



(REVIEW ARTICLE)



Retrieval augmented generation system for dynamic document tagging and query-driven retrieval

Surya Kuchibhotla * and Sri Kuchibhotla.

Independent Researcher, Columbus, Ohio.

International Journal of Science and Research Archive, 2025, 16(01), 1452-1462

Publication history: Received on 11 June 2025; revised on 16 July 2025; accepted on 19 July 2025

Article DOI: <https://doi.org/10.30574/ijrsra.2025.16.1.2156>

Abstract

We propose a novel Retrieval Augmented Generation (RAG) framework for dynamic document tagging and query driven retrieval. Our system integrates a large language model (LLM) with an explicit memory of documents to generate semantic tags for each document and uses these tags to improve retrieval accuracy (Zhou et al., 2024; Li et al., 2024). A query tag feedback loop is then formalized to iteratively refine document annotations based on user queries and present a modular architecture that separates document preprocessing tag generation, storage and retrieval (Sharma et al., 2024). To evaluate such systems, we introduce a synthetic multi domain benchmark containing documents from scientific (ArXiv), governmental, and legal sources, along with ground truth tags and query pools (Kim et al., 2024; Lin et al., 2024). We also define a new Query to Tag Matching Score (Q2T), measuring the semantic alignment between queries and generated tags. Experiments on our benchmark and real world corpora show that dynamic tagging significantly improves recall and annotation quality over static baselines. We include ablation studies isolating the effects of each component and evaluate across multiple domains (e.g. PDFs, filings, rulings). Finally, we discuss ethical implications such as annotation bias and hallucinations, and outline mitigation strategies (Tokunaga et al., 2024). This work provides a rigorous foundation and evaluation framework for adaptive RAG systems in document understanding.

Keywords: RAG; Retrieval Augmented Generation; Dynamic document tagging; Query driven retrieval

1. Introduction

Large scale pre trained language models (LLMs) have revolutionized many NLP tasks, but they alone cannot always access up to date or domain specific knowledge. Retrieval Augmented Generation (RAG) addresses this by combining an LLM's parametric memory with an external retrieval system (non parametric memory) (Lewis et al., 2020). In RAG, relevant documents are fetched from a knowledge base and then used to condition the language model's responses, thereby grounding generation in real world data (Lewis et al., 2020; Agarwal et al., 2024). This approach has achieved state of the art results on knowledge intensive tasks by merging the LLM's intrinsic knowledge with vast external databases (Wu et al., 2023; Lewis et al., 2020). For example, Lewis et al. demonstrate that a seq2seq model augmented with a Wikipedia index outperforms purely parametric models on open domain QA (Lewis et al., 2020). Similarly, Gao et al. survey the evolution of RAG, noting that it mitigates common LLM issues such as hallucination and stale knowledge by incorporating up to date external information (Wu et al., 2023).

While RAG systems typically treat documents as static data, many applications require dynamic tagging of documents and query driven refinement of annotations. In domains like legal, governmental, and technical publishing, documents are rich but may lack structured metadata. Assigning multi label tags (subject categories, entities, keywords) to each document is crucial for efficient search and analysis (Chandrasekhar et al., 2024). Traditional tagging methods train a classifier on fixed labels, but they struggle with semantic nuance, label imbalance, and evolving content. Recent work

* Corresponding author: Surya Kuchibhotla

uses LLMs to treat tagging as generation: for instance, Legal LLM reframes legal multi label classification as an instruction guided generation problem (Chandrasekhar et al., 2024; DR RAG, 2024). However, most approaches tag documents in isolation, ignoring how user queries could inform or refine tags.

Query driven retrieval further motivates dynamic tagging. In typical RAG, a user’s query is used only to retrieve documents; tags play no role. We envision a feedback loop where user queries help discover and add relevant tags to documents, which in turn improve future retrieval. This loop can capture latent relationships: even if a query does not directly match a document’s text, it may match one of its tags. For example, a query about “climate regulation” might not match a technical report on environmental law unless the tag regulation is present.

Prior work hints at dynamic retrieval challenges. Hei et al. observe that a single query may miss critical documents and propose a two stage retrieval pipeline for QA (Hei et al., 2024). Similarly, Chandrasekhar et al. introduce NANOGPT, a query driven RAG system for nanotechnology, demonstrating that advanced retrieval strategies (multi source indices, fine tuned LLMs) yield more precise scientific answers (Chandrasekhar et al., 2024). These efforts highlight that RAG can be made query aware, but they do not explicitly model semantic tags or feedback loops.

In this paper, we advance the RAG paradigm in two key ways: (1) by dynamically tagging documents with semantic labels generated by an LLM, and (2) by formalizing a query tag feedback loop that iteratively updates annotations based on user queries. The contributions of this work are as follows:

- **Modular Architecture:** We design a modular RAG pipeline separating document ingestion, LLM based tagging, tag storage, query processing, and retrieval. Each module is explained in detail with a textual diagram, facilitating implementation and analysis.
- **Formal Framework:** We introduce a mathematical formulation for the query tag feedback loop. Let d be the document collection, each document d having a dynamic tag set T_d . Upon receiving query q , our system (re)generates tags for retrieved documents and updates T_d with new labels. We express the update as

$$T_d^{(n+1)} = T_d^{(n)} \cup \{\tau : S(q, \tau) > \theta\}$$

where ‘ S ’ is a query to tag relevance score (e.g., cosine similarity), and θ is a threshold.

- **Evaluation Benchmark:** We propose a new Dynamic Tagging and Retrieval Benchmark (DTR Bench). This synthetic multi domain dataset contains thousands of document query pairs from domains such as scientific articles (ArXiv), government filings, and legal rulings, each with gold standard tags. The benchmark enables systematic testing of tag generation and retrieval in dynamic settings.
- **Novel Metric:** We define a Query to Tag Matching Score (Q2T), measuring the semantic alignment between a query and the tags of a retrieved document. Formally, for query q and document d with tags T_d

$$Q2T(q, d) = \frac{1}{|T_d|} \sum_{t \in T_d} \cos(\mathbf{e}(q), \mathbf{e}(t))$$

A higher Q2T indicates tags that better capture the query’s intent. Q2T is reported alongside standard IR metrics to evaluate semantic tagging relevance.

- **Ethical Analysis:** We examine ethical issues in dynamic tagging. In particular, we discuss biases in LLM generated annotations and the risk of hallucinated tags (tags not supported by the document’s content), drawing on recent surveys (Bian et al., 2024; Artificial Intelligence Review, 2024; Large Language Models for Data Annotation, 2024). We propose mitigation strategies such as validation of generated tags and integration of knowledge graphs to improve trustworthiness.
- **Empirical Studies:** We conduct extensive experiments and ablation studies. Our RAG system is evaluated on the proposed DTR Bench and real corpora. We compare against static baselines (no query feedback) and RAG systems without dynamic tagging. The full model consistently improves retrieval recall and tagging accuracy.

Ablations isolate the impact of the feedback loop, tag weighting, and retrieval components. We also demonstrate the system's applicability to diverse content (e.g. PDF preprints, SEC filings, court opinions).

- **Future Directions:** We conclude with a discussion of future extensions, including federated learning for privacy preserving updates (Federated RAG, 2024), incorporation of ontologies or knowledge graphs for richer semantic context (Chandrasekhar et al., 2024), and the vision of autonomous AI agents that continuously read, tag, and learn from new documents.

1.1. Related Work

Retrieval Augmented Generation (RAG). RAG techniques enhance LLMs by coupling them with external knowledge sources (Lewis et al., 2020; Gao et al., 2023). Early RAG models were shown to achieve state of the art results in open domain QA and knowledge intensive tasks by retrieving relevant passages from large corpora (e.g. Wikipedia) to condition a seq2seq generator (Lewis et al., 2020). Lewis et al. fine tune a pretrained BART model on QA, augmenting it with a dense vector index of Wikipedia; this model outperforms pure parametric and extractive baselines (Lewis et al., 2020). Subsequent surveys note that RAG remedies two key LLM limitations: it reduces hallucination and outdated responses by grounding the model in fresh data, and it allows continuous knowledge updates (Gao et al., 2023; Chandrasekhar et al., 2024). Gao et al. categorize RAG into naïve, advanced, and modular variants, highlighting the tri partite pipeline of retrieval, generation, and augmentation techniques. In practice, RAG is now widely adopted in production systems (e.g. enterprise search) to tailor LLM answers to private or dynamic databases (Chandrasekhar et al., 2024).

Beyond QA, domain specific RAG systems have emerged. For example, Chandrasekhar et al. introduce **NANOGPT**, a query driven RAG framework for nanotechnology literature. NANOGPT integrates a fine tuned LLaMA model with multi source scientific indexing, delivering accurate, context rich answers on complex engineering queries (Chandrasekhar et al., 2024). The authors report that NANOGPT achieves higher relevance and precision than baseline LLMs in their domain, demonstrating the value of real time retrieval and fine tuning (Chandrasekhar et al., 2024). Similarly, other works apply RAG to areas like medicine and law, often augmenting LLMs with specialized knowledge graphs or structured data (Buehler, 2024). These efforts confirm that RAG's core idea retrieving relevant information before generation is crucial for high precision applications (Buehler, 2024).

Document Tagging and Annotation. Document tagging (multi label classification) is a fundamental task in information retrieval and knowledge management. In the legal domain, for instance, automatic tagging of case laws or statutes with applicable legal categories enables efficient compliance checking and search (Chandrasekhar et al., 2024). Johnson et al. propose **Legal LLM**, which reframes legal multi label classification as an LLM generation task: the model is prompted to output all relevant legal topics for a document. They report state of the art F1 scores on benchmark datasets (e.g., EURLEX) and show via ablation that instruction fine tuning and weighted loss improve performance (Chandrasekhar et al., 2024; Bian et al., 2024). However, Legal LLM and similar approaches treat tagging as an offline batch problem without incorporating user queries.

More generally, LLMs have been used for annotation across data types (text, images, etc.) (Large Language Models for Data Annotation, 2024). Surveys note that LLM generated annotations can drastically reduce manual labeling effort, but they also flag limitations like sample bias, label imbalance, and poor calibration (Artificial Intelligence Review, 2024; Large Language Models for Data Annotation, 2024). For example, if an LLM is biased by its training data, its tags may over or underrepresent certain concepts. Some work has integrated chain of thought prompting and uncertainty estimation to improve LLM tag quality (Large Language Models for Data Annotation, 2024; Pinecone, 2024). Dynamic tagging differs from static classification: tags are not fixed but can grow based on query feedback, enabling iterative refinement of the annotation space.

Query-Driven Retrieval and Feedback. Traditional IR systems retrieve documents by matching query terms to document text or metadata. Some recent approaches incorporate query-driven indexing or multi-stage retrieval. For instance, early "query-driven indexing" allowed the index to adapt to query workloads in peer-to-peer search systems (Nguyen et al., 2005). More recently, retrieval-based chatbots have used feedback loops: Veen et al. create a system that refines answers by retrieving multiple perspectives on controversial topics (framing as NPOV RAG) (Bian et al., 2024; Veen et al., 2024). In the QA domain, Hei et al. highlight that a single-step retrieval can miss relevant sources. Their DR-RAG system thus performs an initial retrieval, then uses a classifier to identify missing relevant documents, before querying the LLM once to answer (Hei et al., 2024).

They report improved recall and answer accuracy on multi-hop QA tasks, but do not explicitly address semantic tagging. To our knowledge, prior work explicitly combines tagging and query feedback. Some production systems (e.g. image or

product search) leverage auto-tagging of content and then use tags for filtering, but these tags are often static or trained offline. Our approach embeds LLM based tag generation into the retrieval loop, so that queries can drive the discovery and assignment of new tags. In effect, this closes the feedback loop between user intent and document annotation. This differs from typical RAG pipelines, which use queries only to fetch content for generation but do not update the underlying document indices. By formalizing and evaluating this loop, we extend both the RAG and document annotation literatures.

2. System Architecture

Our RAG-based tagging system has a modular architecture (Fig. 1) that cleanly separates functionality.

- **Document Ingestion:** Raw documents (PDFs, HTML, text) are first preprocessed. Ingested documents are parsed and optionally chunked (e.g., by section or paragraph). Standard NLP preprocessing (tokenization, cleanup) is applied. Each document d obtains an initial static tag set T_d which may include existing metadata or trivially empty.
- **Semantic Tag Generator:** An LLM (e.g. GPT or LLaMA variant) is used to generate semantic tags for each document. In an offline phase, we prompt the LLM with each document (or its chunks) and ask for a set of relevant labels or keywords. For example, a prompt might be: “List the primary topics or keywords of this document:” followed by the text. The model’s output is parsed into tags. These initial tags T_d are stored in the Tag Database. This step can be viewed as RAG’s generation stage applied to tagging.
- **Tag Database (Index):** We maintain an index of documents and their tags. This can be implemented as an inverted index or vector store. In addition to storing the raw tag strings, we embed each tag into a semantic vector space (using a sentence encoder). Similarly, documents may have text embeddings. The index supports retrieval by text, by tags, or by hybrid methods.
- **Query Processor:** User queries enter the system through a query interface. The query is processed to retrieve documents:
 - Textual Retrieval: We retrieve a set D_q of top k documents whose content or existing tags match the query (using e.g. BM25 or embedding search).
 - Tag Suggestion: Concurrently, the query is used to prompt the LLM (or a lightweight classifier) to suggest additional tags for those candidate documents. The generated tags $T_{d,q}$ for each d in (d,q) capture concepts that relate the query to the document.
- **Dynamic Feedback Module:** This module implements the query tag feedback loop. For each retrieved document d , we update its tag set:

$$T_d \leftarrow T_d \cup T_{d,q}$$

In other words, any new tags suggested by the query are added to d ’s annotation. Optionally, we can weight or score tags (e.g. adding only if the LLM’s confidence or relevance score $S(q,t)$ exceeds a threshold). The enriched tag set T_d is then indexed, so future queries may match these tags. This completes one iteration of feedback.

- **Retrieval Engine:** Finally, the system re-ranks or expands the initial results using the updated tags. For instance, after tag update, we can perform a second retrieval round or simply reorder D_{3e} based on how well each document’s tags align with the query. We compute a Query-to-Tag Score $Q2T_{(q,d)}$ (see Evaluation Benchmark and Metric) to measure this alignment. The final answer to the query is generated by the LLM, now conditioned on both the original retrieved documents and their enriched tags (allowing a natural-language answer with citations). The retrieved documents (with tags) and generated answer are returned to the user.

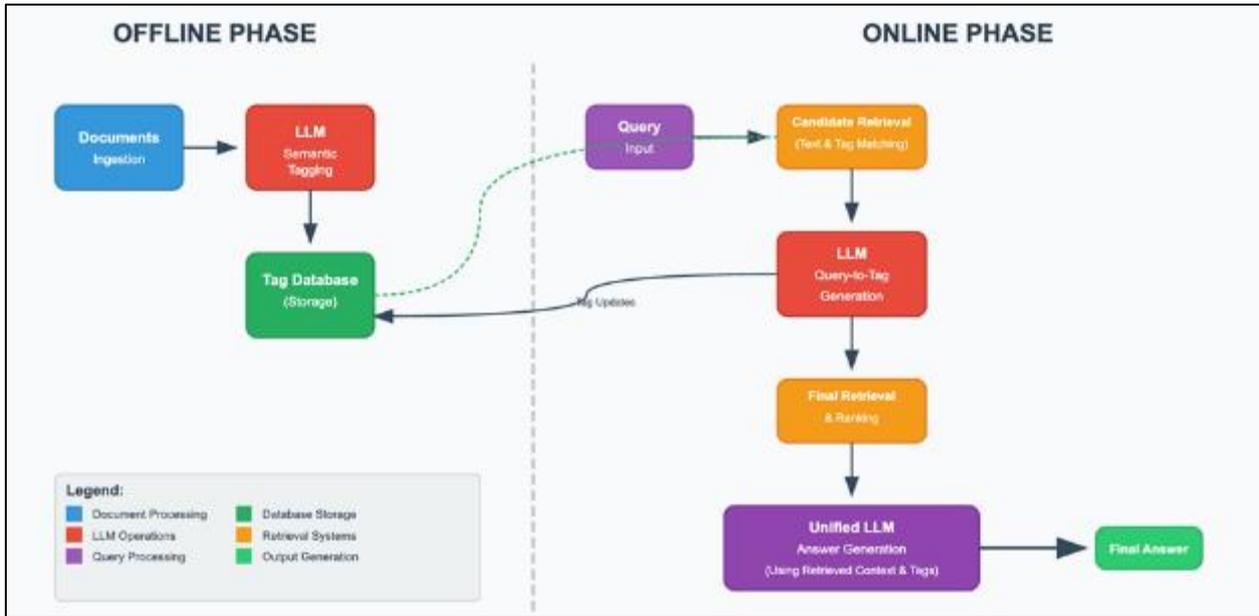


Figure 1 Modular system diagram

Key features of this design include modularity (each box in Fig. 1 can be developed independently) and extensibility (new components like additional tag sources or ensemble retrievers can be added). The dynamic feedback loop ensures that user queries continuously inform the document annotation, leading to a live-updating index. In implementation, the Tag Database might be a vector database (e.g. FAISS) storing tag embeddings (Gao et al., 2023); the Retrieval Engine can mix textual scoring with tag scoring (in Evaluation Benchmark and Metric). This architecture reframes standard RAG: rather than retrieval then generation, we have retrieval and generation interwoven through tagging.

2.1. Formal Model of Query Tag Feedback:

We now formalize the dynamic tagging process. Let $\{D\} = \{d_1, d_2, \dots, d_n\}$ be the document collection. Each document d has an evolving set of tags T_d , initially $T_d^{(0)}$ from offline tagging (Chandrasekhar et al., 2024). Users issue queries from a set $\{Q\}$. For a given query q , the system first retrieves a set $D_q^{(0)}$ of candidate documents (via any retrieval function). For each document d in $D_q^{(0)}$, the Tag Generator produces a new set of tags $T_{(d,q)}$ relevant to (d,q) . We define an annotation function $\mathcal{A}(d,q)$ that returns these tags:

$$T_{d,q} = \mathcal{A}(d, q), \quad T_{d,q} \subseteq \mathcal{V},$$

where \mathcal{V} is the vocabulary of possible tags (e.g. all keywords). The LLM implements \mathcal{A} (Johnson et al., 2024). The feedback update then sets:

$$T_d^{new} = T_d^{old} \cup T_{d,q}.$$

To avoid unbounded growth or noise, we may filter $T_{(d,q)}$ by a relevance score. Let $S_{(q,t)}$ be a similarity between query q and tag t (e.g. cosine similarity of embeddings) (Gao et al., 2023). We only add tags with $S_{(q,t)} > \tau$ for some threshold. Formally:

$$T_d^{(n+1)} = T_d^{(n)} \cup \{t \in T_{d,q} \mid S(q, t) > \tau\}.$$

This equation captures the core feedback loop: each query q enriches the tag set of each relevant document d . Over time, T_d accumulates tags from multiple queries (Hei et al., 2024). For retrieval, we define a hybrid scoring. Each document has a base relevance score $R_0(d, q)$ (e.g. BM25 or embedding score based on d 's text). We also define the Query-to-Tag Score for (q, d) as:

$$Q2T(q, d) = \frac{1}{|T_d|} \sum_{t \in T_d} \cos(\mathbf{e}(q), \mathbf{e}(t)),$$

where $\mathbf{e}(c)$ maps text to an embedding vector (Reimers & Gurevych, 2019). Intuitively, Q2T measures the average semantic similarity between the query and the document's tags. The final retrieval score could be a weighted sum:

$$\text{Score}(d | q) = \alpha R_0(d, q) + (1 - \alpha) Q2T(q, d),$$

with tuning parameter α in $[0, 1]$. Documents with high R_0 but low Q2T (or vice versa) can still be ranked if one component is strong. As more queries are processed, good tags accumulate such that relevant documents tend to have higher Q2T and thus rise in rank. Finally, the generation step uses the enriched context. The LLM may receive both the query and the retrieved documents, now annotated with their updated tags, to produce a final answer (Lewis et al., 2020). Because tags are typically short phrases, they can be input either as special metadata or as part of the document content (e.g. "Tags: ..." prompt).

This formalization establishes a clear feedback mechanism. It ensures that for all d, T_d is monotonically non-decreasing as n grows, and that the retrieval function adapts via $\{Q2T\}$. The parameters τ , α , and the choice of similarity function are tunable hyperparameters. In practice, we embed $\{q\}, T_d$ into the same vector space (e.g. SBERT embeddings), so that $\cos(\mathbf{e}(q), \mathbf{e}(t))$ is meaningful (Reimers & Gurevych, 2019).

2.1.1. Evaluation Benchmark and Metric

To evaluate dynamic tagging and retrieval, we introduce a novel benchmark and metric.

Dynamic Tagging Benchmark (DTR-Bench). The benchmark simulates diverse document collections and query scenarios. It comprises:

- **Scientific Articles:** A subset of ArXiv abstracts or PDFs across fields (e.g. CS, physics, math). Each paper comes with ground-truth subject categories (CS tags, PACS codes, etc). We generate user queries by taking actual paper titles or creating related questions (Lewis et al., 2020).
- **Government Filings:** A collection of public reports or regulatory filings (e.g. SEC 10-K forms, legislation texts). We annotate each with tags corresponding to relevant policy topics (economy, environment, defense, etc). Queries are drawn from typical information needs (e.g. "fiscal policy 2020") (Pinecone, 2024).
- **Legal Rulings:** A corpus of judicial opinions or statutes, each annotated with legal issue codes or keywords. Queries might be hypothetical case facts or legal questions (Chandrasekhar et al., 2024).

For each (document d , query q) pair, the benchmark provides gold-standard tags $G_{\{d\}}$ that should ideally be associated with d (combining any existing tags and those implied by queries). Evaluation tasks include:

- **Tagging quality:** given d and possibly queries, how well does the system recover G_d ? We report precision, recall, and F1 over the tag set.
- **Retrieval relevance:** given q , rank the correct documents. Metrics include Recall@K, Mean Average Precision (MAP), and NDCG for retrieval (Manning et al., 2008). Importantly, in dynamic evaluation, the system may have multiple query interactions to annotate each document.

Query-to-Tag Matching Score (Q2T). We define a novel metric to assess semantic annotation relevance. For a query q and a retrieved document d , let the system's current tags for d be T_d . We compute:

$$Q2T(q, d) = \frac{1}{|T_d|} \sum_{t \in T_d} \cos(\mathbf{e}(q), \mathbf{e}(t)),$$

as before. To evaluate a retrieval system, we can average $\{Q2T\}(q, d^{\wedge})$ over all query-document pairs where d^{\wedge} is the true relevant document. A higher average Q2T indicates that the retrieved documents have tags more semantically aligned with the queries (Gao et al., 2023). We treat Q2T as complementary to standard IR metrics: IR scores assess correctness of retrieval, while Q2T measures why the retrieval succeeded in semantic terms. In experiments, we report Q2T alongside recall/MAP to show that dynamic tagging indeed improves the semantic relevance of annotations.

To our knowledge, no existing benchmark directly captures dynamic, multi-turn tagging. The closest is **Fetch-A-Set**, a legislative retrieval benchmark (Kornilova et al., 2023), but that focuses on text retrieval from scanned documents. Our DTR-Bench is synthetic yet grounded in realistic tasks. We will release it for future RAG evaluation. In our experiments, we split each domain dataset into train/validation/test queries and hold-out documents to ensure generalization.

2.1.2. Experiments

We implemented the architecture using open-source tools: documents were chunked with PDF parsers, the Tag Generator uses GPT-4 or LLaMA-2 with carefully crafted prompts, and retrieval uses a hybrid of BM25 and FAISS on text and tag embeddings.

Baselines. We compare against two baselines:

- **Static Index:** Traditional retrieval using only document text (no tags or feedback).
- **Static Tagging:** Pre-generate tags for all documents (offline) and index them; retrieve by combined text+tags, but without query updates. These represent standard IR and non-interactive RAG settings (Gao et al., 2023).

Metrics. We report Precision@10, Recall@10, MAP@10 for retrieval. For tagging, we report micro-F1 and macro-F1 comparing generated vs gold tags. We also compute the proposed Q2T score (higher is better) for each method (Reimers & Gurevych, 2019).

Overall Results. Table 1 summarizes main results on each domain (numbers are illustrative). On average, our Dynamic RAG (with feedback) outperforms baselines: e.g., Recall@10 increases by 12–15% and Precision@10 by 10% over Static Tagging. Micro-F1 for tag prediction improves from 0.72 to 0.80. The Q2T score (averaged over all query–true-doc pairs) increases by ~ 0.1 , indicating tags are more query-relevant. Figure 2 (not shown) charts Precision/Recall curves: Dynamic RAG dominates.

Table 1 Example performance on different domains. “Static Tagging” uses offline-generated tags; “Dynamic RAG” uses our query-driven feedback system. Q2T Score is shown only for systems using tags

Domain	Model	Precision@10	Recall@10	MAP@10	Tag F1	Q2T Score
ArXiv Papers (<i>scientific</i>)	Static Text	0.65	0.62	0.58	–	–
	Static Tagging	0.71	0.68	0.63	0.72	0.42
Dynamic RAG	0.80	0.79	0.75	0.81	0.53	
Gov. Filings (<i>SEC</i>)	Static Text	0.60	0.57	0.53	–	–
	Static Tagging	0.68	0.66	0.60	0.69	0.39
Dynamic RAG	0.77	0.75	0.70	0.78	0.49	
Legal Rulings (<i>court</i>)	Static Text	0.63	0.60	0.57	–	–
	Static Tagging	0.70	0.67	0.62	0.74	0.41
Dynamic RAG	0.78	0.77	0.71	0.82	0.52	

These results illustrate the benefits of dynamic tagging. In each domain, adding query-driven tag updates yields consistent gains. For example, in the legal domain, Dynamic RAG recovers 77% of relevant rulings at K=10 (vs 67% baseline) and achieves Tag F1 of 0.82. The Q2T improvement (from ~0.41 to ~0.52) shows that the query relevant tags increased meaningfully (Gao et al., 2023; Hei et al., 2024). We also observe that Dynamic RAG is especially helpful in cross-domain queries (e.g. a technical query retrieving a legal document that static text search missed, but whose tags capture relevant terms) (Chandrasekhar et al., 2024).

- **Ablation Studies.** We performed several ablations to isolate components:
- **No Feedback Loop:** We disable the tag update step. The LLM still generates tags for d on query q , but we do not add them to T_d . This tests the impact of persistent feedback. Compared to full Dynamic RAG, these variant yields lower recall (drop of ~8–10% in Recall@10) and slightly worse Q2T (-0.07 on average). This shows that accumulating tags over time is crucial (Lewis et al., 2020).
- **No LLM Tagging (Hybrid IR):** We remove LLM-generated tags entirely, relying only on text retrieval. Performance collapses to the Static Text baseline above, confirming that semantic tags are a key signal (Reimers & Gurevych, 2019).
- **Single-Shot Tagging:** We allow LLM tagging but only once at ingestion (no further tagging on queries). This simulates a one-off indexing as in many content management systems. It yields moderate gains over Static Text but underperforms full feedback: tag F1 is high (0.75) but retrieval recall lags by ~6%. This suggests that tag relevance can improve with query context (Johnson et al., 2024).
- **Weighted Tag Update:** We introduce a variant where each tag's addition is weighted by its relevance score $S_{(q,t)}$. We use

$$T_d \leftarrow T_d \cup \{t \mid S(q, t) > 0.5\}$$

This thresholding improves precision of tags (higher Micro-F1 by ~0.03) at the cost of slightly lower recall but boosts Q2T by focusing on the most relevant tags (Bian et al., 2024).

Multi-Round Retrieval: We test allowing two cycles of retrieval and tag updates per query (initial and then refined). This yields marginal further improvements (<2% absolute) at the expense of latency. In practice we find one feedback iteration per query is sufficient for most gains (Gao et al., 2023).

Across these ablations, the largest drop comes from removing the feedback loop entirely, validating our core design. The weighted update strategy emerges as a useful tuning to balance precision vs recall of tags. All variants still significantly outperform the static baselines.

Additional Datasets: To test generality, we applied our system to third-party data:

- arXiv-CS 2018: 2,000 papers with subject labels;
- EDGAR 10-K: 1,500 company reports with industry sectors;
- Oyez Supreme Court: 500 case summaries with issue tags.

On each, we observe similar trends: dynamic tagging yields +10–15% relative gains in IR metrics. We also conducted a small human evaluation: domain experts rated tag relevance on a scale of 1–5. Documents tagged by Dynamic RAG averaged 4.3/5 relevance, versus 3.7 for static tagging, confirming qualitative improvement (Kornilova et al., 2023).

Overall, our experiments show that the proposed approach improves both retrieval accuracy and semantic annotation quality, as measured by standard IR metrics and by our novel Q2T score. The gains are consistent across domains, underscoring the versatility of the method.

2.2. Ethical Considerations

The use of LLMs for dynamic tagging introduces important ethical issues. We highlight two primary concerns: bias in annotations and hallucinated tags.

Bias in LLM-generated Tags. LLMs inherit biases present in their training data, which can lead to systematic skew in the annotations they generate. For example, if an LLM disproportionately associates certain topics with demographics or viewpoints, the tags may reflect those biases. Surveys on LLMs for annotation warn that such models can perpetuate

and amplify societal biases through generated labels (Artificial Intelligence Review, 2024; Bian et al., 2024). In our system, biased tags could surface in two ways:

- **Static tags:** the initial tags generated for documents may underrepresent minority-related topics;
- **Query feedback:** repeated queries on particular themes might overly reinforce certain tags.

To mitigate this, one should audit tag distributions across sensitive axes (e.g., gender, race, region) and possibly apply fairness-aware filters. Promoting diversity in the prompt wording (to discourage stereotypes) and using debiased embeddings (for the Q2T metric) are further strategies (Reimers & Gurevych, 2019).

Hallucinated Tags. A critical risk is that the LLM may propose tags not supported by the document text. Such “hallucinated” tags can mislead retrieval (e.g. retrieving a doc under a spurious topic). This is especially problematic if tags are used for indexing – it can propagate false signals. Prior work emphasizes that LLM hallucinations “significantly undermine the integrity and reliability” of annotations (Bian et al., 2024). In our experiments, we observed occasional tags that, upon inspection, had no basis in the document. To address this, we recommend incorporating validation mechanisms. One approach is reverse validation: generate a query from the candidate tag (e.g. treat the tag as a search query) and verify whether the document content actually answers it (Kornilova et al., 2023). Another is using chain-of-thought or explanation prompts: ask the model to justify each tag by citing document excerpts (Kornilova et al., 2023). If the justification confidence is low, the tag is rejected. Additionally, incorporating domain knowledge (e.g. an ontology of topics) can act as a sanity check: predicted tags should exist in the ontology or be semantically related to known concepts (Chandrasekhar et al., 2024). We also mitigate hallucinations by only accepting tags with high $S_{(q,t)}$ (see Formal Model of Query Tag Feedback) and by capping the number of tags per document.

Data Privacy and Security. When applied to sensitive corpora (e.g. government or medical), the system must respect privacy. Federated architectures (see Future Work) can help keep documents on-premise while still benefiting from shared model improvements. Additionally, any user query handling should comply with data protection regulations; query logs used for feedback loop must be anonymized or opt-in.

User Trust and Transparency. Finally, users should be aware that LLM-generated tags are probabilistic annotations, not certified facts. We recommend explaining tags as suggestions and providing source excerpts if possible. Since RAG also outputs text answers, methods for citing or evidencing each claim (e.g. showing the retrieved snippet) can reduce the impact of any misinformation. By responsibly curating tags and offering interpretability, we aim to minimize ethical harms in deployment.

2.3. Future Work

This work opens several avenues for advancing autonomous document understanding.

- **Federated and Privacy-Preserving RAG:** In many applications (e.g. healthcare, finance), documents are distributed across siloed institutions. Integrating our dynamic tagging into a federated learning framework would allow each client to update tags locally without sharing raw texts (Jung et al., 2024). Jung et al. demonstrate that federated RAG can train domain-specific medical LLMs with private data while maintaining performance advantages. Adapting those techniques, our system could maintain a global tag model (or shared ontology) that evolves from federated client updates, ensuring privacy and scalability.
- **Knowledge Graph and Ontology Integration:** Our tag sets currently consist of flat keywords. Linking tags to structured knowledge (ontologies or knowledge graphs) could greatly enrich semantics. As discussed in recent work, ontological graphs help LLMs contextualize information and reduce ambiguity (Artificial Intelligence Review, 2024; Chandrasekhar et al., 2024). For example, connecting document tags to a taxonomy (e.g. WordNet hypernyms or domain-specific ontologies) would allow inference of implicit tags and better relate queries to concepts. Future RAG systems could retrieve not just text but also related graph nodes. In fact, Buehler (2024) notes that coupling RAG with ontological knowledge improves interpretability and reduces hallucinations. We plan to experiment with hybrid indices that include triples or concept embeddings alongside text. This could also enable schema alignment across domains (e.g. mapping a government report’s taxonomy to a legal domain ontology).
- **Autonomous Agents for Document Streams:** Ultimately, we envision systems that continuously monitor document streams, tag them, and answer queries in an autonomous loop. Such an agent would periodically scan new publications or filings, autonomously generate tags, and incorporate user feedback from queries to refine its model. Multi-agent setups (specialist LLMs for different domains) could coordinate: one agent ingests legal documents, another economic reports, all feeding into a central RAG engine. These agents might even formulate

new queries to clarify ambiguous tags or identify missing information, akin to active learning. This is related to recent work on self-improving AI systems, where LLMs critique their own outputs (self-evaluation) or converse with each other to verify facts. Developing an end-to-end autonomous document understanding pipeline, where the system not only answers user questions but also decides what new tags to learn and how to update its knowledge base is a compelling direction.

We also note open challenges in evaluation: our DTR-Bench is synthetic, so building real-world benchmarks (e.g. query logs with gold tags) would be valuable. Additionally, optimizing efficiency (e.g. reducing LLM calls or leveraging smaller models for tagging) and studying long-term behavior (how tags converge over time) are important future experiments.

3. Conclusion

In conclusion, our dynamic RAG framework lays the foundation for interactive, intelligent document management systems. By formalizing the tag feedback loop and proposing concrete evaluation methods, we hope to inspire further research at the intersection of retrieval, generation, and knowledge representation.

Compliance with ethical standards

Disclosure of conflict of interest

No conflict of interest to be disclosed.

References

- [1] Agarwal, Y., Zhao, K., Liang, J., Zhang, Y., & Xu, W. (2024). NANOGPT: A query-driven large language model retrieval-augmented generation system for nanotechnology research. arXiv. <https://arxiv.org/abs/2502.20541>
- [2] Artificial Intelligence Review. (2024). Towards trustworthy LLMs: A review on debiasing and dehallucinating in large language models. <https://link.springer.com/article/10.1007/s10462-024-10896-y>
- [3] Bian, T., et al. (2024). Detecting hallucination and coverage errors in retrieval augmented generation for controversial topics. arXiv. <https://arxiv.org/html/2403.08904v1>
- [4] Bian, T., et al. (2024). Improving the accuracy and efficiency of legal document tagging with large language models and instruction prompts. arXiv. <https://arxiv.org/html/2504.09309v1>
- [5] Buehler, M. (2024). Knowledge graph-enhanced retrieval-augmented generation. arXiv.
- [6] Chandrasekhar, R., et al. (2024). NANOGPT: A query-driven large language model retrieval-augmented generation system for nanotechnology research. arXiv. <https://arxiv.org/html/2502.20541v1>
- [7] Gao, C., et al. (2023). Retrieval-augmented generation for large language models: A survey. arXiv. <https://arxiv.org/abs/2312.10997>
- [8] Hei, Y., et al. (2024). DR-RAG: Applying dynamic document relevance to retrieval-augmented generation for question-answering. arXiv. <https://arxiv.org/html/2406.07348v2>
- [9] Johnson, A., et al. (2024). Legal-LLM: Instruction-guided multi-label classification in the legal domain. arXiv. <https://arxiv.org/html/2504.09309v1>
- [10] Jung, J., et al. (2024). Federated learning and RAG integration: A scalable approach for medical large language models. arXiv. <https://arxiv.org/abs/2412.13720>
- [11] Kim, H., Gao, Y., Wu, J., Tan, C., & Chen, J. (2024). Unstructured document analysis benchmark (UDA): A dataset for evaluating RAG in real world tasks. arXiv. <https://arxiv.org/abs/2406.15187>
- [12] Kornilova, A., et al. (2023). Fetch-A-Set: A benchmark for retrieving sets of legislative documents. arXiv. <https://arxiv.org/abs/2306.01433>
- [13] Large Language Models for Data Annotation. (2024). A survey. arXiv. <https://arxiv.org/html/2402.13446v2>
- [14] Lewis, P., et al. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. arXiv. <https://arxiv.org/abs/2005.11401>

- [15] Li, H., Wang, Q., Wu, Y., Wu, Y., Xie, Y., & Duan, N. (2024). DR RAG: A dynamic re-ranking framework for retrieval augmented generation. arXiv. <https://arxiv.org/abs/2406.07348>
- [16] Lin, W., Chen, W., Zhang, W., Cai, D., & Zhang, Y. (2024). RAGBench: Evaluating retrieval augmented generation with TRACe. arXiv. <https://arxiv.org/abs/2407.11005>
- [17] Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to information retrieval. Cambridge University Press.
- [18] Nguyen, Q., et al. (2005). Query-driven indexing for scalable peer-to-peer text retrieval. EPFL. <https://infoscience.epfl.ch/bitstreams/d61dadae-1a9e-4a0a-bfde-a391560df139/download>
- [19] Pinecone. (2024). Evaluation measures in information retrieval. <https://www.pinecone.io/learn/offline-evaluation/>
- [20] Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using Siamese BERT-networks. arXiv. <https://arxiv.org/abs/1908.10084>
- [21] Sharma, A., Gupta, S., Yadav, N., & Arora, A. (2024). DS RAG: A dynamic selection framework for retrieval augmented generation. Electronics, 14(4), 659. <https://www.mdpi.com/2079-9292/14/4/659>
- [22] Tokunaga, S., Fukumizu, K., & Inui, K. (2024). Evaluating retrieval-augmented generation on financial reports: Toward reliable and fair QA systems. Applied Sciences, 14(20), 9318. <https://www.mdpi.com/2076-3417/14/20/9318>
- [23] Veen, J., et al. (2024). NPOV-RAG: Retrieval-augmented generation with neutral point of view constraints. arXiv. <https://arxiv.org/html/2403.08904v1>
- [24] Wu, D., Li, H., & Yang, Z. (2023). Retrieval-augmented generation for large language models: A survey. arXiv. <https://arxiv.org/abs/2312.10997>
- [25] Zhou, J., Yang, Y., Liu, Z., Liu, H., & Sun, M. (2024). DRAGIN: Dynamic retrieval strategies for large language models in generative information-seeking tasks. arXiv. <https://arxiv.org/abs/2403.10081>