

Incorporating Q&A Nuggets into Retrieval-Augmented Generation

Laura Dietz¹✉, Bryan Li², Gabrielle Liu³, Jia-Huei Ju⁴, Eugene Yang⁵,
Dawn Lawrie⁵, William Walden⁵, and James Mayfield⁵

¹ University of New Hampshire, Durham, New Hampshire, USA dietz@cs.unh.edu

² University of Pennsylvania, Philadelphia, Pennsylvania, USA

bryanli@seas.upenn.edu

³ Yale University, New Haven, Connecticut, USA kaili.liu@yale.edu

⁴ University of Amsterdam, Amsterdam, Netherlands j.ju@uva.nl

⁵ Human Language Technology Center of Excellence, Johns Hopkins University,
Baltimore, Maryland, USA

{eugene.yang,lawrie,wwalden1,mayfield}@jhu.edu

Abstract. RAGE systems integrate ideas from automatic evaluation (E) into Retrieval-augmented Generation (RAG). As one such example, we present CRUCIBLE, a Nugget-Augmented Generation System that preserves explicit citation provenance by constructing a bank of Q&A nuggets from retrieved documents and uses them to guide extraction, selection, and report generation. Reasoning on nuggets avoids repeated information through clear and interpretable Q&A semantics—instead of opaque cluster abstractions—while maintaining citation provenance throughout the entire generation process. Evaluated on the TREC NeuCLIR 2024 collection, our CRUCIBLE system substantially outperforms GINGER, a recent nugget-based RAG system, in nugget recall, density, and citation grounding.⁶

Keywords: RAG · LLM judge · nugget-based evaluation.

1 Introduction

Retrieval-Augmented Generation (RAG) has become the dominant framework for grounding LLM outputs in evidence [13, 16]. At the same time, nugget-based evaluation methods have emerged as the standard for measuring relevance of long-form responses [10, 18]. A *nugget*, i.e., short Q&A pair, fact, or claim, is a fine-grained reusable unit for assessing whether key information is covered. By encoding the information that must appear in a system answer, nuggets enable evaluation metrics such as “nugget recall” that directly quantify the amount of useful information given by the information system. We argue that nuggets are valuable not only for evaluation but also for guiding retrieval and generation, especially given that LLMs can produce high-quality nuggets automatically [4, 9].

⁶ Appendix at <https://github.com/hltcoe/ecir26-crucible-system-appendix/>

Contributions. To support this argument, we present CRUCIBLE, a nugget-centric RAG system that automatically constructs its own nugget bank and uses it as a control signal throughout the pipeline. For example, nuggets form the basis for controlling redundancy without content-based clustering, while naturally preserving citation provenance. On the TREC NeuCLIR 2024 test collection, we provide a direct comparison to a recent nugget-oriented RAG system GINGER [14]. We find that CRUCIBLE decisively outperforms GINGER across several evaluation metrics, from nugget recall to citation support.

2 Related Work

Summarize one-document-at-a-time. Early RAG systems either produce a summary per retrieved passage [16] or encode each passage separately [13]. This retains exact citation provenance, but the model must attend to many independent inputs, which often harms coherence and fluency.

Joint representation. Later work uses a single joint representation to represent the retrieved document set. REALM [12] and EMDR2 [24] learn a unified encoding; xRAG [3] pushes this to the extreme by compressing all content into a single token representation. These approaches improve synthesis and fluency, yet the compact latent representation discards explicit links to source documents, making citation grounding difficult and the pipeline opaque to developers. A recommended remedy is to first hallucinate an answer and then retrieve supporting evidence, as in HyDE [11], which raises lingering concerns about trustworthiness.

Clustering-based RAG. GINGER [14] is a nugget-based RAG system that defines nuggets as verbatim text spans copied directly from retrieved passages. These spans are first clustered into topical facets using BERTopic, after which the clusters are reranked to identify the most relevant facets. Summaries of the top clusters form the system’s output, which is then rewritten for fluency.

Although the pipeline has access to evidence for generation, the clustering summarization step abstracts away from the original document extractions, impacting the faithful citation grounding.

Agentic RAG. Agentic frameworks decompose the pipeline into subtasks and let a planner invoke them as needed [1, 27]. WEBGPT [19] embeds search instructions in the prompt, letting the LLM choose passages to cite. A prominent agentic system is GPTRESEARCHER [7, 8], which orchestrates query generation, retrieval, extraction, and report writing, with optional verification and trace-back.

3 Nugget-first RAG Approach: CRUCIBLE

Figure 1 illustrates the workflow of CRUCIBLE, which builds its pipeline around structured Q&A nuggets that guide retrieval, extraction, and assembly. This contrasts with GINGER, which operates on verbatim spans clustered into facets.

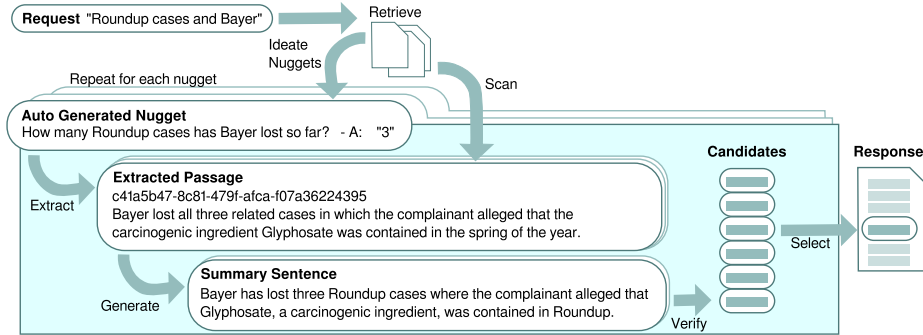


Fig. 1. For each generated nugget, CRUCIBLE extracts candidate sentences and adds the best k sentences to the final response. No content clustering is needed.

Nugget ideation. We begin by retrieving an initial pool of documents; in our study we use the PLAID-X dense retriever [28]. For each document we generate a short query-focused summary with an LLM, then prompt the same LLM (using LLaMA-3.3-70B-Instruct) to produce Q&A nuggets conditioned on the summary and the user request.

To reduce redundancy among nuggets, we first detect paraphrases with an LLM prompt, then successively merge most confident paraphrases until the desired number of canonical nuggets is obtained.

Nugget ranking. The resulting nugget bank is reranked with a Support-Vector Classifier (SVC) trained on 19 quality features, which utilize LLM judge prompts to measure how well each nugget addresses the information need (task statement, background, role, communication style, and scope), degree of vitality [22], and several features from researchy questions [23]. This is combined with basic readability metrics to quantify reading level [2] and sentence complexity [17]. Using the number of paraphrases as an indicator for nugget quality, we fuse this feature-based ranking with a popularity-based ranking. Obtaining the top 20 nuggets for each request will form the nugget bank that drives the remainder of the response generation.

Retrieval. Our base system uses documents of which nuggets were extracted. Additionally, the pipeline is evaluated with the top 100 of several retrieval models with their default configurations: PLAID-X [28] a dense multilingual retriever, using top 25 per language, dense retrieval with Qwen3 [30], and Milco [20].

Scanning and generation. Using the nugget bank, we scan retrieved documents⁷ for passages that directly answer each nugget. Using the prompt in Fig. 2⁸ we (1) locate a supporting passage, (2) extract a concise self-contained sentence, and (3) record the LLM’s token-likelihood as the extraction confidence.

⁷ Segmented into 1000 character chunks, split at sentence boundary.

⁸ The system message was omitted due to a bug. Updated results in online appendix.

Given the following question and answer, please find the sections in the provided source document that support and validate the answer to the question.
 If the source document does not support the given answer, then set confidence to 0.0 and all fields to None.
 Provide the section from the document that supports the answer, with complete sentences directly from the document. Please include context around each supporting segment, making sure that there is enough context to support why the answer is a correct response to the question. Your response should include just the extracted text segment, and nothing else. Then condense the extracted text into one concise sentence that clearly demonstrates how the question is answered, without referring to the source document.

- Emit **each field exactly once**, in the order shown.
- Do NOT repeat headers. Do NOT include any field more than once.
- If a field is unknown, omit it entirely - do NOT write 'None' or 'null'.

Field name	Description
In: nugget.text	Input question
In: answer	Valid answer or list of answers
In: source.document	Input document chunk
In: title	Title query
In: background	User background
In: problem.statement	User's problem statement
Out: extracted.text.segment	Passage returned by the model
Out: summary	Single concise sentence drawn from the passage
Out: reasoning	Internal chain of thought (ignored downstream)
Out: confidence	Model reported confidence value (float)

Fig. 2. Prompt for scanning, extraction, and sentence candidate generation.

Verification (optional). Extractions may be double-checked for nugget coverage and citation support. This step invokes prompts used by the AUTOARGUE system [26] albeit with the automatically-generated nuggets identified after the Nugget Ranking step; the LLM is prompted with a binary “YES/NO” question asking whether a candidate sentence is truly supported by the cited span or whether the nugget is truly covered by the candidate sentence and extracted passage. As this step risks leaking evaluator knowledge into the system, we also evaluate the system with this step skipped as “CRUCIBLE-BASE”.

Sentence selection. We rank the remaining candidates for each nugget by the extraction confidence and choose the top k sentences ($k = 1$ herein; results for larger k are in the online appendix). This ensures that each sentence is tied to exactly one citation.

Assembly. Selected sentences are concatenated into a report, ordered by nugget quality ranking. Sentences with same stopped/stemmed text are omitted. Because every sentence is self-contained and atomic, the order of sentences does not affect the readability. Every sentence cites exactly one document.

4 Evaluation Setup

Dataset. We evaluate on the TREC NeuCLIR 2024 Report Generation Pilot [15], which requires drafting topic-focused reports supported by citations drawn from a multilingual corpus, where all language subsets are machine-translated to En-

glish. The topics are broad, e.g., “Roundup cases and Bayer”, yet each report must be tailored to the specific problem statement and user background. The collection’s evaluation tool, AUTOARGUE [26], uses held-out, manually curated gold-nugget banks.

Compared systems. Using **LLaMA-3.3-70B-Instruct** and PlaidX retriever, we compare two variants of CRUCIBLE with state-of-the-art systems:

CRUCIBLE-BASE (**ours**), as described in Section 3 without the verification step. CRUCIBLE-VERIFIED (**ours**), full pipeline with check for coverage/support.

GINGER [14], a nugget-informed RAG system that generates responses based on clusters built on the extracted nuggets from the retrieved documents. While the released implementation does not produce citations in its response, we assign citations to the document that contains the nuggets closest to the cluster centroid. This is an alternative system that is based on ideas similar to those in CRUCIBLE. We report the best variant, where initial set of documents is retrieved by Qwen3-Embedding bi-encoder [30], underlying LLM is GPT-4o.

GINGER-LLAMA: Using the GINGER implementation with the same LLaMA LLM to allow a fair comparison to CRUCIBLE.

GPTRESEARCHER [7, 8] a simple agentic system with subqueries, retrieval, and writing. Uses LLaMA as LLM (despite the name).

BULLETPOINTS [29], an extractive pipeline (best in TREC NeuCLIR 2024, submitted as **hltcoe-eugene**). Uses LLaMA as LLM.

Evaluation metrics. We report metrics from the AUTOARGUE framework [26].

Nugget Recall = $\frac{\text{covered nuggets}}{\text{gold nuggets}}$ measures the recall of gold nuggets.

Nugget Density = $\frac{\text{covered nuggets}}{\text{sentences}}$ balance of coverage vs. concise summaries.

Sentence Novelty = $\frac{\text{sentences mention a new nugget}}{\text{all sentences}}$, fraction of sentences that introduce new relevant information, according to gold nuggets.

Relevant Sentences = $\frac{\text{sentence with nuggets}}{\text{all sentences}}$ measures how many sentences mention relevant nuggets.

Citation Support = $\frac{\text{supported citations}}{\text{citations}}$, measures how many citations actually support their summary sentence.

5 Results

Table 1 reports results averaged across NeuCLIR topics. CRUCIBLE consistently outperforms GINGER and GINGER-LLaMA on all nugget-oriented metrics. Relative gains are especially pronounced for *nugget recall* (+42 to +65%) and *nugget density* (+21 to +25%).

Table 1. Comparison of RAG systems (each with PlaidX retriever) on TREC NeuCLIR 2024. Best results in bold, paired t-test \blacktriangle/∇ with reference BASE.

System	Nugget Recall	Nugget Density	Sentence Novelty	Relevant Sent.	Citation Support
CRUCIBLE-BASE	0.429	0.448	0.255	0.703	0.902
CRUCIBLE-VERIFIED	0.438	0.457	0.267	\blacktriangle0.733	\blacktriangle0.961
GPTRESEARCHER	∇ 0.177	∇ 0.131	∇ 0.083	∇ 0.265	∇ 0.571
GINGER	∇ 0.244	∇ 0.264	∇ 0.162	∇ 0.285	∇ 0.436
GINGER-LLaMA	∇ 0.241	∇ 0.134	∇ 0.097	∇ 0.136	∇ 0.476
BULLETPOINTS	0.508	∇ 0.340	0.243	∇ 0.468	0.835

Table 2. Downstream effects of the document retrieval method and optional verification; using the same nugget bank across variations; reference CRUCIBLE.

	Retrieval	Nugget Recall	Nugget Density	Sentence Novelty	Relevant Sent.	Citation Support
CRUCIBLE-BASE	From nuggets	0.429	0.255	0.448	0.703	0.902
	Milco	0.467	\blacktriangle 0.300	\blacktriangle 0.501	0.708	∇ 0.839
	Qwen3	0.446	\blacktriangle 0.292	\blacktriangle0.491	\blacktriangle0.758	0.892
	PlaidX	0.424	\blacktriangle0.324	0.468	0.728	∇ 0.807
CRUCIBLE-VERIFIED	From nuggets	0.438	0.267	0.457	0.733	0.961
	Milco	0.464	\blacktriangle0.337	\blacktriangle0.511	0.706	∇ 0.931
	Qwen3	0.434	\blacktriangle 0.308	0.481	0.751	0.959
	PlaidX	0.412	\blacktriangle 0.305	0.463	0.695	∇ 0.921

Overall performance. We attribute CRUCIBLE’s gains to its nugget-first design: explicit ideation yields systematic coverage; per-nugget extraction enforces grounding; and fingerprint duplicate checks preserve density. In contrast, GINGER’s cluster-based summarization risks citation traceability, while not providing the required information. We remark that GINGER was designed for the TREC RAG 24 [25] task and the Autonuggetizer [21] for evaluation.

Ranking model. CRUCIBLE obtains its great performance under all explored retrieval models. Retrievers that offer broad coverage and emphasize recall, such as Milco and Qwen3 consistently yield improvements on nugget metrics, but may hurt citation support. (These other retrievers do not change results on GINGER.)

Verification step. Using prompts to verify the nugget coverage separately from sentence extraction obtains small improvements—especially for the sentence novelty metric. However, independent verification of citation support leads to a 60% error reduction. After manual inspection, we believe that it is not that original citations were incorrect, but during verification we prefer sentences where the support is clearly articulated in the extracted document passage, which is ultimately helpful for the user to trust the system.

More experiments. CRUCIBLE runs were submitted to several TREC 2025 tracks: DRAGUN, RAG, RAGTIME [5]. While many relevance judgments are not available at the time of writing, we have some preliminary corroborating evidence.

Costs. This system was designed to study possible quality improvements, not to obtain a fast system ready for deployment. The runtime cost is dominated by LLM calls. One of the most expensive steps is the paraphrasing detection of the nugget ideation with $\Theta(D^2)$ prompts, which could be scaled with fingerprinting and SimHash approaches. Other expensive steps are nugget scanning with $\Theta(N \cdot D)$ LLM calls and verification with $\Theta(S) \approx \Theta(N \cdot D)$ LLM prompts, these could be made more efficient by tracking the quality of extractions, and stopping when k sentences of sufficiently high quality have been found.

Limitations. CRUCIBLE is developed working closely with organizers of TREC NeuCLIR and developers of the AUTOARGUE evaluation system. As shown in Dietz et al. [6], this knowledge can lead to a distortion of empirical evaluation results. Since, the verification phase, represents a circularity with AUTOARGUE, we offer results with and without this phase. We further study this potential vulnerability in the upcoming TREC 2026 Auto-Judge track.

6 Conclusion

We demonstrate one example of how RAGE systems, which incorporate ideas from automatic evaluation (E) into RAG systems, can lead to significantly better performing systems. To substantiate this claim we introduce CRUCIBLE, a nugget-first pipeline that uses a bank of Q&A nuggets to guide retrieval, extraction, and final assembly. Optionally, CRUCIBLE uses LLM-judge prompts as a verification steps. On the NeuCLIR benchmark CRUCIBLE decisively outperforms GINGER, demonstrating the practical benefits controlling the generation.

The key insight behind CRUCIBLE is that citation provenance can be lost whenever a system condenses raw text into latent or clustered summaries. Most contemporary RAG approaches either (1) struggle to avoid redundancy (e.g. FiD [13]) or (2) sever the explicit link between a generated sentence and its source document. CRUCIBLE avoids both pitfalls by first constructing a set of canonical Q&A nuggets (any required clustering occurs at this stage) and only then extracting and summarizing sentences that address each nugget. The generated report will contain the best sentence for each nugget, each sentence is supported with exactly one citation. Obtaining near-perfect citation support score of 0.9 (CRUCIBLE-BASE) and 0.96 (CRUCIBLE-VERIFIED) supports that either this solves the citation-support problem, or that we need more precise evaluation methods for citation support.

We hope that future work will explore a greater variety of RAGE systems across the gamut of RAG and LLM-as-a-Judge paradigms.

Disclosure of Interests. Authors have no competing interests.

Bibliography

- [1] Asai, A., Wu, Z., Wang, Y., Sil, A., Hajishirzi, H.: Self-rag: Learning to retrieve, generate, and critique through self-reflection. In: Proceedings of the International Conference on Learning Representations (ICLR) (2024)
- [2] Bansal, S.: textstat (Jul 2025), URL <https://github.com/textstat/textstat>
- [3] Cheng, X., Wang, X., Zhang, X., Ge, T., Chen, S.Q., Wei, F., Zhang, H., Zhao, D.: xrag: Extreme context compression for retrieval-augmented generation with one token. In: Advances in Neural Information Processing Systems (NeurIPS) (2024)
- [4] Dietz, L.: A workbench for autograding retrieve/generate systems. In: Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 1963–1972 (2024)
- [5] Dietz, L., Li, B., Mayfield, J., Lawrie, D., Yang, E., Walden, W.: [hltcoe] hltcoe evaluation team at trec 2025: Rag, ragtime, and dragun. In: The Thirty-Fourth Text REtrieval Conference Proceedings (TREC2025) (2025)
- [6] Dietz, L., Li, B., Yang, E., Lawrie, D., Walden, W., Mayfield, J.: Insider knowledge: How much can rag systems gain from evaluation secrets? In: Proceedings of the 48th European Conference on Information Retrieval (ECIR 2026) (2026)
- [7] Duh, K., Yang, E., Weller, O., Yates, A., Lawrie, D.: HLTCOE at LiveRAG: GPT-Researcher using ColBERT retrieval. arXiv preprint arXiv:2506.22356 (2025)
- [8] Elovic, A.: gpt-researcher (Jul 2023), URL <https://github.com/assafelovic/gpt-researcher>
- [9] Farzi, N., Dietz, L.: Exam++: Llm-based answerability metrics for ir evaluation. In: Proceedings of LLM4Eval: The First Workshop on Large Language Models for Evaluation in Information Retrieval (2024)
- [10] Farzi, N., Dietz, L.: Pencils down! automatic rubric-based evaluation of retrieve/generate systems. In: Proceedings of the 2024 ACM SIGIR International Conference on Theory of Information Retrieval, pp. 175–184 (2024)
- [11] Gao, L., Ma, X., Lin, J., Callan, J.: 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), pp. 1762–1777 (2023)
- [12] Guu, K., Lee, K., Tung, Z., Pasupat, P., Chang, M.W.: Retrieval augmented language model pre-training. In: Proceedings of the 37th International Conference on Machine Learning (ICML), pp. 3929–3938, PMLR (2020)
- [13] Izacard, G., Grave, É.: Leveraging passage retrieval with generative models for open domain question answering. In: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pp. 874–880 (2021)

- [14] Lajewska, W., Balog, K.: Ginger: Grounded information nugget-based generation of responses. In: Proceedings of the 48th International ACM SIGIR Conference (SIGIR '25) (2025), URL <https://krisztianbalog.com/files/sigir2025-ginger.pdf>, sIGIR 2025 paper
- [15] Lawrie, D., MacAvaney, S., Mayfield, J., McNamee, P., Oard, D.W., Soldaini, L., Yang, E.: Overview of the TREC 2024 neuclir track. arXiv preprint arXiv:2509.14355 (2025)
- [16] 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 knowledge-intensive NLP tasks. *Advances in neural information processing systems* **33**, 9459–9474 (2020)
- [17] Liu, W., Zeng, W., He, K., Jiang, Y., He, J.: What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. arXiv preprint arXiv:2312.15685 (2023)
- [18] Mayfield, J., Yang, E., Lawrie, D., MacAvaney, S., McNamee, P., Oard, D.W., Soldaini, L., Soboroff, I., Weller, O., Kayi, E., Sanders, K., Mason, M., Hibbler, N.: On the evaluation of machine-generated reports. In: Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 1904–1915 (2024)
- [19] Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., Jiang, X., Cobbe, K., Eloundou, T., Krueger, G., Button, K., Knight, M., Chess, B., Schulman, J.: Webgpt: Browser-assisted question-answering with human feedback. arXiv preprint arXiv:2112.09332 (2021)
- [20] Nguyen, T., Lei, Y., Ju, J.H., Yang, E., Yates, A.: Milco: Learned sparse retrieval across languages via a multilingual connector. arXiv [cs.IR] (2025)
- [21] Pradeep, R., Thakur, N., Upadhyay, S., Campos, D., Craswell, N., Lin, J.: Initial nugget evaluation results for the TREC 2024 rag track with the autonuggetizer framework (2024), URL <https://arxiv.org/abs/2411.09607>, arXiv preprint
- [22] Pradeep, R., Thakur, N., Upadhyay, S., Campos, D., Craswell, N., Soboroff, I., Dang, H.T., Lin, J.: The great nugget recall: Automating fact extraction and rag evaluation with large language models. In: Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, p. 180–190, SIGIR '25, Association for Computing Machinery, New York, NY, USA (2025), ISBN 9798400715921, <https://doi.org/10.1145/3726302.3730090>, URL <https://doi.org/10.1145/3726302.3730090>
- [23] Rosset, C., Chung, H.L., Qin, G., Chau, E.C., Feng, Z., Awadallah, A., Neville, J., Rao, N.: Researchy questions: A dataset of multi-perspective, decompositional questions for llm web agents. arXiv preprint arXiv:2402.17896 (2024)
- [24] Sachan, D.S., Reddy, S., Hamilton, W.L., Dyer, C., Yogatama, D.: End-to-end training of multi-document reader and retriever for open-domain question answering. In: *Advances in Neural Information Processing Systems* (NeurIPS) 2021, pp. 25968–25981 (2021)

- [25] Upadhyay, S., Pradeep, R., Thakur, N., Campos, D., Craswell, N., Soboroff, I., Dang, H.T., Lin, J.: A large-scale study of relevance assessments with large language models: An initial look. arXiv preprint arXiv:2411.08275 (2024)
- [26] Walden, W., Weller, O., Dietz, L., Li, B., Liu, G.K.M., Hou, Y., Yang, E.: Auto-ARGUE: LLM-based report generation evaluation. arXiv preprint arXiv:2509.26184 (2025)
- [27] Yang, D., Rao, J., Chen, K., Guo, X., Zhang, Y., Yang, J., Zhang, Y.: Im-rag: Multi-round retrieval-augmented generation through learning inner monologues. In: Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 730–740 (2024)
- [28] Yang, E., Lawrie, D., Mayfield, J., Oard, D.W., Miller, S.: Translate-distill: learning cross-language dense retrieval by translation and distillation. In: European Conference on Information Retrieval, pp. 50–65, Springer (2024)
- [29] Yang, E., Lawrie, D., Weller, O., Mayfield, J.: HLTCOE at TREC 2024 NeuCLIR track (2025), URL <https://arxiv.org/abs/2510.00143>
- [30] Zhang, Y., Li, M., Long, D., Zhang, X., Lin, H., Yang, B., Xie, P., Yang, A., Liu, D., Lin, J., et al.: Qwen3 embedding: Advancing text embedding and reranking through foundation models. arXiv preprint arXiv:2506.05176 (2025)