# Is Agentic RAG worth it? An experimental comparison of RAG approaches

**Pietro Ferrazzi[1,2], Milica Cvjeticanin[3], Alessio Piraccini [4], Davide Giannuzzi[5],**

[1]Fondazione Bruno Kessler, Trento, Italy , [2]Univerità di Padova, Padova, Italy,
[3]Cargill Geneve, Switzerland, [4] Alkemy, Milan, Italy, [5]Affiliation 5
**Correspondence:** pferrazzi@fbk.eu

## Abstract

Retrieval-Augmented Generation (RAG) systems are usually defined by the combination of a generator and a retrieval component that extracts textual context from a knowledge base to answer user queries. However, such basic implementations exhibit several limitations, including noisy or suboptimal retrieval, misuse of retrieval for out-of-scope queries, weak query–document matching, and variability or cost associated with the generator. These shortcomings have motivated the development of "Enhanced" RAG, where dedicated modules are introduced to address specific weaknesses in the workflow. More recently, the growing self-reflective capabilities of Large Language Models (LLMs) have enabled a new paradigm, which we refer to as "Agentic" RAG. In this approach, the LLM orchestrates the entire process—deciding which actions to perform, when to perform them, and whether to iterate—thereby reducing reliance on fixed, manually engineered modules. Despite the rapid adoption of both paradigms, it remains unclear which approach is preferable under which conditions. In this work, we conduct an extensive, empirically driven evaluation of Enhanced and Agentic RAG across multiple scenarios and dimensions. Our results provide practical insights into the trade-offs between the two paradigms, offering guidance on selecting the most effective RAG design for real-world applications, considering both costs and performance.

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable results in natural language understanding and generation tasks. However, their capabilities remain limited by their fixed training data, raising the need for integrating private enterprise knowledge into LLM-based systems. Retrieval-Augmented Generation (RAG) has emerged as a practical solution: instead of relying solely on parametric knowledge, the model retrieves relevant documents from an external knowledge base and conditions its generation on this evidence (Lewis et al., 2020). This method has garnered significant attention in both the research community (Wang et al., 2024; Fan et al., 2024; Wang et al., 2024) and industry, with cloud providers offering their own RAG solutions (IBM, 2025; Amazon Web Services, 2025; Microsoft Azure, 2025). Since this first initial definition RAG workflows have been expanded to what we define as **Enhanced RAG** (Figure 1, left). Such systems add to the retrieval and generation blocks components that perform further refinement, such as query rewriting, document re-ranking, and query routing. Recently, LLMs' increasing self-reflective capabilities have enabled a shift towards **Agentic RAG** (Figure 1, right), where the LLM acts as an orchestrator, deciding which actions to perform, being able to utilize different tools for different purposes. These systems are no longer fixed pipelines, but rather iterative loops with no predefined order, where the model is in charge of all the decisions. Although initial work on identifying theoretical distinctions between Enhanced and Agentic RAG systems has been proposed (Neha and Bhati, 2025), it remains unclear what the performance differences are between the two systems. To this end, we conducted an experiment-driven comparison of the two in different domains.

Our **first contribution** consists of the evaluation of the performance of the two systems in four dimensions selected from the literature (Table 1). We identified weaknesses of base RAG systems, analyzed how they have been addressed in previous work, and proposed an experimental setting to compare Enhanced and Agentic implementations. The selected dimensions are:

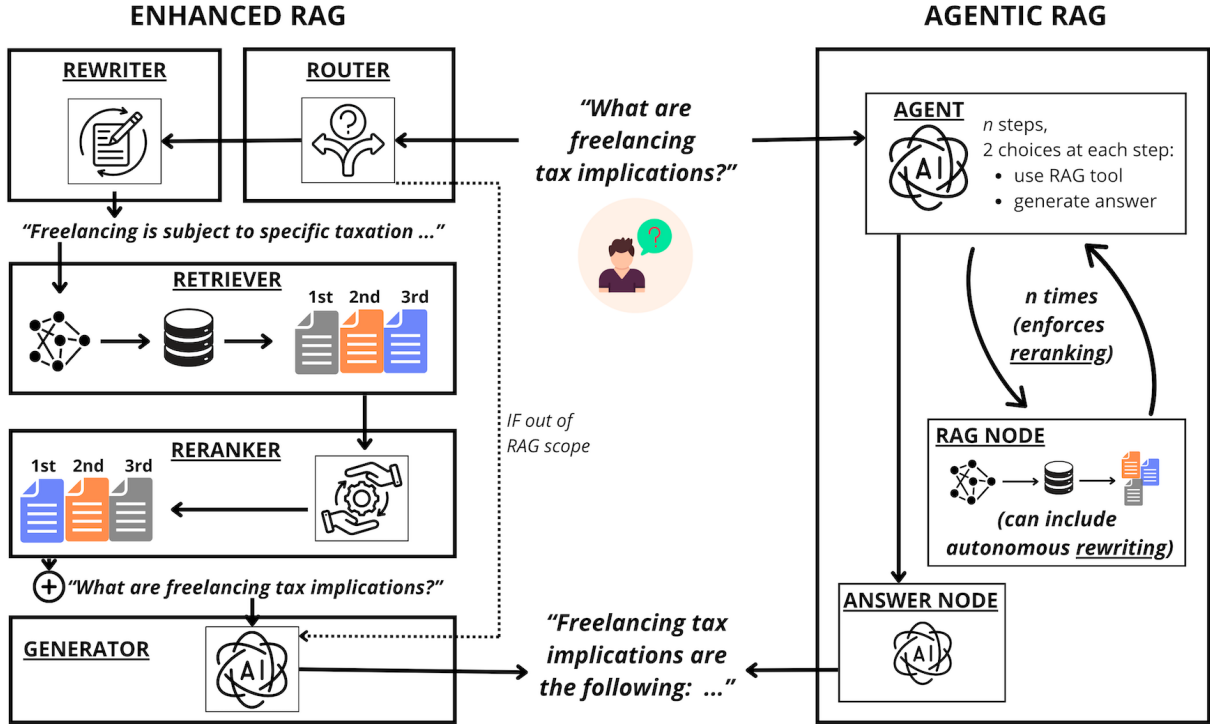1. **User intent handling**, measuring how sys-

Figure 1: **Left — Enhanced RAG**. The system is composed by a sequence of modules, each responsible for improving a specific stage of the RAG pipeline. A router determines whether a query should trigger retrieval; a rewriter reformulates the query; a retriever selects candidate chunks from the knowledge base; and a reranker orders the retrieved context before passing it to the generator. The workflow is fixed: information flows through predefined blocks intended to mitigate known weaknesses of **naïve RAG** systems (defined by the simple composition of the *retriever* and *generator* blocks). **Right — Agentic RAG**. The LLM acts as an agent that orchestrates the entire process. At each step, it can choose whether to call a RAG tool or proceed directly to answer generation. Retrieval and context refinement can be repeated, and the agent autonomously selects and sequences operations based on the evolving state of the task. This yields an iterative pipeline without predefined intermediate modules.

tems discriminate between queries that do or do not require external knowledge;

2. **Query-documents alignment**, measuring how systems can align queries to the format of knowledge base documents;

3. **Retrieved documents adjustment**, measuring the ability of systems to further tune the selection of documents after retrieval;

4. **Impact of LLM quality**, measuring how robust systems are to changes in the capability of the underlying models.

Our **second contribution** consists of a detailed analysis of costs and computational time required by the two systems under several scenarios. Finally, we propose a short summary of our findings, hoping to help researchers and practitioners select the most appropriate system for their needs.

## 2 Retrieval Augmented Generation

The term Retrieval-Augmented Generation (*RAG*) is often inconsistently used across the literature. Here we introduce the definitions used in this work.

**RAG definition** We define a RAG system as the integration of *i)* a knowledge base (KB) of textual chunks; *ii)* a retrieval method that extracts relevant chunks from the knowledge base given a query; and *iii)* a generator that produces an answer conditioned on both the query and the retrieved chunks.

**Naïve RAG** We refer to *Naïve Rag* (Figure 1) as the simplest instantiation of the RAG paradigm, where the generator is a Large Language Model and the workflow comprises four sequential steps: when *i)* a query is received, *ii)* the retriever deterministically selects a fixed number of chunks from the KB (*retrieval*); *iii)* the retrieved chunks are prepended to the query; *iv)* and passed to the generator, which produces an answer (*generation*).

| Naïve RAG shortcoming | What we evaluate | Implementations | |
|---|---|---|---|
| | | **Enhanced** | **Agentic** |
| Retrieval is performed even for queries that do not require it | Accuracy of RAG usage for in-scope and out-of-scope queries | Semantic routing system | Agent decides whether to do retrieval or not |
| Queries and documents in the KB differ in format or semantics, causing weak retrieval | Impact of query rewriting techniques | Hyde-based query rewriting | Agent rewrites query as it wishes |
| Noisy or suboptimal retrieval | Impact of retrieved document list refinement techniques | Encoder-based re-ranker | Agent can redo retrieval multiple times |
| The underlying LLM is too weak / takes too much time / is too expensive | Impact of selecting more / less powerful models | Test different LLMs | Test different LLMs |

Table 1: Summary of the evaluation dimensions we select. For each shortcoming in Naïve RAG, we define an evaluation dimension ("What we evaluate") and an implementation to test how Enhanced and Agentic RAG overcome such a limitation.

**Enhanced RAG** Following the taxonomy outlined in Huang and Huang (2024), we define Enhanced RAG (Figure 1, left) as any Naïve RAG pipeline augmented with additional steps designed to improve its performance. These enhancements can occur at different stages of the workflow and implemented in various ways, which we discuss in detail in the related sections.

**Agentic RAG** We define Agentic RAG (Figure 1, right) as a system in which the LLM assumes control over the workflow, being able to dynamically decide to perform actions. This agentic control allows the system to adapt the retrieval–generation loop to the specific context and task at hand. Unlike Enhanced RAG, no extra components are introduced with respect to the basic knowledge base-retriever-generator setting.

## 3 Related work

**RAG** The concept of RAG, first introduced by (Lewis et al., 2020), has undergone intensive research, rapidly evolving to highly sophisticated systems. A comprehensive overview is presented by the surveys proposed by Gao et al. (2023b); Fan et al. (2024); Wang et al. (2024). As large language models (LLMs) have acquired the capacity to multi-step reasoning and reflection, their consistency has enabled a paradigm shift toward Agentic RAG solutions (Shinn et al., 2023; Madaan et al., 2023). While the definition of the properties that characterize AI agents has evolved (Masterman et al., 2024), this emerging research direction has not yet been comprehensively categorized within a unified taxonomy, with initial attempts made by Singh et al. (2025); de Aquino e Aquino et al. (2025).

**Agents design and implementation** Several open-source frameworks have recently emerged to support the development of Agentic RAG systems, reflecting the rapid growth and experimentation in this area. SmolAgents (Roucher et al., 2025), LangGraph (LangChain Inc., 2025), LlamaIndex (Liu, 2022), CrewA (Contributors, 2025a), AutoGen (Microsoft, 2025), PydanticAI (Contributors, 2025b), and Atomic Agents (BrainBlend-AI, 2025) offer different advantages and reflect different design choices. We use PocketFlow, a lightweight and minimal framework that offers a simple graph-based abstraction and avoids the complexity of larger libraries, ensuring a clean and controlled evaluation. Table 8 in the Appendix summarizes our framework choice.

**The Need for Empirical Comparison** Despite rapid conceptual and architectural progress in Retrieval-Augmented Generation (RAG), the literature still lacks comprehensive experimental comparisons between Agentic RAG and Enhanced RAG systems. A preliminary effort in this direction is presented by Neha and Bhati (2025), who propose a useful set of definitions and evaluation dimensions but stop short of conducting a full empirical study.

## 4 Data

To conduct our experiments, we require datasets that are representative of common RAG applications, consisting of queries paired with a knowledge base. Following the taxonomy proposed by Arslan et al. (2024), which categorizes RAG use cases by application area, we focus on the most prominent, natural language based scenar-

ios. Specifically, we consider two major categories: i) Question Answering (QA), cases where RAG is utilized to ground answers in factual knowledge, and ii) Information Retrieval and Extraction (IR/E), cases where RAG is intended as a tool to get knowledge from data through natural language queries.

For QA, we selected FIQA (Maia et al., 2018) and NQ (Kwiatkowski et al., 2019) in the version proposed by Thakur et al. (2021). For IR/E, we used FEVER (Thorne et al., 2018) and CQADupStack-English (Hoogeveen et al., 2015). Each dataset is chosen as it represents a different real-world scenario, and allows to implement a different testing task (Table 2):

- FIQA (Financial QA) represents the RAG use case where queries are domain-specific questions to be answered by grounding responses on expert knowledge. The task consists of providing an answer to the user's query.

- NQ (Natural Questions) captures broad QA use cases, where users can ask any type of question. The task consists of providing the most appropriate answer to the user's query.

- FEVER tests the claim verification task, common in legal and biomedical applications, where users seek evidence for or against a statement. The task consists of finding documents that support or contradict the user query, returning a final assessment with a summary of the references.

- CQADupStack-English represents settings where RAG is used to find previously resolved tickets and discussions. The task consists of finding existing blog posts that address the same question posed by the user, providing an user-friendly summarisation.

# 5 Evaluation

The evaluation of RAG systems has been frequently decomposed into assessments of individual subcomponents (Es et al., 2024; Chen et al., 2020), each corresponding to a dimension that influences overall effectiveness. We design our evaluation following the same approach. We first identify a list of limitations of Naïve RAG pipeline based on the work done by Huang and Huang (2024), then construct an experimental setting to compare two implementations of Enhanced and Agentic RAG

| Task | Domain | Dataset | #Query | #Doc | Avg D/Q |
|------|--------|---------|--------|------|---------|
| QA | General | NQ | 3,452 | 2,681,468 | 1.2 |
| | Finance | FiQA | 648 | 57,638 | 2.6 |
| IE/R | Grammar forum | CQAD-EN | 1,570 | 40,221 | 1.4 |
| | Wikipedia | FEVER | 6,666 | 5,416,568 | 1.2 |

Table 2: Overview of the four selected datasets used for the experimental settings. For both Question Answering (QA) and Information Extraction and Retrieval (IE/R), 2 datasets are selected. Each query has a labelled list of relevant and irrelevant documents. Avg. D/Q indicates the average number of relevant documents per query.

systems. Table 1 summarizes our formulation of the key shortcomings of Naïve RAG and illustrates how Enhanced and Agentic RAG address them. For each of the four identified dimensions, we *i)* define it, *ii)* detail our implementation choices to enable both systems to address it, *iii)* present the evaluation setting, *iv)* define the evaluation metrics.

## 5.1 User intent handling

**Definition** In the context of this work, user intent handling refers to the need to determine whether a certain query requires the usage of retrieval or not. While prior surveys on RAG (Huang and Huang, 2024; Gao et al., 2023b; Fan et al., 2024) do not explicitly address nor mention it, Wang et al. (2024) highlight its importance by proposing a dedicated classifier for this task. We argue that intent detection is crucial in real-world RAG systems, as it prevents unnecessary or inappropriate retrieval calls and ensures that RAG is invoked only when it contributes meaningfully to answering the query.

**Enhanced Implementation** We implement an Enhanced RAG routing system using the semantic-router framework[1]. A router is defined by two sets of example queries, labelled as *valid* and *invalid* respectively. At inference time, the user query is compared to these groups and it is classified as *valid* or *invalid* accordingly. The system uses RAG to answer valid queries and avoid it for invalid ones. For our experiments, we utilize OpenAI's `text-embedding-3-small` as embedder. More details on the structure of the routing system are reported in Appendix A.2.

**Agentic Implementation** An Agentic RAG system embeds the ability to discriminate between

---

[1] https://github.com/aurelio-labs/semantic-router.git

| Setting | QA | | IR/E | | | |
|---|---|---|---|---|---|---|
| | **FIQA** | | **FEVER** | | **CQA-EN** | |
| | rec | F1 | rec | F1 | rec | F1 |
| naïve | 100 | 66.7 | 100 | <u>66.7</u> | 100 | 66.7 |
| enhanced | 95.1 | <u>95.7</u> | 84.4 | **87.9** | 94.7 | <u>96.6</u> |
| agentic | 97.7 | **98.8** | 49.3 | 64.6 | 100 | **99.8** |

Table 3: **User intent handling** recall and F1 score on 500 valid and invalid queries per dataset. The baseline is represented by Naïve RAG, where retrieval is performed for each user query. The Enhanced settings are based on the semantic router approach, while the agent autonomously decided whether to use the RAG tool.

| Setting | QA | | IR/E | | |
|---|---|---|---|---|---|
| | **FIQA** | **NQ** | **FEVER** | **CQAD** | **AVG** |
| naïve | 45.3 | 43.7 | 66.2 | 45.8 | 50.3 |
| enhanced | 43.5 | 43.9 | 81.1 | 42.8 | <u>52.8</u> |
| agentic | 43.2 | 51.7 | 83.1 | 44.3 | **55.6** |

Table 4: **Query rewriting** performances in terms of NDCG@10. The Baseline is represented by Naïve RAG, where the user query is directly embedded without any further transformation.

queries that require retrieval by design. When a query is received, the agent can freely decide whether to utilise the RAG node or directly answer. We use *gpt-4o* as underlying LLM.

**Experimental setting**   We tested performances on a dataset composed by an equal number of valid and invalid queries (500 for each of the four datasets). We selected the valid queries from the train splits of each dataset, while we generated the invalid ones prompting *gpt-4o* via 5-shot. We make publicly available the datasets of valid and invalid queries[2]. We excluded the NQ dataset from this evaluation stage as it handles by design any type of query, preventing the definition of *invalid* ones.

**Evaluation metric**   To evaluate if systems correctly handle queries, we utilized F1 score and recall, to take into account performances on both the valid and invalid classes.

## 5.2   Query rewriting

**Definition**   Much attention has been given to query rewriting techniques, first introduced by (Ma et al., 2023). The idea is that when the user query is tested against the knowledge base for a similarity search, the comparison is often performed among fairly different texts: the query is usually a short

---

and dense question, while chunks in the KB can be long and complex. Query rewriting techniques aim to reduce this delta by rewriting the query into a text with a structure more similar to the target chunks. Hyde (Gao et al., 2023a), consisting in substituting the query with a short paragraph that answers it, has emerged as one of the best performing techniques (Wang et al., 2024).

**Enhanced Implementation**   We forced the enhanced RAG system to perform Hyde query rewriting. Each user query is automatically rewritten before performing the retrieval step by prompting gpt-4o with the following: *Please write a passage to answer the question. \n Question: {user_query}\n Passage:"*.

**Agentic Implementation**   We design a prompt to make aware the agent that rewriting might help. Once the Agent has chosen to use the RAG tool, it can decide to use perform this step: *"Convert the user query into a {type_of_doc}"*, where *type_of_doc* differs based on the dataset (*"longer blog post"* for CQADupStack, *"passage to answer it"* for FiQA and NQ, *"longer factual statement"* for FEVER).

**Experimental Settings**   We run all the queries in the test sets of the four datasets against the two systems. To enforce fair comparison, in those cases where the Agentic setting did not perform query rewriting, we calculate the retrieval metric on the original query.

**Evaluation metric**   All queries in the four datasets come with annotations on the ground truth documents they should be linked to. When evaluating the quality of the retrieved chunks, we can then use NDCG@10. NDCG stands for Nomarlized Discounted Cumulative Gain (Järvelin and Kekäläinen, 2002) and is one of the most common metrics used to assess the effectiveness of a ranking model. NDCG at cutoff $K$ is defined as:

$$\text{NDCG@K} \equiv \frac{\text{DCG@K}}{\text{maxDCG@K}} \quad (1)$$

where maxDCG@K is the maximum DCG@K that can be obtained from the given relevance labels, and where DCG@K is defined as:

$$\text{DCG@K} \equiv \sum_{i=1}^{K} \frac{2^{l_i} - 1}{\log(1 + i)} \quad (2)$$

|  | QA | IR/E |  |
|---|---|---|---|
|  | **FIQA** | **CQA-EN** | **AVG** |
| naïve | 45.0 | 46.0 | 45.5 |
| enhanced w/o rewriting | 49.0 | 47.0 | <u>48.0</u> |
| enhanced with rewriting | 51.0 | 48.0 | **49.5** |
| agent | 43.4 | 44.4 | 43.9 |

Table 5: **Document list refinement** performances in terms of NDCG@10. The Agentic setting performs this task by iterating retrieval as many times as needed, whereas the Enhanced setting achieves it through a re-ranker model. Since Agentic retrieval indirectly performs query rewriting, we report results with and without rewriting for the enhanced setting, allowing for a fair comparison.

where $l_i$ is the relevance label (the ground truth label) of the document in position $i$ in the rank. Since Equation 5.2 is always positive, Equation 5.2 is a number bounded between 0 and 1, where NDCG@K equal to 1 means that we have a perfect ranked list.

## 5.3 Document list refinement

**Definition** The retrieval step selects a number of chunks that are not necessarily the most appropriate ones. Previous work has shown how initial retrieval may include partially irrelevant or noisy results, and proposed approaches to improve the selection process via reranking strategies (Sachan et al., 2022; Sun et al., 2023; Qin et al., 2024). Reranking consists in sorting the retrieved chunks, selecting a subset of highly relevant ones with respect to the user query.

**Enhanced Implementation** We experiment with this dimension by using an ELECTRA-based reranker (Déjean et al., 2024)[3] on top of the list of most similar 20 documents for each user query.

**Agentic Implementation** The Agentic RAG system can inherently attempt to consider a more suitable context when needed. Specifically, the agent may trigger additional retrieval rounds and adapt the query formulation as it deems appropriate, allowing it to iteratively obtain more relevant context. We calculate the metric on the last reformulation of the query that the Agent uses for the RAG tool, which directly precedes answer generation.

---

[3] https://huggingface.co/naver/trecdl22-crossencoder-electra

**Experimental Settings** We run all the queries in the test sets of FIQA and CDQStack-English against the two systems to assess performance on both Qa and IR/E tasks, respectively. We did not consider NQ and FEVER due to their size. As detailed in Section 6.1, Agentic RAG would take >7 days on each of them.

**Evaluation metric** To evaluate reranking, we utilize NDCG@10, comparing the ground-truth links between queries and documents with the selected ones. The metric is the same described above for query rewriting.

## 5.4 Underlying LLM

**Definition** Both Agentic and Enhanced settings are highly impacted by the choice of the underlying LLM, as different models produce different answers even when provided with the same query and retrieved context. Furthermore, the role of the generator is particularly critical in Agentic RAG, where the model must not only produce the final answer but also make decisions at each stage of the workflow. We are interested in understanding and quantifying the impact that a weaker generator has on the system, compared to what a stronger one would have.

**Enhanced and Agentic Implementation** To assess this effect, we tested four generators of varying capability, namely Qwen3-0.6B (without thinking), Qwen3-4B, Qwen3-8B, and Qwen3-32B (Yang et al., 2025). We define generator capability based on the benchmark presented in Table 6. We do not aim to assess which model performs best as a generator, but how the overall RAG performances are impacted by generators of different power. We utilized the Enhanced RAG settings with query rewriting and reranking, and the Agentic settings with the system prompt defined in Figure 3, Appendix A.1.

**Experimental Settings** We run all queries in the test sets of FIQA and CDQStack-English against the two systems to assess performance on both QA and IR/E tasks, respectively. We did not consider NQ and FEVER due to their size.

**Evaluation metric** We evaluate the quality of the final answers generated by both systems using automatic metrics, employing the LLM-as-a-judge paradigm. Out of the many approaches proposed in the literature (Kim et al., 2024; Wang et al., 2025; Es et al., 2024) we select Seleni-70B (Alexandru

| Model | General | Align | Reason | AVG |
|---|---|---|---|---|
| Qwen3 | GPQA-D | IFEVAL | AIME '24 | |
| 0.6B | 22.9 | 54.5 | 3.4 | 26.9 |
| 4B | 55.9 | 81.9 | 73.0 | 70.3 |
| 8B | 62.0 | 85.0 | 76.0 | 74.3 |
| 32B | 68.4 | 85 | 81.4 | 78.3 |

Table 6: Performances of the models selected to analyse the impact of the underlying LLM as reported by Yang et al. (2025). The selected benchmarks are Graduate-Level Google-Proof QA Diamond (Rein et al., 2023), Instruction Following Eval (Zhou et al., 2023), and the American Invitational Mathematics Examination (AIME, 2024).

et al., 2025), a fine-tune model based on Llama-3.3-70B-Instruct, as its smaller version scores among the top ones in the LLM-as-a-Judge Arena[4]. The authors also show that these models' judgments closely match human evaluations on financial QA tasks, which is relevant given our use of the FIQA dataset. For this dataset, we therefore adopt the binomial 0–1 classification metric—its most human-aligned option—and report average scores.

On the other hand, for CQADupStack-English, such alignment has not been demonstrated. Therefore, we assessed it by performing manual annotation on a subset of the testing data ($5\%$, 312 answer pairs in total), and calculated the agreement rate between the two human annotators and automatic metrics. We select the pairwise metric (given two answers of two different models, select the best one). Inter-annotator agreement rate (ratio of times they both chose A or B) is $71.9\%$, while the human-model agreement is $65.4\%$. Manual annotation took an average of $1.5$ minutes per pair, resulting in a total of $15.5$ hours overall. The annotation guidelines are reported in the Appendix. To summarize the impact of underlying model changes, we calculate the ratio of times when the larger model wins over the smaller counterpart.

## 6 Results

**User intent handling**  Table 3 reports the results for user intent handling. We found that Agentic slightly overperforms Enhanced settings in the FIQA and CDQADupStack-EN tasks. In the case of FEVER, the latter outperforms the former by a margin ($+28.8$ F1 points), due to a very low recall ($49.3$). Such low recall is due to the fact that the system often uses retrieval even in cases it should

---

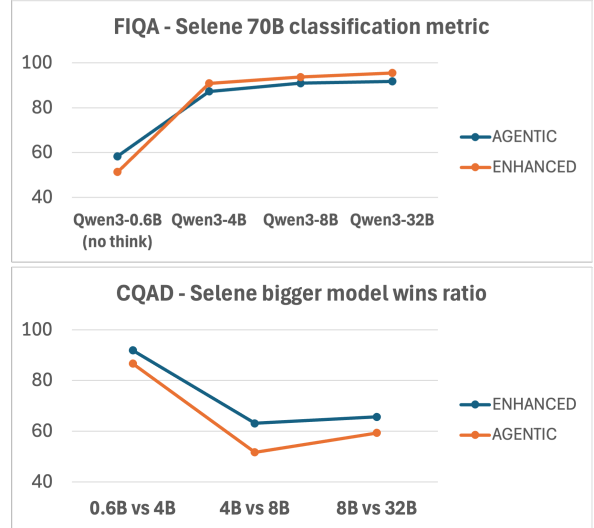[4]https://huggingface.co/blog/arena-atla

Figure 2: Performances of Enhanced and Agentic settings when changing the underlying LLM. For FIQA, the metric is the overall ratio of the classification metric (1 if the answer follows the instruction, 0 otherwise). For CQADupStack-EN, the metric is based on a pairwise analysis, calculated as the ratio of times when the larger model's answer is better than that of the smaller counterpart. Both metrics are calculated via LLM-as-a-Judge (Selene-70B).

not. We attribute this to the first two datasets having a very clear domain definition (finance, English grammar), whereas the FEVER task is much less restrictive by design, as it aims to verify user queries on factual information, which makes it harder for the agent to understand exactly what kind of requests require retrieval.

**Query rewriting**  Results in Table 4 show that the Agentic setting performs better than the enhanced one, with an average gain of $2.8$ NDCG@10 points. We attribute this to the flexibility of the former, which can dynamically decide whether to perform rewriting, and how. This suggests that query rewriting is always beneficial, and having an adaptive approach to decide what to do on a case-by-case basis is the most reliable approach.

**Document list refinement**  In the Enhanced setting, re-ranking has a substantial positive impact on performance. In contrast, the Agentic RAG setting gains no benefit from iterating the retrieval step, as it appears unable to identify better documents than those initially retrieved.

**Underlying LLM**  We aim to assess whether Enhanced and Agentic systems exhibit distinct performance patterns across model scales. Figure 2 re-

ports the resulting average scores for both datasets, which show that the two systems do not present significant differences in patterns when changing the underlying LLM. In fact, the performance increase in FIQA follows the same distribution both Enhanced and Agentic RAG, and the same is true for the ratio of times in which the bigger model is preferred over the smaller one in CQADupStack-En settings. Full results on the evaluation metrics are reported in Appendix A.3.

## 6.1 Cost and Time

When comparing Enhanced and Agentic RAG systems, an important dimension is their computational and monetary cost. In many real-world deployments, practitioners may reasonably accept slightly lower performance if it results in reductions in latency or resource usage.

**Fixed time and cost for Enhanced and Agentic settings** The two systems share some fixed costs due to hardware requirements. In this work, we utilized for both settings a `t3.large` ec2 AWS instance (0.09 \$/h per on-demand usage) to instantiate a relational database with vector search capabilities (`pgvector`[5]). We implement the RAG application backend on a `t2.medium` (0.05 \$/h). We tested open LLMs in a proprietary $8\times$a40 cluster (46GB). We run Qwen3 0.6B, 4B, and 8B on a single gpu, while the 32B version on $4\times$ a40. A similar setting on AWS can be `g4ad.8xlarge`, which costs 1.9 \$/h. The time and cost related to retrieval are the same in both settings. We used OpenAI `text-embedding-3-small` as default embedder model with cosine similarity. Enhanced RAG does re-ranking with a 300M-parameter model on the same cluster as the LLMs, so the added cost is negligible.

**Runtime costs and number of tokens** We approximate the runtime RAG cost with the number of processed input and output tokens. This metric is hardware-agnostic, allowing cost to be estimated for any deployment given the model's throughput and hourly pricing. We analyse token usage for both GPT and Qwen models to quantify cost differences between Enhanced and Agentic settings. We also report end-to-end latency—the time from receiving a query to returning the final answer. For open LLMs, latency is heavily influenced by the underlying hardware. Table 9 (Appendix) reports a

[5]https://github.com/pgvector/pgvector

summary of time and tokens per approach.

We find that, for FIQA, the Agentic setting requires on average $2.7\times$ more input tokens and $1.7\times$ more output tokens than Enhanced RAG. For CQADupStack-En, the increases are even larger-$3.9\times$ and $2.0\times$, respectively. In both scenarios, time increase by $1.5\times$ on average. These differences translate directly into higher financial costs. For valid queries in Enhanced RAG, we find that roughly 45–50% of the total time is spent generating the answer, a similar proportion is spent on query rewriting, 0–5% on retrieval, and 0–2% on document re-ranking. The dominant factor in latency is the LLM calls; therefore, any performance optimization should focus primarily on that.

## 7 Conclusion

Our experimental comparison reveals that neither Enhanced nor Agentic RAG is universally superior. First, we observe that in narrow, well-defined domains with highly structured user behavior, Agentic RAG excels at **handling user intent,** thanks to its ability to understand the user query. However, in broader or noisier domains, our Enhanced RAG routing system proves more reliable. Second, with respect to **alignment of the query to** the structure and semantics of the **documents** in the KB, Agentic RAG outperforms Enhanced RAG retrieval quality. Its dynamic use of query rewriting allows for retrieval of more relevant context. Third, we found that when Agentic RAG selects certain documents, it is not as good as the **re-ranking** done by Enhanced RAG at selecting just the most meaningful docs. Our results suggest that integrating an explicit re-ranking step into Agentic pipelines could provide substantial gains. Fourth, we observe that changing the **underlying LLM** produce the same changes in performance in both settings. The impact of model size follows a similar trend for both paradigms: as the underlying LLM becomes larger, performance improves at comparable rates. Our **cost** analysis highlights that Agentic RAG is systematically more expensive—up to 3.6 times more—due to additional reasoning steps and repeated tool calls. This cost difference should not be overlooked: a well-optimized Enhanced RAG can match or exceed Agentic performance while remaining more efficient.

## Limitations

This study has several limitations that should be considered when interpreting the results. First, we do not release the implementation of our PocketFlow-based agent, which may limit full reproducibility. Second, although our evaluation covers key dimensions of RAG behavior, additional aspects—such as document summarization, document repacking (re-sorting documents in the context according to their importance), and output refinement have not been considered. Furthermore, our agent is equipped with only a single tool, whereas richer agentic setups might exhibit different behaviors. Overall, each block of Enhanced RAG is implemented following previous work but lacks of compairison of different appraoches, which might perform better than the ones we selected.

## References

AIME. 2024. Aime problems and solutions. https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions. Accessed 2024.

Andrei Alexandru, Antonia Calvi, Henry Broomfield, Jackson Golden, Kyle Dai, Mathias Leys, Maurice Burger, Max Bartolo, Roman Engeler, Sashank Pisupati, Toby Drane, and Young Sun Park. 2025. Atla selene mini: A general purpose evaluation model. *Preprint*, arXiv:2501.17195.

Amazon Web Services. 2025. Amazon web services — bedrock. https://aws.amazon.com/bedrock/.

Muhammad Arslan, Hussam Ghanem, Saba Munawar, and Christophe Cruz. 2024. A survey on rag with llms. *Procedia Computer Science*, 246:3781–3790. 28th International Conference on Knowledge Based and Intelligent information and Engineering Systems (KES 2024).

BrainBlend-AI. 2025. Brainblend-ai/atomic-agents: Modular agents framework for building agents from atomic components. https://github.com/BrainBlend-AI/atomic-agents. MIT License; accessed: 2025-11-18.

Andrew Chen, Andy Chow, Aaron Davidson, Arjun DCunha, Ali Ghodsi, Sue Ann Hong, Andy Konwinski, Clemens Mewald, Siddharth Murching, Tomas Nykodym, Paul Ogilvie, Mani Parkhe, Avesh Singh, Fen Xie, Matei Zaharia, Richard Zang, Juntai Zheng, and Corey Zumar. 2020. Developments in mlflow: A system to accelerate the machine learning lifecycle. In *Proceedings of the Fourth International Workshop on Data Management for End-to-End Machine Learning*, DEEM '20, New York, NY, USA. Association for Computing Machinery.

CrewAIInc Contributors. 2025a. crewAIInc/crewAI: Framework for orchestrating role-playing, autonomous ai agents. https://github.com/crewAIInc/crewAI. Accessed: 2025-11-18.

Pydantic Contributors. 2025b. pydantic/pydantic-ai: Genai agent framework, the pydantic way. https://github.com/pydantic/pydantic-ai. Accessed: 2025-11-18.

Gustavo de Aquino e Aquino, Nádila da Silva de Azevedo, Leandro Youiti Silva Okimoto, Leonardo Yuto Suzuki Camelo, Hendrio Luis de Souza Bragança, Rubens Fernandes, Andre Printes, Fábio Cardoso, Raimundo Gomes, and Israel Gondres Torné. 2025. From rag to multi-agent systems: A survey of modern approaches in llm development. *Preprints*.

Hervé Déjean, Stéphane Clinchant, and Thibault Formal. 2024. A thorough comparison of cross-encoders and llms for reranking splade. *Preprint*, arXiv:2403.10407.

Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. 2024. RAGAs: Automated evaluation of retrieval augmented generation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 150–158, St. Julians, Malta. Association for Computational Linguistics.

Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '24, page 6491–6501, New York, NY, USA. Association for Computing Machinery.

Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023a. Precise zero-shot dense retrieval without relevance labels. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1762–1777, Toronto, Canada. Association for Computational Linguistics.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2023b. Retrieval-augmented generation for large language models: A survey. *Preprint*, arXiv:2312.10997.

Doris Hoogeveen, Karin M. Verspoor, and Timothy Baldwin. 2015. Cqadupstack: A benchmark data set for community question-answering research. In *Proceedings of the 20th Australasian Document Computing Symposium*, ADCS '15, New York, NY, USA. Association for Computing Machinery.

Yizheng Huang and Jimmy Huang. 2024. A survey on retrieval-augmented text generation for large language models. *Preprint*, arXiv:2404.10981.

IBM. 2025. Ibm watsonx — rag development. https://www.ibm.com/it-it/products/watsonx-ai/rag-development.

Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of ir techniques. *ACM Trans. Inf. Syst.*, 20(4):422–446.

Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024. Prometheus: Inducing fine-grained evaluation capability in language models. *Preprint*, arXiv:2310.08491.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.

LangChain Inc. 2025. langchain-ai/langgraph: Build resilient language agents as graphs. https://github.com/langchain-ai/langgraph. MIT License; accessed: 2025-11-18.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.

Jerry Liu. 2022. LlamaIndex.

Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. Query rewriting in retrieval-augmented large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5303–5315, Singapore. Association for Computational Linguistics.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: iterative refinement with self-feedback. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.

Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. Www'18 open challenge: Financial opinion mining and question answering. In *Companion Proceedings of the The Web Conference 2018*, WWW '18, page 1941–1942, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.

Tula Masterman, Sandi Besen, Mason Sawtell, and Alex Chao. 2024. The landscape of emerging ai agent architectures for reasoning, planning, and tool calling: A survey. *Preprint*, arXiv:2404.11584.

Microsoft. 2025. microsoft/autogen: Autogen – agent orchestration and tool calling. https://github.com/microsoft/autogen. MIT License; accessed: 2025-11-18.

Microsoft Azure. 2025. Azure ai search — rag solution tutorial. https://docs.azure.cn/en-us/search/tutorial-rag-build-solution.

Fnu Neha and Deepshikha Bhati. 2025. Traditional rag vs. agentic rag: A comparative study of retrieval-augmented systems. *Authorea Preprints*.

Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, and Michael Bendersky. 2024. Large language models are effective text rankers with pairwise ranking prompting. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 1504–1518, Mexico City, Mexico. Association for Computational Linguistics.

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. *Preprint*, arXiv:2311.12022.

Aymeric Roucher, Albert Villanova del Moral, Thomas Wolf, Leandro von Werra, and Erik Kaunismäki. 2025. 'smolagents': a smol library to build great agentic systems. https://github.com/huggingface/smolagents.

Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving passage retrieval with zero-shot question generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3781–3797, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.

Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talaei Khoei. 2025. Agentic retrieval-augmented generation: A survey on agentic rag. *Preprint*, arXiv:2501.09136.

Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2023. Is ChatGPT good at search? investigating large language models as re-ranking

agents. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14918–14937, Singapore. Association for Computational Linguistics.

Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. Beir: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1.

James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.

PeiFeng Wang, Austin Xu, Yilun Zhou, Caiming Xiong, and Shafiq Joty. 2025. Direct judgement preference optimization. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 1979–2009, Suzhou, China. Association for Computational Linguistics.

Xiaohua Wang, Zhenghua Wang, Xuan Gao, Feiran Zhang, Yixin Wu, Zhibo Xu, Tianyuan Shi, Zhengyuan Wang, Shizheng Li, Qi Qian, Ruicheng Yin, Changze Lv, Xiaoqing Zheng, and Xuanjing Huang. 2024. Searching for best practices in retrieval-augmented generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17716–17736, Miami, Florida, USA. Association for Computational Linguistics.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. Qwen3 technical report. *Preprint*, arXiv:2505.09388.

Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. *Preprint*, arXiv:2311.07911.

# A  Example Appendix

## A.1  Prompt for Agentic RAG

The Agentic RAG is defined by three nodes: the orchestrator, the answer node, and the RAG node. Here we report the prompts utilized by the orchestrator (Figures 3 and 4) and the answer node (Figure 5), while the RAG node does not have a specific system prompt.

| Qwen | CQAD | | FIQA | |
|---|---|---|---|---|
| | enhanced | agentic | enhanced | agentic |
| **0.6B** | 51,6 | 15,4 | 51,4 | 58,4 |
| **4B** | 95,7 | 81,8 | 90,9 | 87,3 |
| **8B** | 94,9 | 88,8 | 93,7 | 91 |
| **32B** | 94,5 | 91,3 | 95,5 | 91,8 |

Table 7: Classification metric based on Selene-70B for **FIQA**.

## A.2  Routing system details

The schema of the routing system we implement for Enhanced RAG settings is described in 8.
The routing mechanism relies on example-based classification. Queries are compared to two reference sets: valid and invalid. If a query is sufficiently similar to the valid set, it is processed; otherwise, it is rejected. The key challenge is determining what "sufficiently similar" means, i.e., selecting an appropriate similarity threshold.

**Threshold selection**  Threshold selection can be approached in multiple ways (e.g., random search, linear models, classification algorithms). If class definitions are clear and well-separated, threshold tuning becomes less critical, since queries will naturally cluster around their correct class. In practice, however, class boundaries often contain noise, making the threshold an essential safeguard against misclassification.

## A.3  Underlying LLM evaluation metrics

The results for the metrics calculated by means of Selene-70B are reported in Table 7 (FIQA) and Figure 7 (CQADupStack-EN).

| Framework | pros | cons |
|---|---|---|
| smolagents | minimal; codeagent | lack of low-level abstraction; lower control and security issues because of `codeagent` (can be handled) |
| langgraph | graph abstraction; many integrations; many resources for implementing popular patterns | |
| llamaindex | many integrations; many resources for implementing popular patterns | high-complexity codebase |
| pocketflow | minimal; graph abstraction | need to perform implementation |
| crewAI | | focus on multi-agents |
| autogen | mature ecosystem | high complexity |
| pydanticAI | type safety; Python oriented | early stages of development; requires learning some concepts |
| atomic agents | modular | early stages of development, limited documentation |

Table 8: Comparison of frameworks considered for Agentic RAG implementation. "Pros" and "cons" are defined in the scope of this work, i.e. constructing the simplest agent possible with a single RAG tool.

| Model | | FIQA | | | | | | CQAD-EN | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | time | tot tokens | | ratio (Ag/En) | | | time | tot tokens | | ratio (Ag/En) | | |
| | | *s* | *input* | *output* | *time* | *input* | *output* | *s* | *input* | *output* | *time* | *input* | *output* |
| GPT-4.1-nano | En | 9.0 | 1683 | 465 | 1.1 | 2.2 | 0.8 | 7.0 | 856 | 331 | 1.2 | 3.0 | 0.9 |
| | Ag | 10,2 | 3676 | 348 | | | | 8.6 | 2463 | 297 | | | |
| Qwen3-0.6B | En | 8.1 | 1743 | 236 | 2.2 | 2.9 | 4.1 | 8.1 | 862 | 254 | 1.1 | 3.5 | 3.4 |
| | Ag | 22.1 | 4978 | 979 | | | | 8.9 | 3032 | 867 | | | |
| Qwen3-4B | En | 35.5 | 1743 | 1435 | 1.1 | 2.8 | 1.2 | 31.0 | 862 | 1019 | 1.2 | 3.9 | 1.9 |
| | Ag | 38.6 | 4834 | 1704 | | | | 37.2 | 3372 | 1943 | | | |
| Qwen3-8B | En | 58.5 | 1743 | 1490 | 1.2 | 2.8 | 1.2 | 58.4 | 862 | 1339 | 0.9 | 3.9 | 1.5 |
| | Ag | 69.9 | 4943 | 1837 | | | | 54.3 | 3394 | 1983 | | | |
| Qwen3-32B | En | 62.6 | 1743 | 1695 | 1.5 | 2.7 | 1.1 | 43.9 | 862 | 1109 | 2.3 | 3.9 | 1.5 |
| | Ag | 93.8 | 4766 | 1866 | | | | 101.9 | 3359 | 1636 | | | |
| **AVG ratio** | | | | | **1.5** | **2.7** | **1.7** | | | | **1.4** | **3.6** | **1.8** |

Table 9: Analysis on costs (measured by number of input and output tokens) and time (end-to-end latency experienced by the user when running a query). The presented ratios represent the multiplicative factor to go from the Enahnced to the Agentic settings. Qwen3-0.6B is run without thinking mode, resulting in substantially shorter outputs compared to the larger Qwen models. Agentic RAG always performed a maximum of 3 turns. In scenarios requiring more turns, tokens consumed by the agent would increase.

## Context
You are a powerful agentic AI agent. Your task is to {task_description}

## Inputs
• **query: {query}**
• **previous_tool_results: {previous_tool_results}**

## Guidelines
To do complete your task, you can take several actions:
• **answer: {answer_node_description}**
• **document_retrieval:{retrieval_node_description}**

## Important
• **Don't repeat actions unnecessarily.**
• **When composing the query for RAG search you must always write a passage to answer the input query provided by the user.**
• **If you think that the retrieved documents in the previous_tool_results can be further improved, you can rewrite the query and call document_retrieval again.**
• **If you see the relevant retrieved documents already in the previous_tool_results, you should proceed with answering!**

## Metadata:
• **current reasoning step: {current_reasoning_step}**
• **maximum reasoning steps: {max_reasoning_steps}**

Figure 3: System prompt template that defines the agent for all the tasks.

**FIQA:**
*task_description*:
**help users on financial issues using the tools at your disposal.**

*answer_node_description*:
**provide the final answer to the user, utilizing the retrieved documents if necessary. Choose this action when you have gathered enough information from the previous_tool_results.**

*retrieval_node_ description*:
**retrieve documents from a knowledge base of financial information. Choose this to ground answers on external sources.**

**NQ:**
*task_description*:
**help users using the tools at your disposal.**

*answer_node_description*:
**provide the final answer to the user, utilizing the retrieved documents if necessary. Choose this action when you have gathered enough information from the previous_tool_results.**

*retrieval_node_ description*:
**retrieve documents from a knowledge base of potentially useful information. Choose this to ground answers on external sources.**

**FEVER:**
*task_description*:
**help users verifying their queries on factual knowledge. For each user claim, you determine if it is supported or refuted by the world-knowledge on the topic.**

*answer_node_ description*:
**provide the final answer to the user, utilizing the retrieved documents if necessary. Choose this action when you have gathered enough information from the previous_tool_results.**

*retrieval_node_description*:
**retrieve documents from a knowledge base of factual information. Choose this to ground answers on gold-standard knowledge. Be careful, your knowledge might be outdated.**

**CQADUPSTACK-EN:**
*task_description*:
**help users finding blog posts about english grammar issues that are related to the user query. You have to find the relevant documents, summarize them, and inform the user about them.**

*answer_node_description*:
**provide a summary of the retrieved blog posts to the user. Choose this action when you have gathered enough information on grammar blog posts from the previous_tool_results.**

*retrieval_node_ description* :
**retrieve blog posts from a knowledge base of English grammar issues. Choose this to find relevant blog posts. Never use it for queries not related to english grammar issues.**

Figure 4: System prompt filling values that define each of the four selected tasks.

**FIQA:**
## Context

You are a blog posts summarizer. You'll be provided the original user query, a list of tools called and their result, a draft answer, and conversation history. Provide a bullet point list of the relevant blog posts related to the user query. One bullet point per blog post.
Each bullet point is composed of 1) a short title (max 10 words); 2) a one sentence summary (max 15 words) of the blog post; 3) the reference to the document.
You do not answer the question. Instead, you provide an overview of the retrieved blog posts to the user. Your objective is to make the user aware of the relevant blog posts related to their query.

## Inputs
• question: {query}
• previous_tool_results: {previous_tool_results}
• draft answer: {draft_answer}

## Guidelines
• Be consistent with previous responses in the conversation
• Introduce the bullet point list with a short sentence like "Here are some blog posts that might be relevant to your query:"
• Do not answer the question directly. Your goal is to summarize the retrieved blog posts.

**FEVER:**
## Context
You are a question-answering assistant. You'll be provided the original user query, a list of tools called and their result, a draft answer, and conversation history. Provide a summary of the grounding reasons and evidence.

## Inputs
• question: {query}
• previous_tool_results: {previous_tool_results}
• draft answer: {draft_answer}

## Guidelines
• Be consistent with previous responses in the conversation.
• If a follow-up question references a previous conversation, use the conversation history to provide context.
• Summarize the reasons that support or contradict the statement.

**FIQA & NQ:**
## Context

You are a question-answering assistant. You'll be provided the original user query, a list of tools called and their result, a draft answer, and conversation history.
Provide the final answer to the user.

## Inputs
• question: {query}
• previous_tool_results: {previous_tool_results}
• draft answer: {draft_answer}

## Guidelines
• Be consistent with previous responses in the conversation.
• If a follow-up question references a previous conversation, use the conversation history to provide context.

Figure 5: Answer node system prompts that define the answer generation.
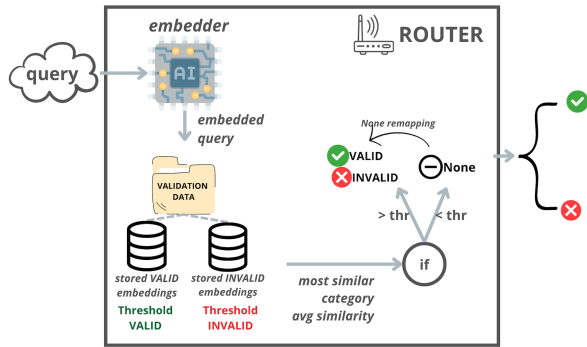
Figure 6: Schema of the routing system utilized for Enhanced We define two query classes—valid and invalid—each represented by embedded examples from the validation set. For each incoming query, the router embeds it, retrieves the top-20 most similar examples via cosine similarity, and selects the class with the highest average similarity. If the mean similarity of the selected class exceeds a predefined threshold, that label is assigned; otherwise, the router returns *None*, which can be further remapped according to business logic.
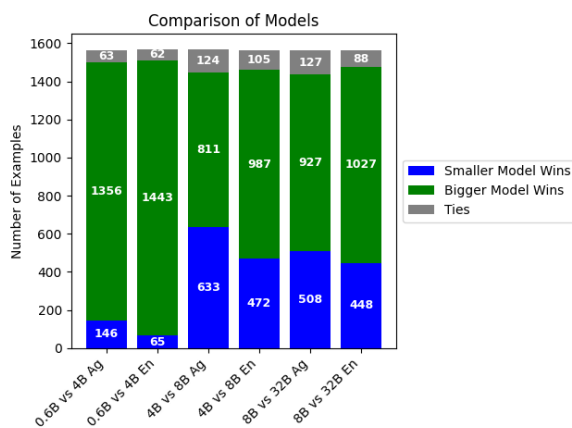


Figure 7: Pairwise metric based on Selene-70B for **CQADupStack-EN**.

**GENERAL DESCRIPTION**
You are evaluating a system that, given a user query, retrieves a list of similar queries previously submitted by other users. These past queries are called documents.
The system's output must be a list of bullet points, each representing one and only one document.
An ideal bullet point contains:
- a title of at most five words,
- a brief description of the retrieved document, and
- a reference to the document number.

Important: The documents do not contain answers to the user's query; they are actual queries submitted in the past. If the system's response includes a document description that explains the grammar question rather than restating the query itself, this is incorrect.

**ANNOTATION**
For each query, there are two system responses, A and B. Your task is to select the better response based on the criteria above.
If you are unsure, you may write "I don't know." If both responses are of equal quality, write "same."

Common situations include:
- The document description contains an answer rather than the original query (serious error).
- One of the listed documents is irrelevant to the query (serious error).
- If the two responses differ in the number of documents and all documents are relevant, always choose the one with more documents.
- The response mixes the document number with the title (error).
- The document description is excessively long (error).
- The document number reference is missing (error).

Figure 8: Guidelines for evaluating CQADupStack-English. They are used by both human annotators and LLM-as-a-Judge (Selene model).
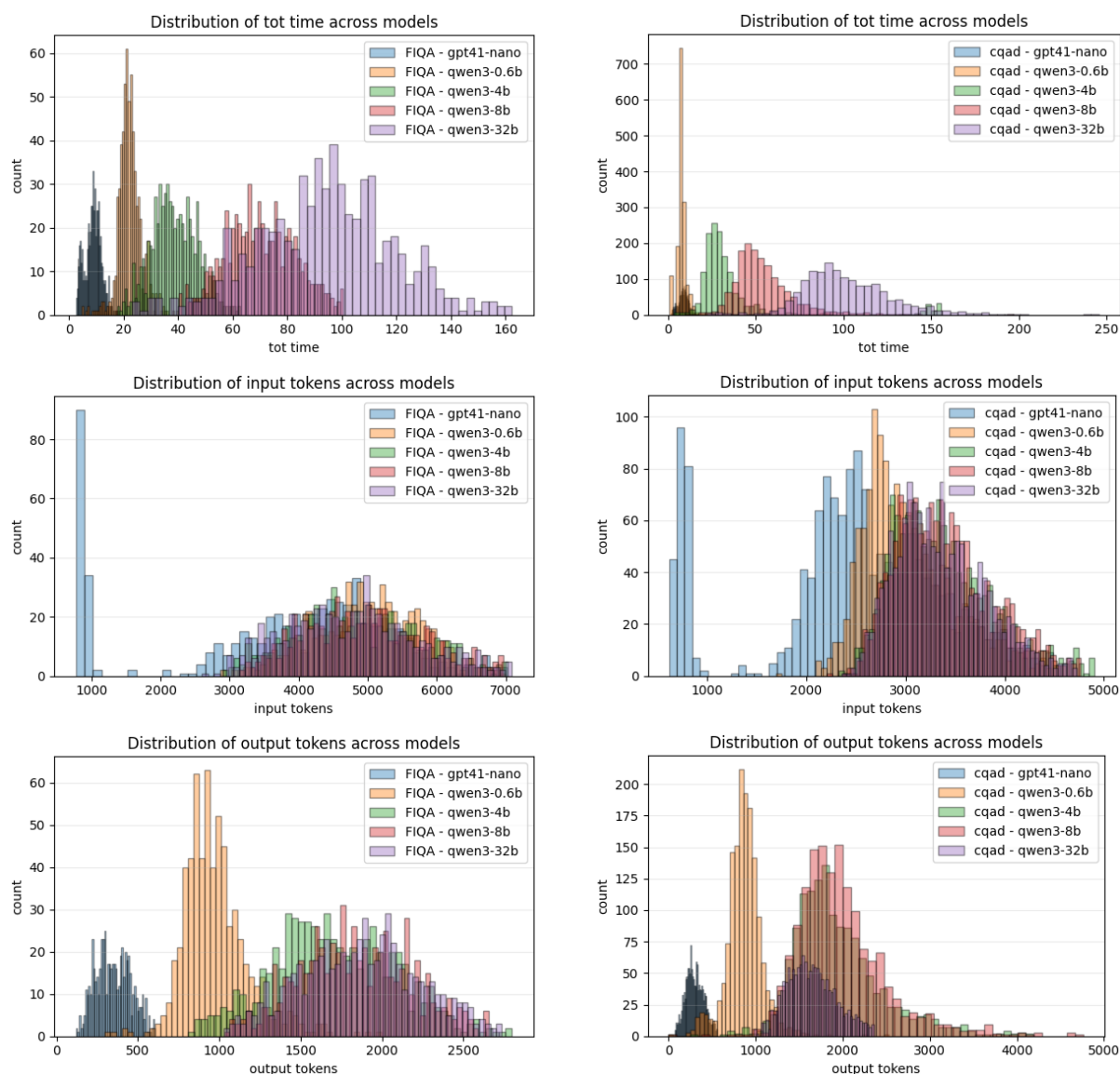
Figure 9: Overall computational cost and token usage for each model in the **Agentic** setting when processing a single user query. Qwen3 4B, 8B, and 32B operate in thinking mode, whereas the 0.6B variant is run without it. Qwen3 0.6B, 4B, and 8B are executed on a single NVIDIA A40, while Qwen3 32B is run on 4×A40 GPUs.

**Top:** Distribution of total latency, measured from the moment the system receives the query to the moment the final answer is produced. Within the Qwen family, Qwen3-0.6B achieves the lowest latency due to its smaller size and the absence of thinking mode. **Middle:** Average number of input tokens. This value increases slightly with model size. The peak of low input-token values for GPT-4.1-nano arises because the model frequently opts not to use the RAG tool, thereby reducing the number of required reasoning steps.

**Bottom:** Average number of output tokens. Here, the largest difference emerges: enabling thinking mode leads Qwen3 4B, 8B, and 32B to produce substantially longer outputs.