# Ragnarök: A Reusable RAG Framework and Baselines for TREC 2024 Retrieval-Augmented Generation Track

Ronak Pradeep*, Nandan Thakur*, Sahel Sharifymoghaddam, Eric Zhang, Ryan Nguyen, Daniel Campos, Nick Craswell, Jimmy Lin

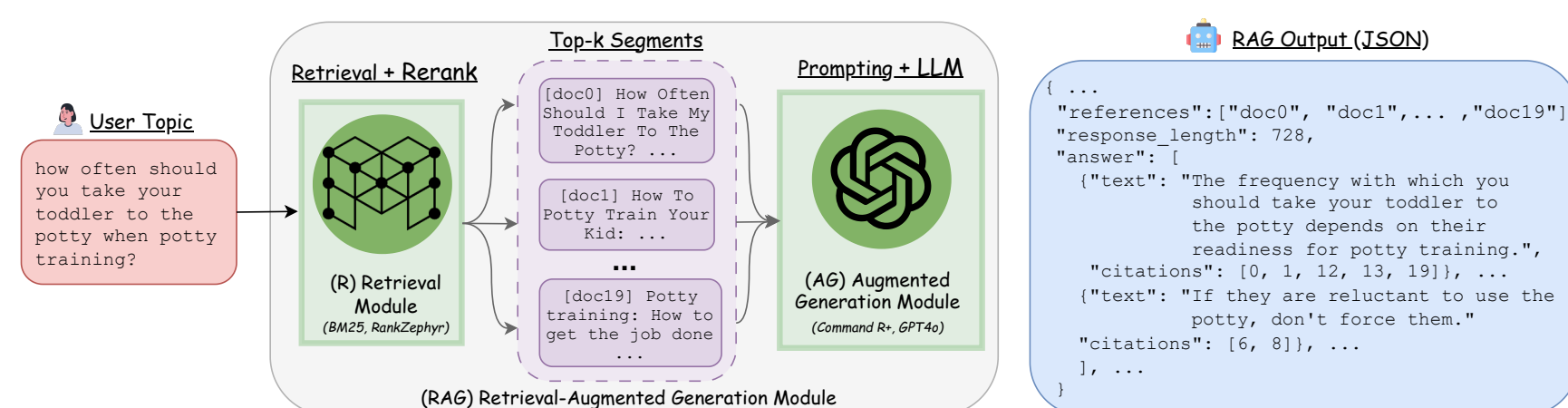University of Waterloo    Snowflake Inc.    Microsoft

## The Challenge: Standardizing RAG Evaluation

Retrieval-Augmented Generation (RAG) powers modern search (Bing, ChatGPT, Google AI Overviews), but lacks standardized frameworks for research and evaluation. The RAG community needs tools!

We introduce **Ragnarök**: an open-source RAG inference framework built with reproducibility in mind. Ragnarök has successfully been used as a SOTA pipeline in several TREC tracks — RAG, NeuCLIR, and BioGen.

Additionally, we present the MS MARCO V2.1 Collection and two topic sets to power the TREC 2024/2025 RAG Tracks.

## The Ragnarök Framework



Ragnarök provides an end-to-end RAG pipeline:

- **(R) Retrieval**: Integrates `pyserini` for first-stage retrieval & `rank_llm` for reranking (monoT5, RankZephyr, RankGPT).
- **(AG) Augmented Generation**: Supports several proprietary LLM endpoints (OpenAI, Cohere, Gemini) and open-weight LLMs (via vLLM) with configurable prompting for cited answer generation.
- **Standardized I/O**: I/O specs with sentence-level citations.
- **Tools**: REST APIs & WebUI for easy use and evaluation.

## Foundation: Corpus & Topics for TREC RAG

**MS MARCO V2.1 Corpora**

- Deduplicated MS MARCO V2 documents (reduced by 8%).
- Segmented corpus using a sliding window (10-sentence window, 5-sentence stride) resulting in **>113 million segments**.
- Results in a retrieval chunk size suitable for both eval & RAG.

**Topic Sets:** Designed to challenge RAG systems beyond factoid QA.

- **TREC-RAGgy 2024**: 120 topics filtered from TREC DL '21-'23. Focus on long-form answers and information aggregation (*24% aggregation topics*). Includes mapped relevance judgments.
- **TREC-Researchy 2024**: Fresh, diverse, non-factoid topics sampled from Researchy Questions. High emphasis on knowledge intensity queries (*80%*) and multi-facetedness (*76%*). Includes relevance judgments from TREC 2024 RAG.

*These resources form the basis for rigorous retrieval and RAG benchmarking in the TREC RAG Track.*

## Retrieval Baselines

We provide several state-of-the-art retrieval pipelines:

- **First-stage**: BM25, Dense (GTE, ArcticEmbed), Hybrid (RRF).
- **Reranking**: Pointwise (monoT5), Listwise (RankZephyr, RankGPT).

| Model | nDCG@10 | MAP@100 | Recall@100 |
|---|---|---|---|
| *Lexical & Dual Encoders* | | | |
| (1a) BM25 | 0.4227 | 0.1561 | 0.2807 |
| (1b) Hybrid | 0.6064 | 0.2592 | 0.3990 |
| *Rerankers* | | | |
| (2a) RRF(1b, monoT5-3B) | 0.6175 | 0.2708 | 0.4208 |
| (2b) RRF(2a, RankZephyr) | 0.6357 | 0.2770 | 0.4208 |

Table 1. Baseline results on Document Ranking of the TREC-RAGgy 2024 set
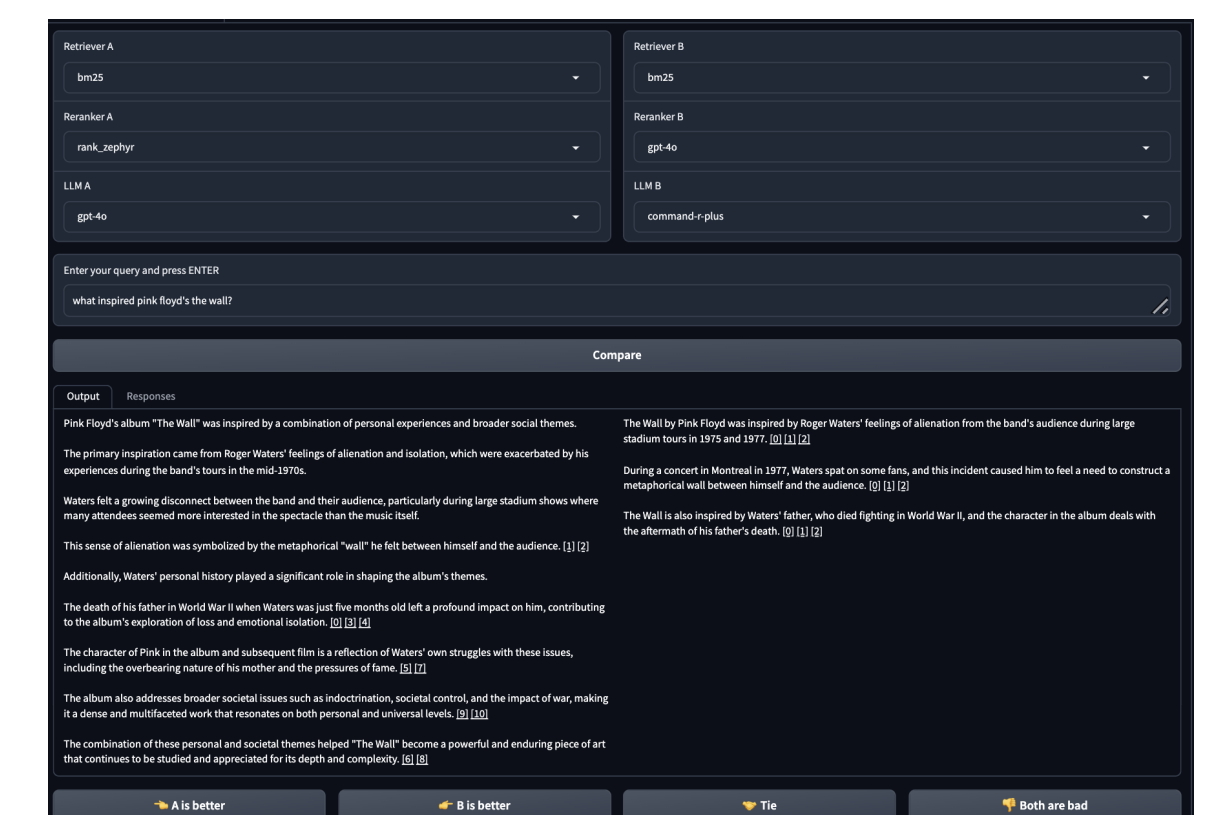
## Augmented Generation Baselines

Using the retrieved context, Ragnarök generates cited answers with leading LLMs like GPT-4o, Command R+, LLaMA3.3-70B.

Configurable prompting (e.g., ChatQA-style) enables quick iterations to get consistent, well-cited answers following the output requirements of popular TREC 2024/2025 Tracks.

*Qualitative analysis indicates GPT-4o is more detailed, with fewer yet more relevant citations compared to Command R+.*

## Ragnarök System Arena: Interactive Evaluation

Evaluate and compare RAG pipelines head-to-head!



- **Pairwise Comparison**: Judge outputs from two pipelines side-by-side (blinded or unblinded).
- **ELO Leaderboards**: Rank individual modules (Retrieval, Generation) or full RAG pipelines.

## Check it Out!