel, and Jiawei Han. ORAG:
831 Ontology-Guided Retrieval-Augmented Generation for Theme-Specific Entity Typing. In *Pro-*
832 *ceedings of the First Conference on Language Modeling, COLM '24*, Philadelphia, PA, USA,
833 October 2024. OpenReview. URL <https://openreview.net/forum?id=cKBmZ2PZ6c>.
- 834
835 Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. Benchmarking Retrieval-Augmented
836 Generation for Medicine. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of*
837 *the Association for Computational Linguistics: ACL 2024*, pp. 6233–6251, Bangkok, Thailand,
838 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.37
839 [2](https://aclanthology.org/2024.findings-acl.372/). URL <https://aclanthology.org/2024.findings-acl.372/>.
- 840
841 Zhipeng Xu, Zhenghao Liu, Yibin Liu, Chenyan Xiong, Yukun Yan, Shuo Wang, Shi Yu, Zhiyuan
842 Liu, and Ge Yu. ActiveRAG: Revealing the Treasures of Knowledge via Active Learning, October
2024. URL <https://arxiv.org/abs/2402.13547>. arXiv:2402.13547.
- 843
844 Huiyin Xue and Nikolaos Aletras. Pit One Against Many: Leveraging Attention-Head Embeddings
845 for Parameter-Efficient Multi-Head Attention. In Houda Bouamor, Juan Pino, and Kalika Bali
846 (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 10355–
847 10373, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v
848 [1/2023.findings-emnlp.695](https://aclanthology.org/2023.findings-emnlp.695/). URL <https://aclanthology.org/2023.findings-emnlp.695/>.
- 849
850 Biwei Yan, Kun Li, Minghui Xu, Yueyan Dong, Yue Zhang, Zhaochun Ren, and Xiuzhen Cheng. On
851 Protecting the Data Privacy of Large Language Models (LLMs) and LLM Agents: A Literature
852 Review. *High-Confidence Computing*, 5(2):100300:1–100300:21, June 2025. ISSN 2667-2952.
853 doi: <https://doi.org/10.1016/j.hcc.2025.100300>. URL [https://www.sciencedirect.com/scie](https://www.sciencedirect.com/science/article/pii/S2667295225000042)
854 [nce/article/pii/S2667295225000042](https://www.sciencedirect.com/science/article/pii/S2667295225000042).
- 855
856 Chening Yang, Duy-Khanh Vu, Minh-Tien Nguyen, Xuan-Quang Nguyen, Linh Nguyen, and Hung
857 Le. SuperRAG: Beyond RAG with Layout-Aware Graph Modeling. In Weizhu Chen, Yi Yang,
858 Mohammad Kachuee, and Xue-Yong Fu (eds.), *Proceedings of the 2025 Conference of the Na-*
859 *tions of the Americas Chapter of the Association for Computational Linguistics: Human Lan-*
860 *guage Technologies (Volume 3: Industry Track), NAACL '25*, pp. 544–557, Albuquerque, NM,
861 USA, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-194-0. URL
862 <https://aclanthology.org/2025.naacl-industry.45/>.
- 863
864 Wenhao Yu, Hongming Zhang, Xiaoman Pan, Peixin Cao, Kaixin Ma, Jian Li, Hongwei Wang, and
865 Dong Yu. Chain-of-Note: Enhancing Robustness in Retrieval-Augmented Language Models. In
866 Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference*
on Empirical Methods in Natural Language Processing, EMNLP '24, pp. 14672–14685, Miami,

864 FL, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.e
865 mnlp-main.813. URL <https://aclanthology.org/2024.emnlp-main.813/>.
866

867 Jihao Zhao, Zhiyuan Ji, Zhaoxin Fan, Hanyu Wang, Simin Niu, Bo Tang, Feiyu Xiong, and Zhiyu
868 Li. MoC: Mixtures of Text Chunking Learners for Retrieval-Augmented Generation System. In
869 Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Pro-
870 ceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume
871 1: Long Papers)*, ACL '25, pp. 5172–5189, Vienna, Austria, July 2025a. Association for Com-
872 putational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.258. URL
873 <https://aclanthology.org/2025.acl-long.258/>.

874 Jihao Zhao, Zhiyuan Ji, Yuchen Feng, Pengnian Qi, Simin Niu, Bo Tang, Feiyu Xiong, and Zhiyu Li.
875 Meta-Chunking: Learning Text Segmentation and Semantic Completion via Logical Perception,
876 May 2025b. URL <https://arxiv.org/abs/2410.12788>. arXiv:2410.12788.

877 Konstantinos G. Zografos, Michael A. Madas, and Yiannis Salouras. A Decision Support System
878 for Total Airport Operations Management and Planning. *Journal of Advanced Transportation*, 47
879 (2):170–189, March 2013. doi: 10.1002/atr.154. URL [https://onlinelibrary.wiley.com/
880 doi/ftr/10.1002/atr.154](https://onlinelibrary.wiley.com/doi/ftr/10.1002/atr.154).
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

APPENDIX

A REVIEW OF MULTI-ASPECTUALITY IN INDUSTRY

Real-world decision-making and analysis tasks often require a *holistic integration of multiple, semantically distinct information sources*. Literature from diverse domains—from safety investigations to diagnostics and enterprise analytics—consistently emphasizes that no single data stream suffices for complex tasks. Instead, success comes from combining heterogeneous elements (design details, human factors, environmental context, etc.) into a unified understanding. Below we highlight several examples that support this multi-aspectual reasoning framework.

In **safety-critical environments**, accident and incident investigations rely on aggregating information from disparate sources. For example, a U.S. Federal Railroad Administration report stresses that collecting multiple sources of information is essential for effective event reconstruction (Office of Railroad Safety, 2014). The Columbia Accident Investigation Board similarly combined telemetry, debris forensics, environmental data, and even autopsy reports to reconstruct the chain of events behind the Space Shuttle Columbia disaster (Columbia Accident Investigation Board, 2003). These investigations exemplify the need to integrate technical, procedural, and human-centered aspects of data to obtain actionable conclusions.

In **medicine**, particularly oncology, the concept of *integrated diagnostics* exemplifies multi-aspect decision-making. Physicians routinely combine patient histories, radiological scans, lab values, pathology slides, and genetic tests to form a diagnosis. This is no longer optional: as modern datasets become more complex, integrating semantically distinct modalities has become necessary to reach accurate, personalized outcomes (Messiou et al., 2023; Kalia, 2013).

The same applies to **enterprise settings**. For example, airport operations management integrates foot traffic sensors, weather feeds, and gate schedules – among others – to optimize personnel allocation and prevent congestion (Zografos et al., 2013). Similarly, in legal and financial firms, cross-silo systems integrate internal metrics (e.g., billing and staffing) with external sources (e.g., news, contracts, social media) to guide decision-making and strategic planning (OpenText, 2022). These examples show that modern organizational workflows demand the integration of multiple semantically distinct data sources.

In summary, across domains like industrial safety, healthcare, and business intelligence, it is widely recognized that *multi-aspectuality*—combining diverse, independently relevant information fragments—is essential for accurate and effective decision-making (Office of Railroad Safety, 2014; Columbia Accident Investigation Board, 2003; Reason et al., 2006; Kalia, 2013; Packham, 2017; Messiou et al., 2023; Multer et al., 2013; Coury et al., 2010; Harle, 1997; O’Hare, 2000; Gordon et al., 2005; Bridger, 2021; OpenText, 2022). In all these cases, the embeddings of documents from such divergent subdomains would be far away from one another in the embedding space when using standard RAG pipelines (as also confirmed by our own datasets in legal and plant accident use cases, see Section 5), underlying the relevance of MRAG.

B MATHEMATICAL FORMULATION: ADDITIONAL DETAILS

We omit, for clarity, unnecessary details such as layer normalizations. The output of attention head h for the i th token x_i is defined as follows (Vaswani et al., 2017):

$$\text{head}^h(\mathbf{x}_i) = \sum_j w_{ij} \mathbf{v}_j^h, \quad \text{where} \quad (1)$$

$$w_{ij} = \text{softmax} \left((\mathbf{q}_i^h)^T \mathbf{k}_j^h \right), \quad \mathbf{q}_i^h = \mathbf{W}_q^h \mathbf{x}_i, \quad \mathbf{k}_j^h = \mathbf{W}_k^h \mathbf{x}_j, \quad \mathbf{v}_j^h = \mathbf{W}_v^h \mathbf{x}_j \quad (2)$$

where $\mathbf{W}_q^h, \mathbf{W}_k^h, \mathbf{W}_v^h$ are, respectively, learnable query, key, and value projections associated with head h , and \mathbf{x}_j is the vector embedding of the j th token x_j . These outputs get combined to form the output of the i th multi-head attention block as follows:

$$\text{MHA}(\mathbf{x}_i) = \mathbf{W}_o \text{concat}(\text{head}^1(\mathbf{x}_i), \dots, \text{head}^h(\mathbf{x}_i))^T \quad (3)$$

where matrix \mathbf{W}_o is the linear layer that combines the outcomes of all attention heads.

B.1 STANDARD RAG EMBEDDING

In standard RAG, a single embedding vector $\mathbf{e}_{\text{std}} \in \mathbb{R}^d$ is computed for a chunk by extracting the decoder output for the final token x_n after the final feed-forward layer:

$$\mathbf{e}_{\text{std}} = \text{FFN}(\text{MHA}(x_n)) \quad (4)$$

B.2 MULTI-HEAD RAG EMBEDDING

Instead of compressing all head outputs into a single embedding, MRAG leverages the individual output of each head on the final token x_n :

$$\mathcal{S} = \{\mathbf{e}_k = \text{head}^k(\mathbf{x}_n) \in \mathbb{R}^{d/h} \mid k = 1, \dots, h\} \quad (5)$$

This results in h head-wise embeddings per chunk, capturing diverse semantic aspects. Crucially, this design avoids any increase in memory or compute during inference, as head-level vectors are computed as part of the standard MHA process.

1026 C SYSTEM DESIGN & IMPLEMENTATION: ADDITIONAL DETAILS

1027

1028 We provide addition details on system design and implementation.

1029

1029 C.1 RANKING STRATEGY DETAILS

1030

1030 The scoring of embedding spaces is detailed in Algorithm 1.

1031

1032

1032 **Algorithm 1** Importance scores for heads.

1033

1033 **for** each head h_i **do**

1034

1034 $a_i, b_i, count_a_i, count_b_i \leftarrow 0$

1035

1035 **for** each embedding e_{ij} in h_i **do**

1036

1036 $a_i \leftarrow a_i + \|e_{ij}\|$

1037

1037 $count_a_i \leftarrow count_a_i + 1$

1038

1038 **for** subset of m embeddings e_{ih} sampled uniformly at random **do**

1039

1039 $b_i \leftarrow b_i + \text{cosine-distance}(e_{ij}, e_{ih})$

1040

1040 $count_b_i \leftarrow count_b_i + 1$

1041

1041 **end for**

1042

1042 **end for**

1043

1043 $a_i \leftarrow a_i / count_a_i; b_i \leftarrow b_i / count_b_i$

1044

1044 $s_i \leftarrow a_i \cdot b_i$

1045

1045 **end for**

1046

1047

1047 C.2 RANKING STRATEGY DETAILS

1048

1048 The voting strategy used by MRAG in its reranker is pictured in Algorithm 2.

1049

1050

1050 **Algorithm 2** Voting strategy.

1051

1051 $l \leftarrow []$

1052

1052 **for** each head h_i and its score s_i **do**

1053

1053 find best matching c text chunks

1054

1054 **for** each chunk $d_{i,p}$ with index p in top c **do**

1055

1055 $w_{i,p} \leftarrow s_i \cdot 2^{-p}$

1056

1056 add tuple $(d_{i,p}, w_{i,p})$ to l

1057

1057 **end for**

1058

1058 **end for**

1059

1059 sort l using weights $w_{i,p}$

1060

1060 return top k elements of l

1061

1062

1063

1063 C.3 INTEGRATION WITH DATA STORES

1064

1064 MRAG can be seamlessly used with different classes of data stores  and nearest neighbor (NN) search approaches. It can be combined with both the exact and the approximate NN to find the matching (embedding, chunk)-pairs. These two parts of the broader RAG processing pipeline are orthogonal to MRAG.

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

D SPECIFICATION OF PROMPTS

D.1 PROMPT TEMPLATE FOR READER

Prompt Template for Reader

You are presented with a series of articles, each potentially addressing different aspects of a topic.

1. <article title>
 [<metadata>]
 <body>

2. <article title>
 [<metadata>]
 <body>

<...>

Task: Carefully analyze the articles. When formulating your response to the question below, identify the relevant aspects or claims made within each article. Construct your answer by comparing, contrasting, or synthesizing these points in a coherent and logically structured manner. Your response should be supported by specific references to the content of the articles. Where applicable, acknowledge differences in perspectives, data, or assumptions across the sources. Aim for clarity, precision, and concise reasoning grounded in the evidence provided.

Please answer the following question: <query text>

D.2 PROMPT TEMPLATE FOR THE SYNTHETIC DATASET GENERATION

Prompt Template for Query Generation

Please create a story about the attached <number of articles> articles on the topics <list of titles>.

It is very important that each of the attached articles is relevant to the story, in a way that references the content of the article, not just its title. But please also mention each title at least once. Please make sure that all of the attached articles are relevant to your story, and that each article is referenced in at least two sentences! They do not necessarily have to be referenced in the same order, but make sure no article is forgotten.

Important: Output only the story, no additional text. And do not use bullet points, or paragraphs.

Articles:

Article <title>:
 <body>

<...>

Again, make sure that you reference all the following topics in your story: <list of titles>

E MULTI-ASPECTUALITY WITH ATTENTION HEADS WITHOUT ADDITIONAL TRAINING

In MRAG, we extract embeddings from the hidden representations immediately after the attention block in the last decoder layer, avoiding any fine-tuning. This decision is based on an existing hypothesis that attention heads in Transformer models naturally differentiate during training, each attending to distinct aspects of the input data distribution.

E.1 LITERATURE SURVEY

This hypothesis has been substantiated across various Transformer families. Wang et al. (2025a) introduce the *Local Learning Coefficient* (LLC), a measure of training dynamics at the head level. They show that attention heads begin with similar behavior but quickly diverge into functionally distinct clusters, each specializing in different patterns, ranging from local structure to multigram token groups. Olsson et al. (2022) identified “induction heads” in GPT-style models: specific heads that attend to earlier repeated sequences (e.g., in patterns like “X ... Y ... X”), enabling the model to learn simple in-context repetition *without additional supervision*. Further studies observed heads that consistently attend to names, suppress repetition, or shift attention predictably to structurally aligned tokens (McDougall et al., 2024; Gould et al., 2024).

Research in BERT-style models further supports these trends. Clark et al. (2019) demonstrate that different heads attend to direct objects, nominal modifiers, or punctuation tokens. Kovaleva et al. (2019) identify broad attention patterns like “vertical” heads (focusing on special tokens) and “diagonal” heads (attending to adjacent words). Htut et al. (2019) show that many heads correlate with syntactic dependency arcs.

Notably, while many heads specialize in meaningful ways, others appear to contribute little to model performance. Voita et al. (2019) and Michel et al. (2019) show that a large fraction of heads can be pruned without substantial performance loss, suggesting that specialization tends to concentrate in a smaller subset of effective heads.

E.2 ANALYSIS OF MULTI-HEAD PATTERNS

We investigated the attention heads of two models in detail: LLaMA-2 7B and SFR-Embedding-Mistral. We selected these two models for a detailed investigation because the former represents models that are not fine-tuned for text embeddings, while the latter is specifically the text embedding model that we used for our experiments. For each model, we looked specifically at the attention scores within each attention head, i.e., how much attention each head pays to each input token during the inference. Knowing the semantics of the input tokens enables then deriving certain conclusions about multi-aspectuality and attention heads.

We plot selected results in Figure 7. Each heatmap shows the dot-product between key- and value-projections inside a given specified attention head, where line i of a heatmap for attention head h indicates the dot-products between the query-projection of token i and the key-projections of all previous tokens $j < i$ (both models use causal attention).

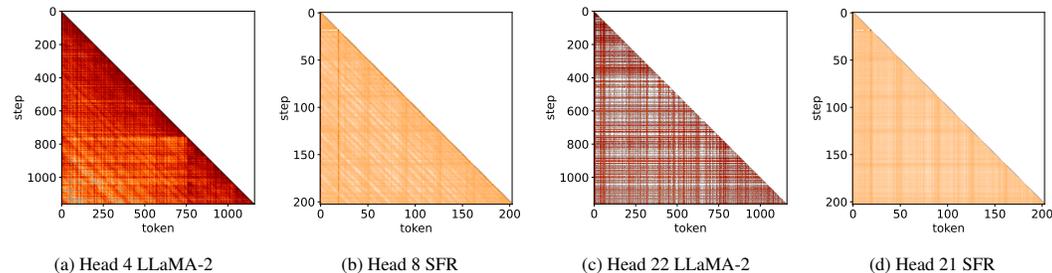


Figure 7: Heatmap plots for selected attention heads of the LLaMA-2 7B and SFR-Embedding-Mistral models.

For both models, we found out that the attention patterns vary significantly between the different attention heads. Still, we encountered two distinct patterns. First, the diagonal lines in Figures 7a and 7b indicate that, when processing a certain input token x , elevated attention is paid to some tokens that came a constant numbers of steps before x . We postulate that this pattern is likely beneficial to understanding the overall rhythm of a natural language, allowing the model to better

identify which words are semantically connected, and which parts of the input text refer to each other. Second, horizontal and vertical lines in Figures 7c and 7d show that these heads learned to pay attention to specific tokens, regardless of how far apart they are within the input sequence. An intuitive justification for such patterns is the focus on certain semantic aspects of the input sequence.

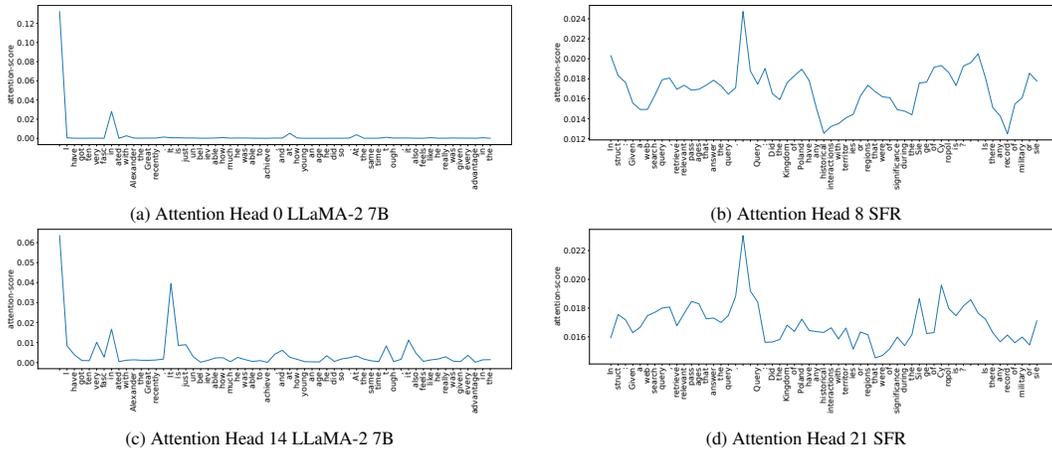


Figure 8: Attention scores for selected attention heads of the LLaMA-2 7B and SFR-Embedding-Mistral models.

We also detail attention scores (after applying softmax) of selected heads in Figures 8 and 9, when the model is processing the last token of its input. We see that some tokens gather a lot of attention from most heads, yet there is always a plethora of passages which are attended differently by any two attention heads. An interesting pattern we encountered was that for the SFR-Embedding-Mistral model (see Figure 9), all heads’ attention spiked significantly on the first line-break in the input sequence - either positively or negatively. We conjecture that this is a consequence of how the embedding model was fine-tuned and its intended usage pattern: embedding queries are usually prepended with a retrieval instruction, which is terminated by a line-break. The model likely learnt to summarise the necessary information about this instruction inside the terminating line-break.

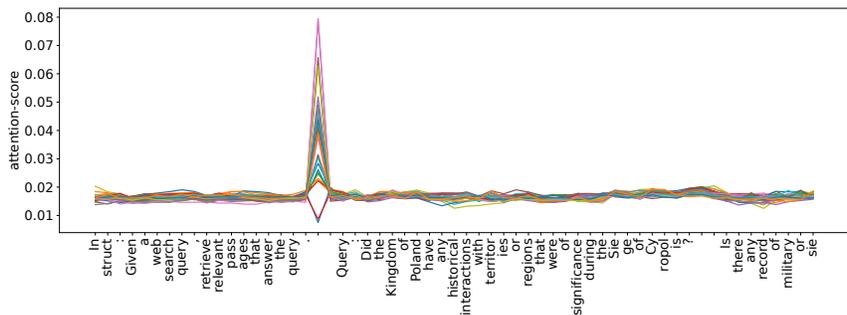


Figure 9: Attention scores for all attention heads of the SFR-Embedding-Mistral model.

F COMPLEXITY ANALYSIS: ADDITIONAL DETAILS

We now discuss additional details for our complexity analysis of the RAG schemes.

In Poly-encoders (Humeau et al., 2020) retrieval, we use $s(n)$ to represent additional attention evaluations needed to prepare the m code embeddings and their final context embeddings with all n documents.

In Lewis et al. (2020) retrieval, the work and depth are increased by the additional generator model evaluation.

In ColBERT (Khattab & Zaharia, 2020) retrieval, the work is increased linearly and depth logarithmically with the average query and document size due to the `maxsim` evaluations. For storage, the overhead scales linearly with the average size of the document, due to per-token embeddings.

In EMDR (Singh et al., 2021) retrieval, the work is increased $k + 1$ times due to T5 encoders run on the top k documents. This corresponds to a D_m depth increase as the evaluations can be parallelized.

In Self-RAG (Asai et al., 2024), once the model decides to conduct retrieval using the retrieval token, it embeds the query, runs nearest k search, and evaluates the model twice for each best- k document to predict the `issup`, `isrel`, `isuse` labels. The depth increases only by D_m as the evaluations can be parallelized.

In Chain-of-Note (Yu et al., 2024), during retrieval the k best documents are fetched and then each is summarized sequentially by generating notes, resulting in a $ks(n)$ work and depth increase, with $s(n)$ representing the cost to generate an average note.

In RAPTOR (Sarathi et al., 2024), preprocessing involves creating a document tree, with leaves representing documents, and parents containing summaries of children, a cost we represent with $s(n)$ both for work and storage. Retrieval involves traversing the tree, which we encapsulate in a different $s(n)$.

RAGraph (Jiang et al., 2024) preprocessing involves embedding the documents by a graph model and creating the toy graphs, which we account for using W_e and D_e . The toy graphs also create additional space requirements, which we account for using $s(n)$. For retrieval, additional key information such as environment or position-aware codes is used, represented by another $s(n)$ in our notation.

In RQ-RAG (Chan et al., 2024), the retrieval may be decomposed into multiple separate queries that are generated by the model and then evaluated in sequence using standard techniques, resulting in increases in work and depth that we denote using $s(n)$.

In ActiveRAG (Xu et al., 2024), retrieval is conducted using standard techniques, followed by three separate model evaluations: the Knowledge Assimilation (KA), Self-Inquiry (Q), and Thought Accommodation (TA) agents. As TA is similar to the work of an LLM answering the query, we omit its cost, like in all other frameworks, resulting in triple the work. Q and KA can be evaluated in parallel, increasing the depth only twice.

HiQA (Chen et al., 2024c) combines multiple retrieval strategies, using vector similarity matching, elastic search with BM25, and keyword matching. We represent these with $s(n)$ for work. Note that as these are independent, the depth is not increased. For preprocessing, HiQA uses a Hierarchical Contextual Augmentor, which creates a data hierarchy and introduces an overhead we denote by $s(n)$. HiQA also stores more information alongside vectors, such as keywords, increasing requirements by $s(n)$.

GraphRAG (Edge et al., 2025) creates a knowledge graph during preprocessing, which we denote with W_e and D_e . It then extracts communities and their summaries, which we account for using $s(n)$. These are stored and require additional space also denoted by a different $s(n)$. During retrieval, GraphRAG uses an LLM to rank all communities in parallel by how useful they are which we estimate using $s(n)$.

In Fusion RAG (Rackauckas, 2024) retrieval, k queries are generated and for each a standard RAG is evaluated, together with a reranking achieved by another model evaluation, which we estimate as $s(n)$.

1296 In Meta-chunking (Zhao et al., 2025b), the documents are preprocessed by splitting them into chunks
1297 based on perplexity or margin sampling. As the decision whether to split or combine sentences in
1298 a document is based on an LLM, the preprocessing requires l_d evaluations with some $s(n)$ postpro-
1299 cessing. The storage requirements are also increased by a different $s(n)$ as documents are stored not
1300 as a single vector but multiple vectors based on chunks. For the same reason, retrieval requires more
1301 work, as the number of vectors is increased by $s(n)$.

1302 In MoC (Zhao et al., 2025a), chunks are created based on different granularities. Similarly to Meta-
1303 chunking, this means a larger number of vectors resulting in increased storage requirements, and
1304 retrieval work by $s(n)$. As MoC also includes routing and meta-chunkers, the preprocessing cost is
1305 increased by $s(n)$.

1306 In Parametric RAG (Su et al., 2025), documents are represented as deltas of parameters that can be
1307 applied to the model. During preprocessing, the model is fine-tuned on a given document, resulting
1308 in a considerable increase of $s(n)$ in work and depth. As parameters are considerably larger than
1309 the hidden dimension, storage is also increased by a different $s(n)$. During retrieval, standard RAG
1310 fetches the documents, and the original model needs to be updated with their parameters, increasing
1311 work and depth by $s(n)$.

1312 SuperRAG (Yang et al., 2025) first embeds the documents in a knowledge graph, which we represent
1313 by W_e and D_e , and then uses this knowledge graph in retrieval to index it using W_i and D_i . These
1314 also include any reranking SuperRAG might do.

1315 HiRAG (Huang et al., 2025a) creates a hierarchical knowledge graph used for indexing in retrieval.
1316 Similarly, to SuperRAG, this increases the preprocessing and retrieval costs. We include in these the
1317 additional description and report generation that HiRAG conducts.
1318

1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

G FULL SPECIFICATION OF BENCHMARKING MULTI-ASPECTUALITY

We provide more details on how to benchmark multi-aspectuality in RAG. Figure 10 shows an example query and metrics usage. Each query requires retrieving a specific number of documents and the corresponding non-overlapping categories which define the ground truth. We fetch the top k documents from a database, where k is the “total number of documents fetched for a tested RAG scheme” (including potentially mismatches). Among these k documents, we search for matches with the ground truth.

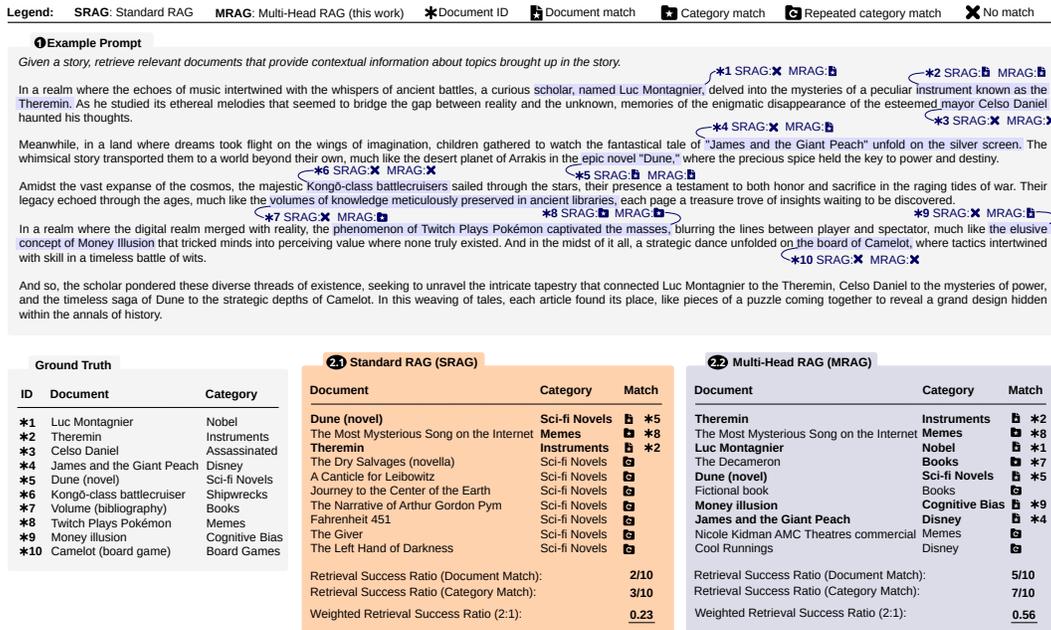


Figure 10: An example query used to evaluate different RAG strategies. We mention the documents to be fetched in the text and then assess the success ratio of different RAG strategies in finding these documents and their categories. We mark exact document matches 📄, category matches 📁, documents that match a category multiple times 🔄, and text segments with no matching document ✖. Finally, we show the weighted success ratio for each strategy, taking a 2:1 weighting (prioritizing the exact article matches).

G.1 MULTI-ASPECT DATASETS

We first select conceptually different categories of documents for a synthetic dataset. Here, we harness publicly available Wikipedia articles. In the dataset construction pipeline, the user selects a given number of categories (e.g., countries, board games, historical swords, shipwrecks, etc.) and then, for each category, they sample a specified number of documents. The first part of the document (overview) is used as a text chunk to be embedded. We enforce that each overview must have at least 800 characters, matching commonly used chunk sizes in RAG schemes. We also use multi-aspect **real-world inspired datasets** consisting of NDAs and reports describing industry accidents in chemical processing plants. We ensure the usefulness of these datasets by working directly with tech leaders from 3 corporations that rely on RAG in their in-house LLM-driven report generation and analytics frameworks. Example categories of the legal documents are legal areas (energy law, family law, criminal law, etc.) or document language style (aggressive, mild, neutral, etc.). Examples of accident causes are natural disasters, human mistakes, or lack of proper training. We fully release these datasets to propel RAG research. Details on all three datasets can be found in the Appendix H.2. In our evaluation, we use a total of 16,500 documents.

G.2 MULTI-ASPECT QUERY GENERATION

We also require queries that touch upon a given *number of n aspects*. For example, a query with 10 aspects must contain a question about 10 different documents from 10 different categories. We create such queries by selecting n categories, sampling a document from each selected category (ensuring there are no duplicates overall), and then generating a story that combines these documents, using an LLM (GPT-4o). We construct 160 queries with 1, 2, 3, 4, 5, 6, 10, 15, 20 and 25 aspects (1600 queries in total). An example multi-aspect query sent to the LLM that requires retrieving 10 documents from 10 different categories, is pictured in the top part of Figure 10.

1404 G.3 METRICS

1405 We also design novel metrics to assess how well a given RAG scheme supports multi-aspectuality.
 1406 For a query Q , a used Reranker scheme S (detailed in Section 2.3), and n documents from n categories
 1407 to retrieve, Q_{rel} denotes the *ideal* set of documents that should be retrieved for Q . Then,
 1408 $S(Q, n)$ is the set of the *actually* retrieved documents. We define the *Retrieval Success Ratio* as
 1409 $\Xi(Q, n) = \frac{|S(Q, n) \cap Q_{rel}|}{|Q_{rel}|}$, i.e., the ratio of successfully retrieved relevant documents. Moreover,
 1410 there is a case when a RAG scheme does not retrieve the *exact* desired document, but it still retrieves
 1411 successfully *some other document* from *the same* category. While less desired, it still increases
 1412 chances for a more accurate LLM answer following the retrieval. For example, when asking the
 1413 LLM to determine the cause of an industry accident, fetching the documents in the same category as
 1414 the accident being queried about, improves the chances for the LLM to give a more relevant answer.
 1415 To consider such cases, we use another measure, the **Category Retrieval Success Ratio** or Ξ_c . It
 1416 has the same form as $\Xi(Q, n)$ above, with one difference: $S(Q, n)$ is now the set of all the retrieved
 1417 documents that belong to categories of the ideal desired documents. Finally, to combine these two
 1418 metrics, we use the **Weighted Retrieval Success Ratio** Ξ_w as $\Xi_w = \frac{w \cdot \Xi + \Xi_c}{w + 1}$.

1419 An example of using these metrics to assess how well MRAG and Standard RAG capture multi-
 1420 aspectuality is pictured in the bottom part of Figure 10.
 1421

1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

1458 H EVALUATION SETUP: ADDITIONAL DETAILS

1459 H.1 COMPUTE RESOURCES

1461 Our experiments were executed with compute nodes containing 4x NVIDIA GH200 and a total
 1462 memory of 800 GB. In general one GPU with at least 40GB of memory should suffice. We used
 1463 at most 50GB of storage and the OpenAI API as an external resource. The full experiments took
 1464 at most three hours of GPU time and the cost for the OpenAI API were at most \$15. We carried
 1465 out additional experiments, which amounted to around 20 hours of GPU time and cost of \$25 for
 1466 the OpenAI API. Additional evaluation was executed with a mix of compute resources including
 1467 NVIDIA A100 and V100 GPUs.

1468 H.2 DATASET DETAILS

1469 Table 4: Overview of the structure and the number of documents in the respective datasets.
 1470

| 1471 dataset | 1472 #categories | 1473 #topics | 1474 #documents | 1475 total #documents |
|-----------------------|------------------|-----------------------------------|-----------------|-----------------------|
| 1476 Wikipedia | 1477 80 | 1478 50 documents per category | 1479 | 1480 4000 |
| 1481 Legal Documents | 1482 25 | 1483 25 per category 10 per topic | 1484 | 1485 6250 |
| 1486 Accident Reports | 1487 25 | 1488 25 per category 10 per topic | 1489 | 1490 6250 |

I EVALUATION: ADDITIONAL RESULTS

We provide additional empirical evaluation results.

I.1 HARNESSING DIFFERENT DECODER BLOCKS

We analyze the impact of using embeddings from **different decoder blocks** for MRAG (instead of the last one). Here, we consider taking multi-aspect embeddings from three different layers of the embedding model: after the first multi-head attention block, after multi-head attention block 16 (in the middle of the decoder architecture), and the final multi-head attention. We discover that the last multi-head attention performs the best when compared with the Standard RAG.

I.2 ANALYZING DIFFERENT VOTING STRATEGIES

We also illustrate selected representative data from a long investigation into two **additional voting strategies** for MRAG. We compare **MRAG (1)** where only the exponential lowering of significance of selected chunks is applied ($w_{i,p} = 2^{-p}$), and **MRAG (2)** which assigns the weight for each text chunk based on the distance between the particular text chunk ($d_{i,p}$) and the query (q) ($w_i = \frac{1}{\text{distance}(d_{i,p}, q)}$). Figure 11 shows that these voting strategies perform worse on average than our selected strategy for MRAG, justifying its design and selection (described in Section 2.3).

We also consider two voting strategies for Split RAG, to further deepen the empirical evaluation. **Split (1)** only uses the exponential lowering of significance ($w_{i,p} = 2^{-p}$) and **Split (2)** which uses the same strategy as MRAG ($w_{i,p} = s_i \cdot 2^{-p}$). Figure 11 (on the right) shows that these voting strategies are on-par with each other while being worse than MRAG, further showcasing the advantages of MRAG.

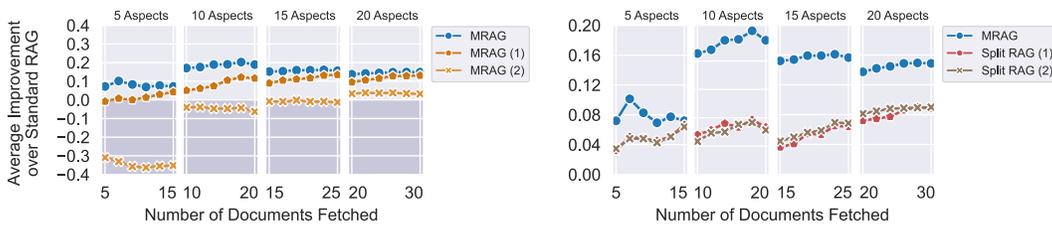


Figure 11: Evaluation of different voting strategies for MRAG and Split RAG.

I.3 ANALYZING PREPROCESSING OVERHEAD

One-time head importance scoring in MRAG introduces minimal preprocessing overhead on top of the standard embedding scheme. The scoring consists of computing: (i) average L2 norms per embedding space, and (ii) average pairwise cosine distances among embeddings within each space. To assess practical overhead, we analyze six datasets using the SFR-Embedding-Model. The additional time required to compute importance scores is measured as a percentage of the original embedding time. For example, on SciFact, the overhead was just 2.7%, and on NFCorpus, only 1.75%. Even for moderate-scale corpora such as ArguAna (310 seconds total encoding time), the overhead remained under 5%. These results assume full pairwise distance computation across all chunks; in practice, the harnessed sampling makes them even lower.