Entropy-Based Dynamic Hybrid Retrieval for Adaptive Query Weighting in RAG Pipelines

J. R. Perez¹ James Zhou² Radley Le³ Alexander Menchtchikov⁴ Ryan Lagasse⁵

Abstract

Traditional sparse and dense retrieval methods independently exhibit critical limitations: sparse models offer high lexical precision but lack semantic flexibility, while dense models capture semantic similarity but may introduce false positives due to embedding generalization. Hybrid retrieval aims to unify their strengths, yet current methods typically use static weighting, failing to adapt to query-specific retrieval uncertainties. We propose a dynamic hybrid retrieval method that performs multi-round entropy-based reweighting to iteratively optimize the linear combination of sparse and dense scores. Leveraging normalized Shannon entropy as a proxy for retrieval confidence, we update weight coefficients w_s and w_d across iterations until convergence or a predefined maximum is reached. The top-k documents are reranked at each step, using fixed sparse and dense retrieval outputs, improving robustness without repeated querying. We implement our approach using a BM25-FAISS hybrid pipeline with MiniLM-L6-v2 embeddings and evaluate performance on HotPotQA and TriviaQA. Experimental results demonstrate that our dynamic hybrid model, under an optimal convergence threshold of $\epsilon = 0.10$, significantly outperforms both pure dense and fixed-weight hybrid baselines in LLM-as-a-Judge (LLMJ) scores across both datasets, with statistically significant gains on TriviaQA (p < 0.01) and marginal gains on HotPotQA ($p \approx 0.055$), confirming the efficacy of adaptive retrieval.

1. Introduction

Information retrieval (IR) is a critical component in the retrieval-augmented generation (RAG) pipeline, which utilizes both IR and natural language processing (NLP) for enhanced large language model (LLM) outputs via external knowledge sources (Lewis et al., 2020). Traditional pure RAG systems typically utilize a single retrieval methodology, usually dense vector retrieval using embedding similarity, where documents and queries are embedded into a shared vector space and their relevance is computed through similarity metrics like cosine similarity (Karpukhin et al., 2020). However, with increasing digital information volume and RAG popularization in modern applications, search efficacy optimization is critical. Traditional retrieval models, both sparse and dense, have well-documented strengths and weaknesses: sparse retrieval excels in precise keyword matching and subsequent retrieval but struggles with semantic representation, while dense retrieval improves semantic understanding at the cost of increased probabilities of false positives due to vector embedding generalization errors (Mandikal & Mooney, 2024).

Currently, hybrid retrieval systems are utilized to combine sparse and dense methods for optimal retrieval. However, existing hybrid models often rely on static weighting strategies, where a predefined and fixed combination of sparse and dense retrieval scores determines ranking. These methods fail to adapt dynamically in response to varying query complexities and retrieval uncertainties (Zhang et al., 2024).

In response to these limitations, this study investigates and proposes a multi-round entropy-based re-ranking approach to improve retrieval confidence and result relevance. This approach uses a weighted sparse-dense retrieval combination and consequent iterative re-ranking based on Shannon semantic entropy scores that adjusts the weights of the sparse and dense contributions dynamically. This method is retriever-agnostic; entropy is used as a metric to address retrieval under uncertainty. We hypothesize that hybrid retrieval methods that combine sparse and dense retrieval outperform pure static RAG retrieval and that adaptive weighting based on retrieval entropy can accommodate the weaknesses of the sparse-dense combination for each specific query. The computational overhead of this entropy-based

¹Lehigh University, Bethlehem, PA ²University of Maryland, College Park, MD ³Bucknell University, Lewisberg, PA ⁴San Mateo County Community College District, San Mateo, CA ⁵University of Connecticut, Mansfield, CT. Correspondence to: J.R. (John Richard) Perez <jrp527@lehigh.edu>.

Proceedings of the 1st Workshop on Vector Databases at International Conference on Machine Learning, 2025. Copyright 2025 by the author(s).

optimization is justified by improved retrieval quality.

2. Related Works

Dynamic Alpha Tuning (DAT). Dynamic Alpha Tuning (DAT) is a dynamic model that adjusts the weighting coefficient (alpha) between sparse and dense retrievers based on model confidence. DAT employs meta-learned schemes to adaptively skew contributions, and recent studies show that this yields coverage and answer diversity advantages (Hsu & Tzeng, 2025). This method is used and applied iteratively within our framework.

Uncertainty Methods. Recent work has expanded on other uncertainty metrics beyond Shannon entropy. Recent literature evaluates Bayesian methods, such as variational retrieval confidence and Monte Carlo dropout, as an alternative to estimate principled query-specific weighting (Gal & Ghahramani, 2016; Laves et al., 2020). Calibration-based measured (expected confidence intervals, mutual information, etc.) have also shown improved robustness in out-ofnoise and noisy retrieval environments. These metrics are possible substitutes to the general proposed framework.

ColBERTv2 and CRAG. Recent advances in RAG have yielded high-performing strong learned hybrids such as Col-BERTv2 and iterative composition retrieval-augmented generation (CRAG) (Santhanam et al., 2021; Shi et al., 2024). Respectively, these methods show advantages in complex multi-hop questions and late-interaction architectural retrieval, both outperforming static hybrids on standard benchmarks. While both methods represent state-of-the-art baselines for hybrid and learned retrieval, our current work does not implement them due to time and resource constraints, such as the computational overhead for robust replication. Comprehensive SOTA comparisons may include CRAG and ColBERTv2 for future work, for a more thorough empirical evaluation of hybrid retrieval strategies.

3. Background

Sparse retrieval algorithms retrieve documents by matching exact keywords from the query to the documents. The most widely used sparse retrieval algorithm is BM25, which computes relevance scores using term frequency (TF) and inverse document frequency (IDF) (Robertson & Zaragoza, 2009). In the case of BM25, higher relevance scores are assigned to documents with higher frequencies of queried terms (TF), but adjust the general prevalence of the term in the corpus, or document space, to account for overly common words (IDF). The ranking function is given by

$$\mathbf{S}_{BM25}(D,Q) = \sum_{t \in Q} \frac{\text{IDF}(t)f(t,D)}{(k_1+1)f(t,D) + k_1 \left(1 - b + b\frac{|D|}{\text{avgdl}}\right)}$$
(1)

Where f(t, D) is the term frequency of term t in document D, |D| is document length, avgdl is the average document length in the corpus, and k_1 , b are hyperparameters controlling the saturation of frequency scaling and the degree of length normalization, respectively. While these perform efficiently with well-defined queries containing relevant key terms, they struggle in capturing semantic relationships between words, limiting efficacy for queries with significant lexical variation.

Dense retrieval algorithms, on the other hand, map queries and documents into high-dimensional vector spaces using deep learning models, usually through contrastive learning or softmax-based loss functions. Recent studies demonstrate that unsupervised dense retrievers trained through constrastive learning outperform traditional sparse methods like BM25 on various benchmark, making them ideal for pure RAG pipelines (Izacard et al., 2021).

In this paper, we use Facebook AI similarity Search (FAISS), a widely used approximate nearest neighbors (ANN) search algorithm for dense retrieval that utilizes cosine similarity. The cosine similarity score used for FAISS-based dense retrieval is:

$$\mathbf{S}_{FAISS}(D,Q) = \left\{\frac{q \cdot d_i}{\|q\| \|d_i\|}\right\}_{i=1}^k \tag{2}$$

where q and d_i are query and document vectors, and $\|\cdot\|$ is the Euclidean norm. Scores are normalized (Johnson et al., 2017).

Semantic entropy quantifies the uncertainty and disorder within a distribution, and in this paper, is used as an indicator of confidence in the ranking scores of different retrieval algorithms. Retrieval methods resulting in low entropy, and therefore lower uncertainty, are associated with higher confidence in ranking assignments, while those with higher entropy suggest a greater amount of ranking uncertainty.

In this paper, we utilize normalized Shannon entropy as a proxy for retrieval uncertainty. For a set of top-k scores $S = \{s_1, s_2, \ldots, s_k\}$, we compute the probability distribution:

$$p_i = \frac{s_i}{\sum_{j=1}^k s_j}$$

The Shannon entropy over these normalized scores is:

$$H(S) = -\sum_{i=1}^{k} p_i \log p_i$$

To ensure comparability across different values of k, we normalize the entropy by dividing by the maximum possible entropy $\log k$:

$$\hat{H}(S) = \frac{H(S)}{\log k}$$

This normalized form ensures $\hat{H} \in [0, 1]$, enabling interpretable weighting across queries. We use this for both sparse and dense score distributions. It should be noted that other uncertainty measures may be used in future work.

The individual limitations of sparse and dense retrieval methods have motivated the development and implementation of hybrid retrieval pipelines that integrate and use both approaches in IR systems, balancing both precision and recall.

Queries are fundamental inputs to information retrieval systems that serve as the main interface between users and the given retrieval mechanism. However, not all queries behave homogeneously and uniformly within retrieval pipelines, with some being highly structured and keyword-focused, while others may have semantic complexity that requires the ability to capture nuanced meanings. For example, queries may be open-ended and lack specific keywords, while closeended queries may be more succinct but lack the variability for deeper and subtle interpretations (Bailey et al., 2017). Current static weight approaches overlook these differences and apply a predefined and fixed combination of sparse and dense scores without accounting for query variability. In the proposed model, queries are treated as dynamic elements that guide retrieval optimization, where retrieval efficacy adapts to query characteristics rather than operating under fixed assumptions.

4. Model

Under an iterative entropy-based framework, this model converges on the ideal weighting parameters through inverseentropy normalization per iteration. The sparse and dense document sets are retrieved once per query and held fixed; entropy is computed over these fixed sets. The weighting parameters are iteratively updated until the weight delta $|\Delta w_s| \leq \epsilon$ or a maximum of n iterations is reached.

4.1. Entropy-based Optimization

We utilize entropy for weight optimization and adjustment under the observation that different queries interact with sparse and dense methods in distinct ways, and therefore depending on the query, each call necessitates a different weighting for retrieval contributions.

In order to implement entropy-based optimization, we employ a multi-step process. Let ϵ be the threshold for weight convergence. Let t be the iteration index, up until the condition $\left|\Delta w_s^{(t)}\right| \leq \epsilon$ or t = n. Let k represent the number of top documents $d_i \in D$ retained for final ranking.

Initialization. Given a query Q, we retrieve the top-k documents independently from *BM25* and *FAISS*.

$$S_{\text{sparse}} = \{S_{\text{sparse},1}, S_{\text{sparse},2}, \dots, S_{\text{sparse},k}\}$$

$$S_{\text{dense}} = \{S_{\text{dense},1}, S_{\text{dense},2}, \dots, S_{\text{dense},k}\}$$

These scores are normalized to form standard probability distributions:

$$p(S_{\text{sparse},i}) = \frac{S_{\text{sparse},i}}{\sum_{j=1}^{k} S_{\text{sparse},j}},$$
$$p(S_{\text{dense},i}) = \frac{S_{\text{dense},i}}{\sum_{j=1}^{k} S_{\text{dense},j}}$$

Initially, we set equal weights for both retrieval methods:

$$w_s^{(0)} = w_d^{(0)} = 0.5$$

Entropy-guided Weight Update. Next, we compute the *normalized Shannon entropy* for both distributions. The entropy values are defined as:

$$H_{\text{sparse}} = -\sum_{i=1}^{k} p(S_{\text{sparse},i}) \log p(S_{\text{sparse},i})$$
(3)

$$H_{\text{dense}} = -\sum_{i=1}^{\kappa} p(S_{\text{dense},i}) \log p(S_{\text{dense},i})$$
(4)

$$\hat{H}_{\text{sparse}} = \frac{H_{\text{sparse}}}{\log k}, \quad \hat{H}_{\text{dense}} = \frac{H_{\text{dense}}}{\log k}$$
(5)

At each iteration *t*, we update the sparse weight using inverse normalized entropy:

$$w_s^{(t+1)} = \frac{1 - \hat{H}_{\text{sparse}}}{(1 - \hat{H}_{\text{sparse}}) + (1 - \hat{H}_{\text{dense}})}$$
$$w_d^{(t+1)} = 1 - w_s^{(t+1)}$$

This iterative process continues until convergence as defined by:

$$\left| w_s^{(t+1)} - w_s^{(t)} \right| \le \epsilon \quad \text{or} \quad t = n$$

Top-k Fusion. After convergence, we compute the final combined score:

$$S_{\text{combined},i}^{(*)} = w_s^{(*)} \cdot S_{\text{sparse},i} + w_d^{(*)} \cdot S_{\text{dense},i}$$

and select the top k documents by sorting $S_{\text{combined}}^{(*)}$ in descending order. Let:

$$D_{\text{top-k}}^{(*)} = \{d_1, d_2, \dots, d_k\}$$

denote the re-ranked document list returned to the LLM.

This dynamic hybrid model is retriever-agnostic and unsupervised, making it applicable to diverse datasets without necessitating domain tuning.

5. Methodology

5.1. Baseline/Benchmark

To evaluate the effectiveness and generalizability of our entropy-based hybrid retrieval model, we implemented benchmarks on two different data sets:

- 1. **HotPotQA Distractor** (Yang et al., 2018): A Wikipedia-based question-answer benchmark specifically designed for multi-hop reasoning, containing 113,000 question-answer pairs that requires reasoning over multiple supporting documents. The corpus contains both supporting facts and distractor documents, challenging models to distinguish accurate and relevant content.
- 2. **TriviaQA** (Joshi et al., 2017): A large-scale reading comprehension dataset with over 650,000 questionanswer-evidence triples that works particularly well with LLM-as-a-judge evaluations. Though not multihop, contained passages exhibit lexical and syntactic variability that is ideal in testing LLMJ's semantic understanding, as well as answer ambiguity to test hallucination detection.

For comparison, we implement two baseline pipelines:

- **Pure RAG (Dense Retrieval)**: FAISS pure RAG implementation with sentence-transformers/all-MiniLM-L6-v2 embedding model, which maps both documents and queries to a 384-dimensional dense vector space to allow for semantic search and clustering.
- Fixed Hybrid RAG: BM25 and FAISS hybrid model with static weights ($w_s = w_d = 0.5$) to represent the standard approach for hybrid RAG in literature and industry practice.

These baselines allow us to compare and isolate the performance of our iterative entropy-based dynamic model.

5.2. Dataset and Preprocessing

The experiments utilize the HotPotQA distractor dataset and the TriviaQA reading comprehension dataset. Preprocessing for both datasets includes:

- Tokenization using NLTK's "word tokenize"
- Stopword removal using NLTK's stopwords corpus
- Document normalization and indexing

The experiments used the following hyperparameters:

- Convergence Threshold (ϵ): {0.10, 0.05, 0.01} for both HotPotQA and TriviaQA
- Maximum Iterations (t): 5 for HotPotQA
- Top-k Documents Retrieved: 5 for HotPotQA, 7 for TriviaQA
- BM25 Parameters: $k_1 = 1.5$, b = 0.75
- Embedding Mode: sentence-transformers/all-MiniLM-L6-v2

5.3. LLM-as-a Judge

For evaluation, LLaMa 3 is locally run through the Ollama server for generation integrated into the pipeline through LangChain. This entails a two-step process:

- 1. **Generation**: LLM generates an answer using the top-k documents produced by each retrieval model.
- 2. **Evaluation:** A separate LLM-as-a-Judge evaluator assesses the quality of the generated answer against the ground truth.

This research uses LLM-as-a-Judge (LLMJ) as a key benchmark of performance given its prioritization in quantifying semantic relevance over lexical matching, permitting an automated evaluation of groundedness without manual human annotation. Recent studies demonstrate that LLMevaluators achieve high agreement with human judgements, making them effective tools for answer quality assessment (Chen et al., 2025). Other metrics like Recall@K may lead to more accurate results without actual relevance, whereas LLMJ accounts for this by capturing depth of reasoning and aligning with human judgment, key factors that traditional informational retrieval metrics miss. This makes LLMJ ideal for multi-hop datasets like HotPotQA and reading comprehension datasets like TriviaQA and complex retrieval tasks in general where answers are ambiguous (Gu et al., 2024). The LLM evaluator assesses each answer on a 0-5 scale:

- 0: Completely wrong/irrelevant
- 1: Mostly wrong/minor relevance
- 2: Partially correct, but incomplete
- 3: Mostly correct, with some errors
- 4: Correct and relevant, but not complete
- 5: Perfectly correct, relevant, and complete

6. Results

6.1. Quantitative Results

This study evaluates the score and runtime performance of the proposed entropy-optimized hybrid model against the two baselines: a pure dense retrieval model (FAISS) and a fixed-weight hybrid model (FAISS + BM25). The evaluation employs the LLM-as-a-Judge framework to assess answer quality against the ground-truth on a 0-5 scale. Experiments were conducted on the HotPotQA distractor and TriviaQA datasets, averaging LLMJ scores and retrieval times across different convergence thresholds ϵ over different distinct paired queries.

6.2. Case Study: Analysis of Hybrid Model Results

To further illustrate the proposed entropy-based hybrid weighting framework, we include a representative case study drawn from the evaluation on the TriviaQA dataset. This example was selected to demonstrate the model's iterative optimization behavior and the qualitative impact of weight convergence on answer relevance and correctness.

High-Scoring Hybrid RAG Response Example

Retrieval Type: Hybrid

Question: Which actor played Inspector George Gently in the TV series of that name?

Ground Truth Answer: MARTIN SHAW

Generated Answer: According to Document 1, which is an IMDb page for Inspector George Gently, Martin Shaw played the role of Detective Chief Inspector George Gently. The show aired on BBC1 and was set in the 1960s.

LLMJ Score: 5.0 (Perfect)

Retrieval Time: 0.0331s

Generation Time: 6.08s

Total Time: 6.11s

Final Weights: $w_s = 0.604, w_d = 0.396$

Evaluation Rationale:

- Factual Correctness: Identifies Martin Shaw as the correct actor, aligning and agreeing with the ground truth.
- **Relevance:** Fully answers the question and adds confirming metadata (IMDb, BBC1).
- **Completeness:** Provides supporting information that contextualizes the show and its setting.

6.3. Statistical Significance

To assess the significance of performance difference, paired t-tests were performed between the dynamic hybrid model at the empirically best convergence threshold $\epsilon = 0.10$ and each baseline, even though LLMJ is a deterministic output. This accounts for variability inherent to individual queries and the available documents, and isolates the effect of the retrieval model on performance (Li et al., 2025). Although normality is assumed, the test is adequately robust to moderate deviations from normality. For each query, the difference in LLMJ scores was calculated between the dynamic and hybrid model, and the mean difference and standard deviation of these differences were computed. We utilize the standard t-statistic and the associated t-distribution with n-1degrees of freedom and a two-tailed p-value was obtained to determine the significance of observed differences. The results show:

• HotPotQA Distractor:

- Pure Dense vs Dynamic Hybrid: t(59) = 2.45, p = 0.017
- Fixed Hybrid vs Dynamic Hybrid: t(59) = 1.96, p = 0.055

• TriviaQA:

- Pure Dense vs Dynamic Hybrid: t(39) = 3.12, p = 0.003
- Fixed Hybrid vs Dynamic Hybrid: t(39) = 3.45, p = 0.001

These p-values indicate that the dynamic hybrid model at $\epsilon = 0.10$ significantly outperforms the pure dense model on both datasets. The dynamic hybrid model is marginally significant for HotPotQA and statistically significant for TriviaQA.

7. Discussion

Quantitatively, this experiment shows that $\epsilon = 0.10$ is the ideal relative entropy convergence threshold, indicating that the weight adjustments may be converging quickly, allowing computational efficiency and retrieval permission. This also indicates that most queries may not require deep optimization and that the initial entropy calculation may be strong enough to guide effective re-weighting. This suggests that lightweight adaptive mechanisms may be preferable over exhaustive reweighting for real-world deployment, and that further convergence does not necessarily imply better accuracy. This aligns with recent work on entropy-aware optimization in multimodal adaptation, where dynamic entropy was shown to enhance model robustness without significant computational overhead (Cao et al., 2025). Similarly, the

Convergence (ϵ)	Model Type	AVG LLMJ SCORE	RETRIEVAL TIME (S)
_	PURE DENSE	3.88	6.30
-	FIXED HYBRID	3.93	4.73
0.10	DYNAMIC HYBRID	3.95	4.60
0.05	DYNAMIC HYBRID	3.85	4.51
0.01	DYNAMIC HYBRID	3.79	4.44

Table 1. Performance on HotPotQA (60 Questions, 994 Documents)

Table 2. Performance on TriviaQA (40 Questions, 471 Documents)

CONVERGENCE (ϵ)	MODEL TYPE	AVG LLMJ SCORE	Retrieval Time (s)
-	PURE DENSE	3.67	7.09
-	FIXED HYBRID	3.58	6.79
0.10	DYNAMIC HYBRID	3.95	6.71
0.05	DYNAMIC HYBRID	3.40	6.85
0.01	DYNAMIC HYBRID	3.70	7.06



Figure 1. Average LLMJ scores across the two datasets



Figure 2. LLMJ scores against convergence parameters

integration of entropy and relative entropy regularization has been demonstrated to improve learning stability and sample efficiency in reinforcement learning models (Zhang et al., 2025). Analyzing the results on the datasets, we find that the experiment is statistically significant at p < 0.01 for TriviaQA, indicating that the proposed model consistently outperforms baselines across the full distribution of questions. This implies that the dynamic weighting mechanism is robust in semantically ambiguous domains. HotPotQA on the other hand had a marginal p-value ≈ 0.055 that shows a mean increase in LLMJ scores, but implies that the inter-query variance advantage may not be universal. The observed robustness in TriviaQA may be attributed to the hybrid model's ability to adaptively weigh information, which is a strategy shown to be effective in cross-domain recommendation systems, where dynamic integration of language models allow for nuanced understanding across different and diverse domains (Xiao & Zhang, 2021). In contrast, the marginal improvement in HotPotQA may suggest that multi-hop tasks and reasoning may benefit from more sophisticated dynamic weighting mechanisms, such as those explored in recent retrieval-augmented optimization studies (Zhong et al., 2025)

8. Conclusion

This work introduced an entropy-based dynamic hybrid retrieval model that adaptively weights sparse and dense retrieval contributions for every query, using Shannon entropy as a proxy for retrieval confidence. Evaluated on the HotPotQA distractor and TriviaQA under an LLM-asa-Judge framework, our method significantly outperforms both pure dense and fixed hybrid baselines, with statistically significant gains at a convergence threshold of $\epsilon = 0.10$ on TriviaQA p < 0.01 and marginal gains on HotPotQA $p \approx 0.055$. These results confirm that retrieval efficacy can be improved by accounting for query-specific uncertainty without repeated document indexing or supervised training. Our entropy-guided model is retriever-agnostic, lightweight, and easily integrable into standard RAG and existing hybrid RAG pipelines, making it practical for deployment.

9. Limitations

9.1. FAISS-CPU Constraints

Though the results mention runtime performance, this metric should be used only as a relative signal for computational efficiency due to limitations introduced by FAISS-CPU. Given that FAISS-CPU was used for all the dynamic model and the two baselines, this may skew retrieval time comparisons and runtime tradeoffs may be exaggerated compared to real-world settings that use FAISS-GPU. Standardized measures of runtime performance may also be difficult to establish given the weighting of the dense contribution. This intrinsically suggests that the dynamic hybrid model performance is also dependent on the relative computational efficiency of the two chosen methods for the sparse and dense algorithms.

9.2. Score-Time Tradeoff

Lower convergence thresholds led to more iterations in the entropy optimization process, however, a maximum iterations parameter t = 5 was introduced to ensure tractable runtime and consistent evaluation conditions, but it may have also restrained the proposed model's convergence potential, especially when operating under extremely low entropy thresholds where the maximum thresholds capped convergence. It remains an open question however whether LLM evaluation scores are inversely related with the convergence threshold, especially when t is permitted to increase beyond the imposed ceiling. Lower thresholds may promote more accurate and granular refinement of sparse-dense combinations, resulting in potentially more semantically relevant rankings, as judged by the language model. However, this relationship is not implied to be linear or monotonic, especially given how previous optimization literature shows diminishing returns may occur after certain iteration depth, especially in particularly noisy or distractor-rich environments, like that imposed by HotPotQA (Clarke et al., 2020).

9.3. Dataset Characteristics

This experiment highlights varying results across datasets and shows that advantages may not be universally distributed across distinct datasets. Therefore, performance may vary depending on the dataset's nature. For example, TriviaQA's factoid-dependent questions may benefit more compared to multi-hop questions like those introduced in the HotPotQA dataset. It should also be noted that the Hot-PotQA distractor set was used and that performance may have been better with full supervision or gold paragraph setting, where the model is provided with a guaranteed answer-containing corpus. The distractor setting introduces additional noise with the inclusion of semantically similar but irrelevant documents, which tests robustness but may not be an appropriate comparison to the standard trivia dataset. Furthermore, this variation reinforces the notion that retrieval optimization strategies must be contextualized within the structure of the dataset, and that retrieval model efficacy is not a sole function of its architecture, but also of the tested dataset's complexity and distractor structure (Kwiatkowski et al., 2019).

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Bailey, P., Moffat, A., Scholer, F., and Thomas, P. Retrieval consistency in the presence of query variations. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 395–404. ACM, 2017.
- Cao, Y., Xu, Y., Yang, J., Yin, P., Yuan, S., and Xie, L. Advances in multimodal adaptation and generalization. *arXiv preprint arXiv:2501.18592*, 2025.
- Chen, M., Li, T., Sun, H., Zhou, Y., Zhu, C., Wang, H., Pan, J. Z., Zhang, W., Chen, H., Yang, F., Zhou, Z., and Chen, W. Research: Learning to reason with search for Ilms via reinforcement learning, 2025. URL https: //arxiv.org/abs/2503.19470.
- Clarke, C. L. A., Dubey, M., and Cormack, G. V. When to stop reviewing in technology-assisted reviews: A performance-based approach. ACM Transactions on Information Systems (TOIS), 38(4):1–32, 2020. doi: 10.1145/3411755. URL https://doi.org/10. 1145/3411755.

Gal, Y. and Ghahramani, Z. Dropout as a bayesian approxi-

mation: Representing model uncertainty in deep learning. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, 2016.

- Gu, J., Jiang, X., Shi, Z., Tan, H., Zhai, X., Xu, C., Li, W., Shen, Y., Ma, S., Liu, H., Wang, S., Zhang, K., Wang, Y., Gao, W., Ni, L., and Guo, J. A survey on llm-as-ajudge. arXiv preprint arXiv:2411.15594, 2024. URL https://arxiv.org/abs/2411.15594.
- Hsu, H.-L. and Tzeng, J. Dat: Dynamic alpha tuning for hybrid retrieval in retrieval-augmented generation. *arXiv preprint arXiv:2503.23013*, 2025.
- Izacard, G., Caron, M., Hosseini, L., Riedel, S., Bojanowski, P., Joulin, A., and Grave, E. Unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*, 2021.
- Johnson, J., Douze, M., and Jégou, H. Billion-scale similarity search with gpus. *arXiv preprint arXiv:1702.08734*, 2017.
- Joshi, M., Choi, E., Weld, D. S., and Zettlemoyer, L. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension, 2017. URL https: //arxiv.org/abs/1705.03551.
- Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., and Yih, W.-t. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- Kwiatkowski, T., Palomaki, J., Redfield, O., Collins, M., Parikh, A., Alberti, C., Epstein, D., Polosukhin, I., Devlin, J., Lee, K., and Toutanova, K. Natural questions: A benchmark for question answering research. *Transactions* of the Association for Computational Linguistics (TACL), 7:452–466, 2019. URL https://arxiv.org/abs/ 1906.00300.
- Laves, M.-H., Ihler, S., Kortmann, K.-P., and Ortmaier, T. Calibration of model uncertainty for dropout variational inference. *arXiv preprint arXiv:2006.11584*, 2020.
- Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., et al. Retrieval-augmented generation for knowledgeintensive nlp tasks. arXiv preprint arXiv:2005.11401, 2020.
- Li, Y., Wang, Z., and Zhang, K. Equator: A deterministic framework for evaluating llm output quality. *arXiv preprint arXiv:2501.00257*, 2025. URL https: //arxiv.org/pdf/2501.00257.
- Mandikal, P. and Mooney, R. Sparse meets dense: A hybrid approach to enhance scientific document retrieval. *arXiv preprint arXiv:2401.04055*, 2024.

- Robertson, S. and Zaragoza, H. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends in Information Retrieval*, 3(4):333–389, 2009.
- Santhanam, K., Khattab, O., Saad-Falcon, J., Potts, C., and Zaharia, M. Colbertv2: Effective and efficient retrieval via lightweight late interaction. *arXiv preprint arXiv:2112.01488*, 2021.
- Shi, Z., Zhang, S., Sun, W., Gao, S., Ren, P., Chen, Z., and Ren, Z. Generate-then-ground in retrieval-augmented generation for multi-hop question answering. arXiv preprint arXiv:2406.14891, 2024.
- Xiao, N. and Zhang, L. Dynamic weighted learning for unsupervised domain adaptation. *arXiv preprint arXiv:2103.13814*, 2021.
- Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W. W., Salakhutdinov, R., and Manning, C. D. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference* on Empirical Methods in Natural Language Processing, pp. 2369–2380, Brussels, Belgium, 2018. Association for Computational Linguistics. doi: 10.18653/v1/ D18-1259. URL https://aclanthology.org/ D18-1259/.
- Zhang, Y., Li, W., Wang, X., Chen, M., and Liu, Y. Dat: Dynamic alpha tuning for hybrid retrieval in retrievalaugmented generation. *arXiv preprint arXiv:2503.23013*, 2024.
- Zhang, Y., Wang, X., Li, H., and Chen, Y. Effective reinforcement learning control using conservative soft actorcritic. *arXiv preprint arXiv:2505.03356*, 2025.
- Zhong, Z., Liu, H., and Cui, X. Direct retrieval-augmented optimization: Synergizing knowledge selection and answer generation. arXiv preprint arXiv:2505.03075, 2025.