

# LLM-as-MMR: Stepwise Diversification for Redundancy-Aware Document Reranking

Anonymous ACL submission

## Abstract

Modern retrieval pipelines often rely on relevance-only ranking, which can return highly redundant passages and fail to cover complementary evidence needed to answer ambiguous queries. We propose an LLM-driven diversification reranker inspired by Maximal Marginal Relevance (MMR). Given a query and a fixed pool of candidate documents, our method constructs the ranked list iteratively: at each step, the LLM selects the next document by jointly considering (i) relevance to the query and (ii) marginal novelty with respect to the already selected set, explicitly penalizing information overlap while rewarding complementary coverage. This “LLM-as-MMR” formulation leverages the model’s semantic understanding to estimate redundancy beyond surface-level similarity, without relying on hand-crafted similarity functions. Experiments on ambiguous question answering style reranking benchmarks show that our approach improves answer coverage and ranking quality while producing more diverse, less redundant top- $k$  results compared to similarity-based ranking and heuristic diversification baselines.

## 1 Introduction

In modern Retrieval-Augmented Generation (RAG) and open-domain Question Answering (QA), the utility of a retrieval system is not defined solely by the relevance of individual documents, but by the *collective information coverage* of the top- $k$  results (Lewis et al., 2021; Karpukhin et al., 2020). This is particularly critical for **ambiguous queries**—queries that encompass multiple valid interpretations or require multifaceted evidence to answer comprehensively (Min et al., 2020).

Current retrieval pipelines, dominated by bi-encoder embeddings and pointwise cross-encoders, typically score documents independently (Karpukhin et al., 2020; Qu et al., 2021; Nogueira

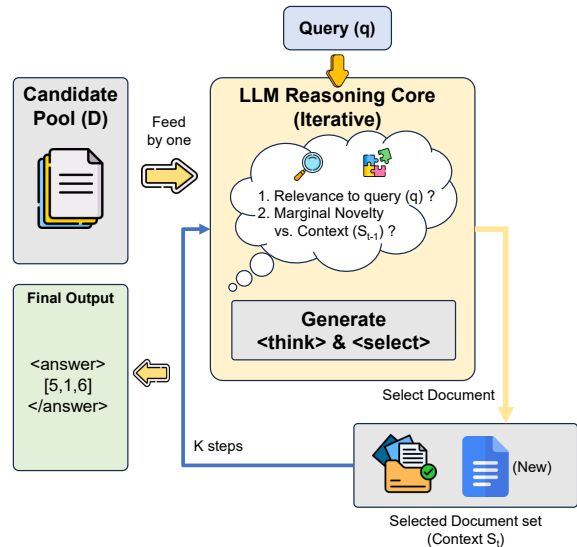


Figure 1: The LLM-as-MMR Framework. The model iteratively selects documents by explicitly reasoning about their relevance to the query ( $q$ ) and marginal novelty against the selected context ( $S_{t-1}$ ), effectively performing stepwise diversification.

and Cho, 2019). In addition to classic dense retrievers, recent open embedding models (e.g., E5 and BGE/FlagEmbedding) have further strengthened this approach (Wang et al., 2024; Xiao et al., 2024). While effective for many queries, independent scoring often leads to the “**Redundancy Trap**”: for ambiguous queries, the top- $k$  list can contain near-duplicates of the dominant interpretation, leaving fewer slots for complementary evidence.

To address this, the classic *Maximal Marginal Relevance* (MMR) algorithm (Carbonell and Goldstein, 1998) introduced a greedy selection criterion that explicitly balances relevance with novelty. More broadly, diversification and subtopic-aware retrieval have been studied extensively (Zhai et al., 2003; Clarke et al., 2008; Agrawal et al., 2009; Gollapudi and Sharma, 2009; Santos et al., 2010). However, classical MMR relies on heuristic vector similarity (e.g., cosine similarity) to estimate

061 redundancy. This approach is brittle: it struggles  
062 to distinguish between *true semantic redundancy*  
063 (which should be penalized) and *lexically similar*  
064 *but complementary details* (which should be pre-  
065 served).

066 In this work, we propose **LLM-as-MMR**, an  
067 LLM-based reranker that applies the MMR idea in  
068 an iterative selection procedure. Instead of using  
069 a fixed similarity function as a redundancy proxy,  
070 we cast reranking as a **sequential decision process**.  
071 At each step, the model conditions on the query  
072 and the *already selected* set, produces a short rationale,  
073 and selects the next document based on its  
074 incremental contribution to the set.

075 We train the reranker in three stages. First, we  
076 perform Supervised Fine-Tuning (SFT) on expert  
077 trajectories to teach the model the required tag format  
078 and the “compare-and-select” behavior, following  
079 standard instruction-tuning practices (Ouyang  
080 et al., 2022; Wang et al., 2023; Rafailov et al.,  
081 2024). Second, because imitation alone does not  
082 directly optimize set-level coverage, we apply **Re-**  
083 **inforcement Learning (RL)** via Group Relative  
084 Policy Optimization (GRPO) to maximize the *Answer Coverage*  
085 objective under format constraints. Third, to reduce  
086 inference cost, we adopt **Coconut** compression  
087 (Hao et al., 2025), distilling each explicit step  
088 into a compact latent representation and achieving  
089 an  $\approx 8\times$  speedup.

090 We evaluate our framework on the **AmbigQA**  
091 benchmark using a **redundancy-based difficulty**  
092 **stratification**. Our experiments show that in highly  
093 redundant pools, RL training improves answer coverage  
094 and achieves performance comparable to strong  
095 pairwise rerankers, while also outperforming  
096 prompted GPT-4.1 baselines.

097 **Contributions.** Our contributions are:

- 098 • We propose *LLM-as-MMR*, an iterative  
099 reranker that implements the relevance–  
100 novelty trade-off with an explicit step-wise  
101 selection format rather than a hand-crafted  
102 similarity heuristic.
- 103 • We show that **RL (GRPO)** improves answer  
104 coverage in high-redundancy pools compared  
105 to supervised baselines.
- 106 • We apply **Coconut** compression to reduce  
107 inference latency by distilling explicit steps  
108 into latent representations.
- 109 • We evaluate on AmbigQA with redundancy-  
110 based stratification and report robustness

gains in the most redundant candidate pools. 111

## 2 Related Work 112

Document reranking has a long history in information retrieval, with most approaches focusing on query–document relevance. Our work is most closely related to (i) diversification-aware reranking, (ii) neural/LLM reranking, and (iii) post-training and efficiency techniques for controllable reranking. 113–119

**Diversification-aware reranking.** Maximal Marginal Relevance (MMR) is a classical greedy criterion that explicitly balances relevance to the query and novelty with respect to already selected items, reducing redundancy in ranked lists (Carbonell and Goldstein, 1998). Beyond MMR, result diversification has been studied through intent/subtopic-aware retrieval and explicit diversification objectives (Zhai et al., 2003; Clarke et al., 2008; Agrawal et al., 2009; Gollapudi and Sharma, 2009; Santos et al., 2010). This line of work motivates our formulation of reranking as iterative set construction; our method differs by using an LLM to make the relevance–novelty decision rather than relying on a fixed similarity function. 120–135

**Neural and LLM reranking.** Neural rerankers based on pretrained language models substantially improved passage ranking by better modeling semantic relevance (Nogueira and Cho, 2019), and late-interaction designs further improve efficiency–effectiveness trade-offs (Khattab and Zaharia, 2020). On the representation side, recent open embedding models such as E5 and BGE/FlagEmbedding have become strong foundations for dense retrieval and reranking pipelines (Wang et al., 2024; Xiao et al., 2024). Beyond these foundations, modern reranking has explored hybrid and collaborative rerankers that leverage multiple retrieval signals (Zhang et al., 2023), and efficiency-oriented reranking for generation pipelines via lattice encoding (Singhal et al., 2023). Related zero-shot reranking approaches based on question generation have also been explored (Sachan et al., 2023). Recent work has also explored listwise reranking with LLMs via prompted permutation generation (Ma et al., 2023; Pradeep et al., 2023; Liu et al., 2025), highlighting both effectiveness and efficiency constraints due to context length and decoding cost. For practical LLM reranking usage and 136–159

reproducible evaluation, we also note recent tooling that standardizes LLM-based reranking pipelines (Sharifymoghaddam et al., 2025). Building on this trend, we use an instruction-following LLM as a controllable reranker that performs step-wise selection, making the relevance–diversity trade-off explicit via a structured output format.

**Post-training and efficiency.** Beyond supervised fine-tuning, reinforcement learning can optimize rankers toward downstream objectives. We use GRPO-style group-relative updates implemented with the VERL training stack (volcengine, 2024) to directly optimize answer coverage under format constraints. Our use of RL is aligned with common policy-optimization recipes used in instruction tuning (Ouyang et al., 2022; Schulman et al., 2017b,a). To reduce the inference cost of explicit step-wise reasoning, we adopt progressive latent-space reasoning and compression following Coconut (Hao et al., 2025).

### 3 Method

We formulate document reranking as a **sequential decision-making process** driven by a Large Language Model (LLM). Unlike traditional pointwise or pairwise rerankers that score documents independently, our approach, *LLM-as-MMR*, iteratively constructs a ranked list by modeling the conditional probability of the next document given the query and the previously selected context. This section details our problem formulation, the three-stage training pipeline (Supervised Fine-Tuning, Reinforcement Learning, and Inference Compression), and the specific mechanisms designed to balance relevance and redundancy.

#### 3.1 Task Formulation: Reranking as Iterative Selection

Let  $q$  be an ambiguous query and  $\mathcal{D} = \{d_i\}_{i=1}^n$  be a candidate pool of retrieved documents. The objective of reranking is to generate an ordered subset  $\mathcal{S}_k = (d_{i_1}, \dots, d_{i_k})$  of size  $k$  (where  $k \leq n$ ) that maximizes the coverage of distinct answers required by  $q$ .

We model this process as an autoregressive generation task. At each step  $t$  ( $1 \leq t \leq k$ ), the model selects a document  $d_{i_t} \in \mathcal{D} \setminus \mathcal{S}_{t-1}$  based on the query  $q$  and the history of selected documents  $\mathcal{S}_{t-1}$ . Inspired by the Maximal Marginal Relevance (MMR) criterion (Carbonell and Gold-

stein, 1998), the selection at step  $t$  aims to maximize the marginal utility:

$$d_{i_t} = \arg \max_{d \in \mathcal{D} \setminus \mathcal{S}_{t-1}} \left[ \lambda \text{Rel}(q, d) - (1 - \lambda) \max_{d' \in \mathcal{S}_{t-1}} \text{Sim}(d, d') \right] \quad (1)$$

While classical MMR relies on heuristic similarity functions (e.g., cosine similarity) for  $\text{Sim}(\cdot, \cdot)$ , our method learns to approximate this trade-off implicitly through natural language reasoning. The LLM acts as a policy  $\pi_\theta$  that outputs a reasoning path (rationale) followed by a selection decision, effectively “soft-computing” the relevance and novelty of candidates in real-time.

**Inference cost.** The step-wise procedure selects  $k$  documents sequentially, so inference cost scales linearly with  $k$  (one generation step per rank). In practice, this makes the approach best suited for reranking a moderate candidate pool (e.g., top-50) rather than producing very deep ranked lists. We therefore focus on  $k = 3$  in our main evaluation.

#### 3.2 Stage I: Supervised Fine-Tuning (SFT) for Stepwise Reasoning

To initialize the policy with the capability of comparative reasoning, we construct a high-quality dataset of reranking trajectories derived from the AmbigQA benchmark.

**Data Synthesis Pipeline.** We develop a five-stage pipeline to synthesize training data that explicitly exposes the “thinking process” of an expert reranker:

- Context Extraction:** We extract pairs of  $(q, \mathcal{D})$  from AmbigQA, where  $\mathcal{D}$  contains both relevant passages answering different sub-intents of  $q$  and distracting passages.
- Oracle Annotation:** Using ground-truth answers, we label each document in  $\mathcal{D}$  with fine-grained tags (e.g., *fully\_relevant*, *partially\_relevant*, *redundant*). Note that these labels are **hidden** from the student model.
- Teacher Self-Decision:** A strong teacher LLM (GPT-4o; Ethics Statement) generates “optimal” reranking trajectories. We prompt the teacher to simulate the MMR process: for each step, it must explain *why* a candidate adds new information compared to  $\mathcal{S}_{t-1}$ .

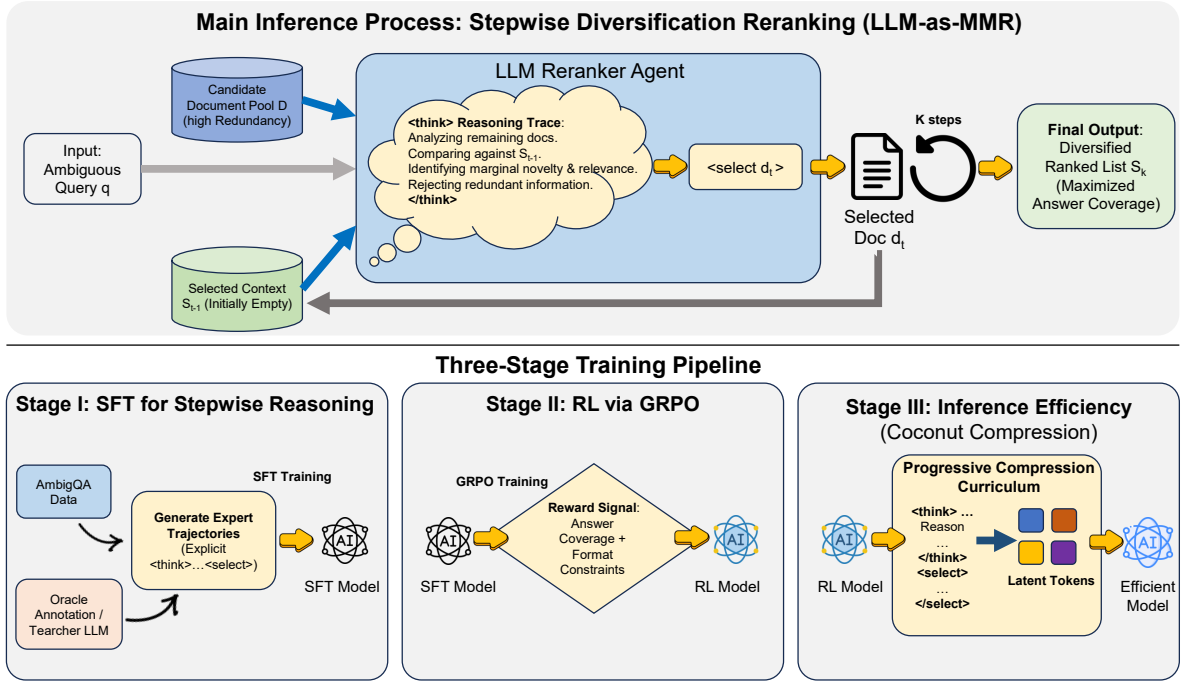


Figure 2: Overview of the LLM-as-MMR framework. (Top) Inference Process: The model functions as a stepwise agent, iteratively selecting documents ( $d_t$ ) by generating explicit reasoning traces to evaluate relevance and penalize redundancy against the selected context ( $S_{t-1}$ ). (Bottom) Three-Stage Training Pipeline: The model is optimized through: (I) SFT on expert trajectories to learn the reasoning format; (II) RL via GRPO to directly maximize answer coverage (DeepSeek-AI et al., 2025); and (III) Progressive Compression (Coconut) to distill explicit reasoning into latent tokens for efficient inference (Hao et al., 2025).

251 4. **Trajectory Formatting:** We curate sequences for multiple target lengths  $k \in \{1, 3, 5\}$  and a **dynamic termination mode**. In dynamic mode, the model learns to output <answer>[ ]</answer> when the remaining pool offers no marginal value.

257 5. **Quality Filtering:** We filter out trajectories where the teacher’s selection contradicts the ground-truth evidence.

260 **Structured Input-Output Format.** The model takes the query and the indexed candidate pool as input. The target output is a structured chain of thought. We introduce special tokens <think>, <select>, and <answer> to enforce a rigid schema:

266 **Input:** <query>... <docs>[1]...[n]...</docs>  
**Output:**  
 <think> Analyzing... Doc [1] covers X. Doc [2] is redundant. Select [3]. </think>  
 <select> 3 </select>  
 ...  
 <answer>[ 1, 3, ... ]</answer>

268 This format forces the model to perform **explicit comparison** before selection, grounding the deci-

270 sion in the document content rather than position bias.

### 272 3.3 Stage II: Reinforcement Learning via Group Relative Policy Optimization (GRPO) 273 274

275 While SFT teaches the *format* of reasoning, it relies on imitation and may not optimally solve the combinatorial optimization of covering all answers. We employ Reinforcement Learning (RL) to directly optimize the end metric (Answer Coverage). 276 277 278 279

280 **Reward Engineering.** We design a composite reward function  $r(y)$  that balances structural validity and semantic utility. For a generated response  $y$ : 281 282

- 283 • **Format Reward ( $r_{\text{fmt}}$ ):** A strict rule-based score that enforces valid XML-like structure and selection constraints. 284 285
- 286 • **Coverage Reward ( $r_{\text{cov}}$ ):** An answer-level coverage score computed by matching gold short answers against the selected documents. 287 288

289 The total reward is a weighted combination:

$$290 r(y) = 2(\alpha \cdot r_{\text{fmt}}(y) + (1 - \alpha) \cdot r_{\text{cov}}(y)) - 1 \quad (2)$$

We empirically set  $\alpha = 0.3$  (i.e., 30% format and 70% coverage). We also include a sparse bonus for perfect termination in dynamic mode.

**Format Reward Details.** The format reward  $r_{\text{fmt}}$  is computed by checking five aspects: (1) **Required tags** (0.25): the output contains `<think>`, `<select>`, and `<answer>`; (2) **Tag pairing** (0.20): tags are properly opened/closed and correctly paired; (3) **Valid indices** (0.25): selected document IDs are within range and contain no duplicates; (4) **Answer format** (0.15): the answer list follows `<answer>[N1,N2,...]</answer>`; (5) **Mode check** (0.15): in fixed mode the list length equals  $k$ , and in dynamic mode the model either selects at least one document or terminates with `<answer>[]</answer>`.

**Coverage Reward Details.** Given gold short answers  $\{a_j\}_{j=1}^m$  from the dataset and the concatenated text of selected documents, we compute an answer-level match score  $s(a_j) \in [0, 1]$ : **exact match**: if  $a_j$  appears as a case-insensitive substring,  $s(a_j) = 1$ ; **soft match**: otherwise, we compute a lexical overlap ratio  $\text{overlap}(a_j, T) \in [0, 1]$  between the tokenized answer  $a_j$  and the tokenized selected text  $T$ , defined as the fraction of answer tokens that appear in  $T$ :

$$\text{overlap}(a_j, T) = \frac{|\text{tok}(a_j) \cap \text{tok}(T)|}{|\text{tok}(a_j)|}$$

If  $\text{overlap} \geq 0.5$ , we set  $s(a_j) = \text{overlap}$ ; otherwise, we down-weight it as  $s(a_j) = 0.5 \cdot \text{overlap}$ . We then define  $r_{\text{cov}}$  as the mean score over all gold answers:  $r_{\text{cov}} = \frac{1}{m} \sum_{j=1}^m s(a_j)$ .

**Relation to the Reported Metric.** The soft-match component is used only for RL reward shaping. The reported  $\text{Cov}@k$  metric in Section 4 is computed using case-insensitive string matching.

**Hard Constraints and Bonuses.** We apply additional shaping terms: in fixed mode, if the output length mismatches  $k$ , we apply a penalty of  $-0.3$ ; if the format is perfect and the length matches  $k$ , we add a bonus of  $+0.1$ ; if the model outputs an empty selection without a proper dynamic termination, we apply a penalty of  $-0.2$ .

**Optimization Objective.** We use the VERL framework (volcengine, 2024) to implement GRPO (DeepSeek-AI et al., 2025). For a query  $x$ , we sample a group of  $N$  outputs  $\{y_i\}_{i=1}^N$  from  $\pi_\theta$ . The advantage  $A_i$  is computed as  $A_i =$

$(r(y_i) - \bar{r}) / (\sigma_r + \epsilon)$ . The objective maximizes the advantage-weighted log-likelihood with KL regularization:

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}} \left[ \frac{1}{N} \sum_{i=1}^N \left( A_i \log \pi_\theta(y_i | x) - \beta \text{KL}(\pi_\theta || \pi_{\text{ref}}) \right) \right] \quad (3)$$

This stage aligns the model’s “internal utility function” with the true goal of diverse retrieval.

### 3.4 Stage III: Inference Efficiency via Progressive Compression (Coconut)

The explicit reasoning steps generated in Stages I and II improve performance but incur high latency costs. To resolve this, we adapt the **Coconut** paradigm (Hao et al., 2025) to compress reasoning into latent space.

**Curriculum Learning Schedule.** We employ a progressive compression schedule. The training proceeds in stages defined by a schedule function  $s(e)$  dependent on the epoch  $e$ :

$$s(e) = \left\lfloor \frac{e}{E_{\text{stage}}} \right\rfloor \quad (4)$$

In stage  $s$ , the first  $s$  explicit reasoning steps are replaced by  $s \cdot c_{\text{thought}}$  latent tokens. As training progresses ( $s \rightarrow k$ ), the model gradually internalizes the entire MMR reasoning loop.

**Latent Token Budget.** In our implementation, we set  $c_{\text{thought}} = 1$ , meaning each step’s `<think>...</think><select>...</select>` block is compressed into a single latent token during training. We set  $E_{\text{stage}} = 3$ , i.e., we train each compression stage for 3 epochs before increasing the number of replaced reasoning steps.

## 4 Experiments

In this section, we empirically validate the effectiveness of *LLM-as-MMR*. We evaluate our approach on the AmbigQA benchmark, focusing on the model’s ability to maximize answer coverage under strict ranking constraints ( $k = 3$ ).

### 4.1 Experimental Setup

#### 4.1.1 Dataset Stratification

To systematically analyze model robustness against redundancy, we stratify the validation set into three

377 tiers based on the **Redundancy Score** ( $\rho$ ) of the  
378 candidate pool. We use Qwen3-embedding-8B  
379 (Zhang et al., 2025) to encode documents and com-  
380 pute  $\rho$  as the average pairwise cosine similarity.

- 381 • **Tier I (Low Difficulty)**: Diverse pools ( $\rho <$   
382  $Q_{0.33}$ ).
- 383 • **Tier II (Medium Difficulty)**: Moderate over-  
384 lap ( $Q_{0.33} \leq \rho \leq Q_{0.67}$ ).
- 385 • **Tier III (High Difficulty)**: Highly redundant  
386 pools ( $\rho > Q_{0.67}$ ).

387 **Candidate Pool Source.** We directly rerank the  
388 per-query candidate pool provided by the Am-  
389 bigQA dataset release (i.e., no external retriever  
390 is used to form  $\mathcal{D}$ ).

### 391 4.1.2 Baselines and Metrics

392 We compare our method against: (1) **Embedding**  
393 **Sim**: Standard cosine similarity ranking using  
394 Qwen3-embedding-8B. (2) **MMR**: Classical diver-  
395 sification with  $\lambda = 0.5$ . (3) **Qwen3-Reranker**: A  
396 state-of-the-art pairwise reranker (8B) (Zhang et al.,  
397 2025). (4) **GPT-4.1 (Prompted)**: A prompted  
398 LLM baseline using GPT-4.1 (OpenAI et al., 2024),  
399 instructed to follow our stepwise reasoning-and-  
400 selection format.

401 We report **Answer Coverage (Cov@3)** and  
402 **NDCG@3** (Järvelin and Kekäläinen, 2002). Note  
403 on AvgRel: While we monitor Average Relevance,  
404 we observe that high AvgRel in Tier III often corre-  
405 lates with low Coverage (the “Redundancy Trap”),  
406 as the model selects repetitive passages that are all  
407 highly similar to the query.

408 **Answer Matching for Coverage.** AmbigQA  
409 provides multiple gold short answers per query. We  
410 compute coverage by performing case-insensitive  
411 string matching of each gold answer against the  
412 concatenated text of the selected top- $k$  documents.  
413 Cov@3 is the mean match rate over all gold an-  
414 swers for each query, averaged across queries.

415 **Relevance Labels for NDCG.** We assign a  
416 graded relevance label to each candidate document  
417 based on its **answer coverage ratio**, defined as the  
418 fraction of gold short answers that appear (case-  
419 insensitive string match) in that document. We use  
420 this ratio as the gain for computing NDCG@3.

421 **Similarity Function for MMR.** For classical  
422 MMR,  $\text{Sim}(\cdot, \cdot)$  is cosine similarity computed over  
423 Qwen3-embedding-8B embeddings.

### 424 4.1.3 Implementation Details

425 Our models are based on Llama-3-8B (Grattafiori  
426 et al., 2024).

- 427 • **SFT**: Fine-tuned on expert MMR trajectories.
- 428 • **RL**: Optimized with Group Relative Policy  
429 Optimization (GRPO) to maximize coverage  
430 reward.
- 431 • **Coconut**: Our efficient inference variant. We  
432 train the model to compress the explicit rea-  
433 soning steps into **1 latent token** per selec-  
434 tion step. We apply a multi-stage curriculum  
435 where each implicit reasoning stage is trained  
436 for **3 epochs** to ensure the latent representa-  
437 tions effectively capture the rejection logic.

438 **Input/output format.** All LLM-based methods  
439 use the same numbered candidate format and the  
440 tag schema from §3.2; prompt templates are pro-  
441 vided in Appendix A. We evaluate on the Am-  
442 bigQA validation candidate pools using consistent  
443 case-insensitive matching for coverage and NDCG  
444 labels.

445 **Fixed- $k$  evaluation protocol.** Unless stated oth-  
446 erwise, we use fixed- $k$  reranking with  $k=3$  for all  
447 methods and report Cov@3 and NDCG@3 on the  
448 resulting top-3 list. This isolates set-level coverage  
449 behavior under a strict budget and makes compar-  
450 isons across embedding, heuristic, and LLM-based  
451 rerankers consistent.

### 452 4.2 Main Results

453 Table 1 presents the performance breakdown across  
454 difficulty tiers.

455 **Overall Performance.** As shown in Table 1, our  
456 methods consistently outperform the Embedding  
457 and MMR baselines. **Ours (Coconut)** achieves  
458 the highest overall coverage (42.6%), outperform-  
459 ing GPT-4.1 (Prompted) (41.4%). This demon-  
460 strates that specialized fine-tuning on redundancy-  
461 aware trajectories is more effective than prompt-  
462 only baselines.

463 **Performance by Difficulty.** The stratified analy-  
464 sis reveals distinct behaviors:

- 465 • **Tier I (Low Difficulty)**: In diverse pools,  
466 Coconut performs exceptionally well (Cov  
467 43.0%). When candidates are distinct, the  
468 model effectively identifies the optimal subset  
469 without "over-thinking."

Table 1: Performance breakdown on AmbigQA ( $k = 3$ ). **Bold** indicates the best performance among methods in that tier. Data is stratified by pool redundancy (Tier I=Low, Tier III=High). Note that in Tier III, while Embedding Sim fails (Cov 36.5%), our RL model maintains high robustness (Cov 41.7%).

Method	Overall (Avg)		Tier I (Low Diff.)		Tier II (Med Diff.)		Tier III (High Diff.)	
	Cov@3	NDCG@3	Cov@3	NDCG@3	Cov@3	NDCG@3	Cov@3	NDCG@3
<i>Baselines</i>								
Embedding Sim	37.1	23.6	39.9	24.6	34.8	22.8	36.5	23.3
MMR ( $\lambda = 0.5$ )	35.8	20.2	37.2	20.1	35.7	19.6	34.3	20.9
Qwen3-Reranker	41.6	32.7	42.1	30.1	40.9	31.4	41.7	<b>36.6</b>
GPT-4.1 (Prompted)	41.4	33.1	39.8	32.2	39.8	31.3	41.0	35.7
<i>Ours</i>								
Ours (SFT)	42.2	32.8	43.0	34.1	41.8	31.2	41.7	33.2
<b>Ours (RL)</b>	41.7	32.9	41.7	<b>34.6</b>	41.6	30.8	<b>41.7</b>	33.2
Ours (Coconut)	<b>42.6</b>	<b>33.6</b>	<b>43.0</b>	33.6	<b>42.2</b>	<b>33.4</b>	39.7	33.8

- **Tier III (High Difficulty):** This tier represents the "Redundancy Trap." Embedding Sim struggles (Cov 36.5%) due to retrieving near-duplicates. **Ours (RL)** maintains robust coverage (41.7%), matching the strong Qwen3-Reranker. This indicates that explicit RL optimization helps the model navigate highly redundant spaces.

- **Coconut in Tier III:** We observe a slight dip in Coconut’s performance in Tier III (39.7%) compared to RL. This suggests that extremely subtle redundancy distinctions might be harder to compress into a single latent token.

## 5 Analysis

In this section, we provide a comprehensive analysis of the model’s reasoning process, efficiency, and limitations.

### 5.1 Reasoning Process and Efficiency

**Reasoning Form.** Our method employs an iterative selection process where the model explicitly reasons about document relevance and redundancy. Unlike black-box scorers (e.g., Qwen3-Reranker) or vector heuristics (MMR), our model generates a structured <think> trace before each selection. This trace evaluates the candidate’s marginal utility against the already selected set  $S_{t-1}$ , enabling interpretable and controllable diversification.

**Efficiency Comparison.** While explicit reasoning improves performance, it incurs latency. Table 2 compares the average inference time per query (batch size=1) on a single A100 (80GB) GPU. **Coconut** compresses the reasoning chain into 1 latent

Method	Latency (ms)
Ours (SFT/RL; vLLM)	$\approx 810$
<b>Ours (Coconut; no vLLM)</b>	<b>106.9</b>

Table 2: Efficiency comparison (A100 80GB, batch size=1). vLLM uses PagedAttention (Kwon et al., 2023). Coconut runs with a vanilla Transformer backend (no vLLM). Coconut achieves an  $\approx 8\times$  speedup over explicit reasoning models.

token per step, achieving a  $7.7\times$  **speedup** (106.9 ms) compared to the explicit models (810 ms), even without vLLM optimization.

### 5.2 Successful Case Study: Disentangling Ambiguity

We present a real-world example from the validation set to illustrate how our model handles high-redundancy scenarios.

**Query:** *Who played george washington in the john adams series?*

**Candidate Pool:** Contains multiple documents about the "John Adams" miniseries (Docs [1], [5], [8]), the biography book (Doc [2]), and actors (David Morse [9], William Daniels [3]).

The baseline Embedding Sim model tends to select Docs [1], [2], and [5] because they all strongly match the phrase "John Adams series", failing to pinpoint the specific actor. In contrast, our model produces the following stepwise trace:

- **Step 1 Selection: [9]** (David Morse).  
*Reasoning:* "The most important aspect... is identifying the actor... Document [9] is the only candidate that directly provides this information... compared to background documents, [9] delivers the answer in the clearest way."

529 • **Step 2 Selection: [1]** (Miniseries Info).  
 530 *Reasoning:* “To build a comprehensive answer, it’s important to include... contextual  
 531 details... While the first selected document  
 532 clearly states who played George Washington,  
 533 it does not explain what the ‘John Adams’  
 534 miniseries is... Among remaining options, [1]  
 535 stands out by providing a concise overview...”  
 536

537 • **Step 3 Selection: [8]** (Broader Context).  
 538 *Reasoning:* “The selected documents so far  
 539 cover the basics... To create a more compre-  
 540 hensive answer, it’s important to add broader  
 541 historical context. Document [8] achieves this  
 542 by situating the miniseries within the wider  
 543 landscape... This information is not present in  
 544 the previously selected documents...”

545 **Outcome:** The model successfully retrieves the  
 546 direct answer (David Morse) first, then supple-  
 547 ments it with essential context (Miniseries info)  
 548 and broader context (Genre info), achieving 100%  
 549 coverage and a perfect ranking order. This demon-  
 550 strates the model’s ability to plan a cohesive infor-  
 551 mation package rather than just picking keywords.

552 **5.3 Failure Analysis**

553 Despite the success, our analysis of Tier III failure  
 554 cases reveals specific limitations.

555 **Case 1: Temporal Ambiguity (The Simpsons).**  
 556 Query: “When did the simpsons first air on televi-  
 557 sion?”. Valid answers include both 1987 (shorts)  
 558 and 1989 (series). Our model failed to retrieve  
 559 the 1987 date. *Reasoning:* The candidate pool  
 560 was dominated by documents discussing the 1989  
 561 premiere. The model’s redundancy detection mech-  
 562 anism correctly identified that multiple 1989 docu-  
 563 ments were redundant, but the semantic signal for  
 564 the 1987 short was too weak compared to the over-  
 565 whelming consensus of the 1989 date. The model  
 566 prioritized the “authoritative” 1989 date and filled  
 567 remaining slots with context about that specific  
 568 event, missing the secondary interpretation.

569 **Case 2: Abstract Concepts (Consubstantial).**  
 570 Query: “What does consubstantial mean?”. Valid  
 571 answers span theological and philosophical defini-  
 572 tions. *Reasoning:* The model successfully retrieved  
 573 the common theological definition (“of the same  
 574 being”) but missed the broader philosophical defi-  
 575 nition (“common humanity”). *Analysis:* The model  
 576 exhibits a bias towards the “dominant” interpreta-  
 577 tion present in the training data. While it avoids

578 verbatim redundancy (choosing different phrasings  
 579 of the theological definition), it fails to recognize  
 580 that the philosophical definition represents a *dis-*  
 581 *tinct valid answer category* rather than just a less  
 582 relevant distractor. 583

584 **5.4 Discussion**

585 The qualitative analysis highlights a critical distinc-  
 586 tion: **Redundancy Avoidance vs. Answer Diver-**  
 587 **sity.** Our model excels at *Redundancy Avoidance*  
 588 (e.g., rejecting two documents that say the exact  
 589 same thing). However, in cases like “The Simp-  
 590 sons” or “Consubstantial”, the challenge is *Answer*  
 591 *Diversity* (recognizing that two factually different  
 592 statements are both valid answers). Future work  
 593 should focus on training objectives that explicitly  
 594 reward retrieving *conflicting* or *multi-faceted* infor-  
 595 mation when the query is inherently ambiguous.

596 **6 Conclusion**

597 In this paper, we introduced *LLM-as-MMR*, an  
 598 iterative reranking method inspired by Maximal  
 599 Marginal Relevance. By casting document se-  
 600 lection as a sequential prediction problem, our  
 601 approach replaces a fixed similarity-based redun-  
 602 dancy proxy with step-wise selection conditioned  
 603 on the already selected set, allowing the model to  
 604 trade off relevance and marginal novelty explicitly.

605 Our empirical evaluation on the AmbigQA  
 606 benchmark demonstrates that our approach effec-  
 607 tively addresses the “Redundancy Trap” inherent  
 608 in traditional retrieval systems. Specifically, our  
 609 RL-optimized model maintains robust answer cov-  
 610 erage in high-difficulty scenarios where standard  
 611 embedding methods collapse due to severe infor-  
 612 mation overlap. We further established that this  
 613 reasoning capability can be efficiently deployed:  
 614 by integrating the Coconut paradigm, we com-  
 615 pressed the explicit reasoning chain into compact  
 616 latent tokens, achieving an  $\approx 8\times$  speedup over  
 617 standard Chain-of-Thought models with minimal  
 618 performance degradation.

619 This work suggests that improving diversity can  
 620 benefit from modeling set construction explicitly,  
 621 rather than relying only on stronger embeddings.  
 622 Future directions include extending the objective  
 623 to multi-criteria reranking (e.g., coverage, faith-  
 624 fulness, coherence) and exploring adaptive com-  
 625 pression strategies that adjust the latent reasoning  
 626 length based on query complexity.

## 626 Limitations

627 Despite the promising results, our approach has in-  
628 herent constraints regarding efficiency and scalabil-  
629 ity. First, while **Coconut** compression significantly  
630 reduces latency to  $\approx 100$  ms, it remains an order  
631 of magnitude slower than lightweight bi-encoders,  
632 restricting its utility to second-stage reranking on  
633 small candidate pools (e.g., top-50). Furthermore,  
634 the fundamental design of *stepwise selection* is au-  
635 toregressive, meaning inference cost scales linearly  
636 with the target list size  $k$ . This prevents paralleliza-  
637 tion across ranks, making the method unsuitable  
638 for generating very deep ranked lists.

639 Second, we observe a trade-off between reason-  
640 ing compression and resolution. In our high-  
641 difficulty experiments (Tier III), the compressed  
642 model exhibited a slight performance degrada-  
643 tion compared to the explicit Chain-of-Thought  
644 model, suggesting that a single latent token may en-  
645 counter an information bottleneck when resolving  
646 extremely subtle semantic redundancies. Addition-  
647 ally, our empirical validation focused exclusively  
648 on AmbigQA to target ambiguity resolution; the  
649 generalization of this redundancy-aware reasoning  
650 to other domains, such as legal or medical retrieval  
651 where "redundancy" may have different definitions,  
652 remains to be verified in future work.

## 653 Ethics Statement

654 We hereby acknowledge that all authors of this  
655 work are aware of the ACL Code of Ethics and  
656 honor the code of conduct.

657 **Data Use and Privacy.** Our experiments are pri-  
658 marily conducted on the public **AmbigQA** bench-  
659 mark (Min et al., 2020), which is derived from  
660 Wikipedia and does not contain private or person-  
661 ally identifiable information (PII). For the construc-  
662 tion of our instruction-tuning dataset, we utilized a  
663 synthetic data generation pipeline where a teacher  
664 LLM (GPT-4o) generated reasoning traces based  
665 on public passages. We have manually verified a  
666 subset of these trajectories to ensure they do not  
667 propagate toxic or harmful content. All data usage  
668 complies with the licenses of the respective source  
669 datasets.

670 **Use of AI Assistants.** We employed large lan-  
671 guage models (specifically GPT-4.1) as a **writing**  
672 **assistant** to refine the grammatical structure and  
673 clarity of the manuscript. We also used a teacher

LLM (GPT-4o) to synthesize the reasoning trajec- 674  
tories for Supervised Fine-Tuning, as described in 675  
Section 3.2. The core scientific ideas, experimental 676  
design, and analysis were conducted solely by the 677  
authors. 678

## References 679

- Rakesh Agrawal, Sreenivas Gollapudi, Alan Halver- 680  
son, and Samuel Jeong. 2009. [Diversifying search 681](#)  
[results](#). In *Proceedings of the Second ACM Interna- 682*  
*tional Conference on Web Search and Data Mining*, 683  
pages 5–14. 684
- Jaime Carbonell and Jade Goldstein. 1998. The use of 685  
MMR, diversity-based reranking for reordering doc- 686  
uments and producing summaries. In *Proceedings 687*  
*of the 21st Annual International ACM SIGIR Confer- 688*  
*ence on Research and Development in Information 689*  
*Retrieval*, pages 335–336. 690
- Charles L.A. Clarke, Maheedhar Kolla, Gordon V. 691  
Cormack, Olga Vechtomova, Azin Ashkan, Stefan 692  
Büttcher, and Ian MacKinnon. 2008. [Novelty and 693](#)  
[diversity in information retrieval evaluation](#). In *Pro- 694*  
*ceedings of the 31st annual international ACM SIGIR 695*  
*conference on Research and development in informa- 696*  
*tion retrieval*, pages 659–666. 697
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, 698  
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, 699  
Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, 700  
Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi- 701  
hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 702  
2025. [Deepseek-r1: Incentivizing reasoning capa- 703](#)  
[bility in llms via reinforcement learning](#). *Preprint*, 704  
arXiv:2501.12948. 705
- Sreenivas Gollapudi and Aneesh Sharma. 2009. [An 706](#)  
[axiomatic approach for result diversification](#). In 707  
*Proceedings of the 18th international conference on 708*  
*World wide web*, pages 381–390. 709
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, 710  
Abhinav Pandey, Abhishek Kadian, Ahmad Al- 711  
Dahle, Aiesha Letman, Akhil Mathur, Alan Schel- 712  
ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh 713  
Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi- 714  
tra, Archie Sravankumar, Artem Korenev, Arthur 715  
Hinsvark, and 545 others. 2024. [The llama 3 herd of 716](#)  
[models](#). *Preprint*, arXiv:2407.21783. 717
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, 718  
Zhiting Hu, Jason Weston, and Yuandong Tian. 2025. 719  
[Training large language models to reason in a contin- 720](#)  
[uous latent space](#). *Preprint*, arXiv:2412.06769. 721
- Kalervo Järvelin and Jaana Kekäläinen. 2002. [Cumu- 722](#)  
[lated gain-based evaluation of ir techniques](#). *ACM 723*  
*Transactions on Information Systems*, 20(4):422– 724  
446. 725

726	Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen tau Yih. 2020. <a href="#">Dense passage retrieval for open-domain question answering</a> . <i>Preprint</i> , arXiv:2004.04906.	783
727		784
728		785
729		786
730		787
731	Omar Khattab and Matei Zaharia. 2020. <a href="#">Colbert: Efficient and effective passage search via contextualized late interaction over bert</a> . <i>Preprint</i> , arXiv:2004.12832.	788
732		789
733		790
734		791
735	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. <a href="#">Efficient memory management for large language model serving with pagedattention</a> . <i>Preprint</i> , arXiv:2309.06180.	792
736		793
737		794
738		795
739		796
740		797
741	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. <a href="#">Retrieval-augmented generation for knowledge-intensive nlp tasks</a> . <i>Preprint</i> , arXiv:2005.11401.	798
742		799
743		800
744		801
745		802
746		803
747	Qi Liu, Bo Wang, Nan Wang, and Jiaxin Mao. 2025. <a href="#">Leveraging passage embeddings for efficient listwise reranking with large language models</a> . <i>Preprint</i> , arXiv:2406.14848.	804
748		805
749		806
750		807
751	Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. 2023. <a href="#">Zero-shot listwise document reranking with a large language model</a> . <i>Preprint</i> , arXiv:2305.02156.	808
752		809
753		810
754		811
755	Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. <a href="#">Ambigqa: Answering ambiguous open-domain questions</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> .	812
756		813
757		814
758		815
759		816
760	Rodrigo Nogueira and Kyunghyun Cho. 2019. <a href="#">Passage re-ranking with BERT</a> . <i>arXiv preprint arXiv:1901.04085</i> .	817
761		818
762		819
763	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> , arXiv:2303.08774.	820
764		821
765		822
766		823
767		824
768		825
769		826
770		827
771	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> . <i>Preprint</i> , arXiv:2203.02155.	828
772		829
773		830
774		831
775		832
776		833
777		834
778		835
779	Ronak Pradeep, Sahel Sharifmoghaddam, and Jimmy Lin. 2023. <a href="#">Rankvicuna: Zero-shot listwise document reranking with open-source large language models</a> . <i>Preprint</i> , arXiv:2309.15088.	836
780		837
781		838
782		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

- 835 Cheng Xiang Zhai, William W. Cohen, and John Laf-  
836 fertility. 2003. [Beyond independent relevance: meth-](#)  
837 [ods and evaluation metrics for subtopic retrieval](#). In  
838 *Proceedings of the 26th annual international ACM*  
839 *SIGIR conference on Research and development in*  
840 *informaion retrieval*, pages 10–17.
- 841 Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang,  
842 Huan Lin, Baosong Yang, Pengjun Xie, An Yang,  
843 Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren  
844 Zhou. 2025. [Qwen3 embedding: Advancing text](#)  
845 [embedding and reranking through foundation models](#).  
846 *Preprint*, arXiv:2506.05176.
- 847 Zongmeng Zhang, Wengang Zhou, Jiaxin Shi, and  
848 Houqiang Li. 2023. [Hybrid and collaborative pas-](#)  
849 [sage reranking](#). *Preprint*, arXiv:2305.09313.

850	<b>A Prompting and Data Construction</b>	
851	This appendix documents our implementation details. Unless stated otherwise, evaluation follows the AmbigQA reranking setup (§4) and our method definition in §3. Importantly, <b>at inference time our reranker only observes the query and candidate documents</b> ; any gold answers or oracle hints are used <b>only</b> during offline data synthesis/analysis.	
858	<b>A.1 Data Construction Pipeline (SFT Trajectories)</b>	
859		
860	<b>Overview.</b> We synthesize step-wise “expert” reranking trajectories to initialize the student model with (i) the rigid output schema and (ii) comparative reasoning behaviors (redundancy checks, complementarity, and multi-intent coverage). Concretely, each training sample corresponds to a query $q$ and its AmbigQA candidate pool $\mathcal{D}$ , and produces a sequence of step decisions ending with a ranked list of indices.	
861		
862		
863		
864		
865		
866		
867		
868		
869	<b>Implementation notes.</b> Two details are important for reproducibility and to avoid confusion with evaluation: (i) we always rerank the <b>dataset-provided</b> candidate pool (no external retriever); (ii) answer matching for both offline analysis and reported $\text{Cov}@k$ is <b>case-insensitive substring matching</b> (§4); any soft matching is used only for RL reward shaping (§3.3).	
870		
871		
872		
873		
874		
875		
876		
877	<b>Pipeline summary.</b> The full five-stage procedure is summarized in Table 3.	
878		
879	<b>A.2 Prompt Templates</b>	
880	<b>Design goals.</b> All prompts enforce a rigid output contract compatible with our training schema (§3.2): each step emits a <code>&lt;think&gt;</code> rationale and a <code>&lt;select&gt;</code> decision; after completing the process, the model outputs a final <code>&lt;answer&gt;</code> list of selected indices.	
881		
882		
883		
884		
885		
886	<b>Why so many templates?</b> We use different templates for (i) trajectory synthesis (teacher), (ii) offline auditing/analysis (oracle-style), and (iii) prompted baselines. To keep the appendix compact and readable, we present the prompts as structured tables that list the key blocks and constraints; this avoids long verbatim blocks that tend to overflow in a two-column layout.	
887		
888		
889		
890		
891		
892		
893		
	<b>A.2.1 Teacher Prompt for Trajectory Synthesis (GPT-4o)</b>	894
		895
	<b>Structure.</b> The teacher prompt is composed of a <i>context block</i> (query and documents), a <i>selection rule block</i> (how to trade off relevance and redundancy), and a strict <i>output contract</i> . In synthesis-only mode, we may include oracle hints to improve coverage; these are never shown to the student at test time. See Table 4.	896
		897
		898
		899
		900
		901
		902
	<b>Selection principles.</b> We explicitly instruct the teacher to reason about marginal utility relative to the already-selected set $\mathcal{S}_{t-1}$ , rather than scoring each document in isolation. This yields trajectories that expose the redundancy-aware decision rule to the student. See Table 5.	903
		904
		905
		906
		907
		908
	<b>Fixed vs. dynamic mode.</b> Fixed- $k$ is used for standard top- $k$ reranking (e.g., $k \in \{1, 3, 5\}$ ). Dynamic mode teaches the model to stop when the remaining pool provides no additional value, which is useful for ablations and for understanding termination behavior. See Table 6.	909
		910
		911
		912
		913
		914
	<b>A.2.2 Oracle-Style Labeling Prompt (Quality Control)</b>	915
		916
	<b>Purpose.</b> This prompt is used only for <b>offline</b> data auditing and debugging during dataset construction. Given $(q, \mathcal{D})$ and the gold short answers, it asks an LLM to summarize each document’s query relevance and which gold answers appear in the document text. These outputs help us sanity-check synthetic trajectories (e.g., whether answer-containing documents exist but were missed), but they are <b>never</b> exposed to the student reranker at inference time.	917
		918
		919
		920
		921
		922
		923
		924
		925
		926
	<b>Template summary.</b> We summarize the input/output structure and label set in Table 7.	927
		928
	<b>A.2.3 Prompted GPT-4.1 Baseline</b>	929
	<b>Purpose.</b> This baseline evaluates how far a strong general LLM can go with prompting alone, without any fine-tuning. We prompt GPT-4.1 to follow the same step-wise selection contract as our models ( <code>&lt;think&gt;/&lt;select&gt;/&lt;answer&gt;</code> ) so that comparisons reflect <b>ranking behavior</b> rather than formatting differences.	930
		931
		932
		933
		934
		935
		936
	<b>Template summary.</b> The baseline prompt specifies the fixed- $k$ constraint, the required tags, and the numbered candidate pool; see Table 7.	937
		938
		939

Stage	Details
1. Context extraction	Extract $(q, \mathcal{D})$ from AmbigQA, where $\mathcal{D}$ is the per-query candidate pool released with the dataset (no external retriever used).
2. Gold normalization	Normalize gold short answers (lowercasing for matching; whitespace/punctuation normalization) and deduplicate answer variants.
3. Oracle analysis (offline)	Using gold answers, compute per-document answer coverage ratio $g(d)$ (case-insensitive substring match) to support quality control and difficulty analysis.
4. Teacher trajectory synthesis	Use a teacher LLM (GPT-4o; Ethics Statement) to produce step-wise selection traces that follow the <code>&lt;think&gt;/&lt;select&gt;/&lt;answer&gt;</code> schema. In synthesis-only mode, oracle hints (e.g., gold answers) may be provided to improve coverage, but are never exposed at test time.
5. Formatting & filtering	Serialize trajectories into instruction-tuning format; discard samples that violate tag format, select invalid indices, or contradict gold evidence (e.g., selecting a clearly irrelevant set when answer-containing docs exist).

Table 3: Five-stage data construction pipeline for SFT trajectories.

#### A.2.4 SFT/RL Student Model Prompts (Llama-3-8B)

**Prompt variants.** We use two prompt variants for the student model (§4): a **standard trace** prompt that elicits interleaved `<think>` and `<select>` tags (useful for interpretability and RL format checking), and an optional **answer-only** prompt that requests only the final `<answer>` list when intermediate traces are not needed.

**Template summary.** Both variants share the same numbered `<docs>` input format; see Table 7.

**Non-teacher prompts (summary).** To reduce the number of separate floats (and the whitespace they can create in a two-column layout), we summarize the oracle-style prompt, the prompted GPT baseline, and the student prompt variants in a single table below. See Table 7.

#### Reference Tables (Appendix A)

We list the corresponding reference tables below.

Field	Content (abridged)
Role	You are an expert document ranker selecting a comprehensive, non-redundant set for an ambiguous query.
Input blocks	Query: {q}; Oracle hints (optional, synthesis-only): {h}; Selected: {S}; Remaining: {D\S}.
Task objective	Maximize collective answer coverage of the top- $k$ set while penalizing redundancy relative to the already-selected set.
Required behavior	Each step must explicitly compare multiple candidates and state whether each is redundant/non-redundant with respect to $\mathcal{S}_{t-1}$ .
Output contract	Repeat per step: <code>&lt;think&gt; ... &lt;/think&gt;</code> then <code>&lt;select&gt;N&lt;/select&gt;</code> . End with <code>&lt;answer&gt;[N1,N2,...]&lt;/answer&gt;</code> .

Table 4: Teacher prompt skeleton for trajectory synthesis (GPT-4o). Oracle hints are used *only* to synthesize higher-coverage trajectories and are never provided at test time.

Selection principles	Checklist (used in synthesis prompts)
Directness	Prefer documents that directly contain missing answer strings or unambiguous facts early (especially in the first step).
Complementarity	Prefer candidates that add new facets (different intent/aspect/definition/date) not covered by $\mathcal{S}_{t-1}$ .
Redundancy control	Mark near-duplicates as redundant when they repeat the same key facts as already selected docs; avoid selecting paraphrases of the same answer unless needed for disambiguation.
Comparative reasoning	Compare at least 3–4 candidates each round and justify why the chosen candidate has the highest marginal utility.

Table 5: Selection checklist embedded in the teacher prompt to encourage explicit marginal-utility reasoning.

Mode	Additional instructions
Fixed- $k$	Must select exactly $k$ documents in total. Each round must output a valid <code>&lt;select&gt;</code> with an in-range, non-duplicate index until $ S  = k$ .
Dynamic	May terminate early only after checking all remaining candidates and concluding that they are irrelevant or redundant. Termination is represented by emitting <code>&lt;answer&gt;[...]</code> (no further selections).

Table 6: Mode-specific constraints used for teacher trajectory synthesis.

Stage	Prompt summary (abridged)
<b>Oracle-style auditing</b>	<b>Input:</b> Query, Gold answers, and each candidate document. <b>Task:</b> summarize relevance and which answers appear; output structured labels (Core-Complete/Core-Partial/Supplementary/Related/Irrelevant). Used only for offline debugging/quality control.
<b>GPT-4.1 (Prompted)</b>	<b>System:</b> select exactly $k$ documents and output only required tags. <b>User:</b> <code>&lt;query&gt;{q}&lt;/query&gt;</code> and numbered <code>&lt;docs&gt;</code> . <b>Output:</b> interleaved <code>&lt;think&gt;</code> and <code>&lt;select&gt;</code> repeated $k$ times, ending with <code>&lt;answer&gt;[...]</code> .
<b>Student (Llama-3-8B)</b>	<b>Standard (trace):</b> same tag format for interpretability and RL format checking. <b>Answer-only:</b> optional fast path that outputs only <code>&lt;answer&gt;[...]</code> when intermediate traces are not needed.

Table 7: Summary of non-teacher prompts used in our pipeline and baselines.

## B Dataset Statistics and Data Formats

### B.1 Dataset Statistics

**Summary.** Key dataset statistics are summarized in Table 8.

**How the statistics are computed.** We report counts over the evaluation split used in our experiments: the number of queries after exclusion for leakage prevention, the per-query candidate pool size  $|\mathcal{D}|$  provided by AmbigQA, and the number of gold short answers per query as released by the dataset. The averages are computed across queries in this evaluation split.

### B.2 Data Formats

**Summary.** The final serialized schema is shown in Table 9.

**Serialization details.** Each training sample contains an instruction (task description, including fixed- $k$  vs. dynamic mode), an input that concatenates `<query>` and a numbered `<docs>` block,

Statistic (AmbigQA validation)	Value
Total queries (evaluated)	500
Avg. candidate pool size (docs/query)	22.932
Avg. number of gold short answers (answers/query)	3.754
Training queries excluded (leakage prevention)	3,000

Table 8: Dataset statistics from our preprocessing pipeline (seed=42).

Field	Example (abridged)
instruction	... select exactly 3 documents step by step ...
input	<code>&lt;query&gt;... &lt;/query&gt; \n&lt;docs&gt; \n[1] ... \n[2] ... \n&lt;/docs&gt;</code>
output	<code>&lt;think&gt;...&lt;/think&gt; \n&lt;select&gt;9&lt;/select&gt; ... \n&lt;answer&gt;[9,1,8]&lt;/answer&gt;</code>

Table 9: Training sample schema (single-example, abridged) used for SFT and RL data.

and an output containing interleaved `<think>` and `<select>` tags followed by the final `<answer>` list. Document indices in `<select>` and `<answer>` refer to the numbering of candidates in the `<docs>` block.

## Reference Tables (Appendix B)

We list the corresponding reference tables below. Tables are allowed to float for better packing in the two-column layout.

## C Metrics and Reward Details

### C.1 Evaluation Metrics (Expanded)

**Answer Coverage (Cov@ $k$ ).** As defined in §4, AmbigQA provides multiple gold short answers per query. We compute  $\text{Cov}@k$  via **case-insensitive substring matching** of each gold answer against the concatenated text of the selected top- $k$  documents, then average the match rate across answers and queries.

**NDCG@ $k$ .** We compute  $\text{NDCG}@k$  using each document’s **answer coverage ratio** (fraction of gold short answers matched inside that document) as the graded gain (§4). Let  $g(d) \in [0, 1]$  be the per-document coverage ratio and  $(d_1, \dots, d_k)$  be the ranked list:

$$\text{DCG}@k = \sum_{i=1}^k \frac{g(d_i)}{\log_2(i+1)}, \quad \text{NDCG}@k = \frac{\text{DCG}@k}{\text{IDCG}@k}. \quad (5)$$

Component	Definition (summary)
Format reward ( $r_{\text{fmt}}$ )	Rule-based checks of (i) required tags, (ii) tag pairing, (iii) valid non-duplicate indices, (iv) answer-list format, and (v) fixed/dynamic mode validity. Weight: 30%.
Coverage reward ( $r_{\text{cov}}$ )	Answer-level scoring using exact case-insensitive substring match; otherwise soft lexical-overlap shaping (threshold 0.5) to provide denser reward. Weight: 70%.
Hard constraints/bonuses	Fixed mode length mismatch: $-0.3$ ; perfect format & correct length: $+0.1$ ; empty selection without proper termination: $-0.2$ .

Table 10: Reward function summary (detailed definitions are in §3.3).

**Redundancy Score.** To stratify difficulty (§4), we embed each document with Qwen3-embedding-8B and compute the mean pairwise cosine similarity within the candidate pool.

## C.2 RL Reward vs. Reported Metrics

For GRPO training (§3), we use a shaped reward that can include soft lexical-overlap components to provide denser learning signals. However, **all reported Cov@ $k$  results in the paper are computed using strict case-insensitive substring matching**, without soft matching. A compact summary is provided in Table 10.

## Reference Tables (Appendix C)

### D Additional Case Analysis

We provide additional qualitative examples complementing §5. In general, our model succeeds when (i) at least one candidate contains a direct answer string and (ii) remaining documents provide complementary facets; failures typically arise when one interpretation dominates the pool and the alternative answer is weakly expressed or paraphrased, making substring-based coverage hard to satisfy.

**What we analyze.** For each query, we examine (i) whether a selected document contains a gold short-answer string under case-insensitive matching, and (ii) whether each additional selected document contributes *new* answer facets rather than repeating the dominant interpretation. We qualitatively diagnose failures by checking whether the missing answer is absent from the candidate pool, present only via paraphrase (making substring matching hard), or dominated by pool bias where one interpretation overwhelms the alternatives.

## D.1 Detailed Examples (from the AmbigQA Validation Set)

### D.1.1 Successful Example: Disentangling Ambiguity via Step-wise Selection

**Query:** *Who played george washington in the john adams series?*

**Candidate pool (summary):** multiple documents about the *John Adams* miniseries (Docs [1], [5], [8]), a biography book (Doc [2]), and actor pages (David Morse [9], William Daniels [3]).

This example illustrates why independent similarity scoring is vulnerable to the redundancy trap. The Embedding Sim baseline tends to select Docs [1], [2], and [5] because they all match the surface phrase “John Adams series” while missing the core intent (the actor identity). In contrast, our model explicitly reasons about marginal utility and produces the following step-wise selections (same trace as §5):

- **Step 1:** select [9] (David Morse). *Rationale:* the only candidate that directly states the actor who portrayed George Washington, satisfying the primary information need.
- **Step 2:** select [1] (miniseries overview). *Rationale:* adds non-redundant context about what the *John Adams* series is; complements the casting fact without repeating it.
- **Step 3:** select [8] (broader contextual/genre list). *Rationale:* provides broader context that is not present in the first two selections, improving the completeness of the package while avoiding near-duplicates.

Overall, the ranked list [9, 1, 8] demonstrates the intended behavior of *LLM-as-MMR*: prioritize the direct answer first, then allocate remaining slots to complementary facets rather than redundant paraphrases.

**Additional failure modes (brief).** Beyond the examples in §5, we repeatedly observe two practical failure modes: (i) **paraphrase-heavy pools**, where the correct answer is present but expressed without the exact short-answer string, which hurts substring-based coverage; and (ii) **pool bias**, where one interpretation dominates the candidate set and the minority interpretation appears only in a weak or noisy document, making it hard for the model to justify allocating a top- $k$  slot.

1086 **D.1.2 Failure Example 1: Temporal**  
1087 **Ambiguity (The Simpsons)**

1088 **Query.** *When did the simpsons first air on televi-*  
1089 *sion?*

1090 **Multiple valid answers.** This query has two com-  
1091 mon interpretations: (i) the debut of *The Simpsons*  
1092 shorts (1987) and (ii) the premiere of the full TV  
1093 series (1989). A coverage-oriented reranker should  
1094 ideally surface evidence for both.

1095 **Observed failure pattern.** In our inspection, the  
1096 candidate pool is often dominated by documents  
1097 discussing the 1989 premiere. The model correctly  
1098 identifies many 1989-focused documents as redund-  
1099 ant with each other, but the evidence for the 1987  
1100 shorts is comparatively weak and easy to miss. As  
1101 a result, the model prioritizes the “authoritative”  
1102 1989 date and uses remaining slots on context about  
1103 that event, leaving the secondary interpretation un-  
1104 covered (as described in §5).

1105 **D.1.3 Failure Example 2: Abstract Concepts**  
1106 **(Consubstantial)**

1107 **Query.** *What does consubstantial mean?*

1108 **Multiple valid answer facets.** Gold short an-  
1109 swers span both theological and philosophical def-  
1110 initions (e.g., “of the same being” vs. a broader  
1111 philosophical notion such as “common humanity”).

1112 **Observed failure pattern.** The model tends to  
1113 retrieve the dominant theological definition first.  
1114 While it avoids verbatim redundancy by selecting  
1115 different phrasings of that interpretation, it may  
1116 fail to recognize that a philosophical definition is  
1117 a *distinct valid answer facet* rather than a weaker  
1118 distractor. This reflects a gap between redundancy  
1119 avoidance and answer diversity (as discussed in  
1120 §5).