

# PROVENCE: EFFICIENT AND ROBUST CONTEXT PRUNING FOR RETRIEVAL-AUGMENTED GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Retrieval-Augmented Generation improves various aspects of Large Language Models (LLMs) generation, but suffers from computational overhead caused by long contexts as well as the propagation of irrelevant retrieved information into generated responses. Context pruning deals with both aspects, by removing irrelevant parts of retrieved contexts before LLM generation. Existing context pruning approaches are however limited, and do not provide a universal model that would be both *efficient* and *robust* in a wide range of scenarios, e.g., when contexts contain a variable amount of relevant information or vary in length, or when evaluated on various domains. In this work, we close this gap and introduce `Provence` (for Pruning and Reranking Of retrieVED relevaNt ContExts), an efficient and robust context pruner for Question Answering, which dynamically sets the needed amount of pruning for a given context and can be used out-of-the-box for various domains. The three key ingredients of `Provence` are formulating the context pruning task as sequence labeling, unifying context pruning capabilities with context reranking, and training on diverse data. Our experimental results show that `Provence` enables context pruning with negligible to no drop in performance, in various domains and settings, at almost no cost in a standard RAG pipeline. We also conduct a deeper analysis alongside various ablations to provide insights into training context pruners for future work.

## 1 INTRODUCTION

Retrieval-Augmented Generation (RAG) has become a widely-used paradigm for improving factuality, attribution, and adaptability of Large Language Models (LLMs) (Das et al., 2019; Asai et al., 2024; Seo et al., 2019; Lewis et al., 2020; Mallen et al., 2023a; Min et al., 2023). Augmenting a given user’s query with retrieved relevant contexts helps to avoid the generation of untruthful information and enables the provision of references used to generate the answer. Furthermore, using a domain-specific datastore may enable access and reasoning over a previously unknown knowledge – without fine-tuning the LLM. One additional advantage of the RAG approach is the easy plug-and-play architecture (LangChain): practitioners may choose components (retrievers, generator LLMs, context granularity etc.) which best suit their particular cases to maximize the final performance.

At the same time, the use of RAG adds *computational overhead* due to both retrieval latency and the increased input length for the LLMs. It may also propagate *irrelevant information* present in retrieved contexts into generated responses. These issues can be solved by developing more efficient and robust LLMs – either by making architectural changes to process long contexts more efficiently (Nawrot et al., 2024; Dao, 2024; Chevalier et al., 2023) or increasing the diversity of the tuning data to improve processing of irrelevant contexts (Lin et al., 2024). However, tuning the LLM can be highly resource-consuming, or even impossible to apply for proprietary (closed) LLMs. An alternative solution consists in *pruning retrieved contexts* by removing context parts irrelevant to the user’s query – which reduces context lengths and therefore speeds up generation. Such context pruning module can be used in a *plug-and-play manner with any generator LLM*, featuring both easy use and better transparency in the RAG pipeline.

Despite initial efforts on developing context pruners for RAG, none of the existing solutions provide a model ready to be used *out-of-the-box* in practice. First, many approaches are designed for a simplified setting, e.g., with the assumption that only one sentence per context is relevant to the

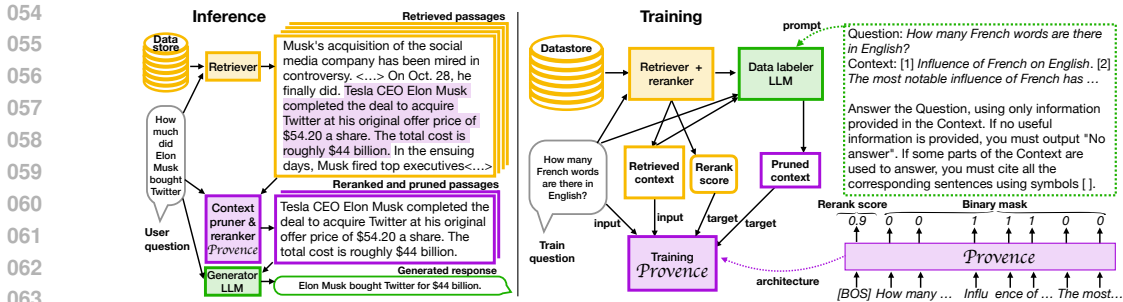


Figure 1: Illustration of inference (left) and training (right) of Provenance.

Table 1: Analysis of existing approaches for context pruning. Violet / Orange highlight practical / less-practical solutions.

Approach	Query-dep.	Granularity	Type	Output	Base arch.	Multi-domain testing	Model re-release
Selective Context	No	token-level	extr.	% of tokens	Llama-7B / GPT2	Yes	Yes
LLMLingua	No	token-level	extr.	% of tokens	Alpaca-7B / GPT2	Yes	Yes
LongLLMLingua	Yes	token-level	extr.	% of tokens	Llama-2-7B-chat	Yes	Yes
LLMLingua2	No	token-level	extr.	% of tokens	RoBERTa / mBERT	Yes	Yes
RECOMP extr.	Yes	sent.-level	extr.	1 sentence	BERT	No	Yes
RECOMP abstr.	Yes	sent.-level	abstr.	$\geq 0$ sentences	T5-L	No	Yes
FilCo	Yes	sent.-level	abstr.	1 sentence	T5-XL / Llama-2-7B	No	No
COMPACT	Yes	sent.-level	abstr.	$\geq 0$ sentences	Mistral-7B	No	Yes
Provenance (ours)	Yes	sent.-level	extr.	$\geq 0$ sentences	DeBERTa	Yes	Yes

input query (Wang et al., 2023; Xu et al., 2024), or that the compression ratio is fixed (Jiang et al., 2023; Pan et al., 2024). However, in practice contexts may contain various portions of relevant information, from empty to full relevant context, and pruners should detect it in an *adaptable* fashion. Second, many works introduce context pruners that are not efficient enough to be used in practice. This includes using billion-sized LLMs as base models for pruners (Jiang et al., 2024; Pan et al., 2024; Wang et al., 2023), or designing abstractive context compressors which require sequential autoregressive generation of the final context (Wang et al., 2023; Xu et al., 2024). We argue that a more practical and *efficient* setting consists in fine-tuning a *small-size model* such as DeBERTa (He et al., 2021b;a), as an *extractive* pruner, i.e., with a lightweight prediction head for selecting relevant context parts. Third, most of the existing works train context pruners for each dataset individually and do not target nor test pruners *robustness* to various data domains.

Table 1 summarizes the properties of various existing methods along specified dimensions and shows that none of them satisfy all listed criteria. The table also includes a dimension of pruning granularity, i.e., token-level vs sentence-level pruning. In this work, we focus on *query-dependent sentence-level* pruning, which prunes out semantic units (sentences) that are deemed not relevant to generate the answer. An alternative approach is token-level pruning which prunes out low-level grammatical units such as articles or interjections, usually in a query-independent fashion. The two approaches are orthogonal and could potentially be combined.

To address listed limitations, we introduce Provenance (Pruning and Reranking Of retrieved relevant Context), an *adaptable*, *efficient* and *robust* sentence-level context pruner for Question Answering, which can be used *out-of-the-box* across various domains and settings. To achieve this, we formulate context pruning as *binary sequence labeling* so that the binary mask predicted by the pruner determines *sentences* (from zero to all) which are relevant to the query, and train our pruner from a lightweight DeBERTa model on diverse data. Furthermore, we notice that context pruning and reranking (i.e., the second step in effective retrieval pipelines) bear a strong resemblance. We therefore propose to **unify these two models into a single one**, completely **eliminating the cost** of context pruning in the RAG pipeline.

More specifically, our contributions are as follows:

- We propose an approach for training an *adaptable*, *robust*, and *efficient* context pruner for QA – and will release our trained models. Three key ingredients of our approach are formulating context pruning as sequence labeling, unifying context pruning and reranking in a single model, and training on diverse data.
- We test `Provenance` on various QA domains and show its *out-of-the-box applicability* to prune contexts with negligible to no drop in performance and at almost no cost, substantially outperforming baseline approaches. We also demonstrate `Provenance` capabilities in detecting the number of relevant sentences at any positions in the context and robustness to various context lengths.
- We conduct multiple ablations to demonstrate which techniques are essential for training robust context pruners, to provide insights for future context pruners development.

**Definitions.** A typical RAG pipeline consists of (0) a user’s question, or query; (1) a *datastore*, i.e., a collection of *documents* (pieces of text) to be retrieved from, (2) an efficient retriever which enables fast retrieval from a large datastore (typically a dual-encoder model, where queries and passages are encoded independently), (3) a more expensive cross-encoder reranker which further reduces and reorders a set of retrieved passages (cross-encoding means encoding a passage together with a query); and (4) a generator LLM which outputs the final response based on the user’s query and the relevant passages. Such a pipeline can be represented as `retrieve >> rerank >> generate`. Context pruning can be incorporated before generation, i.e., `retrieve >> rerank >> prune >> generate`. In our work, we also propose to incorporate context pruning into reranking, an essential and already present component in RAG (Rau et al., 2024a): `retrieve >> rerank+prune >> generate`. This enables **context pruning at almost zero cost**.

## 2 RELATED WORK

**Context pruning.** RECOMP (Xu et al., 2024) focuses on context pruning for RAG and proposes both extractive and abstractive context pruners. The extractive RECOMP approach independently encodes sentences in the context and then selects top sentences with embeddings closest to the query embedding. Such an approach limits context understanding, due to independent processing of both sentences and queries. The method also requires specifying the amount of sentences to keep as a hyperparameter – which is usually unknown in practice and should depend on each particular passage. The abstractive RECOMP summarizes key information from the passage relevant to the query (including zero relevant information) by training on silver summaries generated by GPT-3.5. However, it requires inefficient autoregressive generation of the final context, and can eventually hallucinate facts not present in the input context. FilCo (Wang et al., 2023) similarly proposes to generate contexts autoregressively but is trained on extractive targets, i.e., one sentence from the context selected by one of three criteria. The drawbacks are again inefficiency and the simplified assumption of one relevant sentence per context. A recent approach, COMPACT (Yoon et al., 2024), also proposes to generate filtered contexts autoregressively – hence inefficiently – and introduces an iterative approach for gradually updating the relevant context after processing a new portion of retrieved passages. In contrast to all listed efforts, `Provenance` dynamically detects the amount of relevant information in the context – from zero to all sentences – in an extractive and efficient way. Furthermore, we propose a novel approach of integrating context pruning into a reranker.

Concurrently to our work, DSLR (Hwang et al., 2024) performs extractive sentence-level pruning, by encoding sentences one-by-one, together with the query, using existing rerankers. Similarly to `Provenance`, DSLR keeps sentences with scores higher than a threshold and preserves the original order of sentences. However, in contrast to `Provenance`, DSLR is not capable of keeping groups of semantically connected sentences, due to independent sentence processing.

An orthogonal line of work proposes extractive token-level pruners. LLMLingua (Jiang et al., 2023) and Selective Context (Li et al., 2023) use LLMs to remove tokens with high generation probabilities, independently of the query. LLMLingua2 (Pan et al., 2024) is a small BERT-based model finetuned to eliminate redundant tokens, also independently of the query. LongLLMLingua (Jiang et al., 2024) proposes query-dependent LLM-based token pruning based on contrastive perplexity. Listed approaches remove tokens in a way that it does not break context understanding for the LLM – hence they are not capable of removing semantic parts of the context. LLMLingua models also

162 have many hyperparameters in the interface which are hard to tune in practice. These approaches  
 163 can however also be combined with sentence-level pruning.

164 **Retrieval granularity.** Alternatively to context pruning, one can reformulate datastore content into  
 165 atomic units, e.g., *propositions* as in Dense-X retrieval Chen et al. (2024c) or *decontextualized sen-*  
 166 *tences* (Choi et al., 2021). Such preprocessing is expensive and can lead to some information loss.

167 **Passage filtering.** Another related – and orthogonal – line of works focuses on filtering entire pas-  
 168 sages if they are deemed irrelevant for a given question; such an approach can be straightforwardly  
 169 combined with `Provenance`. A simple method consists in introducing a threshold on the (re)ranking  
 170 score. LongLLMLingua reranks passages based on the probability of a question given the passage.  
 171 (Yoran et al., 2024) use natural language inference models to filter out passages that do not entail  
 172 question-answer pairs, but report that this approach sometimes filters out relevant passages too.

173 **Improving context processing in LLMs.** While context pruners aim to remove context parts ir-  
 174 relevant to the user’s query, another line of work aims to process contexts more efficiently and  
 175 effectively in LLMs. Efficient context processing could be achieved through efficient attention im-  
 176 plementations (Dao, 2024; Anagnostidis et al., 2023), KV cache compression (Nawrot et al., 2024),  
 177 encoding retrieved passages in parallel (Zhu et al., 2024), or compressing contexts into one or more  
 178 context embeddings (Chevalier et al., 2023; Ge et al., 2024; Rau et al., 2024b). Other works aim to  
 179 make LLMs more robust, by exposing them to noisy contexts during training or finetuning (Izac-  
 180 ard et al., 2022; Lin et al., 2024). All such approaches usually require LLM adaptation which may  
 181 complicate application to an arbitrary picked LLM.

### 182 3 PROVENANCE

183 The high-level overview of our proposed approach is illustrated in Figure 1. Our first contribution is  
 184 to pose the context pruning problem as a sequence labeling task. We fine-tune a DeBERTa model to  
 185 encode the query–context pair and output binary masks which are used to filter out irrelevant context  
 186 parts. The labels for training are generated by Llama-3-8B-Instruct (AI@Meta, 2024); **we call them**  
 187 **silver labels since they are generated automatically**. Such an approach solves several limitations of  
 188 existing context pruners: (1) by construction, the model is able to deal with varying noise in contexts  
 189 and select an appropriate pruning ratio; (2) queries are encoded together with context sentences  
 190 (cross-encoding), providing richer representations – compared for instance to extractive RECOMP  
 191 which encodes query and context sentences independently; (3) using a lightweight encoder makes  
 192 our approach more efficient than LLM-based or abstractive methods.

193 Our second contribution consists in unifying reranking and context pruning – instead of considering  
 194 these steps as distinct in the RAG pipeline. In `Provenance`, reranking and pruning can be done *in a*  
 195 *single forward step*, thus eliminating the computational overhead due to context pruning – making  
 196 `Provenance` almost “free”.

197 **Training data.** Our approach requires a set of training questions and a retrieval datastore. Spefi-  
 198 cially, we rely on the train set of the MS MARCO document ranking collection which includes 370k  
 199 queries (Nguyen et al., 2016). The MS MARCO collection is a domain-diverse datastore of 3.2M  
 200 documents crawled from the Web – which is required for the final model’s robustness to various do-  
 201 mains – and is often used to train retrievers and rerankers. We also consider the train set of Natural  
 202 Questions which contains 87k queries Kwiatkowski et al., 2019).

203 **Data processing.** We create a retrieval datastore by splitting MS MARCO documents into passages  
 204 consisting of  $N$  consecutive sentences –  $N$  being a random integer  $\in 1..10$ . This is to enable the  
 205 pruner’s robustness to variable retrieved context lengths. We also prepend page titles to each pas-  
 206 sage. For each question, we retrieve top-5 relevant passages using a strong retrieval pipeline (Rau  
 207 et al., 2024a) consisting of a SPLADE-v3 retriever (Lassance et al., 2024) and a DeBERTa-v3  
 208 reranker (Lassance & Clinchant, 2023). The resulting set of retrieved passages is naturally diverse  
 209 w.r.t. relevance or irrelevance to the question, due to imperfections in retrieval.

210 **Silver labels generation.** Given a question and a retrieved passage (context), we split the passage  
 211 into sentences<sup>1</sup> and prompt Llama-3-8B-Instruct to select sentences relevant to the given question.

212 <sup>1</sup>using the `nltk.sent_tokenize` function: [https://www.nltk.org/api/nltk.tokenize.  
 213 sent\\_tokenize.html](https://www.nltk.org/api/nltk.tokenize.sent_tokenize.html)

One approach would be to use a straightforward prompt such as “*Output indexes of sentences relevant to the given question*”. However, we decided to utilize the strong LLMs’ capabilities of actually *answering* questions while *citing* relevant context sentences. We therefore instruct the LLM to answer the given question using *only* information provided in the given context, and output “*No answer*” in case no relevant information is provided. We also specify the easy-to-parse citation format [i] and number sentences with the same marker in the context. Our prompt can be found in Appendix – Table 6; we use greedy decoding and parse cited sentences using regular expressions. We also compare different prompting strategies in the ablation study.

We found that Llama-3-8B is well capable of answering only based on a given context in most cases and of outputting a citation  $\sim 90\%$  of the time. We filter out cases when no citations are produced and “*No answer*” is not present in the LLM’s output, as these are the cases when the context actually contains relevant information but the LLM “forgot” to cite it. The final labels distribution (number of selected sentences per context, their positions) is shown in Appendix – Figure 5.

**Training of Provence.** Our context pruner receives as input the concatenation of a question and a retrieved context, and outputs *per-token binary labels* denoting whether each token (defined by the pretrained model’s tokenizer) should be included in the selected context. In Section 4.4 (Ablations), we also consider an approach where a special token is inserted at the beginning of each sentence, and labels are predicted per-sentence based on the representations of those tokens. We train `Provence` as a binary per-token classifier with ground truth labels coming from the silver data labeling, and the model can be used as a standalone pruner, i.e., `retrieve >> rerank >> Provence (standalone) >> generate`.

**Unifying compression and reranking.** We note that cross-encoder rerankers (Nogueira & Cho, 2020) share both the same architecture and inputs (pairs of question–passages) as `Provence`. Additionally, the task of context pruning (selecting parts of contexts that are useful for generating the answer to the question) intrinsically bears similarity with re-ranking (estimating the relevance of a context w.r.t. the question) – and we hypothesize the possibility of knowledge transfer between these two related tasks. We therefore propose to *unify* both approaches in a single model, with two different task heads. More specifically, the reranking head outputs a scalar prediction for the BOS token while the pruning head outputs per-token predictions for the passage tokens, as illustrated in Figure 1. To ease training, we propose to further fine-tune a pretrained reranker on our labeling objective, while adding a ranking “regularizer” to preserve initial reranking capabilities. The regularizer is a Mean Squared Error loss on the reranking scores from the initial reranker. This can be viewed as a straightforward pointwise score distillation process, where the initial model serves as the teacher – a method that has demonstrated great effectiveness in Information Retrieval Hofstätter et al. (2021). The final loss function is as follows:

$$\mathcal{L} = \sum_{n=1}^N \left\{ \sum_{k=1}^{L_n} \log P(y_{n,k} | z_{n,k}) + \lambda (s_n - z_{n,0})^2 \right\} \quad z_n = \text{Provence}(x_n) \quad (1)$$

where  $N$  is the number of datapoints (query–passage pairs),  $x_n$  is a sequence of  $L_n + 1$  input tokens (concatenated query, passage and BOS at the 0-th position),  $z_n$  is a sequence of  $L_n + 1$  predictions output by the model,  $y_n$  is a sequence of  $L_n$  target binary labels for context pruning,  $s_n$  is the teacher score (initial reranker),  $z_{n,0}$  is the ranking score predicted from the BOS representation.

In the case of the unified model, re-ranking and context pruning need a single forward step from the encoder, i.e., `retrieve >> Provence (w/ re-ranking) >> generate – making context pruning almost free in terms of execution time`.

**Inference with Provence.** At inference, we feed a concatenation of a question and a retrieved passage through `Provence`, which outputs probabilities of including each token in the final context, as well as the passage score in the case of the unified model. We simply use a threshold  $T$  to binarize the token probabilities (keep or not) – **which has a direct effect on the compression rate**. As shown in the experiments Section, the choice of a threshold is generally transferable across various datasets, making the model flexible to be used out-of-the box in various QA applications<sup>2</sup>.

We note that our model outputs token-level predictions despite the sentence-level labeling task. We found that probabilities of including tokens into the final context are naturally clustered on

<sup>2</sup>Note that tuning the threshold per dataset could of course further improve results.

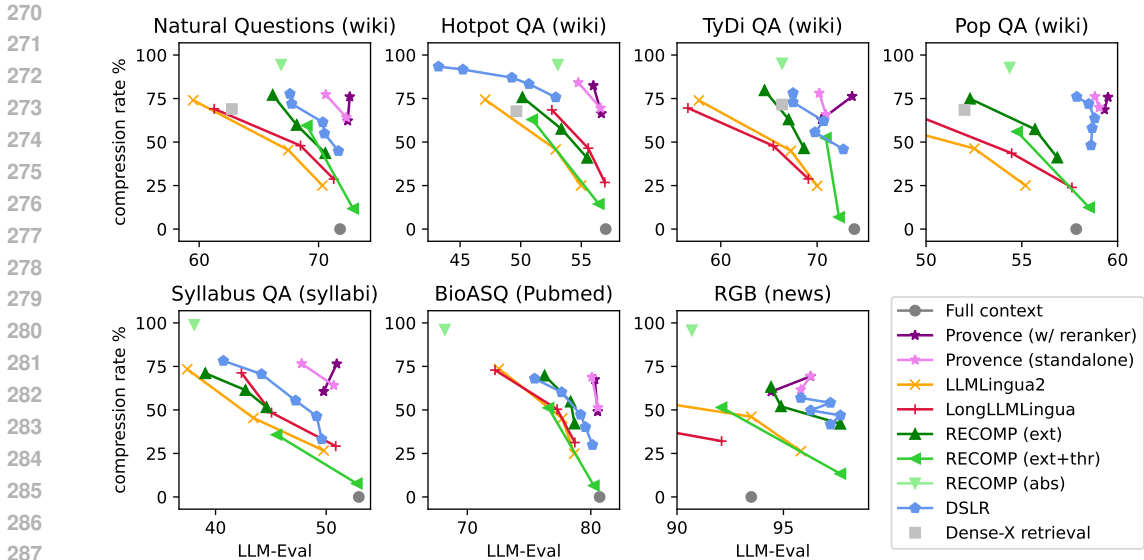


Figure 2: Main results for various QA domains, comparing `Provenance` and baseline models. *Generator*: Llama-2-7B, *retriever*: SPLADE-v3, *reranker*: DeBERTa-v3 (or `Provenance` in the unified setting). Plot titles denote “Dataset name (datastore type)”. *x*-axis denotes QA performance evaluated with LLM-as-a-judge; *y*-axis denotes the context compression ratio. For both metrics, the higher the better: the best model would be closest to the top right corner. Numerical scores are presented in App. Tables 8–9. Main conclusion: `Provenance` consistently lies on the Pareto front.

the sentence level – see example in Appendix Figure 6 – due to the sentence-level targets used in training. However, in rare cases we could still have partial sentences being selected. To avoid this phenomenon, we apply a “sentence rounding” procedure: for each sentence, we check the ratio of kept tokens (predicted label= 1), and select the entire sentence only if it is higher than 0.5.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

**Provenance training details.** We train `Provenance` on the data described in Section 3, using PyTorch (Paszke et al., 2019) and HuggingFace transformers (Wolf et al., 2020). We use DeBERTa-v3 (He et al., 2021a) as our pretrained model for training the standalone `Provenance`. For the unified approach, we start training from an already trained cross-encoder, also based on DeBERTa-v3 (Lassance & Clinchant, 2023). Note that in the latter, we initialize the ranking head from its fine-tuned version, and train the separate pruning head from scratch.

After preliminary experiments, we set the learning rate to  $3 \times 10^{-6}$ , the batch size to 48 and train models for one epoch. For joint training, there is a slight trade-off between pruning and reranking. We set the reranking regularization coefficient  $\lambda$  to 0.05, chosen as the minimal value that does not substantially degrade reranking performance on the MS MARCO development set.

**Evaluation datasets.** We test `Provenance` on a diverse set of QA datasets. First, we consider commonly used datasets relying on Wikipedia datastore: Natural Questions (Kwiatkowski et al., 2019), TyDi QA (Clark et al., 2020), PopQA (Mallen et al., 2023b) (all single-hop questions), and HotpotQA (Yang et al., 2018) (multi-hop questions). Second, we consider datasets with datastores from various domains: BioASQ (Nentidis et al., 2023) (biomedical questions with Pubmed as a datastore), SyllabusQA (Fernandez et al., 2024) (questions about educational course logistics, with courses syllabus as a datastore); and RGB Chen et al. (2024b) (questions about news with Google-searched news as contexts). Further details can be found in Appendix A.

Table 2: Time/MFLOPS required for context pruning. Top-5 retrieved documents, NQ dev set (3k samples).

Pruner	Time (s)	MFLOPS
LongLLMLingua, rate=0.5	2649	$122 \times 10^9$
LLMLingua2, rate=0.5	863	$8 \times 10^9$
RECOMP extr., top=2	351	$1.2 \times 10^9$
RECOMP abstr.	1056	$2.2 \times 10^9$
Providence	471	$4.8 \times 10^9$

Table 3: Speed up in generation due to compression (Providence, 49% compression). Batch sizes 1 or 256.

Generator	bs 1	bs 256
LLama-2-7B	$\times 1.2$	$\times 2$
LLama-2-13B	$\times 1.4$	$\times 2$
SOLAR-10.7B	$\times 1.4$	$\times 1.9$

**Evaluation settings.** We conduct experiments using BERGEN (Rau et al., 2024a), a benchmarking library for RAG, using the recommended experimental setting. For each query, we retrieve top-5 relevant passages using a strong and robust retrieval pipeline: SPLADE-v3 >> DeBERTa-v3 reranker (except for RGB, for which Google-searched passages are already provided). We then pass queries prepended with relevant document (full length or pruned) into LLama-2-7B-chat (Touvron et al., 2023)<sup>3</sup> to generate answers; other RAG settings are further reported in Appendix. Each evaluation dataset comes with short keyword answers, which we use to evaluate responses using LLM-based evaluation (LLMeval in Rau et al., 2024a); match-based metrics are also reported in Appendix. We additionally measure compression as a portion of the context which was pruned out.

We compare Providence to publicly available context pruning models listed in Table 1, except LLMLingua and Selective Context which were shown to underperform LLMLingua2 (Pan et al., 2024). For all context pruners (except abstractive RECOMP for which it is not available), we enforce the selection of the first (title) passage sentence, to avoid ambiguity in understanding the context by the generator. For extractive RECOMP, we use the model trained on NQ, consider using top-1/2/3 sentences, and prepend the passage title to each sentence. For the LLMLingua family, we vary the compression rate in  $\{0.25, 0.5, 0.75\}$  and use code provided on the official repository<sup>4</sup>. We use the XLM-RoBERTa model for LLMLingua2. For Providence, we use  $T = 0.1$  and  $T = 0.5$ . We also compare our method to DSLR based on the same reranker as ours, i.e. DeBERTa-v3.

## 4.2 MAIN RESULTS

Context pruners are often only tested on limited domain data, e.g., with Wikipedia datastore, and an important aspect of our work is evaluating context pruning on a series of QA domains. Figure 2 reports the **trade-off between compression (efficiency) and LLM-evaluated performance (quality)**, for various QA datasets and context pruning methods. We choose to report a figure per dataset to better assess the **Pareto front** of existing solutions, rather than comparing methods with different compression rates in the same table. Figure 7 in Appendix further reports similar results with match-based metric, and Appendix Tables 11–13 show examples of context pruning with various methods.

First, we observe that Providence achieves **the highest performance** across pruning methods, for similar compression ratios. Second, it is noteworthy that Providence outperforms methods requiring more computations such as LLMLingua models, showing that efficiency is not traded for effectiveness. Furthermore, Providence is the only method capable of achieving high compression levels without (or with negligible) performance drops, on all datasets. Moreover, for some datasets, e.g., PopQA, pruning with Providence leads to performance improvements due to noise filtering.

**The effect of threshold.** An important aspect in the out-of-the-box applicability of context pruners is how much effort is needed to select the suitable values of hyperparameters. For Providence, it only consists in setting the pruning threshold  $T$ . In Figure 2 (for which  $T = 0.1$  and  $T = 0.5$ ), we observe that Providence pruning ratio automatically varies from 50% to 80%, depending on the dataset, which demonstrates that the same values for  $T$  work well for all considered domains – making Providence robust to the choice of hyperparameters. If necessary, users can still tune it further for their datasets and/or needs. We note that some models specify the desired compression ratio as a hyperparameter, e.g., LLMLingua models or extractive RECOMP (through top- $N$  sentences).

<sup>3</sup>For main experiments, we chose a “weaker” generator which relies more on contexts, to create a more challenging setting for context pruners; results with stronger generators are reported in Appendix – Figure 8.

<sup>4</sup><https://github.com/microsoft/LLMLingua>

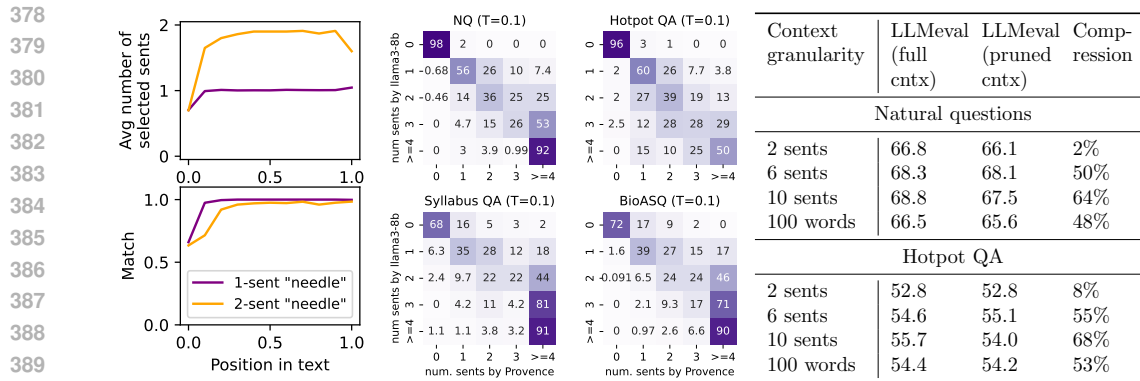


Figure 3: Analyses. (Left) Needle-in-the-haystack test allowing the control of the position of the ground truth sentence(s) in the context. (Middle) Comparison of the number of selected sentences by the silver predictor (LLaMA-3-8B-Instruct) and `Provenance`. Heatmaps are normalized by rows: a cell in position  $(i, j)$  indicates which percentage of contexts that were pruned into  $i$  sentences by the silver predictor, were pruned into  $j$  sentences by `Provenance`. (Right) Testing `Provenance` in settings with different context lengths. All experiments are done with unified `Provenance`,  $T = 0.1$ .

While it may seem convenient to estimate inference cost, the “optimal” compression ratio (without losing performance) is specific to each particular question-context pair. Thus, using a threshold as a hyperparameter is more appropriate for this task. We also experimented with specifying a threshold in extractive RECOMP (shown on the same plot) and found that it often leads to lower performance (compared to top- $N$ ). The reason is that different queries have different ranges of similarity scores.

**Efficiency.** We compare `Provenance` with other pruning methods in terms of efficiency. Table 2 reports compression time and MFLOPS<sup>5</sup> required by different pruning methods. As expected, LongLLMLingua (based on LLaMA-2-7B-chat) is the slowest context pruner. RECOMP abstr. requires less MFLOPS compared to `Provenance`, but its autoregressive nature makes it slower in practice<sup>6</sup>. Note that in the case of the unified model, pruning is almost free – as it’s part of the re-ranking step. Table 3 reports speed-up gains due to compression with `Provenance` model ( $\sim 50\%$  compression rate). All runs were performed on single Tesla V100-SXM2-32GB GPU with vllm Kwon et al. (2023). With large batch sizes, we systematically observe  $2\times$  speed-ups at inference, while smaller batch sizes lead to lower gains (especially for smaller models). We assume this is mostly due to the CPU/GPU communication bottleneck, which masks inference gains due to compression.

### 4.3 ANALYSIS

In this Section, we conduct a more fine-grained evaluation to better understand the properties of `Provenance`.

**Robustness to the position of relevant information in the context.** We design a needle-in-the-haystack experiment which allows us to check the performance of `Provenance` on a simple toy example and to evaluate its robustness w.r.t. the position of the relevant information in the input context. We write 5 questions and answers<sup>7</sup>, and insert answers (“needles”) at random positions between sentences, in a subset of 100 passages sampled from the Wikipedia datastore. Ideally, `Provenance` should only select the “needle” sentences and filter out all other sentences in contexts. We plot the number of selected sentences and percentage of cases when the pruned context contains the “needle” (Figure 3, (Left)). We consider two settings: with 1- and 2-sentence “needles”. We observe that `Provenance` correctly selects “needle” sentence(s) in most cases, except at leftmost and

<sup>5</sup>We use the PyTorch profiler to report FLOPS required by each pruner.

<sup>6</sup>This highlights the fact that MFLOPS do not always align with real inference time, due to different architectural choices.

<sup>7</sup>Example: “Which library was used in the experiments?”, answer: “Experiments were conducted using the Bergen library”. Example reformulation into a 2-sentence answer: “Experiments were conducted using a library. Its name is Bergen.”



Table 4: Effectiveness of reranking top-50 documents retrieved by SPLADE-v3. DeBERTa-v3 is the “baseline” (initialization point for `Providence`, which we aim to preserve performance). We report the R@5 on two RAG datasets (NQ and HotpotQA), MRR@10 on MS MARCO passages (dev set), nDCG@10 on TREC DL’19 (Craswell et al., 2020), and mean nDCG@10 on the 13 open datasets from the BEIR benchmark (Thakur et al., 2021) – Table 7 in Appendix reports the full results.

Model	Dataset				
	NQ	HotpotQA	MS	TREC19	BEIR
DeBERTa-v3	83.0	70.4	40.5	77.4	55.4
<code>Providence</code>	84.4	70.5	40.6	77.2	55.9
<code>Providence</code> (NQ)	84.5	70.3	40.2	77.5	55.1

rightmost positions.<sup>8</sup> In most cases `Providence` does not select any irrelevant sentences. The results are similar for both simpler (1-sentence) and harder (2-sentence) “needles” showing `Providence`’s flexibility in detecting the number of relevant sentences, discussed below in more details.

**Adaptability to the variable number of relevant sentences.** To evaluate the capability of `Providence` to dynamically detect the number of relevant sentences in the context, we compare the number of sentences  $L$  selected by `Providence` and by a silver oracle, for question-context examples from various datasets. A silver oracle is easy to construct for  $L = 0$ , by pairing questions with randomly sampled contexts. For  $L \geq 1$ , we use the labeling produced by Llama-3-8B-Instruct. Figure 3 (*Middle*) shows that the number of relevant sentences detected by `Providence` is close to the silver oracle value in most cases, for all considered datasets. In contrast, extractive RECOMP would always select a prespecified number of sentences.

**Robustness w.r.t. context granularity.** Figure 3 (*Right*) shows `Providence` performance for two datasets, with Wikipedia datstores made of contexts of various granularity. Here, each considered datstore is produced by splitting Wikipedia pages into chunks of  $N$  sentences,  $N \in \{2, 6, 10\}$ , or 100 words, and prepending the page title to each chunk. `Providence` shows high performance in all cases – the performance with pruned contexts being close to the performance obtained using original contexts. As could be expected, the compression ratio is higher for longer contexts.

**Reranking effectiveness.** Table 4 compares **reranking** performance between our reranking baseline and unified `Providence` – whose training starts from the former. We can see that our joint training procedure (on both pruning and ranking tasks) makes it possible to learn a context pruner that preserves initial reranking capabilities. We further include as a comparison point results from a model trained in similar conditions on NQ. Overall, results are similar – further highlighting the robustness of `Providence` w.r.t. training data. We further discuss such aspects in Section 4.4 (Ablations).

**Applicability in different settings.** Figure 8 (App.) demonstrates the applicability of `Providence` in variable retrieval-generator settings – achieving similar results as the ones reported in Figure 2.

#### 4.4 ABLATIONS

In this Section we analyze various design choices made in `Providence` development, to provide insights into training context pruners for future works (results shown in Figure 4). All models in this section are standalone context pruners, trained with the same amount of parameter updates.

**Model size.** We first observe that DeBERTa-large slightly increases the compression rate – when compared to DeBERTa-base. All other ablations are tuned from a DeBERTa-base model, for efficiency reasons. Note that the final `Providence` is trained from a DeBERTa-large model (or its equivalent reranker).

**Data mixtures.** We compare training on NQ (87k queries), MS MARCO downsampled to the same size, and full MS MARCO (370k queries). Despite the observation that using the MS MARCO *type of data* leads to lower results than NQ – with equal number of queries – we also find that using larger

<sup>8</sup>The reason for the drops in the left-most and right-most positions is that training data has little examples of the corresponding types of relevant sentences, see e.g., statistics for the rightmost position in the App. Figure 5, (*Right*). We plan to work on further improving processing of these positions in future work.

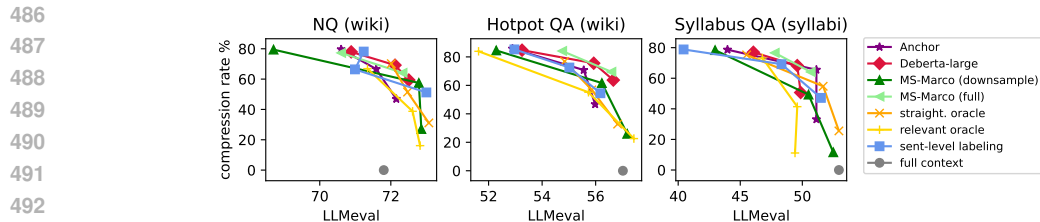


Figure 4: Ablation results. All models are single-component modifications of the anchor model, which is a base-size model, trained on NQ data, with the answer oracle and token-level labeling. Numeric scores for this figure are duplicated in Appendix Table 10, and results with match-based metrics are presented in Appendix – Figure 11.

data (i.e., full MS MARCO) improves results. Our final models are trained on the full MS MARCO – further ablations are conducted on the NQ data, for efficiency reasons.

**Labeling strategies.** As described in Section 3, we can train the pruner either to perform token-level labeling (with sentence rounding at inference) or to perform sentence-level labeling. In the former case sentence representations are richer but the model also needs to learn to output similar predictions for tokens inside one sentence. In the latter case sentence content must be represented in a single embedding which may limit representation expressivity. In practice we observe close performance, with the token-level strategy slightly outperforming the sentence-level one some datasets. In all other experiments we use the token-level strategy.

**Oracle prompts.** We compare three options for prompting an oracle LLM to generate silver labeling: (1) *answer oracle*: asking to answer the given question from the given context, citing corresponding sentences; (2) *relevance oracle*: asking to list any *relevant* information in the context to the question, citing corresponding sentences; (3) *straightforward oracle*: asking to output indexes of sentences which answer the given question. We found that the behavior of the *straightforward oracle* varies on different prompts, while the use of the *answer oracle* makes answers more consistent. The motivation for the *relevance oracle* is that often contexts contain distantly relevant information to the query and it could be reasonable to select the corresponding sentences. Comparing the listed prompts, we observe that the *relevance oracle* underperforms the *answer oracle*, and the *straightforward oracle* performs similarly or slightly lower than the *answer oracle*.

**Unification with reranker.** In Figure 2 we compare `Providence` trained as a standalone model and as a model unified with reranker, and find that both strategies lead to similar results – [although the former relies on two separate inference steps \(re-ranking and pruning\) in a RAG pipeline.](#)

## 5 CONCLUSION

In this work, we present `Providence`, a robust, adaptable, and efficient context pruner for Question Answering – either unified in a single model with reranking capabilities or available as a lightweight standalone model. In contrast to previous extractive approaches, `Providence` dynamically detects the needed pruning ratio for a given context and can be used out-of-the-box for various QA domains. In extensive experiments, we demonstrate that `Providence` prunes contexts with negligible to no drops in performance and in some cases even brings performance improvement due to removing context noise. We also show `Providence` capabilities in correctly detecting the number of relevant sentences in contexts, located at any position, and with contexts of various lengths. Finally, the ablation study highlights the importance of using a large training data and the appropriate prompt in the silver oracle.

**Limitations.** Despite `Providence` being ready to use in various settings, demonstrated in the paper, it is focusing only on QA applications, with a single passage processed at a time, and is trained on English-only data. Future work could consider extending it to other tasks, multi-passage contexts, and languages beyond English.

## REFERENCES

- 540  
541  
542 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md)  
543 [llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 544  
545 Sotiris Anagnostidis, Dario Pavllo, Luca Biggio, Lorenzo Noci, Aurélien Lucchi, and Thomas  
546 Hofmann. Dynamic context pruning for efficient and interpretable autoregressive transform-  
547 ers. *ArXiv*, abs/2305.15805, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:258888224)  
548 [CorpusID:258888224](https://api.semanticscholar.org/CorpusID:258888224).
- 549  
550 Akari Asai, Zexuan Zhong, Danqi Chen, Pang Wei Koh, Luke Zettlemoyer, Hannaneh Hajishirzi,  
551 and Wen-tau Yih. Reliable, adaptable, and attributable language models with retrieval. *arXiv*  
552 *preprint arXiv:2403.03187*, 2024.
- 553  
554 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding:  
555 Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge dis-  
556 tillation, 2024a. URL <https://arxiv.org/abs/2402.03216>.
- 557  
558 Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. Benchmarking large language models in  
559 retrieval-augmented generation. *Proceedings of the AAAI Conference on Artificial Intelligence*,  
560 38(16):17754–17762, Mar. 2024b. doi: 10.1609/aaai.v38i16.29728. URL [https://ojs.](https://ojs.aaai.org/index.php/AAAI/article/view/29728)  
561 [aaai.org/index.php/AAAI/article/view/29728](https://ojs.aaai.org/index.php/AAAI/article/view/29728).
- 562  
563 Tong Chen, Hongwei Wang, Sihao Chen, Wenhao Yu, Kaixin Ma, Xinran Zhao, Hongming Zhang,  
564 and Dong Yu. Dense X retrieval: What retrieval granularity should we use? In Yaser Al-Onaizan,  
565 Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical*  
566 *Methods in Natural Language Processing*, pp. 15159–15177, Miami, Florida, USA, November  
567 2024c. Association for Computational Linguistics. URL [https://aclanthology.org/](https://aclanthology.org/2024.emnlp-main.845)  
568 [2024.emnlp-main.845](https://aclanthology.org/2024.emnlp-main.845).
- 569  
570 Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. Adapting language models  
571 to compress contexts. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of*  
572 *the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 3829–3846,  
573 Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.  
574 [emnlp-main.232](https://aclanthology.org/2023.emnlp-main.232). URL <https://aclanthology.org/2023.emnlp-main.232>.
- 575  
576 Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael  
577 Collins. Decontextualization: Making sentences stand-alone. *Transactions of the Association*  
578 *for Computational Linguistics*, 9:447–461, 2021. doi: 10.1162/tacl.a.00377. URL [https://](https://aclanthology.org/2021.tacl-1.27)  
579 [aclanthology.org/2021.tacl-1.27](https://aclanthology.org/2021.tacl-1.27).
- 580  
581 Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev,  
582 and Jennimaria Palomaki. TyDi QA: A benchmark for information-seeking question answering in  
583 typologically diverse languages. *Transactions of the Association for Computational Linguistics*,  
584 8:454–470, 2020. doi: 10.1162/tacl.a.00317. URL [https://aclanthology.org/2020.](https://aclanthology.org/2020.tacl-1.30)  
585 [tacl-1.30](https://aclanthology.org/2020.tacl-1.30).
- 586  
587 Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. Overview  
588 of the trec 2019 deep learning track, 2020. URL <https://arxiv.org/abs/2003.07820>.
- 589  
590 Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Jimmy Lin. Ms marco: Bench-  
591 marking ranking models in the large-data regime. In *Proceedings of the 44th International*  
592 *ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’21*,  
593 pp. 1566–1576, New York, NY, USA, 2021. Association for Computing Machinery. ISBN  
9781450380379. doi: 10.1145/3404835.3462804. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3404835.3462804)  
[3404835.3462804](https://doi.org/10.1145/3404835.3462804).
- 594  
595 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In  
596 *The Twelfth International Conference on Learning Representations*, 2024. URL [https://](https://openreview.net/forum?id=mZn2Xyh9Ec)  
597 [openreview.net/forum?id=mZn2Xyh9Ec](https://openreview.net/forum?id=mZn2Xyh9Ec).

- 594 Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, and Andrew McCallum. Multi-step retriever-  
595 reader interaction for scalable open-domain question answering. In *International Confer-*  
596 *ence on Learning Representations*, 2019. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=HkfPSh05K7)  
597 [HkfPSh05K7](https://openreview.net/forum?id=HkfPSh05K7).
- 598 Nigel Fernandez, Alexander Scarlato, and Andrew Lan. SyllabusQA: A course logistics ques-  
599 tion answering dataset. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings*  
600 *of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
601 *Papers)*, pp. 10344–10369, Bangkok, Thailand, August 2024. Association for Computational Lin-  
602 guistics. doi: 10.18653/v1/2024.acl-long.557. URL [https://aclanthology.org/2024.](https://aclanthology.org/2024.acl-long.557)  
603 [acl-long.557](https://aclanthology.org/2024.acl-long.557).
- 604 Tao Ge, Hu Jing, Lei Wang, Xun Wang, Si-Qing Chen, and Furu Wei. In-context autoencoder  
605 for context compression in a large language model. In *The Twelfth International Confer-*  
606 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=uREj4ZuGJE)  
607 [uREj4ZuGJE](https://openreview.net/forum?id=uREj4ZuGJE).
- 608 Pengcheng He, Jianfeng Gao, and Weizhu Chen. Debertav3: Improving deberta using electra-style  
609 pre-training with gradient-disentangled embedding sharing, 2021a.
- 610 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. {DEBERTA}: {DECODING}-  
611 {enhanced} {bert} {with} {disentangled} {attention}. In *International Conference on Learning*  
612 *Representations*, 2021b. URL <https://openreview.net/forum?id=XPZIAotutsD>.
- 613 Sebastian Hofstätter, Sophia Althammer, Michael Schröder, Mete Sertkan, and Allan Hanbury.  
614 Improving efficient neural ranking models with cross-architecture knowledge distillation, 2021.  
615 URL <https://arxiv.org/abs/2010.02666>.
- 616 Jennifer Hsia, Afreen Shaikh, Zhiruo Wang, and Graham Neubig. Ragged: Towards informed design  
617 of retrieval augmented generation systems. *arXiv preprint arXiv:2403.09040*, 2024.
- 618 Taeho Hwang, Soyeong Jeong, Sukmin Cho, SeungYoon Han, and Jong Park. DSLR: Document re-  
619 finement with sentence-level re-ranking and reconstruction to enhance retrieval-augmented gen-  
620 eration. In Wenhao Yu, Weijia Shi, Michihiro Yasunaga, Meng Jiang, Chenguang Zhu, Han-  
621 naneh Hajishirzi, Luke Zettlemoyer, and Zhihan Zhang (eds.), *Proceedings of the 3rd Work-*  
622 *shop on Knowledge Augmented Methods for NLP*, pp. 73–92, Bangkok, Thailand, August 2024.  
623 Association for Computational Linguistics. doi: 10.18653/v1/2024.knowledgenlp-1.6. URL  
624 <https://aclanthology.org/2024.knowledgenlp-1.6>.
- 625 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane  
626 Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: Few-shot Learning  
627 with Retrieval Augmented Language Models, November 2022. URL [http://arxiv.org/](http://arxiv.org/abs/2208.03299)  
628 [abs/2208.03299](http://arxiv.org/abs/2208.03299). arXiv:2208.03299 [cs].
- 629 Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. LLMlingua: Com-  
630 pressing prompts for accelerated inference of large language models. In Houda Bouamor,  
631 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Meth-*  
632 *ods in Natural Language Processing*, pp. 13358–13376, Singapore, December 2023. Associa-  
633 tion for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.825. URL <https://aclanthology.org/2023.emnlp-main.825>.
- 634 Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili  
635 Qiu. LongLLMLingua: Accelerating and enhancing LLMs in long context scenarios via prompt  
636 compression. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd*  
637 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.  
638 1658–1677, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL  
639 <https://aclanthology.org/2024.acl-long.91>.
- 640 Dahyun Kim, Chanjun Park, Sanghoon Kim, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo  
641 Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, Changbae Ahn, Seonghoon Yang, Sukyung Lee,  
642 Hyunbyung Park, Gyoungjin Gim, Mikyung Cha, Hwalsuk Lee, and Sunghun Kim. Solar 10.7b:  
643 Scaling large language models with simple yet effective depth up-scaling, 2023.

- 648 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris  
649 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a  
650 benchmark for question answering research. *Transactions of the Association for Computational  
651 Linguistics*, 7:453–466, 2019.
- 652
- 653 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.  
654 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
655 serving with pagedattention, 2023.
- 656
- 657 LangChain. LangChain Documentation. <https://python.langchain.com/>.
- 658
- 659 Carlos Lassance and Stéphane Clinchant. Naver labs europe (splade) @ trec deep learning 2022,  
2023. URL <https://arxiv.org/abs/2302.12574>.
- 660
- 661 Carlos Lassance, Hervé Déjean, Thibault Formal, and Stéphane Clinchant. Splade-v3: New base-  
662 lines for splade, 2024.
- 663
- 664 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman  
665 Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and  
666 Douwe Kiela. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Ad-  
667 vances in Neural Information Processing Systems*, volume 33, pp. 9459–9474. Curran Asso-  
668 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/  
669 6b493230205f780e1bc26945df7481e5-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html).
- 670
- 671 Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. Compressing context to enhance in-  
672 ference efficiency of large language models. In Houda Bouamor, Juan Pino, and Kalika Bali  
(eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Pro-  
673 cessing*, pp. 6342–6353, Singapore, December 2023. Association for Computational Linguis-  
674 tics. doi: 10.18653/v1/2023.emnlp-main.391. URL [https://aclanthology.org/2023.  
675 emnlp-main.391](https://aclanthology.org/2023.emnlp-main.391).
- 676
- 677 Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo  
678 Nogueira. Pyserini: A python toolkit for reproducible information retrieval research with  
679 sparse and dense representations. In *Proceedings of the 44th International ACM SIGIR Con-  
680 ference on Research and Development in Information Retrieval, SIGIR ’21*, pp. 2356–2362, New  
681 York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450380379. doi:  
682 10.1145/3404835.3463238. URL <https://doi.org/10.1145/3404835.3463238>.
- 683
- 684 Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Richard James, Pedro  
685 Rodriguez, Jacob Kahn, Gergely Szilvasy, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih.  
686 RA-DIT: Retrieval-augmented dual instruction tuning. In *The Twelfth International Confer-  
687 ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=  
688 22OTbutug9](https://openreview.net/forum?id=22OTbutug9).
- 689
- 690 Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi.  
691 When not to trust language models: Investigating effectiveness of parametric and non-parametric  
692 memories. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of  
693 the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long  
694 Papers)*, pp. 9802–9822, Toronto, Canada, July 2023a. Association for Computational Linguis-  
695 tics. doi: 10.18653/v1/2023.acl-long.546. URL [https://aclanthology.org/2023.  
696 acl-long.546](https://aclanthology.org/2023.acl-long.546).
- 697
- 698 Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi.  
699 When not to trust language models: Investigating effectiveness of parametric and non-parametric  
700 memories. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of  
701 the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Pa-  
pers)*, pp. 9802–9822, Toronto, Canada, July 2023b. Association for Computational Linguis-  
tics. doi: 10.18653/v1/2023.acl-long.546. URL [https://aclanthology.org/2023.  
acl-long.546](https://aclanthology.org/2023.acl-long.546).

- 702 Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen tau Yih, Pang Wei Koh, Mohit Iyyer,  
703 Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual  
704 precision in long form text generation, 2023. URL <https://arxiv.org/abs/2305.14251>.
- 706 Piotr Nawrot, Adrian La'ncucki, Marcin Chochowski, David Tarjan, and E. Ponti. Dynamic memory  
707 compression: Retrofitting llms for accelerated inference. *ArXiv*, abs/2403.09636, 2024. URL  
708 <https://api.semanticscholar.org/CorpusID:268384862>.
- 710 Anastasios Nentidis, Georgios Katsimpras, Anastasia Krithara, Salvador Lima-López, Eulalia  
711 Farré-Maduell, Luis Gasco Sanchez, Martin Krallinger, and Georgios Paliouras. *Overview  
712 of BioASQ 2023: The Eleventh BioASQ Challenge on Large-Scale Biomedical Semantic In-  
713 dexing and Question Answering*, pp. 227–250. 09 2023. ISBN 978-3-031-42447-2. doi:  
714 10.1007/978-3-031-42448-9.19.
- 715 Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and  
716 Li Deng. Ms marco: A human generated machine reading comprehension dataset. Novem-  
717 ber 2016. URL [https://www.microsoft.com/en-us/research/publication/  
718 ms-marco-human-generated-machine-reading-comprehension-dataset/](https://www.microsoft.com/en-us/research/publication/ms-marco-human-generated-machine-reading-comprehension-dataset/).
- 720 Rodrigo Nogueira and Kyunghyun Cho. Passage re-ranking with bert, 2020.
- 721 Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin,  
722 Victor Ruhle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Dongmei Zhang.  
723 LLMingua-2: Data distillation for efficient and faithful task-agnostic prompt compression. In  
724 Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computa-  
725 tional Linguistics ACL 2024*, pp. 963–981, Bangkok, Thailand and virtual meeting, August 2024.  
726 Association for Computational Linguistics. URL [https://aclanthology.org/2024.  
727 findings-acl.57](https://aclanthology.org/2024.findings-acl.57).
- 728 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
729 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Ed-  
730 ward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,  
731 Lu Fang, Junjie Bai, and Soumith Chintala. *PyTorch: an imperative style, high-performance deep  
732 learning library*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- 734 David Rau, Hervé Déjean, Nadezhda Chirkova, Thibault Formal, Shuai Wang, Vassilina Nikoulina,  
735 and Stéphane Clinchant. Bergen: A benchmarking library for retrieval-augmented generation,  
736 2024a. URL <https://arxiv.org/abs/2407.01102>.
- 737 David Rau, Shuai Wang, Hervé Déjean, and Stéphane Clinchant. Context embeddings for efficient  
738 answer generation in rag, 2024b. URL <https://arxiv.org/abs/2407.09252>.
- 739 Minjoon Seo, Jinhyuk Lee, Tom Kwiatkowski, Ankur Parikh, Ali Farhadi, and Hannaneh Hajishirzi.  
740 Real-time open-domain question answering with dense-sparse phrase index. In Anna Korhonen,  
741 David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Associ-  
742 ation for Computational Linguistics*, pp. 4430–4441, Florence, Italy, July 2019. Association for  
743 Computational Linguistics. doi: 10.18653/v1/P19-1436. URL [https://aclanthology.  
744 org/P19-1436](https://aclanthology.org/P19-1436).
- 746 Xiao Shitao, Liu Zheng, Shao Yingxia, and Cao Zhao. Retromae: Pre-training retrieval-oriented  
747 language models via masked auto-encoder. In *EMNLP, 2022*. URL [https://arxiv.org/  
748 abs/2205.12035](https://arxiv.org/abs/2205.12035).
- 749 Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. BEIR: A  
750 heterogeneous benchmark for zero-shot evaluation of information retrieval models. In *Thirty-fifth  
751 Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round  
752 2)*, 2021. URL <https://openreview.net/forum?id=wCu6T5xFjeJ>.
- 754 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
755 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher,  
Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy

- 756 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,  
757 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel  
758 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,  
759 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,  
760 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,  
761 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh  
762 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen  
763 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,  
764 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,  
765 2023.
- 766 Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. Learning to  
767 filter context for retrieval-augmented generation, 2023. URL [https://arxiv.org/abs/  
768 2311.08377](https://arxiv.org/abs/2311.08377).
- 769 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
770 Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick  
771 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger,  
772 Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural  
773 language processing. In Qun Liu and David Schlangen (eds.), *Proceedings of the 2020 Confer-  
774 ence on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–  
775 45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.  
776 emnlp-demos.6. URL <https://aclanthology.org/2020.emnlp-demos.6>.
- 777 Fangyuan Xu, Weijia Shi, and Eunsol Choi. RECOMP: Improving retrieval-augmented LMs  
778 with context compression and selective augmentation. In *The Twelfth International Confer-  
779 ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=  
780 mLJLVigNHp](https://openreview.net/forum?id=mLJLVigNHp).
- 781 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov,  
782 and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop ques-  
783 tion answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.),  
784 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*,  
785 pp. 2369–2380, Brussels, Belgium, 2018. Association for Computational Linguistics. doi:  
786 10.18653/v1/D18-1259. URL <https://aclanthology.org/D18-1259>.
- 787 Chanwoong Yoon, Taewhoo Lee, Hyeon Hwang, Minbyul Jeong, and Jaewoo Kang. Com-  
788 pact: Compressing retrieved documents actively for question answering, 2024. URL [https:  
789 //arxiv.org/abs/2407.09014](https://arxiv.org/abs/2407.09014).
- 790 Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. Making retrieval-augmented language  
791 models robust to irrelevant context. In *The Twelfth International Conference on Learning Repre-  
792 sentations*, 2024. URL <https://openreview.net/forum?id=ZS4m74kZpH>.
- 793 Yun Zhu, Jia-Chen Gu, Caitlin Sikora, Ho Ko, Yinxiao Liu, Chu-Cheng Lin, Lei Shu, Liangchen  
794 Luo, Lei Meng, Bang Liu, and Jindong Chen. Accelerating inference of retrieval-augmented  
795 generation via sparse context selection. *ArXiv*, abs/2405.16178, 2024. URL [https://api.  
796 semanticscholar.org/CorpusID:270062557](https://api.semanticscholar.org/CorpusID:270062557).
- 797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## 810 A DATA

811  
812 **Evaluation datasets.** We consider the following datasets:

- 813  
814
- Datasets with Wikipedia as a datastore:
    - 815 – Natural Questions (Kwiatkowski et al., 2019). We use a test set of 2.8k ques-  
816 tions, distributed as a part of the KILT collection ([https://huggingface.co/](https://huggingface.co/datasets/facebook/kilt_tasks)  
817 [datasets/facebook/kilt\\_tasks](https://huggingface.co/datasets/facebook/kilt_tasks));
    - 818 – HotpotQA (Yang et al., 2018). We use a test set of 5.6k questions, distributed as a part  
819 of the KILT collection ([https://huggingface.co/datasets/facebook/](https://huggingface.co/datasets/facebook/kilt_tasks)  
820 [kilt\\_tasks](https://huggingface.co/datasets/facebook/kilt_tasks));
    - 821 – PopQA (Mallen et al., 2023b). We use a test set of 14k questions distributed by the  
822 dataset authors.
  - Datasets with individual datastores:
    - 823 – BioASQ (Nentidis et al., 2023). We use a version of the dataset provided by (Hsia  
824 et al., 2024), with 3.8k queries. We only use queries from categories “yes/no”, “fac-  
825 toid”, and “list”.
    - 826 – Syllabus QA (Fernandez et al., 2024). We use the test set of 1.1k questions distributed  
827 by the authors;
    - 828 – RGB (Chen et al., 2024b). We use the test set of 200 questions distributed by the  
829 authors.
- 830  
831

832 All datasets provide short answers (keywords) for each query, which we use to evaluate both match-  
833 based metrics such as Recall and LLM-based metrics Rau et al. (2024a)<sup>9</sup>.

834  
835 **Datastores.** For training `Provence`, we use the MS MARCO document collection (Craswell  
836 et al., 2021). We split each document into overlapping chunks of  $N$  sentences, where  $N$  is random  
837 in  $\in 1..10$  – with a higher probability for longer contexts – to train `Provence` on various context  
838 lengths. Each chunk is prepended with a page title. The resulting datastore contains 34M passages.  
839 We also process the Wikipedia datastore in a similar fashion, for ablation experiments. We download  
840 a 2024 Wikipedia dump and process it using scripts provided by Pyserini (Lin et al., 2021)<sup>10</sup>. We  
841 also prepare versions of this Wikipedia datastore with passages of  $N$  sentences with overlaps of  $N/2$   
842 sentences, for testing `Provence` robustness to various context lengths.

843 All other evaluations on Wikipedia-based datasets – including main evaluations – are conducted on  
844 the Wikipedia datastore provided at [https://huggingface.co/datasets/castorini/](https://huggingface.co/datasets/castorini/odqa-wiki-corpora)  
845 [odqa-wiki-corpora](https://huggingface.co/datasets/castorini/odqa-wiki-corpora). We use a version with passages of 6 sentences with a 3-sentence overlap  
846 – making 9M passages in total.

847 For Pubmed, we use the version of the dataset provided by (Hsia et al., 2024) at [https://](https://huggingface.co/datasets/jenhsia/ragged)  
848 [huggingface.co/datasets/jenhsia/ragged](https://huggingface.co/datasets/jenhsia/ragged). It consists of 58M passages, extracted  
849 from Pubmed abstracts. Each passage (chunk) is prepended with the page’s title.

850 For SyllabusQA, we split each syllabus (provided by the authors) into passages of 100 words. For  
851 RGB, context are provided by the authors.

## 853 B MODELS

854 We list in Table 5 all the main models used to conduct experiments for `Provence`.

855  
856  
857  
858  
859  
860  
861  
862 <sup>9</sup>Using SOLAR-10.7B (Kim et al., 2023).

863 <sup>10</sup>At [https://github.com/castorini/pyserini/blob/master/docs/](https://github.com/castorini/pyserini/blob/master/docs/experiments-wiki-corpora.md)  
experiments-wiki-corpora.md.



864  
865  
866  
867  
868  
869  
870  
871  
872  
873

Model	Checkpoint
SPLADE-v3	naver/splade-v3
RetroMAE	Shitao/RetroMAE-MSMARCO-distill
DeBERTa-v3	microsoft/deberta-large
DeBERTa-v3 (RR)	naver/trecdl22-crossencoder-debertav3
BGE-M3	BAAI/bge-reranker-v2-m3
LLama-2-7B-chat	meta-llama/Llama-2-7b-chat-hf
LLaMA-3-8B-Instruct	meta-llama/Meta-Llama-3-8B-Instruct
Mistral-7B-instruct	mistralai/Mistral-7B-Instruct-v0.2
SOLAR-10.7B-Instruct-v1.0	upstage/SOLAR-10.7B-Instruct-v1.0

874  
875  
876

Table 5: List of all the models used in the experiments with their corresponding HuggingFace checkpoints.

877  
878  
879

Table 6: Prompt used for generating silver labeling with LLaMA-3-8B-Instruct. The sentence citations in the response are parsed using regular expression.

880  
881  
882  
883  
884  
885  
886  
887

Question: {question}  
Context: [1] {sentence1} [2] {sentence2} [3] {sentence3} ...  
Answer the Question, using ONLY information provided in the Context. If no useful information is provided, you MUST output "No answer". If some parts of the Context are used to answer, you MUST cite ALL the corresponding sentences. Use the symbols [ ] to indicate when a fact comes from a sentence in the context, e.g [0] for a fact from sentence 0. You should only answer the given question and should not provide any additional information.

888

Table 7: nDCG@10 on the 13 open BEIR datasets.

889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902

Corpus	DeBERTav3	Provence
TREC-COVID	88.3	88.3
NFCorpus	37.5	37.8
NQ	66.7	66.5
HotpotQA	74.5	74.9
FiQA-2018	47.8	47.6
ArguAna	29.8	33.2
Touché-2020	33.5	33.4
Quora	84.8	85.4
DBpedia	48.9	49.2
SCIDOCS	19.2	19.6
FEVER	86.6	87.9
Climate-FEVER	27.4	28.1
SciFact	75.8	75.3
average	55.4	55.9

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

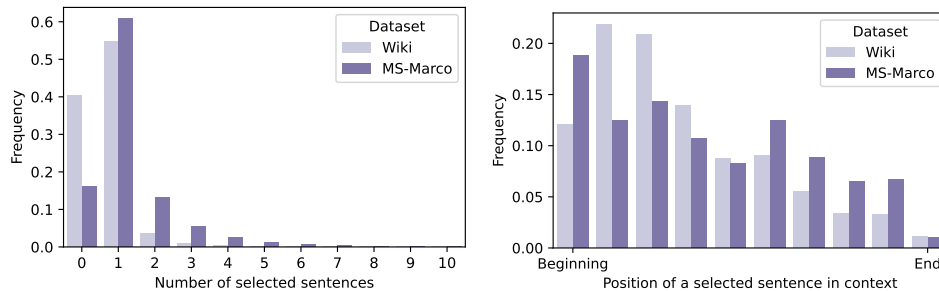


Figure 5: Statistics of the silver contexts labeled by LLaMA-3-8B-Instruct. (Left) the distribution of the number of sentences in silver contexts. (Right) the distribution of the position of the selected sentences in contexts.

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

how many french words are there in the english language **Inf luence of French on English** . The most notable influence of French on English has been its extensive contribution to the English lex icon . It has been estimated that about a third of the words in English are French in origin ; lingu ist Henri ette Walter claims that this total may be as high as two thirds . L ingu ist Anthony Lac oud re has estimated that over 40 , 000 English words come directly from French and may be understood without orth ographical change by French speakers . Albert C . B augh and Thomas Cable note that " although this influx of French words was brought about by the victory of the Conquer or and by the political and social consequences of that victory , it was neither sudden nor immediately apparent . Rather it began slowly and continued with varying tempo for a long time .

COLORS: 0.99 0.1 0

Figure 6: Example visualization of per-token probabilities of being selected in the final context.

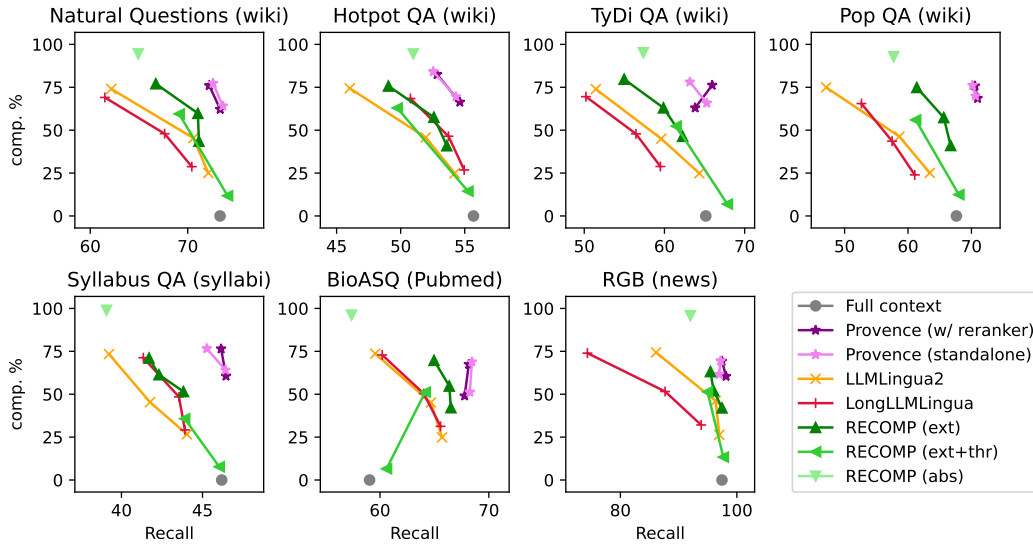


Figure 7: Main results for various QA domains, comparing Provenance and baseline models, metric: Recall. *Generator: Llama-2-7B, retriever: SPLADE-v3, reranker: DeBERTa-v3 (or Provenance in the unified setting)*. Plot titles denote “Dataset name (datastore type)”. *x*-axis denotes QA performance evaluated with Recall; *y*-axis denotes the context compression ratio. For both metrics, the higher the better: the best model would be closest to the top right corner.

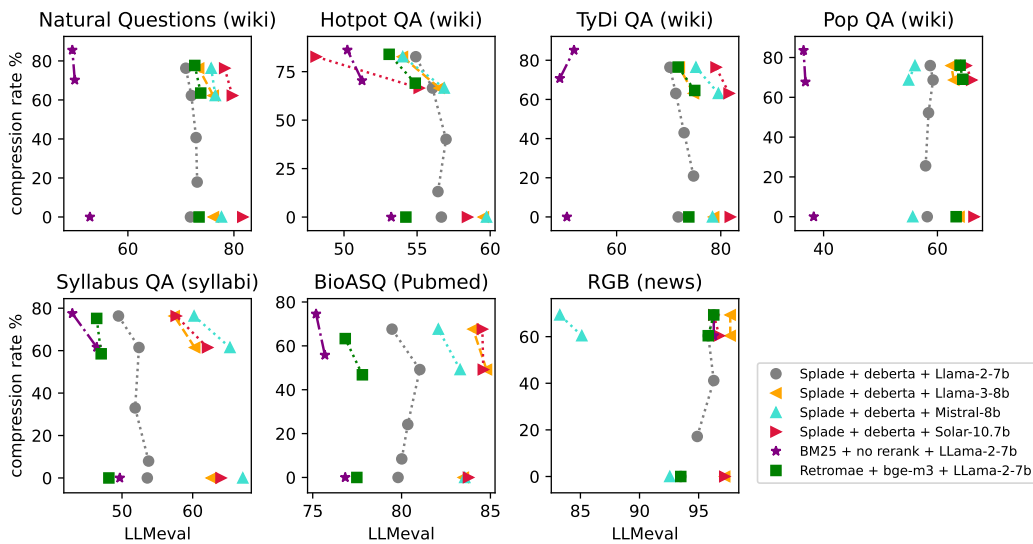


Figure 8: Testing Provenance in various RAG settings (retrieval, re-ranking, generator).

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  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

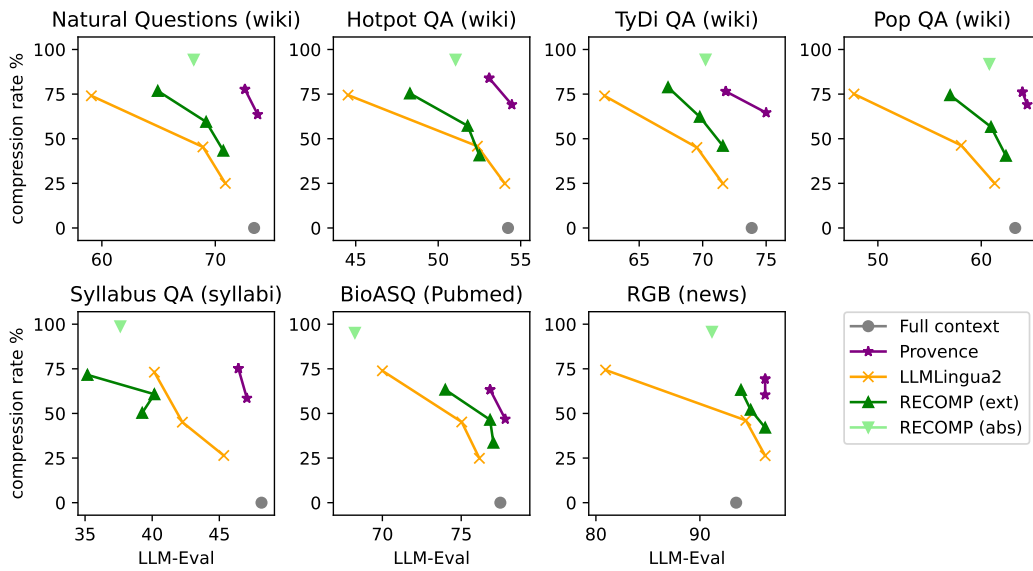


Figure 9: Comparing `Provenance` to a subset of baselines with *retriever*: RetroMAE (Shitao et al., 2022), *reranker*: BGE-M3 (Chen et al., 2024a), *generator*: Llama-2-7B-chat.

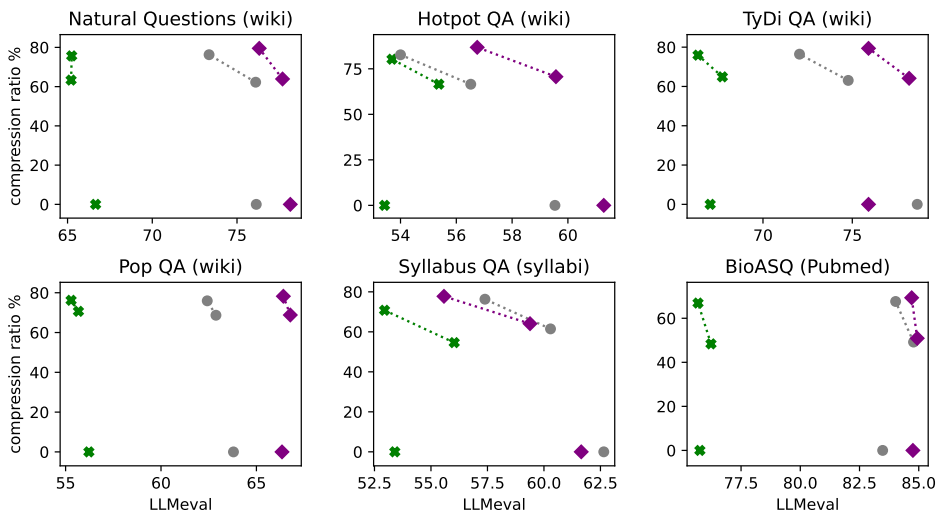


Figure 10: Testing `Provenance` with different top- $k$  documents provided to the generator. The setting is the same as the one in Figure 2.

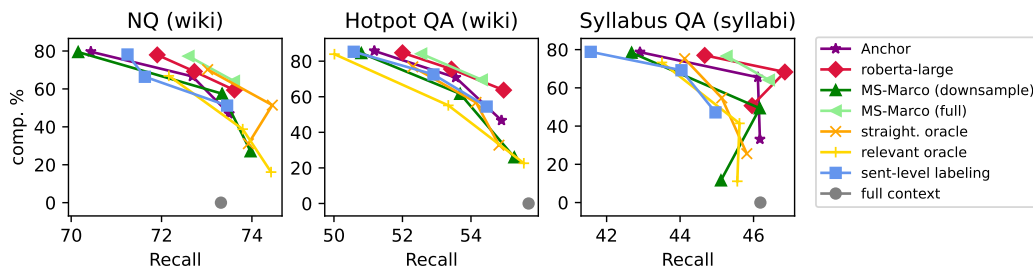


Figure 11: Ablation results with Recall (match-based metric).

1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079

Table 8: Numerical scores corresponding to Figure 2 – NQ, Hotpot QA, Tydi QA, and Pop QA.

	NQ		HotPot QA		Tydi QA		PopQA	
	LLM-Eval	Comp. rate %	LLM-Eval	Comp. rate %	LLM-Eval	Comp. rate %	LLM-Eval	Comp. rate %
Full context	71.8	0.0	57.0	0.0	73.9	0.0	57.8	0.0
Provence (w/ reranker)	72.4	62.2	56.7	66.4	70.5	63.0	59.3	68.6
	72.6	76.0	56.0	82.4	73.6	76.2	59.5	75.8
Provence (standalone)	72.3	64.1	56.6	69.5	70.9	65.8	59.0	69.9
	70.6	77.3	54.8	84.1	70.2	78.1	58.8	76.1
LLMLingua2	59.5	74.0	47.1	74.4	57.7	73.9	42.9	75.0
	67.5	45.4	52.9	45.8	67.3	45.0	52.5	46.3
	70.3	25.0	55.0	24.9	70.0	24.8	55.2	25.1
LongLLMLingua	61.3	69.1	52.6	68.5	56.6	69.5	49.5	65.5
	68.5	47.9	55.6	46.5	65.5	47.8	54.5	43.6
	71.3	28.7	56.9	26.8	69.1	28.8	57.6	23.9
RECOMP (ext)	70.6	43.6	55.5	40.9	68.6	46.4	56.9	41.1
	68.2	59.8	53.4	57.5	67.0	63.0	55.7	57.4
	66.2	77.1	50.1	75.7	64.5	79.7	52.3	74.9
RECOMP (ext+thr)	69.0	59.5	50.9	62.9	70.9	52.4	54.8	56.0
	72.9	11.8	56.4	14.4	72.3	6.9	58.5	12.4
RECOMP (abs)	66.9	94.5	53.1	94.4	66.4	95.2	54.4	92.8
DSLRL	71.7	44.9	52.9	75.7	72.7	45.8	58.6	48.1
	70.5	54.9	50.7	83.4	69.8	55.6	58.7	58.1
	70.4	61.4	49.3	87.0	70.7	62.0	58.8	63.7
	67.7	72.0	45.2	91.7	67.5	72.9	58.5	71.9
	67.6	77.7	43.2	93.4	67.5	78.1	57.9	76.0
Dense-X retrieval	62.7	69.0	49.6	67.7	66.4	71.5	52.0	68.5

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

Table 9: Numerical scores corresponding to Figure 2: Syllabus QA, BioASQ, and RGB.

	Syllabus QA		BioASQ		RGB	
	LLM-Eval	Comp. rate %	LLM-Eval	Comp. rate %	LLM-Eval	Comp. rate %
Full context	52.9	0.0	80.7	0.0	93.5	0.0
Provence	49.8	60.6	80.6	49.0	94.4	60.5
(w/ reranker)	51.0	76.5	80.3	67.4	96.3	69.3
Provence	50.7	64.1	80.6	51.3	95.8	61.6
(standalone)	47.8	76.6	80.1	68.9	96.3	69.4
LLMLingua2	37.4	73.4	72.6	73.6	78.6	74.3
	43.4	45.4	77.7	45.2	93.5	46.1
	49.8	26.6	78.7	24.8	95.8	26.3
LongLLMLingua	42.3	71.3	72.2	72.9	71.6	73.9
	45.1	48.5	77.3	50.4	83.3	51.6
	50.9	29.2	78.7	31.3	92.1	32.1
RECOMP	44.6	51.5	78.7	42.2	97.7	42.0
(ext)	42.7	61.4	78.4	54.8	94.9	52.1
	39.1	71.1	76.3	69.7	94.4	63.2
RECOMP	45.5	35.7	76.6	51.2	92.1	51.4
(ext+thr)	52.8	7.7	80.2	6.5	97.7	13.4
RECOMP (abs)	38.1	98.9	68.2	96.1	90.7	95.7
DSLRL	49.6	33.2	80.1	29.9	97.2	41.6
	49.1	46.4	79.6	40.1	97.7	46.9
	47.2	55.4	79.2	47.3	96.3	49.7
	44.2	70.6	77.6	60.2	97.2	54.1
	40.7	78.2	75.4	68.0	95.8	56.9

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

Table 10: Numerical scores corresponding to Figure 4.

	NQ			HotPot QA		Syllabus QA	
	Thresh.	LLM-Eval	Comp. rate %	LLM-Eval	Comp. rate %	LLM-Eval	Comp. rate %
Anchor	0.01	72.2	46.9	56.0	46.7	51.1	33.1
	0.1	71.6	66.6	55.6	70.8	51.1	65.5
	0.5	70.6	79.6	52.9	85.8	44.0	78.7
Deberta-large	0.01	72.5	59.4	56.7	63.7	49.9	50.7
	0.1	72.1	69.2	55.9	75.6	49.6	68.3
	0.5	70.9	78.0	53.2	84.7	46.1	77.0
MS-Marco (downsample)	0.01	72.9	27.1	57.2	25.9	52.5	11.6
	0.1	72.8	57.5	56.2	61.6	50.5	49.3
	0.5	68.7	79.4	52.3	84.5	43.0	78.4
MS-Marco (full)	0.1	72.3	64.1	56.6	69.5	50.7	64.1
	0.5	70.6	77.3	54.8	84.1	47.8	76.6
straight. oracle	0.01	73.1	31.1	56.8	32.8	52.9	25.6
	0.1	72.5	51.4	55.9	56.8	51.7	54.7
	0.5	72.0	70.2	54.8	76.8	45.5	75.3
relevant oracle	0.01	72.8	16.1	57.4	22.6	49.4	11.1
	0.1	72.6	38.8	55.7	55.1	49.6	41.4
	0.5	71.3	66.9	51.6	83.9	46.8	73.3
sent-level labeling	0.01	73.0	51.2	56.2	54.5	51.5	47.2
	0.1	71.0	66.4	55.0	72.5	48.3	69.2
	0.5	71.2	78.2	53.0	85.3	40.4	78.8
full context	0.01	71.8	0.0	57.0	0.0	52.9	0.0

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

Table 11: Example of context pruning with various approaches. Provenance selects one sentence about the Shepard’s pie and removes sentences about other similar dishes, which is RECOMP (ext) is not capable of by design. RECOMP (abs) correctly generates a summary; LongLLMLingua removes the part relevant to the Shepard’s pie, and LLMLingua2 uniformly removes no-informative tokens.

Question	what goes on the bottom of shepherd’s pie
Original context	Shepherd’s pie. History. In early cookery books, the dish was a means of using leftover roasted meat of any kind, and the pie dish was lined on the sides and bottom with mashed potato, as well as having a mashed potato crust on top. Variations and similar dishes. Other potato-topped pies include: The modern ”Cumberland pie” is a version with either beef or lamb and a layer of bread-crumbs and cheese on top. In medieval times, and modern-day Cumbria, the pastry crust had a filling of meat with fruits and spices.. In Quebec, a variation on the cottage pie is called ”Pâté chinois”. It is made with ground beef on the bottom layer, canned corn in the middle, and mashed potato on top.. The ”shepherdess pie” is a vegetarian version made without meat, or a vegan version made without meat and dairy.. In the Netherlands, a very similar dish called ”philosopher’s stew” () often adds ingredients like beans, apples, prunes, or apple sauce.. In Brazil, a dish called in refers to the fact that a manioc puree hides a layer of sun-dried meat.
Provenance	Shepherd’s pie. In early cookery books, the dish was a means of using leftover roasted meat of any kind, and the pie dish was lined on the sides and bottom with mashed potato, as well as having a mashed potato crust on top.
RECOMP (ext, top-3)	In early cookery books, the dish was a means of using leftover roasted meat of any kind, and the pie dish was lined on the sides and bottom with mashed potato, as well as having a mashed potato crust on top. The ”shepherdess pie” is a vegetarian version made without meat, or a vegan version made without meat and dairy.. In the Netherlands, a very similar dish called ”philosopher’s stew” () often adds ingredients like beans, apples, prunes, or apple sauce..
RECOMP (abs)	Shepherd’s pie is a dish made with leftover roasted meat and mashed potato on the bottom, with variations including a version with beef or lamb and a layer of breadcrumbs and cheese on top. The dish has a history of being used as a means of using leftover roasted meat and has variations on the toppings.
LLMLingua2 (comp.50%)	Shepherd’s pie History early cookery books dish leftover roasted meat lined mashed potato mashed potato crust top Variations similar dishes potato-topped pies include modern ”Cumberland pie” beef or lamb breadcrumbs cheese medieval modern-day Cumbria pastry crust filling meat fruits spices Quebec variation cottage pie ”Pâté ground beef bottom canned corn middle mashed potato top ”shepherdess pie” vegetarian without meat vegan version without meat dairy Netherlands similar dish ”philosopher’s stew” adds ingredients beans apples prunes apple sauce Brazil dish manioc puree hides sun-dried meat
LongLLMLingua (comp. 50%). Processes all passages together and dynamically decides on the compression ratio of each passage.	Shepherd’s pie. Other potato-topped pies include: The modern ”Cumberland pie” is a version with either beef or lamb and a layer of breadcrumbs and cheese on top. In medieval times, and modern-day Cumbria, the pastry crust had a filling of meat with fruits and spices.. In Quebec, a variation on the cottage pie is called ”Pâté chinois”. It is made with ground beef on the bottom layer, canned corn in the middle, and mashed potato on top.. The ”shepherdess pie” is a vegetarian version made without meat, or a vegan version made without meat and dairy.. In the Netherlands, a very similar dish called ”philosopher’s stew” () often adds ingredients like beans, apples, prunes, or apple sauce.. In Brazil, a dish called in refers to the fact that a manioc puree hides a layer of sun-dried meat.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

Table 12: Example of context pruning with various approaches. `Provence` correctly detects that the entire passage is relevant to the query, same as `LongLLMLingua`, while `RECOMP (ext)` is by design not capable of making such a decision.

Question		where does the sweetness of fruit come from
Original context		Sweetness. A number of plant species produce glycosides that are sweet at concentrations much lower than sugar. The most well-known example is glycyrrhizin, the sweet component of licorice root, which is about 30 times sweeter than sucrose. Another commercially important example is stevioside, from the South American shrub "Stevia rebaudiana". It is roughly 250 times sweeter than sucrose. Another class of potent natural sweeteners are the sweet proteins such as thaumatin, found in the West African katemfe fruit. Hen egg lysozyme, an antibiotic protein found in chicken eggs, is also sweet.
Provence		Sweetness. A number of plant species produce glycosides that are sweet at concentrations much lower than sugar. The most well-known example is glycyrrhizin, the sweet component of licorice root, which is about 30 times sweeter than sucrose. Another commercially important example is stevioside, from the South American shrub "Stevia rebaudiana". It is roughly 250 times sweeter than sucrose. Another class of potent natural sweeteners are the sweet proteins such as thaumatin, found in the West African katemfe fruit. Hen egg lysozyme, an antibiotic protein found in chicken eggs, is also sweet.
RECOMP (ext, top-3)		It is roughly 250 times sweeter than sucrose. Another commercially important example is stevioside, from the South American shrub "Stevia rebaudiana". A number of plant species produce glycosides that are sweet at concentrations much lower than sugar.
RECOMP (abs)		[empty context]
LLMLingua2 (comp.50%)		Sweetness plant species produce glycosides sweet lower sugar glycyrrhizin sweet licorice root 30 times sweeter sucrose stevioside South American shrub "Stevia 250 times sweeter sucrose sweeteners sweet proteins thaumatin West African katemfe fruit Hen egg lysozyme antibiotic protein chicken eggs sweet
LongLLMLingua (comp. 50%)		Sweetness. A number of plant species produce glycosides that are sweet at concentrations much lower than sugar. The most well-known example is glycyrrhizin, the sweet component of licorice root, which is about 30 times sweeter than sucrose. Another commercially important example is stevioside, from the South American shrub "Stevia rebaudiana". It is roughly 250 times sweeter than sucrose. Another class of potent natural sweeteners are the sweet proteins such as thaumatin, found in the West African katemfe fruit. Hen egg lysozyme, an antibiotic protein found in chicken eggs, is also sweet.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

Table 13: Example of context pruning with various approaches. Provenance selects one most relevant sentence, which is also ranked first by RECOMP (ext). RECOMP (abs) decides that no information is relevant to the query, while LongLLMLingua on the contrary keeps the entire input, dropping some punctuation marks. LLMLingua2 removes too many tokens which makes text hardly understandable.

Question	what was the tower of london originally used for
Original context	Tower of London. In the 16th century, the Tower acquired an enduring reputation as a grim, forbidding prison. This had not always been the case. As a royal castle, it was used by the monarch to imprison people for various reasons, however these were usually high-status individuals for short periods rather than common citizenry as there were plenty of prisons elsewhere for such people. Contrary to the popular image of the Tower, prisoners were able to make their life easier by purchasing amenities such as better food or tapestries through the Lieutenant of the Tower. As holding prisoners was originally an incidental role of the Tower – as would have been the case for any castle – there was no purpose-built accommodation for prisoners until 1687 when a brick shed, a "Prison for Soldiers", was built to the north-west of the White Tower. The Tower's reputation for torture and imprisonment derives largely from 16th-century religious propagandists and 19th-century romanticists.
Provenance	Tower of London. As a royal castle, it was used by the monarch to imprison people for various reasons, however these were usually high-status individuals for short periods rather than common citizenry as there were plenty of prisons elsewhere for such people.
RECOMP (ext, sorted top-3 sents)	As a royal castle, it was used by the monarch to imprison people for various reasons, however these were usually high-status individuals for short periods rather than common citizenry as there were plenty of prisons elsewhere for such people. This had not always been the case. The Tower's reputation for torture and imprisonment derives largely from 16th-century religious propagandists and 19th-century romanticists.
RECOMP (abs)	[empty context]
LLMLingua2 (comp.25%)	Tower London 16th century grim prison royal castle monarch high-status common citizenry prisoners amenities food Lieutenant Tower no-built accommodation until 1687 "Prison for north-west White Tower reputation torture imprisonment 16th-century propagandists 19th-century romanticists
LongLLMLingua (comp. 25%)	Tower of London In the 6th century, the acquired an enduring reputation as grim, forbidd prison. This had always been the case As a royal castle, it was by the to imprison people for various reasons however these were usually high-status individuals for short rather than common citizenry as there were plenty of prisons elsewhere for such people. Contrary popular of the Tower, prisoners were able to make their life easier purchasing amenities such better food or tapestries through Lieutenant of the Tower. holding prisoners was originally incident role of the– would have been the case for any – was purpose- accommodation for prisoners until 167 a "Prison for Sold", was to thewest of White Tower. The's reputation torture imprisonment derives largely from 6th- religious propagandists and 19th-century romanticists.