

Enhancing Contextual Understanding in Large Language Models through Contrastive Decoding

Anonymous ACL submission

Abstract

Large language models (LLMs) tend to inadequately integrate input context during text generation, relying excessively on encoded prior knowledge in model parameters, potentially resulting in generated text with factual inconsistencies or contextually unfaithful content. LLMs utilize two primary knowledge sources: 1) prior (parametric) knowledge from pretraining, and 2) contextual (non-parametric) knowledge from input prompts. The study addresses the open question of how LLMs effectively balance these knowledge sources during the generation process, specifically in the context of open-domain question answering. To address this issue, we introduce a novel approach integrating contrastive decoding with adversarial irrelevant passages as negative samples to enhance robust context grounding during generation. Notably, our method operates at inference time without requiring further training. We conduct comprehensive experiments to demonstrate its applicability and effectiveness, providing empirical evidence showcasing its superiority over existing methodologies.

1 Introduction

Improving large language models (LLMs) has been a primary focus in natural language processing research. Recent strides have incorporated retrieval mechanisms to enhance LLMs (Lewis et al., 2020; Guu et al., 2020; Izacard and Grave, 2021; Izacard et al., 2023), augmenting their ability to produce contextually relevant and precise responses (Min et al., 2023; Mallen et al., 2023). Retrieval-augmented LLMs, which leverage both *parametric* knowledge acquired during training and *non-parametric* knowledge retrieved during inference, exhibit potential in addressing challenges such as limited memorization (Kandpal et al., 2023), knowledge conflicts (Longpre et al., 2021), and outdated information (Kasai et al., 2022).

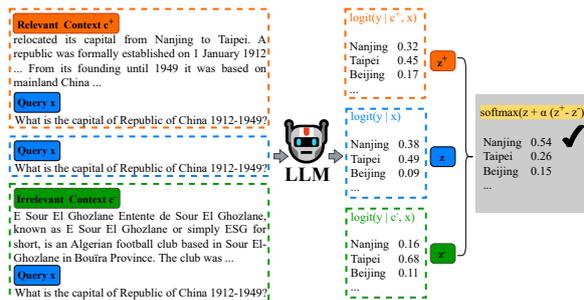


Figure 1: An illustration of our proposed decoding method. Despite the relevant context suggesting the answer as “Nanjing”, it contradicts the LLM’s prior knowledge. After reconciling different knowledge sources, the model correctly predicted the answer.

An ongoing question pertains to how LLMs ought to balance these two knowledge sources during generation. Previous research suggests that LLMs can falter in adequately attending to newly introduced information within the contextual knowledge. To tackle this issue, context-aware decoding (CAD; Shi et al., 2023a) has been proposed. By employing a contrastive output distribution, CAD highlights discrepancies in output probabilities when the model operates with and without context. Their experiments illustrate CAD’s effectiveness in overriding the model’s parametric knowledge in cases of conflict with provided context. However, while prior works often assert context as inherently reliable, our perspective argues that LLMs should possess the capacity to navigate and reconcile both parametric and non-parametric knowledge, ultimately refining their ability to strike a judicious balance. This paper undertakes the development and assessment of a novel decoding strategy tailored for retrieval-augmented LLMs, seeking equilibrium in utilizing parametric and non-parametric knowledge sources. The proposed method involves a contrastive decoding approach (Li et al., 2023), integrating both relevant and irrelevant contexts, wherein the irrelevant context can

067 be adversarially crafted retrieval or bottom-ranked
068 retrieved text. Notably, we emphasize the criticality
069 of leveraging irrelevant contexts, a distinguishing
070 feature of our approach.

071 Through comprehensive experiments span-
072 ning diverse datasets such as Natural Questions
073 (Kwiatkowski et al., 2019), TriviaQA (Joshi et al.,
074 2017), and PopQA (Mallen et al., 2023), and mod-
075 els encompassing various vanilla LLMs including
076 OPT (Zhang et al., 2022), Falcon (Almazrouei
077 et al., 2023), LLaMA families (Touvron et al.,
078 2023a,b), and instruction-tuned Flan-T5 (Chung
079 et al., 2022), we provide empirical evidence sup-
080 porting the superiority of incorporating irrelevant
081 contexts in assisting LLMs to manage knowledge
082 conflicts and seamlessly integrate contexts for gen-
083 erating responses in open-domain question answer-
084 ing against conventional decoding approaches with-
085 out necessitating further fine-tuning. The investiga-
086 tion also delves into the impact of different retrieval
087 sources on the decoding strategy, emphasizing the
088 importance of refining retrieval mechanisms for
089 further enhancements in performance.

090 Additionally, the paper explores different facets
091 of the proposed decoding approach, encompass-
092 ing the influence of various hyperparameters, the
093 effect of scaling model sizes, and the selection
094 of irrelevant contexts. This exploration provides
095 deeper insights into leveraging parametric and non-
096 parametric knowledge sources. We demonstrate
097 that although our approach outperforms regular
098 decoding across most model sizes, it particularly
099 excels with larger models. Moreover, we show our
100 method’s effectiveness even with simple fixed irrel-
101 evant contexts. Additionally, our approach exhibits
102 consistent performance improvements in answering
103 questions with knowledge across varying levels of
104 popularity. Beyond benchmarking against existing
105 methods, this study also explores practical impli-
106 cations and constraints of the proposed decoding
107 strategy, delineating pathways for future research
108 in generative tasks beyond question answering.

109 2 Related Works

110 **Retrieval-augmented LLMs** While LLMs re-
111 lying solely on their parameters can capture ex-
112 tensive world knowledge, they exhibit limited
113 memorization for less frequent entities (Kandpal
114 et al., 2023), susceptibility to hallucinations (Shus-
115 ter et al., 2021), and temporal degradation (Luu
116 et al., 2022; Jang et al., 2022). Furthermore, the

117 acquired parametric knowledge swiftly becomes
118 outdated (Kasai et al., 2022). Recent research
119 emphasizes the enhancement of LLMs with non-
120 parametric memories, referred to as retrieved text
121 chunks, enabling smaller models to match the per-
122 formance of larger counterparts (Izacard et al.,
123 2023). Studies exploring the integration of re-
124 trieved non-parametric memories within intermedi-
125 ate states or output spaces have shown effectiveness
126 in overcoming LLM limitations in memorization
127 and knowledge updating (Zhong et al., 2022; Min
128 et al., 2023). Mallen et al. (2023) extensively an-
129alyze the circumstances favoring the benefits of
130 retrieval augmentation. They demonstrate its effi-
131 cacy in less frequent occurrences but caution about
132 potential misguidance for LLMs. Building upon
133 these insights, they introduce adaptive retrieval and
134 empirically showcase its promising effectiveness.

135 **Knowledge Conflicts** In cases of conflicting
136 knowledge in updated documents, language mod-
137 els are expected to generate responses based on
138 provided contexts rather than relying solely on out-
139 dated parametric knowledge. Retrieval-augmented
140 LLMs (Min et al., 2023; Shi et al., 2023b; Izacard
141 et al., 2023) particularly benefit from this scenario
142 by employing externally retrieved documents to
143 enrich their knowledge. However, the mere addi-
144 tion of documents doesn’t consistently influence
145 model predictions, as current LLMs often over-
146 look contexts and heavily rely on prior parametric
147 knowledge (Longpre et al., 2021; Chen et al., 2022).
148 Existing approaches aiming to improve a model’s
149 fidelity to context, such as prompting-based meth-
150 ods (Zhou et al., 2023), are constrained to large-
151 scale instruction-finetuned LLMs like OpenAI’s
152 text-davinci-003. In contrast, our work investigates
153 a decoding strategy applicable to any LLMs.

154 **Contrastive Decoding** The exploration of con-
155 trastive decoding methods extensively addresses
156 text generation. MMI-based decoding (Li et al.,
157 2016) utilizes a contrastive formulation to enhance
158 output diversity in dialog generation. DExperts
159 (Liu et al., 2021) dampens the output distribution of
160 an anti-expert (e.g., exposed to toxic language) to
161 guide generations away from undesired attributes.
162 Contrastive decoding (Li et al., 2023) demotes an
163 amateur model (e.g., models with minimal param-
164 eters) to distill expert knowledge from larger, com-
165 petitive models. Pozzobon et al. (2023) introduce
166 an innovative toxicity mitigation approach that con-
167 trasts and ensembles the next token probabilities

obtained from a LLM using both toxic and non-toxic retrievals. Context-aware decoding (Shi et al., 2023a) emphasizes output probability differences using a contrastive ensemble between model predictions with and without non-parametric knowledge. It effectively overrides a model’s parametric knowledge when it conflicts with the provided non-parametric information.

3 Methodology

3.1 Problem Statement

We consider decoding approaches for open-domain question answering, where the large language model θ receives an input query \mathbf{x} and aim to generate a faithful answer \mathbf{y} . During the generation of \mathbf{y}_t at each time step t , the language model computes the logits $\mathbf{z}_t \in \mathbb{R}^{|V|}$ for the t -th token, where V represents the vocabulary. The probability distribution over the vocabulary is derived by normalizing and exponentiating \mathbf{z}_t as follows:

$$p_\theta(y_t|\mathbf{x}, \mathbf{y}_{<t}) = \text{softmax}(\mathbf{z}_t).$$

Prompting the model for its parametric knowledge involves sampling the response from the probability distribution conditioned on the query \mathbf{x} and the previously generated response $\mathbf{y}_{<t}$:

$$y_t \sim p_\theta(y_t|\mathbf{x}, \mathbf{y}_{<t}).$$

Similarly, when incorporating additional context \mathbf{c} , containing external knowledge beyond the model’s parametric knowledge, our model θ generates a response \mathbf{y} considering the query, context, and the previously generated response:

$$y_t \sim p_\theta(y_t|\mathbf{c}, \mathbf{x}, \mathbf{y}_{<t}).$$

We observe two sources of knowledge (parametric vs. non-parametric) contributing to model responses, which may sometimes conflict (Longpre et al., 2021; Neeman et al., 2023). While some argue for prioritizing non-parametric knowledge over potentially outdated parametric knowledge (Shi et al., 2023a), we propose the importance of striking a balance between these sources as non-parametric knowledge, derived from external retrievers, may also contain inaccuracies.

3.2 Multi-Input Contrastive Decoding

Context can be both beneficial and problematic. Thus, we segregate context \mathbf{c} into relevant \mathbf{c}^+ and irrelevant \mathbf{c}^- . At each decoding time step t , our

approach combines the model’s prediction based on its parametric knowledge (\mathbf{z}_t) with predictions utilizing relevant (\mathbf{z}_t^+) and irrelevant (\mathbf{z}_t^-) contexts:

$$y_t \sim \text{softmax}(\mathbf{z}_t + \alpha(\mathbf{z}_t^+ - \mathbf{z}_t^-)),$$

where α is a hyperparameter that governs the extent of modification to the parametric answer (\mathbf{z}_t). Equivalently,

$$y_t \sim \tilde{p}_\theta(y_t|\mathbf{c}^+, \mathbf{c}^-, \mathbf{x}, \mathbf{y}_{<t}) \propto p_\theta(y_t|\mathbf{x}, \mathbf{y}_{<t}) \left(\frac{p_\theta(y_t|\mathbf{c}^+, \mathbf{x}, \mathbf{y}_{<t})}{p_\theta(y_t|\mathbf{c}^-, \mathbf{x}, \mathbf{y}_{<t})} \right)^\alpha.$$

In essence, a response will exhibit high probability only if it holds high likelihood under both learned parametric knowledge and relevant non-parametric knowledge, while demonstrating low probability under irrelevant non-parametric knowledge. The ratio $\frac{p_\theta(y_t|\mathbf{c}^+, \mathbf{x}, \mathbf{y}_{<t})}{p_\theta(y_t|\mathbf{c}^-, \mathbf{x}, \mathbf{y}_{<t})}$ functions as a scaling factor used to modify the parametric answer for the given input query. A larger α implies a greater modification, with $\alpha = 0$ resulting in no modification, indicating regular decoding using solely parametric knowledge without additional context.

Fundamentally, our proposed decoding operates as an ensemble involving the logits \mathbf{z}_t , \mathbf{z}_t^+ , and \mathbf{z}_t^- . A similar ensemble approach has been explored in Liu et al. (2021) and Li et al. (2023) for controllable and open-ended text generation, though their ensembles are based on predictions from different models. Another similar work to ours is CAD (Shi et al., 2023a), which examines scenarios where the model’s parametric knowledge contradicts non-parametric knowledge. CAD essentially constitutes a contrastive ensemble between \mathbf{z}_t and \mathbf{z}_t^+ . In this study, we concentrate on the general case of open-domain question answering, proposing a dynamic adjustment of α , controlling the degree of modification without treating it as a fixed hyperparameter. We provide an illustration of our method in Figure 1.

Dynamic α In prior logit adjustment methods (Liu et al., 2021; Malkin et al., 2022; O’Neill et al., 2023; Shi et al., 2023a; Pozzobon et al., 2023), α remains a fixed hyperparameter, requiring exhaustive search within the parameter space. Our innovation lies in dynamically setting α at each time step t without supervision, enabling fine-grained token-level adjustments. We estimate LLM confidence

258 following Jiang et al. (2021) by computing the high-
259 est probability from the normalized predicted token
260 probabilities at each step:

$$261 \quad C = \max_{y' \in V} P_\theta(y' | \mathbf{x}, \mathbf{y}_{<t}).$$

262 Similarly, we estimate LLM confidence using
263 relevant non-parametric knowledge c^+ :

$$264 \quad C_R = \max_{y' \in V} P_\theta(y' | c^+, \mathbf{x}, \mathbf{y}_{<t}).$$

265 At each time step, the value of α is determined
266 as follows:

$$267 \quad \alpha = \begin{cases} 1 - C_R, & \text{if } C > C_R, \\ C_R, & \text{otherwise.} \end{cases}$$

268 Our rationale is that higher LLM confidence in
269 parametric knowledge warrants minor adjustments,
270 while greater confidence in relevant non-parametric
271 knowledge necessitates more substantial modifica-
272 tions to the parametric answer. Note that we use
273 $1 - C_R$ instead of using $C - C_R$ to avoid the case
274 where both C and C_R are low. In such case, a
275 larger modification is still desired.

276 **Selection of c^+ and c^-** Choosing relevant con-
277 text c^+ is straightforward and we follow the
278 retrieval-augmented LLM literature where we use
279 the top retrieved texts from a retrieval module by
280 running our input query over an external knowl-
281 edge base. However, selecting irrelevant context
282 c^- is not trivial. Potential methods include using
283 lower-ranked retrievals, random text, or even delib-
284 erately crafted adversarial text. The primary aim
285 of c^- is to provide adversarial knowledge to elicit
286 incorrect predictions that can be disregarded from
287 the final token distribution. We explore various
288 strategies for selecting c^- in Section 5.3.

289 4 Experimental Setup

290 The present study revolves around open-domain
291 question answering, which involves tasking models
292 to generate responses to factual questions in natural
293 language. Specifically, we concentrate on the *open-*
294 *book* QA setting (Roberts et al., 2020), where we
295 harness non-parametric knowledge by supplying
296 relevant contexts along with the question itself to
297 the model during inference. Consistent with prior
298 investigations, we utilize prompting techniques to
299 assess the models’ performance.

4.1 Datasets and Metrics 300

Datasets Our method undergoes evaluation us-
301 ing three popular QA benchmarks: TriviaQA (Joshi
302 et al., 2017), Natural Questions (NQ; Kwiatkowski
303 et al. 2019), and PopQA (Mallen et al., 2023). Triv-
304 iaQA comprises trivia questions sourced from the
305 Web, whereas NQ consists of questions derived
306 from actual Google search queries, with answer
307 spans located in Wikipedia articles identified by
308 annotators. PopQA is a novel entity-centric open-
309 domain QA dataset covering factual information
310 about entities across a spectrum of popularity, in-
311 cluding *long-tail* knowledge often overlooked in
312 other popular QA datasets. 313

Metrics In line with prior research, our primary
314 metric for evaluating performance is the exact
315 match (EM), which determines whether the pre-
316 dicted sequence matches precisely with one of the
317 correct answers provided within the dataset. 318

4.2 Baselines and Models 319

Baselines Baseline approaches include regular
320 decoding with greedy decoding, following prior
321 work (Izacard and Grave, 2021). We prompt the
322 model for an answer by providing contextual infor-
323 mation. While our primary focus remains on the
324 *open-book* QA setting, we also present a baseline
325 employing the *closed-book* QA setting, where the
326 prompt consists solely of questions. This explo-
327 ration aims to scrutinize the parametric knowledge
328 of LLM. Additionally, we compare our method to
329 CAD, which accentuates the difference in output
330 probabilities when employing a model with and
331 without context. 332

Models Our decoding method undergoes eval-
333 uation across models varying in scale: Flan-T5
334 (XL-3B, XXL-11B; Chung et al. 2022), Falcon
335 (7B, 40B; Almazrouei et al. 2023), OPT (6.7B,
336 13B, 30B, 66B; Zhang et al. 2022), Llama (7B,
337 13B, 33B, 65B; Touvron et al. 2023a), and Llama
338 2 (7B, 13B, 70B; Touvron et al. 2023b), without
339 additional fine-tuning. 340

Instructions We employ a straightforward tem-
341 plate, i.e., “Answer the following question.
342 Question: <question> Answer:”, to for-
343 mat all questions for generative prediction
344 in the closed-book setting. For the open-
345 book setting, the template becomes “Answer
346 the question based on the context
347 below. Context: <context> Question:
348

<question> Answer:”. Although more sophisticated prompts were trialed in preliminary experiments, their marginal improvement over the simple template did not warrant their use, especially considering the risk of overfitting the model. In alignment with prior work (Chung et al., 2022), we employ 5-shot prompting for all models.

Retrieval models As previously mentioned, we explore a retrieval-augmented LLM approach in the open-book setting. This involves running an off-the-shelf retrieval system offline to obtain relevant context from Wikipedia for each query¹, which is then concatenated with the original query. We utilize two widely-used retrieval systems: BM25 (Robertson and Zaragoza, 2009) and Contriever (Izacard et al., 2022). BM25 operates as a static term-based retriever without training, while Contriever is pre-trained on extensive unlabeled corpora. In this study, we leverage Contriever-MS MARCO, a Contriever fine-tuned on MS MARCO (Bajaj et al., 2018). Consistent with Mullen et al. (2023), we utilize the top one retrieved paragraph. Additionally, TriviaQA and NQ datasets provide gold contexts, which we employ to measure the theoretical upper bound of our proposed decoding method. We also investigate the impact of using different retrieval methods in Section 5.2.

Setting alpha Our approach introduces a hyperparameter α to govern the degree of modification atop LLM’s parametric knowledge. For CAD, after a grid search using the validation set, we set $\alpha = 0.5$. In fixed alpha experiments for our method, we set $\alpha = 1.0$. In dynamic alpha, we do not have to set alpha values explicitly. We delve into the effect of α on our method in Section 5.4.

5 Results

We present the results of models featuring the largest variants in Table 1. Notably, employing regular decoding within an open-book setting consistently outperforms the closed-book setting across most models. This inclination suggests that LLM systems require non-parametric knowledge to excel in tasks demanding substantial knowledge assimilation. Interestingly, the performance of Llama 65B and Llama 2 70B in the closed-book setting surpasses that in the open-book setting concerning TriviaQA, indicating these models’ proficiency in factual knowledge retention without resorting

¹We utilize the Wikipedia dump from 2018.

Model	Decoding	NQ	TQA	PopQA
Flan-T5 11B	Reg.-Cl.	14.82	40.5	13.98
	Reg.-Op.	57.84	79.36	31.16
	CAD	47.56	66.08	26.28
	Ours-F	59.58	76.75	31.37
	Ours-D	63.16	80.09	34.64
Falcon 40B	Reg.-Cl.	28.56	71.74	28.79
	Reg.-Op.	53.32	72.05	39.16
	CAD	49.36	20.72	35.31
	Ours-F	53.77	79.56	39.87
	Ours-D	50.53	80.73	38.28
OPT 66B	Reg.-Cl.	13.71	39.65	15.62
	Reg.-Op.	48.73	62.38	34.77
	CAD	45.93	24.51	33.45
	Ours-F	51.97	68.11	34.83
	Ours-D	44.41	63.89	33.44
Llama 65B	Reg.-Cl.	34.13	75.72	35.9
	Reg.-Op.	55.32	74.76	40.31
	CAD	48.03	24.51	31.97
	Ours-F	57.01	76.61	39.9
	Ours-D	52.35	80.28	40.58
Llama-2 70B	Reg.-Cl.	37.87	79.69	40.98
	Reg.-Op.	56.07	76.07	42.7
	CAD	47.53	31.36	33.05
	Ours-F	58.86	78.38	42.59
	Ours-D	55.24	81.7	44.3

Table 1: Results pretrained model using gold retrieval (NQ, TriviaQA), and Contriever retrieval (PopQA). Reg.-Cl. refers to regular decoding with *closed-book* setting (i.e. no retrieval). Reg.-Op. refers to regular decoding with *open-book* setting (i.e. with retrieval). Ours-F refers to our method utilizing a fixed alpha, while Ours-D designates our method incorporating a dynamic alpha.

to non-parametric knowledge. This finding possibly implies that TriviaQA, being the oldest dataset among the three, potentially overlaps with the training data of these LLMs.

Crucially, our proposed decoding approach demonstrates superior performance across all three datasets compared to both regular decoding and CAD. Noteworthy variations exist in the efficacy of employing either the fixed alpha strategy or the dynamic alpha strategy; while in certain instances the fixed alpha approach exhibits better performance, the dynamic alpha approach outperforms in others. In subsequent references within this paper, when mentioning our method, we refer to the setting that delivers superior performance based on Table 1, without explicitly specifying whether it involves dynamic or fixed alpha.

5.1 Effect of Model Scaling

Thus far, our study has elucidated the efficacy of our proposed decoding approach across diverse model families. This segment aims to examine

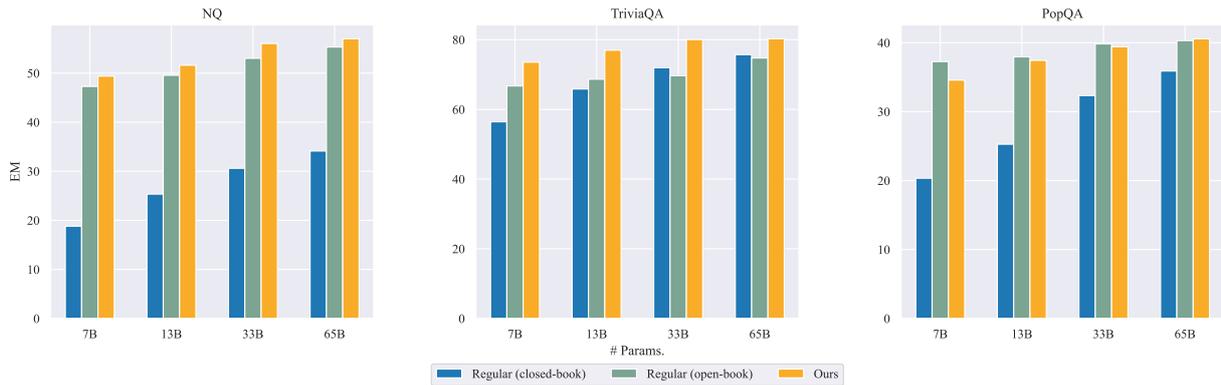


Figure 2: Performance comparison of our method against regular decoding across various sizes of Llama 1 models.

the impact of scaling the model’s parameter count on our methods. The results pertaining to Llama variants—specifically, Llama 7B, 13B, 33B, and 65B—are illustrated in Figure 2. We provide the results of scaling for other models in Appendix A. An observable trend emerges wherein, with an increase in model size, the disparity between closed-book and open-book performance diminishes, indicating that larger models possess greater potential for assimilating parametric knowledge. Furthermore, our decoding method consistently outperforms regular decoding across all model sizes, save for a few instances in the case of PopQA with smaller model variants. We posit this discrepancy to the absence of gold context within the PopQA dataset, leading to reliance on Contriever’s retrieval, which may occasionally introduce inaccuracies.

5.2 Using Different Retrievals

As previously highlighted, our investigation centers on retrieval-augmented LLMs, involving the implementation of retrieval modules over a knowledge base concerning a user query. Subsequently, the retrieved relevant passage supplements the prompt to facilitate the generation of answers by the LLM. In earlier experiments, we utilized the provided gold context by NQ and TriviaQA to establish the theoretical upper bound of our proposed decoding method. This segment aims to examine whether the utilization of off-the-shelf retrieval mechanisms would influence the efficacy of our proposed methods. In Figure 3, we present a comparative analysis between closed-book regular decoding and our decoding method, utilizing retrieval passages from BM25, Contriever, or the provided gold context.

It is pertinent to note that the PopQA dataset lacks gold context. The comparative analysis indicates that results derived from Contriever ex-

hibit superiority over those derived from BM25. Moreover, a substantial disparity exists between outcomes obtained through retrieval and those derived from leveraging gold context. It is essential to underscore that while these observations do not negate the efficacy of our proposed decoding method, they do suggest that enhancements to the retrieval module could yield improved outcomes.

Irr. Passage	NQ	TQA	PopQA
Random	56.74	81.28	43.23
Fixed	57.95	80.84	43.82
Fixed (rand. perm.)	57.17	80.68	42.98
Most distant	58.86	81.7	44.3

Table 2: Comparison of performance on Llama 2 70B across various methods for selecting irrelevant c^- : random selection, fixed adversarially constructed contexts, fixed context with random word permutation, and passages with the most distance from the relevant context.

5.3 Selection of Irrelevant Context

An essential aspect of our decoding method involves the incorporation of the c^- irrelevant context. Here, we investigate various strategies for selecting c^- and its impact on our methods. Initially, we propose employing a random selection of c^- from the complete pool of available contexts (ensuring that the randomly selected c^- differs from c^+). Subsequently, we manually construct an adversarial c^- devoid of meaningful or useful information, details of which are provided in Appendix B. Additionally, we experiment with shuffling the word order within this fixed c^- . Another approach for determining c^- involves using lower-ranked retrievals. However, increasing the retrieval size arbitrarily is computationally inefficient, and

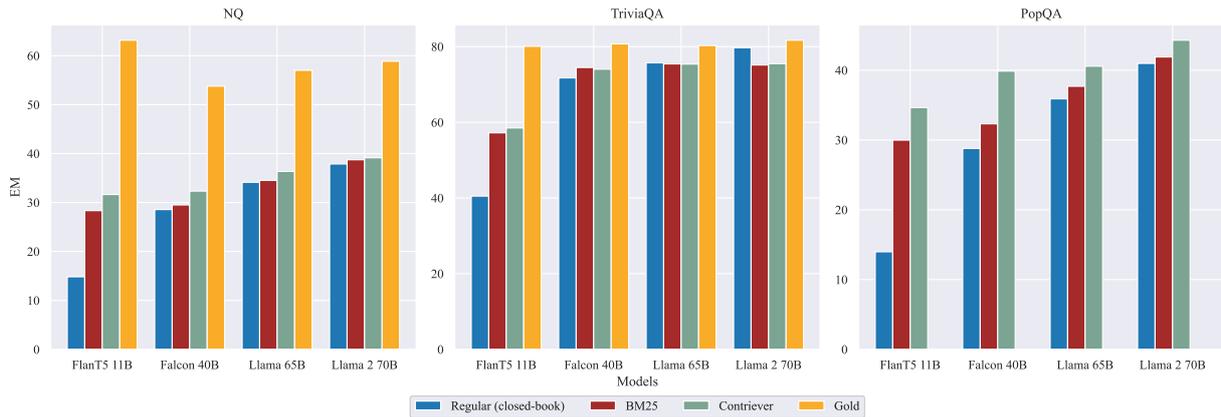


Figure 3: Performance comparison between regular decoding and our method using different sources of retrievals.

even within the top-100 retrievals, relevant information can be present. Therefore, we approximate the bottom-ranked retrieval by selecting the c^- that exhibits the most distance from c^+ , based on the cosine distance of their embeddings in the retrieval module. The comparison results using Llama 2 70B are presented in Table 2. It is evident that c^- with the most distance yields the best performance. Throughout our experiments detailed in this study, if not explicitly specified, we employ the most distant option for selecting c^- . However, if computing distance proves computationally expensive, the use of a fixed adversarial c^- , as demonstrated in our results, remains a viable alternative.

5.4 Adjusting the Knowledge Modification

Our proposed decoding method introduces the hyperparameter α , regulating the degree of modification applied to the parametric answer for a given input query. A larger α signifies a more substantial modification, while $\alpha = 0$ denotes no alteration, thereby reducing decoding to a regular decoding scenario. Despite outlining a strategy to dynamically set this alpha value, we remain interested in assessing the impact of different alpha values on the efficacy of our method. We conducted experiments involving the adjustment of α levels and present the outcomes obtained from Llama models in Figure 5. Our analysis reveals that as the alpha values increase, the effectiveness of the method diminishes substantially. Conversely, setting $\alpha = 1.0$ consistently yields substantial and robust improvements over regular decoding across all three datasets.

5.5 Answering across Knowledge Popularity

The utility of retrieval mechanisms becomes evident in addressing less prevalent factual knowl-

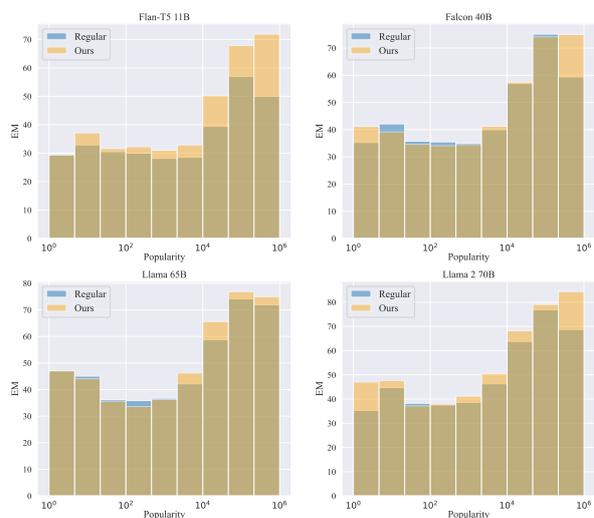


Figure 4: Comparison of performance between regular decoding (open-book) and our method on questions with varying levels of knowledge popularity.

edge, an area where LLMs often exhibit limitations. Therefore, we conducted an analysis to evaluate the efficacy of our proposed decoding approach in facilitating LLMs to accurately respond to factual questions across a spectrum of popularity. Following [Mallen et al. \(2023\)](#), we utilized the popularity of entities gauged by Wikipedia’s monthly page views as an indicator of their frequency in web discussions. Our findings, presented in Figure 4, juxtapose the performance of models employing regular decoding within an open-book setting against those employing our proposed method. The results manifest a consistent trend wherein our proposed method consistently outperforms regular decoding under an open-book setting across varying levels of popularity. This observation underscores the efficacy of our decoding strategy in assisting LLMs to generate more accurate responses to factual queries across a diverse range of entity popularities.

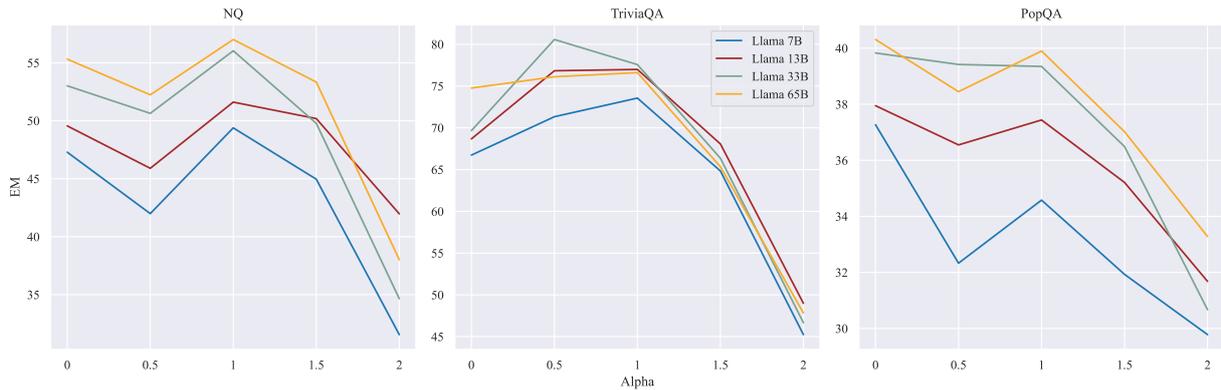


Figure 5: The impact of α on our decoding method across different sizes of Llama 1 models.

Model	Decoding	NQ-SUB
Flan-T5 11B	Reg.-Cl.	0.19
	Reg.-Op.	56.4
	CAD	51.9
	Ours	57.55
Falcon 40B	Reg.-Cl.	0.13
	Reg.-Op.	46.78
	CAD	45.79
	Ours	48.34
Llama 65B	Reg.-Cl.	0.08
	Reg.-Op.	59.25
	CAD	60.41
	Ours	61.65
Llama-2 70B	Reg.-Cl.	0.02
	Reg.-Op.	57.63
	CAD	53.23
	Ours	58.34

Table 3: Comparison of decoding methods on the knowledge conflict dataset. Reg.-Cl. and Reg.-Op. denote regular decoding in closed-book and open-book settings.

5.6 Resolving Knowledge Conflicts

As previously highlighted in the manuscript, tasks reliant on knowledge typically draw from two knowledge sources: parametric knowledge, acquired during training, and non-parametric knowledge, accessed via retrieval modules during inference. The issue of knowledge conflicts, wherein the contextual (non-parametric) information contradicts learned knowledge, has been formally addressed by Longpre et al. (2021) to understand how models utilize these dual sources of knowledge.

To generate question-answer pairs manifesting knowledge conflicts, we followed the methodology proposed by Longpre et al. (2021). Initially, we identified questions in the NQ dataset that contained named entity answers. Subsequently, we obtained the relevant context for each question and replaced the gold answer entity in the context with a random entity. In this setup, an accurate LLM

should produce the substituted entity as the answer when provided with the question and the modified context, disregarding its pre-learned parametric answer. This resulting dataset, termed NQ-SUB, was created for assessing models in scenarios involving knowledge conflicts. The performance results on NQ-SUB are presented in Table 3. Remarkably, all models exhibited poor performance in the regular closed-book setting, given that the task requires the model to disregard its parametric knowledge. However, our proposed decoding method demonstrated superior performance compared to both regular decoding and CAD on this knowledge conflict task.² The comparative results emphasize the effectiveness of our proposed decoding approach in addressing knowledge conflicts, particularly in scenarios where models encounter contradictions between their learned and contextual knowledge.

6 Conclusion

This study introduces a novel decoding strategy, employing contrastive decoding to incorporate relevant and irrelevant context. Through diverse experiments and analyses across datasets and model scales, our approach consistently outperforms regular decoding methods. Notably, it excels in managing knowledge conflicts, surpassing both regular decoding and CAD. Moreover, our exploration of retrieval sources underscores the need for refining these modules to enhance efficacy. The demonstration of the method’s effectiveness in open-domain question answering also sets the stage for future research. The method’s versatility suggests potential applications in various generative tasks, motivating our future exploration in tasks like summarization.

²The CAD results were based on our implementation, due to the unavailability of the CAD at the time of our study.

586 Limitations

587 Our study acknowledges several limitations that
588 warrant consideration. First, we acknowledge
589 the restriction imposed by employing a singular
590 prompt template. The computational complexity
591 inherent in our method limited the scope of experi-
592 ments conducted within this framework. However,
593 this constraint was pivotal in maintaining consis-
594 tency across our comparisons, ensuring the reliabil-
595 ity and robustness of the obtained results despite
596 the limitation in the number of explored templates.

597 Secondly, while our decoding method was
598 specifically showcased in the context of Question
599 Answering (QA) using greedy decoding, it’s es-
600 sential to note that our approach is designed as a
601 general decoding framework applicable to various
602 generative tasks. Although our focus was on QA,
603 expanding this work to encompass other domains
604 like summarization and mitigating language hal-
605 lucination remains a promising avenue for future
606 exploration.

607 Furthermore, it’s imperative to recognize that
608 the scalability and generalizability of our method
609 across different problem domains and decoding
610 strategies might present further challenges and con-
611 siderations. Extending our investigation to encom-
612 pass a broader array of prompt templates and de-
613 coding strategies (such as nucleus sampling) could
614 potentially reveal nuanced insights into the adapt-
615 ability and effectiveness of our proposed method.

616 Additionally, it is crucial to note that the de-
617 coding time required for our method is longer than
618 regular decoding, approximately three times longer,
619 owing to decoding using three logits distributions
620 simultaneously. However, there exists potential
621 for mitigating the time complexity by distributing
622 the decoding of different distributions across mul-
623 tiple GPU machines, thereby enabling paralleliza-
624 tion and potentially reducing the computational
625 overhead. This approach might alleviate the time
626 constraints associated with our decoding method,
627 rendering it more feasible for practical application
628 in real-time scenarios.

629 References

630 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Al-
631 shamsi, Alessandro Cappelli, Ruxandra Cojocaru,
632 Merouane Debbah, Etienne Goffinet, Daniel Heslow,
633 Julien Launay, Quentin Malartic, et al. 2023. Falcon-
634 40b: an open large language model with state-of-

the-art performance. Technical report, Technology 635
Innovation Institute. 636

Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, 637
Jianfeng Gao, Xiaodong Liu, Rangan Majumder, An- 638
drew McNamara, Bhaskar Mitra, Tri Nguyen, Mir 639
Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, 640
and Tong Wang. 2018. *Ms marco: A human gener- 641
ated machine reading comprehension dataset.* 642

Hung-Ting Chen, Michael Zhang, and Eunsol Choi. 643
2022. *Rich knowledge sources bring complex knowl- 644
edge conflicts: Recalibrating models to reflect con- 645
flicting evidence.* In *Proceedings of the 2022 Con- 646
ference on Empirical Methods in Natural Language 647
Processing*, pages 2292–2307, Abu Dhabi, United 648
Arab Emirates. Association for Computational Lin- 649
guistics. 650

Hyung Won Chung, Le Hou, Shayne Longpre, Barret 651
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi 652
Wang, Mostafa Dehghani, Siddhartha Brahma, Al- 653
bert Webson, Shixiang Shane Gu, Zhuyun Dai, 654
Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh- 655
ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, 656
Dasha Valter, Sharan Narang, Gaurav Mishra, Adams 657
Yu, Vincent Zhao, Yanping Huang, Andrew Dai, 658
Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja- 659
cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, 660
and Jason Wei. 2022. *Scaling instruction-finetuned 661
language models.* 662

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu- 663
pat, and Ming-Wei Chang. 2020. *Realm: Retrieval- 664
augmented language model pre-training.* In *Proceed- 665
ings of the 37th International Conference on Machine 666
Learning, ICML’20.* JMLR.org. 667

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Se- 668
bastian Riedel, Piotr Bojanowski, Armand Joulin, 669
and Edouard Grave. 2022. *Unsupervised dense in- 670
formation retrieval with contrastive learning.* *Trans. 671
Mach. Learn. Res.*, 2022. 672

Gautier Izacard and Edouard Grave. 2021. *Leveraging 673
passage retrieval with generative models for open do- 674
main question answering.* In *Proceedings of the 16th 675
Conference of the European Chapter of the Associ- 676
ation for Computational Linguistics: Main Volume,* 677
pages 874–880, Online. Association for Computa- 678
tional Linguistics. 679

Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas 680
Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi- 681
Yu, Armand Joulin, Sebastian Riedel, and Edouard 682
Grave. 2023. *Atlas: Few-shot learning with retrieval 683
augmented language models.* *Journal of Machine 684
Learning Research*, 24(251):1–43. 685

Joel Jang, Seonghyeon Ye, Changho Lee, Sohee Yang, 686
Joongbo Shin, Janghoon Han, Gyeonghun Kim, and 687
Minjoon Seo. 2022. *TemporalWiki: A lifelong 688
benchmark for training and evaluating ever-evolving 689
language models.* In *Proceedings of the 2022 Con- 690
ference on Empirical Methods in Natural Language 691*

692		<i>Processing</i> , pages 6237–6250, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
693			
694			
695	Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. 2021. How can we know when language models know? on the calibration of language models for question answering . <i>Transactions of the Association for Computational Linguistics</i> , 9:962–977.		
696			
697			
698			
699			
700	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.		
701			
702			
703			
704			
705			
706			
707	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In <i>Proceedings of the 40th International Conference on Machine Learning, ICML’23</i> . JMLR.org.		
708			
709			
710			
711			
712	Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui. 2022. Realtime qa: What’s the answer right now?		
713			
714			
715			
716	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.		
717			
718			
719			
720			
721			
722			
723			
724			
725	Patrick S. H. 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. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks . In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .		
726			
727			
728			
729			
730			
731			
732			
733			
734	Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 110–119, San Diego, California. Association for Computational Linguistics.		
735			
736			
737			
738			
739			
740			
741			
742	Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023. Contrastive decoding: Open-ended text generation as optimization . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12286–12312, Toronto, Canada. Association for Computational Linguistics.		
743			
744			
745			
746			
747			
748			
749			
	Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. DExperts: Decoding-time controlled text generation with experts and anti-experts . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6691–6706, Online. Association for Computational Linguistics.		750
			751
			752
			753
			754
			755
			756
			757
			758
			759
	Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7052–7063, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		760
			761
			762
			763
			764
			765
			766
			767
	Kelvin Luu, Daniel Khashabi, Suchin Gururangan, Karishma Mandyam, and Noah A. Smith. 2022. Time waits for no one! analysis and challenges of temporal misalignment . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5944–5958, Seattle, United States. Association for Computational Linguistics.		768
			769
			770
			771
			772
			773
			774
			775
			776
	Nikolay Malkin, Zhen Wang, and Nebojsa Jojic. 2022. Coherence boosting: When your pretrained language model is not paying enough attention . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8214–8236, Dublin, Ireland. Association for Computational Linguistics.		777
			778
			779
			780
			781
			782
			783
	Alex Mullen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.		784
			785
			786
			787
			788
			789
			790
			791
	Sewon Min, Weijia Shi, Mike Lewis, Xilun Chen, Wen-tau Yih, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2023. Nonparametric masked language modeling . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 2097–2118, Toronto, Canada. Association for Computational Linguistics.		792
			793
			794
			795
			796
			797
	Ella Neeman, Roei Aharoni, Or Honovich, Leshem Choshen, Idan Szpektor, and Omri Abend. 2023. DisentQA: Disentangling parametric and contextual knowledge with counterfactual question answering . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10056–10070, Toronto, Canada. Association for Computational Linguistics.		798
			799
			800
			801
			802
			803
			804
			805
	Charles O’Neill, Yuan-Sen Ting, Ioana Ciuca, Jack Miller, and Thang Bui. 2023. Steering language generation: Harnessing contrastive expert guidance and		806
			807
			808

809	negative prompting for coherent and diverse synthetic data generation.		
810			
811	Luiza Pozzobon, Beyza Ermis, Patrick Lewis, and Sara Hooker. 2023. Goodtriever: Adaptive toxicity mitigation with retrieval-augmented models.		
812			
813			
814	Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5418–5426, Online. Association for Computational Linguistics.		
815			
816			
817			
818			
819			
820	Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. <i>Found. Trends Inf. Retr.</i> , 3(4):333–389.		
821			
822			
823	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen tau Yih. 2023a. Trusting your evidence: Hallucinate less with context-aware decoding.		
824			
825			
826			
827	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023b. Replug: Retrieval-augmented black-box language models.		
828			
829			
830			
831	Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3784–3803, Punta Cana, Dominican Republic. Association for Computational Linguistics.		
832			
833			
834			
835			
836			
837			
838	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models.		
839			
840			
841			
842			
843			
844	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models.		
845			
846			
847			
848			
849			
850			
851			
852			
853			
854			
855			
856			
857			
858			
859			
860			
861			
862			
863			
864			
865			
866			
		Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. Opt: Open pre-trained transformer language models.	867
			868
			869
			870
			871
			872
			873
		Zexuan Zhong, Tao Lei, and Danqi Chen. 2022. Training language models with memory augmentation. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5657–5673, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	874
			875
			876
			877
			878
			879
		Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhao Chen. 2023. Context-faithful prompting for large language models.	880
			881
			882

A Additional Results on Scaling Experiments

We present additional scaling experiment results for different model variants. Specifically, we illustrate the outcomes for Flan-T5 variants, 3B and 11B, in Figure 6. The results for Falcon variants, particularly Falcon 7B and 40B, are depicted in Figure 7. Moreover, we showcase the results for OPT variants, encompassing OPT 6.7B, 13B, 30B, and 66B, in Figure 8. Additionally, the findings pertaining to Llama 2 variants, including Llama 2 7B, 13B, and 70B, are illustrated in Figure 9. We can see that our proposed decoding method outperforms regular decoding with open-book setting in most settings across different datasets and model sizes.

B Additional Details on Irrelevant Context

Here we provide the meticulously designed adversarial c^- irrelevant context that is used as the fixed c^- for every query:

“It was a pleasant weather day, with seasonally average temperatures. The local legislative and academic governing bodies held routine meetings regarding budgets and policies. Students focused on their studies while athletes practiced for upcoming competitions. Residents tended to their jobs and daily tasks around their neighborhood. Nothing particularly eventful occurred in the community. It was an ordinary midweek day. The weather was typical for the time of year without any extreme events. Overall it was an average day in the community with people pursuing their regular daily activities.”

Here is the same fixed c^- but with word order permuted:

“an routine Overall was of community. average focused for The around tended upcoming their was policies. their budgets and Residents to eventful held competitions. It particularly extreme with academic temperatures. was day. weather local The their studies events. it meetings average pleasant typical Nothing ordinary time seasonally legislative people an the daily the Students in a neighborhood. activities. community pursuing weather and while in midweek regarding athletes occurred tasks the daily jobs It governing year bodies regular with their for day and practiced on day, was without any”

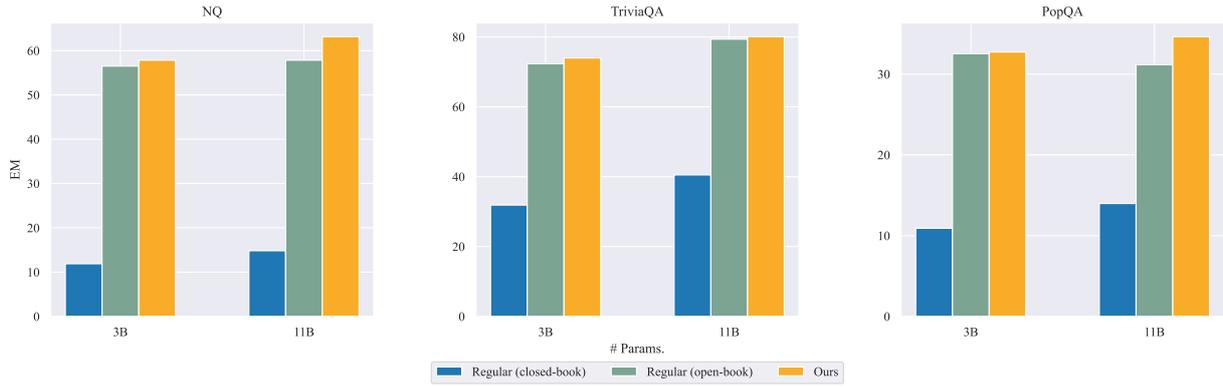


Figure 6: Performance comparison of our method against regular decoding across various sizes of Flan-T5 models.

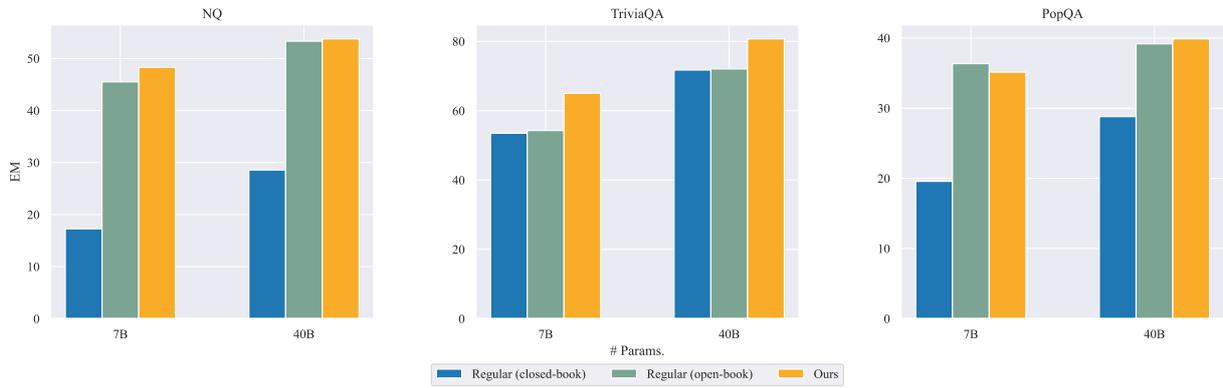


Figure 7: Performance comparison of our method against regular decoding across various sizes of Falcon models.

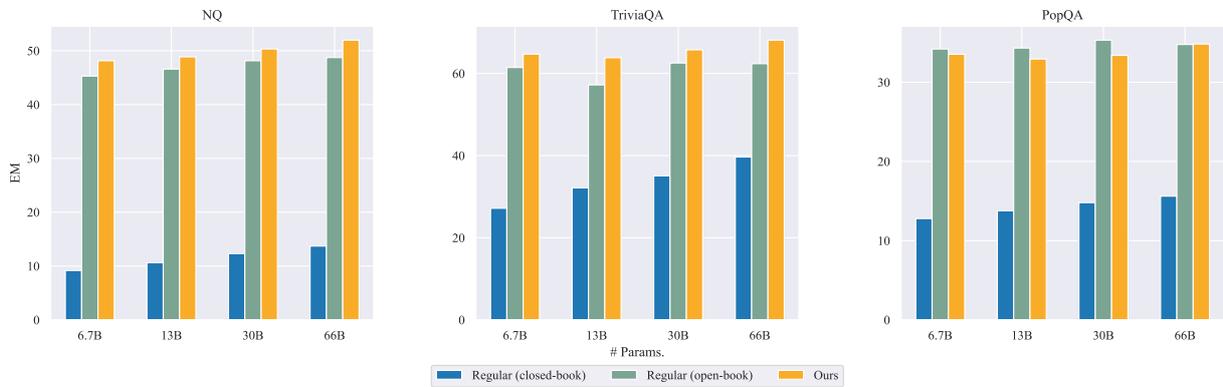


Figure 8: Performance comparison of our method against regular decoding across various sizes of OPT models.

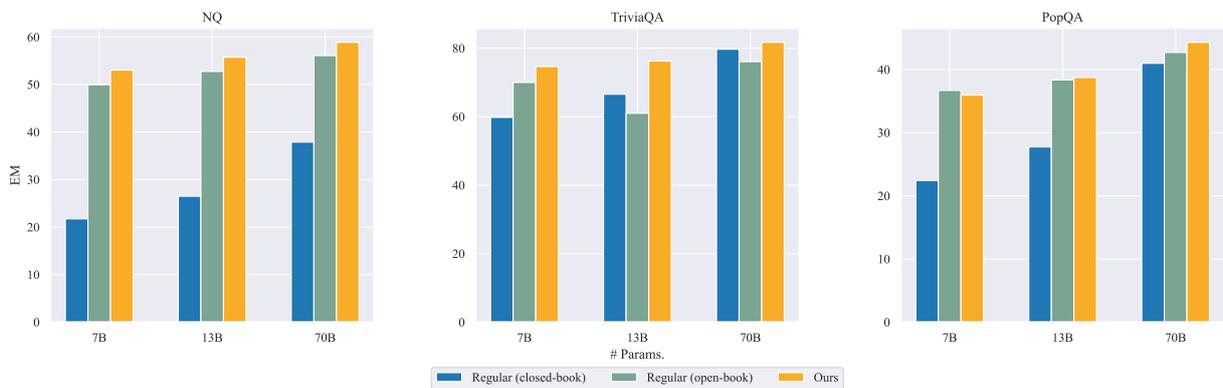


Figure 9: Performance comparison of our method against regular decoding across various sizes of Llama 2 models.