

ZEMI: LEARNING ZERO-SHOT SEMI-PARAMETRIC LANGUAGE MODELS FROM MULTIPLE TASKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Although large language models have achieved impressive zero-shot ability, the huge model size generally incurs high cost. Recently, semi-parametric language models, which augment a smaller language model with an external retriever, have demonstrated promising language modeling capabilities. However, it remains unclear whether such semi-parametric language models can perform competitively well as their fully-parametric counterparts on zero-shot generalization to downstream tasks. In this work, we introduce **Zemi**, a zero-shot semi-parametric language model. To our best knowledge, this is the first semi-parametric language model that can demonstrate strong zero-shot performance on a wide range of held-out unseen tasks. We train Zemi with a novel **semi-parametric multitask prompted training** paradigm, which shows significant improvement compared with the parametric multitask training as proposed by T0 (Sanh et al., 2021). Specifically, we augment the multitask training and zero-shot evaluation with retrieval from a large-scale task-agnostic unlabeled corpus. In order to incorporate multiple potentially noisy retrieved augmentations, we further propose a novel **augmentation fusion** module leveraging perceiver resampler and gated cross-attention. Notably, our proposed Zemi_{LARGE} outperforms T0-3B by 16% on all seven evaluation tasks while being 3.9x smaller in model size.¹

1 INTRODUCTION

Achieving strong generalization ability on unseen tasks while maintaining a reasonably small parameter size is a long-lasting challenge for natural language processing (NLP) models. Although large language models (Brown et al., 2020; Lieber et al., 2021; Rae et al., 2021; Smith et al., 2022; Hoffmann et al., 2022; Zhang et al., 2022; Ouyang et al., 2022; Chowdhery et al., 2022) have shown impressive zero-shot ability on various NLP tasks, the huge model size generally incurs high cost. Alternatively, instead of storing everything in the model parameters, recent work on semi-parametric language models (Grave et al., 2016; Khandelwal et al., 2019; Yogatama et al., 2021; Borgeaud et al., 2021; Zhong et al., 2022) demonstrated competitive **language modeling performance** compared with much larger fully-parametric language models. The intuition is to use a relatively small language model as a reasoning module and augment it with a retrieval system based on some external corpora. Making the model semi-parametric alleviates the need to keep increasing the model size for aligning with the growing data size. However, it is still unclear whether such semi-parametric language models can have strong **zero-shot ability** on unseen tasks as their fully-parametric counterparts such as T0 (Sanh et al., 2021) and FLAN (Wei et al., 2021). Furthermore, improvements in language modeling metrics such as perplexity may not guarantee better performance on downstream tasks especially in low-shot settings (Wei et al., 2022). Thus, in this work, we aim to investigate this unexplored research question, *can semi-parametric language models exhibit strong zero-shot generalization abilities on various types of downstream tasks?*

To this end, we introduce *Zemi*, a semi-parametric language model for zero-shot task generalization. To the best of our knowledge, this is the first semi-parametric language model that shows strong zero-shot performance on a wide range of downstream tasks. In order to effectively train Zemi, we propose a novel training paradigm, *semi-parametric multitask prompted training* (Section 3.1), inspired by T0. Specifically, we augment the multitask training and zero-shot evaluation with retrieved

¹Code and data are available in the supplementary material.

plain text documents. In order to cover a wider coverage of unseen tasks, instead of retrieving from specific corpora for certain tasks, such as exploiting Wikipedia for open-domain question answering (Lee et al., 2019; Karpukhin et al., 2020; Izacard & Grave, 2020), we retrieve from a large-scale task-agnostic corpus, C4 (Raffel et al., 2020) (Section 3.2). Notably, C4 is the unlabeled pre-training corpus of our backbone model (Raffel et al., 2020), which means that we do not require any annotated or curated resources for a specific task. In our preliminary experiments, we find that existing methods (Izacard & Grave, 2020; Brown et al., 2020) that incorporate textual augmentations cannot effectively handle the noise inevitably introduced by information retrieval from large-scale corpora. To address this challenge, we further propose a novel semi-parametric model architecture that can pay attention to salient information of the retrieved augmentations and selectively ignore the noisy ones. Particularly, we borrow ideas from recent advancements in *vision-language fusion* (Alayrac et al., 2022) and introduce a light-weight *augmentation fusion* module between the pre-trained text encoder and decoder (Section 3.3).

Experimental results show that semi-parametric multitask prompted training can effectively improve zero-shot task generalization ability of our proposed Zemi models (Section 4.2). Comparison with state-of-the-art methods shows that Zemi_{LARGE} significantly outperforms its fully-parametric counterpart T0-3B by 16% on all seven evaluation tasks while being 3.9x smaller in scale (Section 4.3).

To sum up, the main contributions of this paper are threefold:

- We introduce **Zemi**, the first semi-parametric language model that demonstrates strong zero-shot performance on a wide range of unseen downstream tasks.
- We train Zemi with a novel training paradigm, **semi-parametric multitask prompted training**, which significantly improves zero-shot task generalization compared with parametric multitask training.
- We propose a novel **augmentation fusion** module to effectively handle multiple potentially noisy retrieved documents.

2 RELATED WORK

2.1 SEMI-PARAMETRIC MODELS

Traditionally, semi-parametric models have been widely applied to specific knowledge-intensive NLP tasks such as open-domain question answering (Chen et al., 2017; Lee et al., 2019; Guu et al., 2020; Wang et al., 2019; Karpukhin et al., 2020; Yang et al., 2019; Lewis et al., 2020; Izacard & Grave, 2020). One representative paradigm is the retriever-reader framework where the retriever first retrieves top k (e.g., 100) relevant passages from a large corpus, and then the reader either generates or extracts the answer conditioned on the query and the retrieved passages. One limitation of those traditional retriever-reader models is that they are designed to solve specific tasks, for example, retrieving from Wikipedia for answering trivia questions (Joshi et al., 2017). In this work, we retrieve from a large-scale **task-agnostic** corpus, C4 (Raffel et al., 2020), which is the unlabeled pre-training corpus of our backbone model.

Recent advancements in semi-parametric language models (Khandelwal et al., 2019; Yogatama et al., 2021; Borgeaud et al., 2021; Zhong et al., 2022) focus on improving language modeling with a relatively small language model and a retrieval system based on a large-scale corpus. Khandelwal et al. (2019) and Yogatama et al. (2021) share the idea of extending a pretrained language model by linearly interpolating it with a k -nearest neighbor model. Borgeaud et al. (2021) proposed to improve language model pretraining with a retrieval-enhanced transformer with chunked cross-attention layers. Although the aforementioned semi-parametric language models have shown competitive performance on language modeling compared with fully-parametric counterparts such as GPT-3 (Brown et al., 2020), they are rarely evaluated on unseen downstream tasks, leaving it unexplored whether semi-parametric language models can perform competitively well on zero-shot task generalization. While concurrent work (Izacard et al., 2022) showed initial success in few-shot settings relying on Fusion-in-Decoder (FiD) (Izacard & Grave, 2020) framework, this work focus on the more challenging **zero-shot** settings. Furthermore, instead of reusing FiD framework as in Izacard et al. (2022), we show that our newly proposed semi-parametric architecture is more effective than FiD due to the **gated mechanism**.

2.2 MASSIVE MULTITASK PROMPTED TRAINING

Based on the assumption that the reasonable zero-shot ability of large language models may come from implicit multitask learning during pretraining, recent studies (Sanh et al., 2021; Wei et al., 2021; Ye et al., 2021; Wang et al., 2022c) have demonstrated that explicitly training a language model on a mixture of diverse tasks can effectively improve its zero-shot performance on unseen tasks. For example, a multitask fine-tuned 11B T0 model (Sanh et al., 2021) can outperform the 175B GPT-3 on a variety of held-out unseen tasks. Followup work (Wang et al., 2022b) shows that encoder-decoder model architecture is particularly suitable for multitask finetuning. In this work, we extend T0’s multitask prompted training to **semi-parametric multitask prompted training**, where we further augment the training and evaluation instances with retrieved documents. Notably, our work is distinguished from previous work ReCross (Lin et al., 2022), which uses upstream training data for augmentation, in twofold. First, we retrieve documents from a much larger task-agnostic corpus instead of clean upstream training instances. Second, in addition to directly concatenating the augmentation with the input just as FiD (Izacard & Grave, 2020), we further propose a novel augmentation fusion module to handle a larger number of potentially noisy retrieved augmentations.

2.3 FUSION OF RETRIEVED AUGMENTATIONS

In this work, the main challenge in designing the semi-parametric language model architecture is to attain the ability to handle multiple potentially noisy retrieved documents. Existing methods on incorporating external texts fall in two categories, *concatenation based* (Lin et al., 2022; Brown et al., 2020; Liu et al., 2021; Lewis et al., 2020; Wang et al., 2022a) and *cross-attention based* (Izacard & Grave, 2020; Prabhumoye et al., 2021; Borgeaud et al., 2021). Concatenation-based methods directly concatenate the retrieved texts with the input and use a unified self-attention to incorporate the augmentations. This method is highly expressive but the number of texts that can be added is strictly limited to the model’s maximum input length. Cross-attention-based methods can take into account a larger number of texts, but they still lack an explicit design for preventing the model from paying attention to the noisy ones. Inspired by recent visual language models Alayrac et al. (2022); Yu et al. (2022); Li et al. (2022); Jiang et al. (2022), we propose to leverage perceiver resamplers (Jaegle et al., 2021) and gated cross-attention layers for incorporating external texts. Although originally proposed for vision-language fusion in Flamingo (Alayrac et al., 2022), we find that the idea can be effectively transferred to do augmentation fusion with potentially noisy retrieved documents. Furthermore, we identify two key differences from Flamingo: first, we use a much smaller encoder-decoder model that is jointly trained with the newly initialized layers instead of frozen layers. Second, instead of inserting the gated cross-attention module into a large frozen language model (Hoffmann et al., 2022), we add only one layer of gated cross-attention on top of the encoder to alleviate the need for extensive pre-training.

3 METHOD

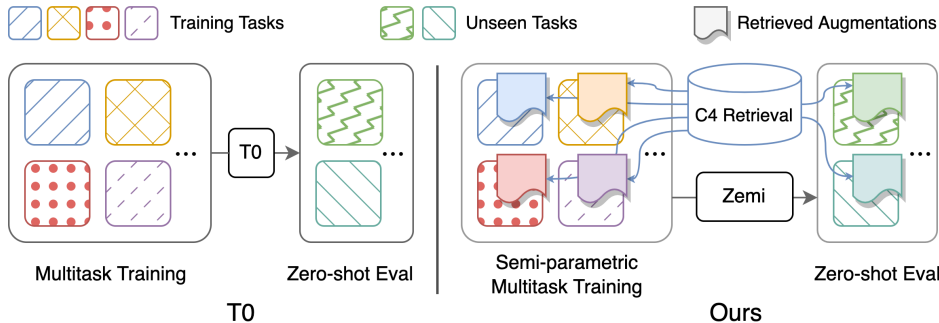


Figure 1: Overview of the semi-parametric multitask prompted training. Each training and evaluation instance is formatted with unified text-to-text prompt templates (Sanh et al., 2021; Bach et al., 2022). In this work, we further augment the prompted instances with retrieved passages from a large-scale task-agnostic corpus, C4 (Sanh et al., 2021). We use a novel semi-parametric language model to incorporate the retrieved augmentations.

3.1 SEMI-PARAMETRIC MULTITASK PROMPTED TRAINING

In this section, we introduce the novel training paradigm for **Zemi**, semi-parametric multitask prompted training. We follow the overall text-to-text framework proposed by the previous parametric multitask prompted training (Sanh et al., 2021) where each input-output pair of a certain task is converted into a prompted text input and a generated text output via human-written templates (Bach et al., 2022)². For Zemi, as illustrated in Figure 1, we further augment the multitask training and zero-shot evaluation with a retrieval system. Instead of using specific corpora for different tasks, such as Wikipedia for open-domain question answering (Chen et al., 2017; Karpukhin et al., 2020; Izacard & Grave, 2020) and textbooks for science question answering (Mihaylov et al., 2018), we retrieve texts from a large-scale task-agnostic corpus, **C4** (Raffel et al., 2020) (Section 3.2). Retrieving from a larger corpus brings wider coverage but also more noisy augmentations. To address this problem we further propose a novel semi-parametric architecture for Zemi that specializes in handling a variable number of potentially noisy augmentations (Section 3.3). After semi-parametric multitask prompted training, we perform zero-shot evaluation on seven diverse held-out unseen tasks (Section 4).

3.2 C4 RETRIEVAL

To build a universal semi-parametric language model that can generalize to various types of NLP tasks, we retrieve documents from Colossal Clean Crawled Corpus (C4). The C4 corpus (750GB in size) contains more than 364 million documents. Performing dense retrieval (Karpukhin et al., 2020) on such a wide-coverage corpus is very expensive. Thus, for efficiency consideration, we perform document-level indexing and retrieval based on BM25 (Robertson et al., 1995) with Elasticsearch (ElasticSearch) and Huggingface Datasets (Lhoest et al., 2021). Despite its simplicity, recent work (Wang et al., 2022a) has demonstrated the effectiveness of using BM25 for retrieving clean training data as augmentations. To further improve the retrieval efficiency, we use 5% of the entire C4 corpus, which is still 3x larger than the Wikipedia corpus (Foundation), as our retrieval corpus in our experiments. For each query, we truncate the query length at 20 tokens and truncate each retrieved document at 256 tokens. See details on the query fields for each dataset in Appendix A.4.

3.3 ZEMI MODEL ARCHITECTURE

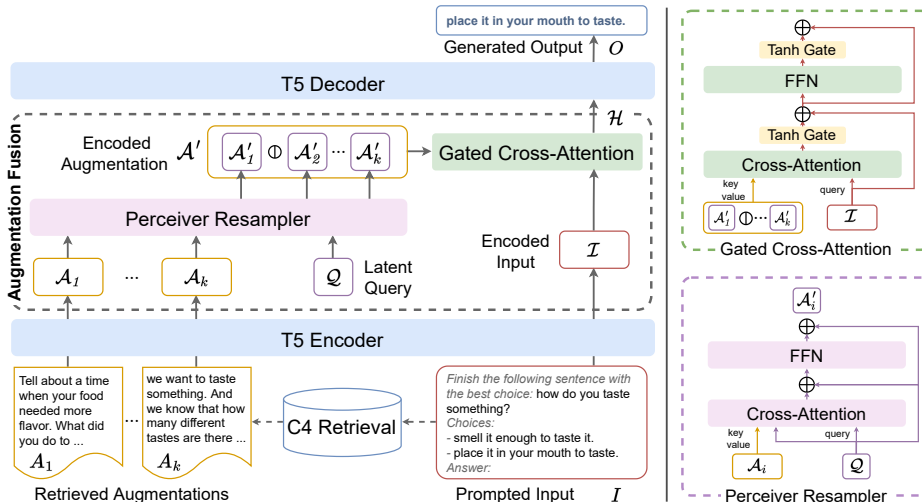


Figure 2: Zemi model architecture with an example of a prompted input and a generated output from the Piqa (Bisk et al., 2020) task. The *italic text* in the prompted input I indicates the prompt template. A_1 and A_k shows two examples of the corresponding retrieved augmentations (documents) from the C4 corpus. To incorporate the potentially noisy retrieved augmentations, we introduce a light-weight newly initialized augmentation fusion module that contains two major components, a single layer perceiver resampler and a single layer gated cross-attention (detailed on the right).

²<https://github.com/bigscience-workshop/promptsources>.

One major challenge of retrieving from a large-scale task-agnostic corpus is that the retrieved augmentations (documents) can be noisy and inaccurately ranked. Examples of good and noisy retrieved documents can be found in Appendix A.1. To address this problem, intuitively, we want the model to have the following two properties: (1) be able to simultaneously pay attention to multiple retrieved augmentations instead of only the top-1 result. (2) be able to identify salient information from the retrieved augmentations and selectively ignore uninformative ones.

To this end, we propose **Zemi**, a semi-parametric language model capable of selectively incorporating multiple potentially noisy retrieved augmentations. The main idea is to add a light-weight **Augmentation Fusion** module between the encoder and decoder, which contains two major components, a *perceiver resampler* and a *gated cross-attention*. The perceiver resampler and the gated cross-attention are adopted from Flamingo (Alayrac et al., 2022), which were originally proposed for vision-language fusion (detailed discussion can be found in Section 2.3).

Figure 2 shows an illustration of the Zemi model architecture. We consider a prompted text input I and a few retrieved augmentations A_1, A_2, \dots, A_k . Let l_I, l_A^i be the length of the prompted input and the i th augmentation, d be the hidden dimension of our backbone model. We first independently encode I and A_1, A_2, \dots, A_k with a shared T5 (Raffel et al., 2020) encoder Enc. We then feed the latent representation of the augmentations $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_k$ through the perceiver resampler.

$$\mathcal{I} = \text{Enc}(I); \quad \mathcal{I} \in \mathbb{R}^{l_I \times d} \quad (1)$$

$$\mathcal{A}_i = \text{Enc}(A_i), \quad \forall i \in \{1, \dots, k\}; \quad \mathcal{A}_i \in \mathbb{R}^{l_A^i \times d} \quad (2)$$

$$\mathcal{A}'_i = \text{PerceiverResampler}(\mathcal{A}_i, \mathcal{Q}), \quad \forall i \in \{1, \dots, k\}; \quad \mathcal{A}'_i \in \mathbb{R}^{l_Q \times d} \quad (3)$$

As shown on the bottom right of Figure 2, the perceiver resampler (Alayrac et al., 2022) is a variant of Perceiver IO (Jaegle et al., 2021), where a cross-attention is performed between a variable-length latent representation of an augmentation \mathcal{A}_i and a fixed-length learnable latent query vector \mathcal{Q} . Let l_Q be the predefined length of the latent query, which is typically smaller than the original length of an augmentation l_A^i . Thus, the output of the perceiver resampler is a compressed fixed-length latent representation of each augmentation. This resampling mechanism not only allows the model to address longer and a larger number of augmentations but also encourages the model to select salient information from the original augmentations.

After the resampler, we concatenate the encoded augmentations $\mathcal{A}' = [\mathcal{A}'_1, \dots, \mathcal{A}'_k] \in \mathbb{R}^{k \times l_Q \times d}$ and perform gated cross-attention with the encoded prompted input \mathcal{I} . As shown in the top right of Figure 2, the gated cross-attention module (Alayrac et al., 2022) contains two Tanh gates controlling the information flow from the cross-attention layer and the feed-forward layer before the addition with the skip connections, i.e., the original encoded input \mathcal{I} . Finally, the hidden states from the gated cross-attention module \mathcal{H} is fed into the T5 decoder Dec to generate the output sequence O .

$$\mathcal{H} = \text{GatedCrossAttn}([\mathcal{A}'_1, \dots, \mathcal{A}'_k], \mathcal{I}); \quad \mathcal{H} \in \mathbb{R}^{l_I \times d} \quad (4)$$

$$O = \text{Dec}(\mathcal{H}) \quad (5)$$

Following Alayrac et al. (2022), we initialize the parameter of the Tanh gate to be 0, allowing the forward pass of the prompted input through the pre-trained T5 encoder-decoder to be intact at the beginning of the training process. With the gated mechanism, the model can learn to gate out noisy augmentations and rely more on the skip connections during semi-parametric multitask training.

4 EXPERIMENTS

4.1 DATASETS

Following Sanh et al. (2021) we partition various types of NLP tasks into two groups, training tasks and held-out unseen tasks. In this work, we are particularly interested in investigating the impact of retrieving plain texts from a large corpus for multitask prompted training and zero-shot task generalization. Thus, when choosing the training and evaluation tasks, we favor knowledge-intensive tasks over extractive tasks such as extractive question answering (QA) or summarization, where most knowledge for solving the problem is already self-contained in the input. Furthermore, we avoid including large datasets, such as DBPedia (Lehmann et al., 2015) (630K instances) and QQP (Shankar Iyer, 2017) (400K instances), due to limited computational resources for retrieval and semi-parametric multitask prompted training.

We use a subset of T0’s (Sanh et al., 2021) training mixture for our semi-parametric multitask prompted training. Specifically, our training mixture contains *eight* multiple-choice QA datasets, including, **CommonsenseQA** (Rajani et al., 2019), **CosmosQA** (Huang et al., 2019), **DREAM** (Sun et al., 2019), **QASC** (Khot et al., 2020), **QUARTZ** (Tafjord et al., 2019), **SciQ** (Johannes Welbl, 2017), **Social IQa** (Sap et al., 2019), and **WIQA** (Tandon et al., 2019). We choose the subset in multiple-choice QA tasks because they are diverse in domains and overall task formats. Ablation studies on including more types of tasks can be found in Section 4.4.

For evaluation tasks, we consider *seven* datasets from *five* diverse categories following the task taxonomy of T0, including, two sentence completion tasks, **COPA** (Roemmele et al., 2011) and **Hel-laSwag** (HSwag) (Zellers et al., 2019), two multiple-choice QA tasks, **OpenbookQA** (OBQA) (Mihaylov et al., 2018) and **Piqa** (Bisk et al., 2020), one word sense disambiguation task, **WiC** (Pilehvar & Camacho-Collados, 2018), one sentiment task, **Rotten Tomatoes** (RT) (Pang & Lee, 2005), and one natural language inference task, **CB** (De Marneffe et al., 2019).

We use *PromptSource* (Bach et al., 2022) with *Huggingface Datasets* (Lhoest et al., 2021) to construct prompted inputs for each training and evaluation instance. During training, we randomly select two templates for each dataset. During evaluation, we follow the exact evaluation procedure as in T0 (Sanh et al., 2021) and report the mean accuracy across all available templates. All scores are reported on the validation set of each dataset. The detailed templates used for training and evaluation can be found in Appendix A.3.

4.2 MULTITASK TRAINING: SEMI-PARAMETRIC V.S. PARAMETRIC

In this section, we show the zero-shot results on the aforementioned seven held-out unseen tasks with four different multitask training settings. The main baseline we are comparing with is *No Aug*, where we perform parametric multitask prompted training without retrieval just as T0 (Sanh et al., 2021). In order to investigate the best way to incorporate the retrieved documents and to test the effectiveness of our proposed *Zemi* architecture, we consider two widely used alternatives to do fusion of retrieved augmentations, including *Concat* and *FiD*.

Concat In *Concat*, we directly concatenate all retrieved augmentations with the prompted input text. The concatenated input is then truncated to the maximum acceptable length of 1024 tokens and fed to our backbone model.

FiD In *FiD*, we implement the model following Izacard & Grave (2020) where we first independently encode each pair of retrieved augmentation and the prompted input text. Then we concatenate the encoder outputs and feed them to the decoder.

We consider two variants of each model with a different pre-trained backbone, i.e., T5-base and T5-large (Raffel et al., 2020). Following T0 (Sanh et al., 2021), we use the language modeling adapted³ checkpoint, which is trained for an additional 100k steps on a language modeling objective. By default, we use five retrieved passages as augmentations for each instance. For *Zemi*, unless otherwise specified, we use one layer of gated cross-attention and one layer of perceiver resampler with a latent size of 64. A comprehensive ablation study on the impact of different aspects of our model design such as the Tanh Gate can be found in Section 4.4. All models are trained on the same training mixture as mentioned in Section 4.1 for ten epochs with a batch size of 32 and a learning rate of 1e-4. We report results from the checkpoint that achieved the best overall performance across all tasks. All experiments are done on eight NVIDIA-V100 32GB GPUs.

As shown in Table 1, we report the mean zero-shot accuracy across all templates on the validation set of each task. We also show the averaged performance across all tasks, i.e., *Avg*. Since *Avg* might be dominated by particularly high or low absolute scores of certain tasks, we further report the **averaged relative gain** against the *No Aug* baselines, i.e., *Rel Δ Avg*, which is calculated as:

$$Rel \Delta Avg = \frac{1}{n} \sum_{i=1}^n ((acc^i - acc_{NoAug}^i) / acc_{NoAug}^i)$$

where acc^i indicates the accuracy of the i^{th} task, and acc_{NoAug}^i indicates the *No Aug* baseline. The major observations are as follows.

³<https://huggingface.co/google/t5-base-lm-adapt>.

Method	Tasks							Avg	Rel Δ Avg
	OBQA	Piqa	RT	CB	COPA	WiC	HSwag		
No Aug _{BASE}	36.64	60.24	64.11	41.46	68.46	49.92	27.97	49.83	0.00
Concat _{BASE}	35.97	60.18	60.06	32.40	68.65	50.91	28.16	48.05	-3.88
FiD _{BASE}	36.94	59.73	68.38	32.71	68.65	50.86	27.42	49.24	-2.03
Zemi _{BASE} (Ours)	35.55	59.23	68.58	50.05	63.57	49.56	29.68	50.89	+3.04
No Aug _{LARGE}	50.50	65.50	82.23	52.40	80.02	50.16	34.12	59.28	0.00
Concat _{LARGE}	48.81	65.87	74.91	44.64	82.69	50.03	30.46	56.77	-4.88
FiD _{LARGE}	50.99	66.71	67.07	60.71	86.30	50.16	32.91	59.26	+0.65
Zemi _{LARGE} (Ours)	51.49	67.85	84.13	62.14	84.50	50.44	35.79	62.33	+5.36

Table 1: Zero-shot performance on held-out unseen tasks. The scores are the mean accuracy (in %) averaged over all templates on the validation set of each task. By default, the backbone model for all methods is T5-base. Subscripts _{BASE} and _{LARGE} indicate using T5-base and T5-large as backbone model respectively. Detailed descriptions of each task can be found in Section 4.1.

Semi-parametric multitask prompted training improves zero-shot task generalization. We find that semi-parametric multitask training with our proposed architecture (*Zemi*) generally outperforms parametric multitask training (*No Aug*). Notably, with a larger backbone model, i.e., *Zemi*_{LARGE}, we not only observe more consistent improvements on all seven tasks but also a larger overall gain of **+5.36%**, as shown in the bottom block of Table 1. This result shows the potential of *Zemi* to achieve even better performance by scaling up to a larger pre-trained backbone model. In Appendix A.1, we further show examples of both good and noisy retrieved documents for each evaluation task.

Zemi is effective on incorporating retrieved augmentations. Comparing the results from different model architectures, we find that for both T5-base and T5-large backbones, *Concat* and *FiD* only achieve marginal to none positive gain against *No Aug* baseline, while *Zemi* consistently shows notable gains. The fact that *Concat* and *FiD* fail to achieve as good performance as *Zemi* demonstrates that the way of doing augmentation fusion has a significant impact on the effectiveness of the semi-parametric multitask training. In Section 4.4, we show that the gated cross-attention is an important factor contributing to the effectiveness of the *Zemi* architecture. Besides, we encounter out-of-memory problem on NVIDIA-V100 during the evaluation of *FiD*_{LARGE}⁴, however, we can successfully train *Zemi*_{LARGE} on NVIDIA-V100 with the same number of augmentations since the perceiver resampler reduces computational complexity.

4.3 COMPARISON TO STATE-OF-THE-ART

Method	# train tasks	# param	Tasks							Avg ₅	Avg ₇
			OBQA	Piqa	RT	CB	COPA	WiC	HSwag		
BART0	36	0.4B	34.40	36.10	-	39.64	-	46.70	39.40	39.25	-
ReCross	36	0.4B	39.58	41.42	-	44.79	-	50.58	47.28	44.73	-
T0-3B	36	3B	42.83	59.28	73.55*	45.46	75.88	50.03	27.29	44.98	53.47
T0-11B	36	11B	59.11	72.49	81.80*	70.12	91.50	55.20	33.51	58.09	66.25
GPT-3	-	175B	57.60	81.00*	83.16	46.40	91.00	49.40 [†]	78.90	62.66	69.64
Zemi _{BASE}	8	0.23B	35.55	59.23	68.58	50.05	63.57	49.56	29.68	44.81	50.89
Zemi _{LARGE}	8	0.77B	51.49	67.85	84.13	62.14	84.50	50.44	35.79	53.54	62.33

Table 2: Comparison to state-of-the-art. The scores are mean zero-shot accuracy (in %) across all templates for each task. Avg₇ indicates averaged performance across all seven tasks. Avg₅ indicates averaged performance on five tasks excluding *RT* and *COPA* due to unreported baseline results. * indicates the task is seen during training. [†] indicates few-shot results with 32 examples.

We compare both the base and large variants of *Zemi* with current state-of-the-arts, i.e., *T0* (Sanh et al., 2021), *BART0* (Lin et al., 2022), *ReCross* (Lin et al., 2022), and *GPT-3* (Brown et al., 2020) on zero-shot task generalization. Table 2 shows the mean zero-shot accuracy across all templates

⁴We end up with using NVIDIA-A100 40GB for evaluating *FiD*_{LARGE}.

for each task. The last two columns of Table 2 show the averaged performance across different sets of tasks, where Avg_7 is averaged across all seven tasks, and Avg_5 considers five tasks excluding RT and COPA due to their unavailable baseline results.

For *BART0*, *ReCross*, and *GPT-3*, we copy the reported scores directly from their original papers. For the missing score of RT on *GPT-3*, we run the latest ⁵ text completion API to get the generated outputs which is then mapped to the most similar answer choice using SentenceBert (Reimers & Gurevych, 2019). For *T0* models, there are some tasks such as OBQA and Piqa that are not evaluated in the original paper (Sanh et al., 2021), and some tasks such as CB and WiC are evaluated with slightly different templates. Thus, for fair comparison, we re-evaluate all seven tasks on T0-3B and T0-11B using the official implementation and checkpoints⁶ with the exact same set of templates as our model. See details on the templates used for each task in Appendix A.3.

We find that **Zemi_{BASE} outperforms ReCross** on the average of five tasks (Avg_5) while being **1.7x smaller in scale**. Notably, **Zemi_{LARGE}, significantly outperforms T0-3B** on all seven evaluation tasks (Avg_7) by **16%** with **3.9x fewer** parameters. We also observe that although trained with **4.5x fewer** training tasks (8 v.s. 36), our model effectively achieves state-of-the-art zero-shot performance. In Section 4.4, we show that adding more tasks into multitask training does not necessarily improve the performance. And the training mixture with multiple-choice QA tasks seems to be highly effective in generalizing to various kinds of unseen tasks.

4.4 ABLATION STUDIES

Ablated setting	Zemi value	Changed value	Tasks							Avg	Rel Δ Avg
			OBQA	Piqa	RT	CB	COPA	WiC	HSwag		
No Augmentation (No Aug)			36.64	60.24	64.11	41.46	68.46	49.92	27.97	49.83	0.00
Zemi_{BASE}			35.55	59.23	68.58	50.05	63.57	49.56	29.68	50.89	+3.04
(i) Tanh Gate	✓	✗	35.02	57.76	55.79	49.90	71.48	51.59	27.94	49.93	+0.93
(ii) Num of Augs	5	1	35.60	58.92	67.05	47.60	65.23	49.48	30.09	50.57	+2.34
		10	35.32	59.42	62.11	46.30	64.55	51.38	29.36	49.78	+0.83
		20*	34.74	58.67	60.34	46.46	61.62	50.13	28.38	48.62	-1.38
		30*	35.10	60.52	58.65	48.19	67.19	50.69	28.51	49.84	+0.81
(iii) Latent size	64	32	34.85	58.84	64.29	44.48	67.87	51.20	28.55	50.01	+0.59
		128	34.74	57.61	63.14	47.81	69.82	50.39	28.13	50.23	+1.11
(iv) Aug length	256 ⁶⁴	512 ⁶⁴	35.32	58.78	58.64	52.50	68.85	50.34	28.82	50.47	+2.37
		512 ¹²⁸	35.38	59.53	69.16	45.42	70.02	49.40	28.71	51.09	+2.39
(v) Training mixture	Zemi	No Aug+	37.58	58.44	-	43.33	-	50.68	28.01	-	+1.16
		Zemi+	34.46	58.67	-	42.81	-	50.13	29.30	-	-0.02

Table 3: Ablation study. Each ablated setting should be compared with the first two rows, i.e., the original *No Augmentation (No Aug)* setting and *Zemi_{BASE}*. The scores are the mean accuracy (in %) across all templates for each task. The last column shows the averaged relative gain across seven tasks against *No Aug* baseline. The superscripted “*” in ablated setting (ii) indicates using the model variant with a frozen augmentation encoder. The superscripted numbers, i.e., 64, in ablated setting (iv) indicates the latent size of the perceiver resampler. See descriptions of each setting in Section 4.4.

In this section, we conduct comprehensive ablation studies on different aspects of our model design. As shown in Table 3, we consider the following five categories of ablated settings on *Zemi_{BASE}*:

(i) **Tanh Gate**, where we replace the gated cross-attention module with vanilla cross-attention in the ablated version. Specifically, we remove the two Tanh Gates as shown in Figure 2. We find that **removing Tanh Gate hurts the zero-shot performance**. Note that the Tanh Gate is also the main difference between *Zemi* and *FiD* (Izacard & Grave, 2020).

⁵We use the most powerful “text-davinci-002” model.

⁶<https://github.com/bigscience-workshop/t-zero>.

(ii) **Number of augmentations**, where we ablated on the number of augmentations. Note that for settings with 20 and 30 augmentations, in order to reduce the computation complexity, we propose another variant of *Zemi* with an additional frozen augmentation encoder. Specifically, we separately encode the augmentations with a frozen pre-trained T5-Encoder instead of sharing the trainable T5-Encoder with the prompted input. We find that increasing the number of augmentations from single to multiple improves the performance. However, further increasing the number to 10 starts to hurt the performance, which indicates that the noise starts to overwhelm the useful signals introduced by the retrieval. As a result, *Zemi* achieves the best performance with 5 augmentations. We also observe that the performance with 30 augmentations outperforms 20 augmentations, we hypothesize that this is due to inaccurate retrieval ranking that leads to some more informative documents being ranked lower. We show an example of this case in Figure 10. Nevertheless, the fact that we are able to achieve positive gain with as many as 30 augmentations shows the **robustness of our model to very noisy augmentations**.

(iii) **Latent size**, where we ablate on the size of the latent query vector in the perceiver resampler. Note that here the latent size is different from the hidden state size of the backbone model. The trade-off of the size of the latent query vector is that, a larger latent size preserves more information from the original augmentation but also includes more noise. A larger latent size can also increase the computational complexity. We find that our model is **relatively robust to the change of the latent size**. We find that *Zemi* achieves the best performance with a latent size of 64.

(iv) **Augmentation length**, where we investigate the impact of different ways of constructing augmentations from the retrieved documents. Specifically, we increase the maximum length of each augmentation from 256 to 512 and fit two retrieved documents into one augmentation. We keep the number of augmentations the same as default, i.e., 5. We then compare this ablated setting with the 10 augmentation variant (row 3). We find that with the same set of retrieved documents, **augmenting the model with longer but fewer augmentations generally outperforms using a larger number of shorter augmentations**.

(v) **Training mixture**, where we investigate the impact of adding new types of training tasks to the original training mixture. We dub the models trained with this new training mixture as *No Aug+* and *Zemi+*, where *No Aug* means training with no retrieval. Specifically, apart from the eight multiple-choice QA tasks, we further include four more tasks: one closed-book QA task **WikiQA** (Yang et al., 2015), one topic classification task **TREC** (Li & Roth, 2002), one sentence completion task **COPA**, and one sentiment task **Rotten Tomatoes** (RT). Note that here we follow T0 to move tasks that are originally in the evaluation split, i.e. COPA and RT, into the training split in this ablated setting. We find that adding new types of tasks does not necessarily increase the performance. Although trained with only 8 tasks (v.s. 36 tasks) we are able to achieve state-of-the-art performance (Section 4.3), which shows that the **multiple-choice QA mixture is highly effective for generalizing to a wide range of held-out unseen tasks**.

5 CONCLUSIONS & LIMITATIONS

Conclusion. In this work, for the first time, we show that we can effectively train a semi-parametric language model with semi-parametric multitask prompted training to achieve strong zero-shot performance on a wide range of downstream tasks. We augment the model with retrieval from a large-scale task-agnostic corpus, C4, which is the unlabeled pre-training corpus of our backbone model. In order to incorporate multiple potentially noisy retrieved augmentations, we further propose a novel model architecture leveraging perceiver resampler and gated cross-attention. We show that our proposed models with fewer parameters, *Zemi_{BASE}* and *Zemi_{LARGE}*, significantly improve zero-shot task generalization compared with current state-of-the-art models.

Limitation. In Section 4.3, we show that our training mixture with multiple-choice QA tasks, although small, is highly effective for multitask training. However, it is still unclear why multiple-choice QA tasks are particularly effective. Future work includes investigating why certain mixtures are more effective than others. Computation overhead is another noticeable limitation of semi-parametric models which is discussed in more details in Appendix A.2. Moreover, although we show that *Zemi* can successfully handle noisy retrieved augmentations, the noise in retrieval can still potentially hurt the effectiveness of semi-parametric multitask training. Future work includes developing efficient and accurate retrieval methods for retrieving from large-scale task-agnostic corpus, such as C4.

REFERENCES

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*, 2022.
- Stephen H. Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-David, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Alan Fries, Maged S. Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Xiangru Tang, Mike Tian-Jian Jiang, and Alexander M. Rush. Promptsource: An integrated development environment and repository for natural language prompts, 2022.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. *arXiv preprint arXiv:2112.04426*, 2021.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading wikipedia to answer open-domain questions. *arXiv preprint arXiv:1704.00051*, 2017.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. The commitmentbank: Investigating projection in naturally occurring discourse. In *proceedings of Sinn und Bedeutung*, volume 23, pp. 107–124, 2019.
- ElasticSearch. Elasticsearch. URL <https://www.elastic.co/>.
- Wikimedia Foundation. Wikimedia downloads. URL <https://dumps.wikimedia.org>.
- Edouard Grave, Armand Joulin, and Nicolas Usunier. Improving neural language models with a continuous cache. *arXiv preprint arXiv:1612.04426*, 2016.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented language model pre-training. In *International Conference on Machine Learning*, pp. 3929–3938. PMLR, 2020.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Cosmos qa: Machine reading comprehension with contextual commonsense reasoning. In *arXiv:1909.00277v2*, 2019.
- Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*, 2020.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*, 2022.
- Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021.

- Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot manipulation with multimodal prompts. *arXiv preprint arXiv:2210.03094*, 2022.
- Matt Gardner Johannes Welbl, Nelson F. Liu. Crowdsourcing multiple choice science questions. 2017.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization through memorization: Nearest neighbor language models. *arXiv preprint arXiv:1911.00172*, 2019.
- Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. Qasc: A dataset for question answering via sentence composition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 8082–8090, 2020.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised open domain question answering. *arXiv preprint arXiv:1906.00300*, 2019.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. Dbpedia—a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic web*, 6(2):167–195, 2015.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander Rush, and Thomas Wolf. Datasets: A community library for natural language processing. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 175–184, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. URL <https://aclanthology.org/2021.emnlp-demo.21>.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. *arXiv preprint arXiv:2201.12086*, 2022.
- Xin Li and Dan Roth. Learning question classifiers. In *COLING 2002: The 19th International Conference on Computational Linguistics*, 2002.
- Opher Lieber, Or Sharir, Barak Lenz, and Yoav Shoham. Jurassic-1: Technical details and evaluation. *White Paper. AI21 Labs*, 2021.
- Bill Yuchen Lin, Kangmin Tan, Chris Miller, Beiwen Tian, and Xiang Ren. Unsupervised cross-task generalization via retrieval augmentation. *ArXiv*, abs/2204.07937, 2022.
- Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hananeh Hajishirzi. Generated knowledge prompting for commonsense reasoning. *arXiv preprint arXiv:2110.08387*, 2021.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.

- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the ACL*, 2005.
- Mohammad Taher Pilehvar and Jose Camacho-Collados. Wic: the word-in-context dataset for evaluating context-sensitive meaning representations. *arXiv preprint arXiv:1808.09121*, 2018.
- Shrimai Prabhunoye, Kazuma Hashimoto, Yingbo Zhou, Alan W Black, and Ruslan Salakhutdinov. Focused attention improves document-grounded generation. *arXiv preprint arXiv:2104.12714*, 2021.
- Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020.
- Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. Explain yourself! leveraging language models for commonsense reasoning. *arXiv preprint arXiv:1906.02361*, 2019.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. Okapi at trec-3. *Nist Special Publication Sp*, 109:109, 1995.
- Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *AAAI spring symposium: logical formalizations of commonsense reasoning*, pp. 90–95, 2011.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*, 2021.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialliqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- Kornél Csernai Shankar Iyer, Nikhil Dandekar. First quora dataset release: Question pairs, 2017. URL <https://quoradata.quora.com/First-Quora-Dataset-Release-Question-Pairs>.
- Shaden Smith, Mostofa Patwary, Brandon Norrick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhunoye, George Zerveas, Vijay Korthikanti, et al. Using deep-speed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. *arXiv preprint arXiv:2201.11990*, 2022.
- Kai Sun, Dian Yu, Jianshu Chen, Dong Yu, Yejin Choi, and Claire Cardie. Dream: A challenge data set and models for dialogue-based reading comprehension. *Transactions of the Association for Computational Linguistics*, 7:217–231, 2019.
- Oyvind Tafjord, Matt Gardner, Kevin Lin, and Peter Clark. Quartz: An open-domain dataset of qualitative relationship questions. *arXiv preprint arXiv:1909.03553*, 2019.
- Niket Tandon, Bhavana Dalvi Mishra, Keisuke Sakaguchi, Antoine Bosselut, and Peter Clark. Wiqua: A dataset for “what if...” reasoning over procedural text. *arXiv preprint arXiv:1909.04739*, 2019.
- Shuohang Wang, Yichong Xu, Yuwei Fang, Yang Liu, Siqi Sun, Ruochen Xu, Chenguang Zhu, and Michael Zeng. Training data is more valuable than you think: A simple and effective method by retrieving from training data. *arXiv preprint arXiv:2203.08773*, 2022a.

- Thomas Wang, Adam Roberts, Daniel Hesslow, Teven Le Scao, Hyung Won Chung, Iz Beltagy, Julien Launay, and Colin Raffel. What language model architecture and pretraining objective work best for zero-shot generalization? *arXiv preprint arXiv:2204.05832*, 2022b.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. Benchmarking generalization via in-context instructions on 1,600+ language tasks. *arXiv preprint arXiv:2204.07705*, 2022c.
- Zhiguo Wang, Patrick Ng, Xiaofei Ma, Ramesh Nallapati, and Bing Xiang. Multi-passage bert: A globally normalized bert model for open-domain question answering. *arXiv preprint arXiv:1908.08167*, 2019.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. End-to-end open-domain question answering with bertserini. *arXiv preprint arXiv:1902.01718*, 2019.
- Yi Yang, Wen-tau Yih, and Christopher Meek. Wikiqa: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pp. 2013–2018, 2015.
- Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. Crossfit: A few-shot learning challenge for cross-task generalization in nlp. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 7163–7189, 2021.
- Dani Yogatama, Cyprien de Masson d’Autume, and Lingpeng Kong. Adaptive semiparametric language models. *Transactions of the Association for Computational Linguistics*, 9:362–373, 2021.
- Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Zexuan Zhong, Tao Lei, and Danqi Chen. Training language models with memory augmentation. *arXiv preprint arXiv:2205.12674*, 2022.

A APPENDIX

A.1 QUALITATIVE ANALYSIS OF THE RETRIEVED DOCUMENTS

Here we visualize one good and one noisy example of the retrieved documents for each evaluation task. A full list of examples for each training and evaluation task can be found in the supplementary material under the “visualization” folder. As shown in Figure 3, 4, 5, 6, 8, 7, and 9, the retrieved augmentations can contain highly correlated information that can be directly helpful for solving a certain task, however, they can also be very noisy. As mentioned in Section 3.2, the retrieved

documents can also be inaccurately ranked, for example in Figure 10, we show that the 21th ranked retrieval result can contain more correlated information than the top ranked ones. Furthermore, as shown in the noisy example of Figure 6, for some tasks such as sentiment analysis, even though the retrieved document is highly correlated with the input text, i.e., with a high BM25 score, the content can steer the prediction into a wrong direction. These observations motivate us to propose the augmentation fusion module with a gated mechanism.

A.2 COMPUTATION OVERHEAD

There are two main computation overhead compared with the fully-parametric counterpart, i.e., the *No Aug* baseline. First, retrieving from a large-scale corpus can be time-consuming. As mentioned in Section 3.2, we apply document-level retrieval with BM25 and truncation on the query to reduce the retrieval time. We also perform the retrieval offline to avoid repeated time commitment. As a result, indexing 5% of the C4 corpus takes 1.03 hours. Offline retrieval for the entire training and evaluation mixture takes 10.98 hours, which is approximately 0.28 seconds per instance. Furthermore, we measure the computation overhead on inference which is caused by the additional retrieved inputs as well as a small amount of newly introduced parameters (+4.6%). The average computation overhead across all evaluation datasets during inference is around 4.36x compared with the *No Aug* baseline.

A.3 FULL LIST OF TASKS AND TEMPLATES

Following T0 (Sanh et al., 2021), we use tasks from Huggingface Datasets (Lhoest et al., 2021) and templates from PromptSource (Bach et al., 2022) marked as “original task” and with “choices.in_prompt”. Specifically, for tasks in the training mixture, we randomly sample two templates per task for semi-parametric multitask prompted training. For tasks in the held-out evaluation mixture, we use all available templates. Table 4, and 5 shows the full list of templates we used for each task during multitask training and zero-shot evaluation.

Mixture	Task	Template Name
Semi-T0 Training	cos_e/v1.11	question_option_description_text description_question_option_id
	cosmos_qa	context_description_question_answer_id description_context_question_answer_text
	dream	baseline read_the_following_conversation_and_answer_the_question
	qasc	qa_with_separated_facts_1 qa_with_separated_facts_4
	quartz	answer_question_below read_passage_below_choose
	sciq	Multiple Choice Multiple Choice Question First
	social_i_qa	Show choices and generate answer Show choices and generate index
	wiqa	effect_with_string_answer effect_with_label_answer
Semi-T0+ Training	wiki_qa	Decide_good_answer found_on_google
	trec	what_category_best_describe trec1
	super_glue/copa	more likely best_option
	rotten_tomatoes	Sentiment with choices Reviewer Opinion bad good choices

Table 4: PromptSource template names used for each task (Part1).

Mixture	Task	Template Name
Semi-T0 Evaluation	openbookqa/main	choose_an_answer_with_options which_correct pick_using_id choices only_options which_correct_inverse pick_answer_with_options
	piqa	what_is_the_correct_ending pick_correct_choice_with_choice_given_before_goal pick_correct_choice_index finish_sentence_with_correct_choice choose the most appropriate solution
	rotten_tomatoes	Reviewer Opinion bad good choices Sentiment with choices
	super_glue/cb	can we infer based on the previous passage claim true/false/inconclusive does it follow that justified in saying always/sometimes/never GPT-3 style consider always/sometimes/never guaranteed true must be true guaranteed/possible/impossible does this imply MNLi crowdsource should assume take the following as truth
	super_glue/copa	exercise ... What could happen next, C1 or C2? i_am_hesitating plausible_alternatives C1 or C2? premise, so/because... ... As a result, C1 or C2? best_option ... which may be caused by more likely cause_effect ... why? C1 or C2 choose
	super_glue/wic	question-context-meaning-with-label grammar_homework affirmation_true_or_false same_sense GPT-3-prompt-with-label polysemous
	hellaswag	complete_first_then Randomized prompts template Predict ending with hint if_begins_how_continues

Table 5: PromptSource template names used for each task (Part2).

A.4 RETRIEVAL QUERY KEY FOR EACH TASK

In order to retrieve most relevant documents for each instance, we specify a certain field for each dataset which will be served as the query to the retrieval system. For example, for most multiple-choice QA tasks, we use the “question” string as our query. Table 6 shows a full list of field names

Task	Query Key
cos_e/v1.11	question
cosmos_qa	question
dream	question
qasc	question
quartz	question
sciq	question
social_i_qa	context
wiqa	question_stem
openbookqa/main	question_stem
piqa	goal
rotten_tomatoes	text
super_glue/cb	hypothesis
super_glue/copa	premise
super_glue/wic	sentence1
hellaswag	ctx
wiki_qa	question
trec	text

Table 6: Retrieval query key used for each task.

we use as retrieval query keys for each dataset. Note that the field name shown in the table is what appears to be in the corresponding Huggingface Dataset (Lhoest et al., 2021).

Good Retrieval Example	Noisy Retrieval Example
<p>===== Instance Index 6063 =====</p> <p>Input Text: If a description of a situation begins like this: [header] How to set macgo mac blu ray player as default player [title] Download mac blu-ray menu player and install it at once. [step] There will be watermark on your screen if you play blu-ray with the trial version. Only 39.95 dollars for the full version of mac blu-ray menu player for now, please buy mac blu-ray player with discount... Then how does it continue?</p> <p>Ending 1: [title] Click " check file associations " under " tools ". [title] Click and macgo mac blu-ray player will be your default player.</p> <p>Ending 2: [title] Choose your video size and port size from the dropdown menu at the top of mac blu-ray menu. [step] Once you have downloaded the blu-ray menu player and installed it, you have to choose your video size and port size.</p> <p>Ending 3: [title] Run the make app and then the itunes installer. [title] Determine the output type for each file in your mac blu-ray player.</p> <p>Ending 4: [title] Uncheck the sidebar at the bottom of " applications ". [step] These are the files that are currently currently on your mac blu-ray player.</p> <p>Target Text: Ending 1</p> <p>#### Retrieved Documents #### Rank: 0 Score: 136.6041 Retrieved Text: Macgo Mac Blu-ray Player has added itself Auto Play function, which means when you insert a disc into your Blu-ray drive, the player will automatically start and play. In order to make this whole process smoother, you'd better set Mac Blu-ray Player as default player on your Mac. Now I'll tell you how to do it. After installing Mac Blu-ray Player, you can go to "Launchpad" and click on its icon to launch the program. The simplified main interface will reduce certain misoperations. You can see a menu at the top of the interface. Click "Check File Associations" under "Tools". Then it will come up with a pop up window. You can choose some media formats which you want to play with Macgo Blu-ray player, then click "Make Mac Blu-ray player my default player". Click "OK" to continue. Then Macgo Mac Blu-ray Player will be your default player. After you set Macgo Mac Blu-ray Player as your default player, you also need to enable Auto Play function to freely enjoy Blu-ray this player. Open "Preferences" under "Mac Blu-ray Player".Open "Playback" and tick under "Auto play when you insert a disc", and then click "OK". Insert a Blu-ray disc into the drive and wait for the program automatically start and display the Blu-ray Menu. You can make some adjustments there or directly click "Play Movie" to enjoy some Blu-ray time.</p>	<p>===== Instance Index 122 =====</p> <p>Input Text: If a description of a situation begins like this: A group of people are in a house. a man... Then how does it continue?</p> <p>Ending 1: is holding cored soap in his hand as he washes with a bottle.</p> <p>Ending 2: is mopping the floor with a mop.</p> <p>Ending 3: is shown wearing skis as he talks about areas he will like to ski on.</p> <p>Ending 4: uses a paintball gun on his child.</p> <p>Target Text: Ending 2</p> <p>#### Retrieved Documents #### Rank: 0 Score: 17.592176 Retrieved Text: #52704883 - Red lanterns, oriental charm, the Spring Festival atmosphere. #35618548 - Silhouette of a man Happy successful raising arms to the sky.. #93113276 - Backlighting portrait of a joyful mother raising her baby outdoors.. #108745447 - Backlighting portrait of a joyful mother raising her baby outdoors.. #37541508 - Group of cheerleaders performing outdoors - Concept of cheerleading.. #108747918 - Back view of young backlit man looking into the distance on illuminated.. #38536931 - Working man walking near airplane wing at the terminal gate of.. #73301104 - Group of urban friends walking in city skate park with backlighting.. #76682984 - People silhouettes putting puzzle pieces together on city background.. #77013901 - People silhouettes putting puzzle pieces together on abstract.. #104666805 - Stylish light gray kitchen interior with modern cabinets with.. #108748125 - Back view of young backlit man looking into the distance on illuminated.. #86815910 - Best friends taking selfie outdoors with backlighting - Happy.. #86815909 - Best friends taking selfie outdoors with backlighting - Happy.. #117963685 - Back view backlighting silhouette of a man alone on a swing looking.. #86815911 - Best friends taking selfie outdoors with backlighting - Happy.. #118172724 - Back view backlighting silhouette of a man sitting on swing alone.. #96363446...</p>

Figure 3: Example of retrieved documents on HellaSwag.

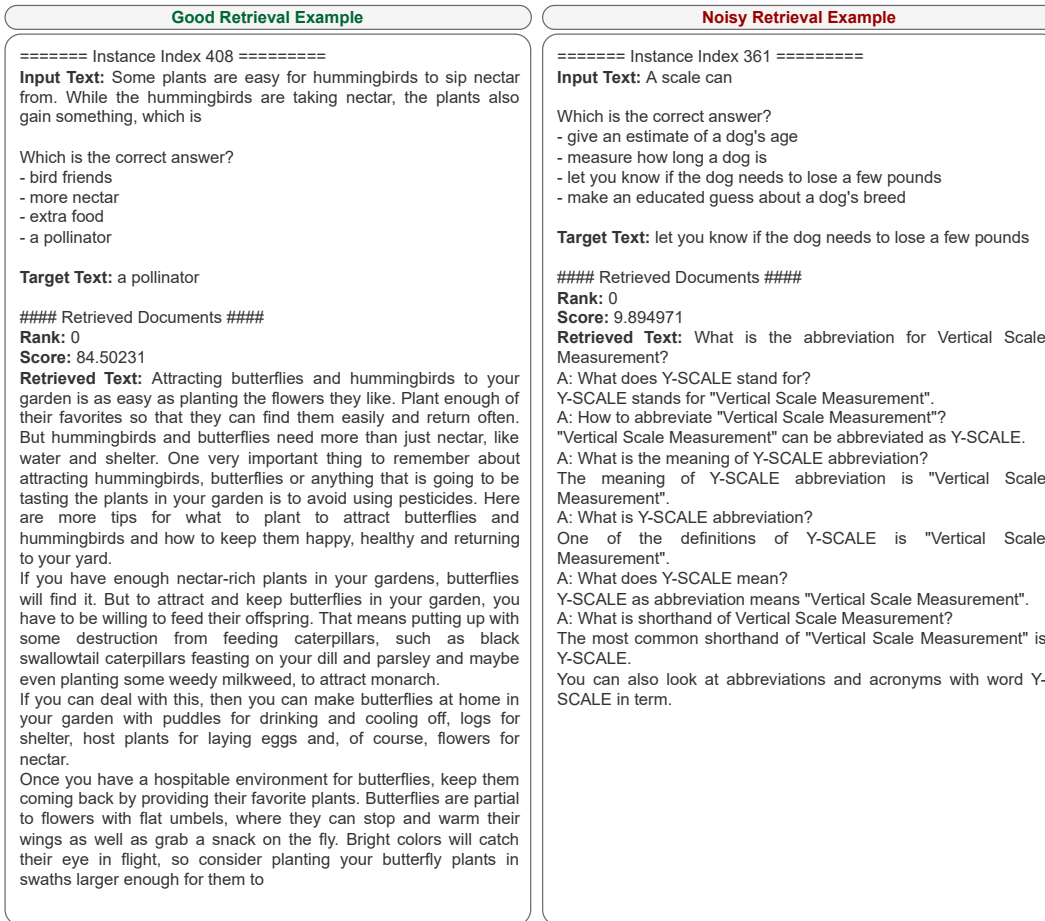


Figure 4: Example of retrieved documents on OpenbookQA.

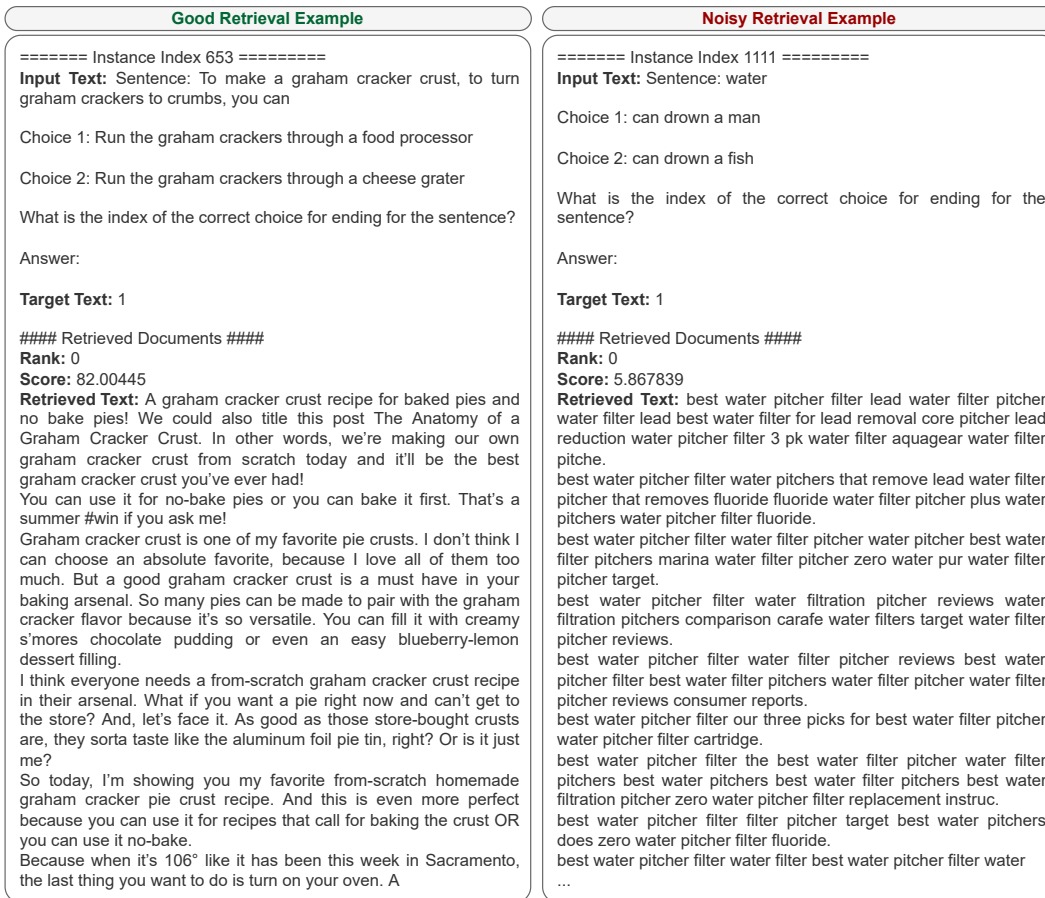


Figure 5: Example of retrieved documents on Piqua.

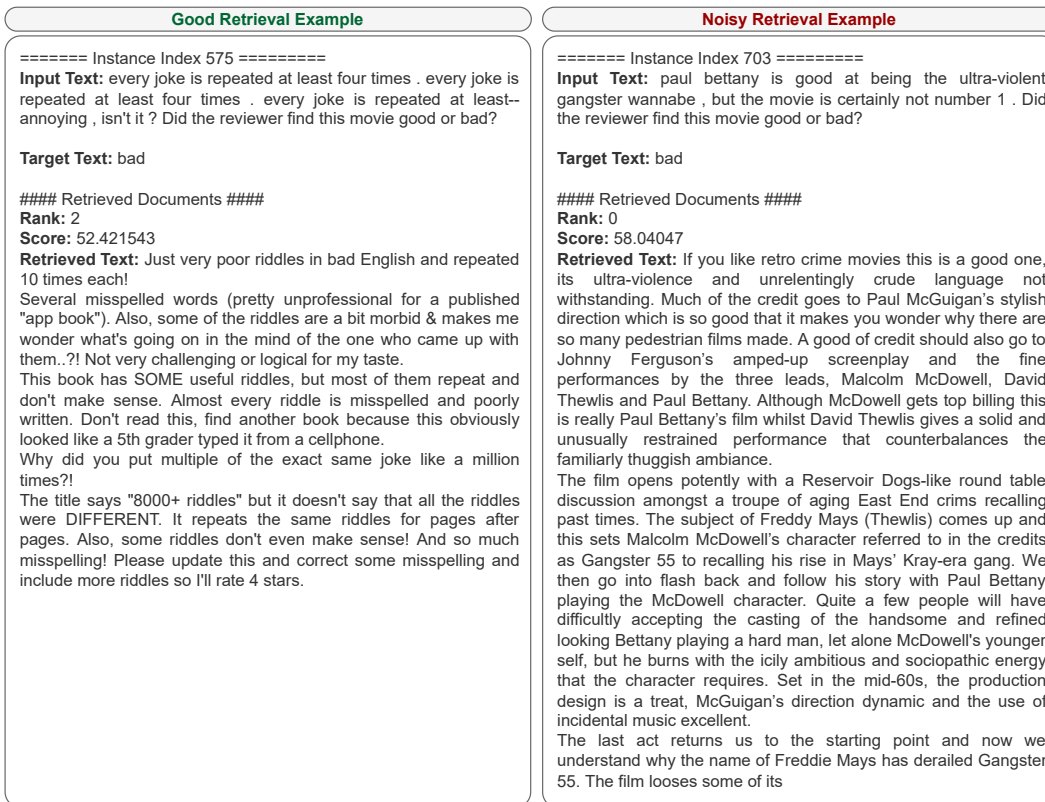


Figure 6: Example of retrieved documents on Rotten Tomatoes.

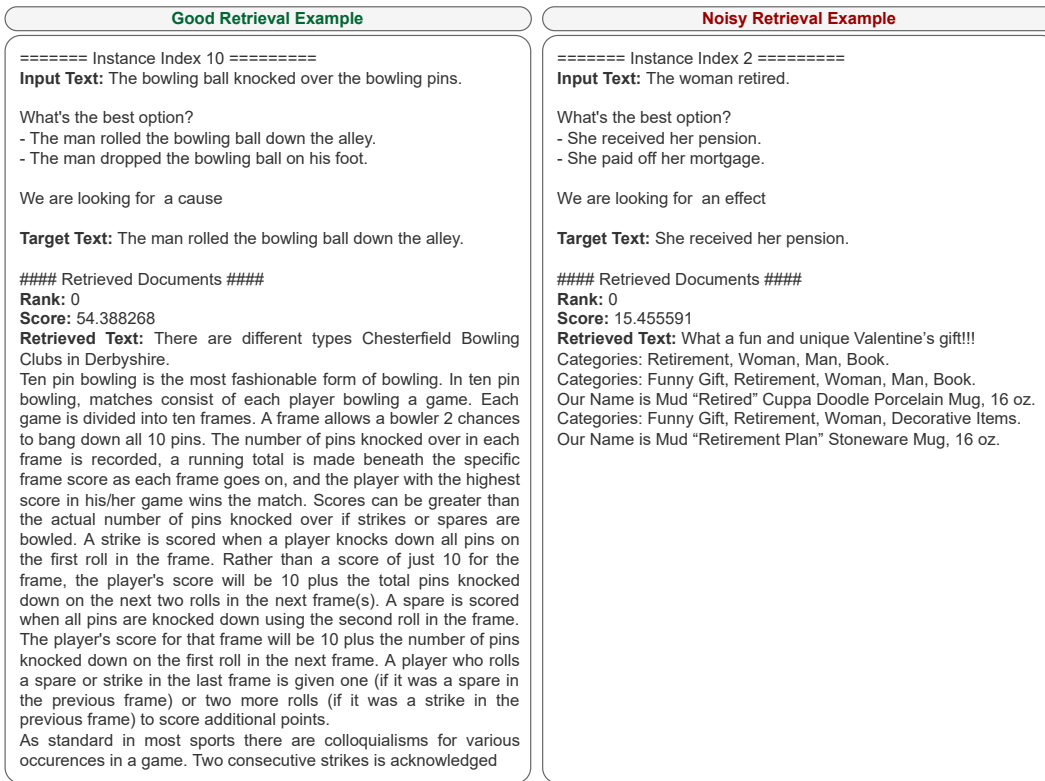


Figure 7: Example of retrieved documents on COPA.

Good Retrieval Example	Noisy Retrieval Example
<p>===== Instance Index 8 =====</p> <p>Input Text: Given that A: And I haven't quite figured that out, if they figure they have got it won or if there's no real hurry because the first three quarters or, uh, uh, if something happens that that adrenalin starts flowing. They say, hey, we got to do something now. And then start playing the game the way the game should be played toward the last few minutes. B: Yeah. A: So, I don't know I'm looking for a good year. I guess we're always looking for a good year. B: So, obviously though, do you think they're going to do anything in the playoffs to make it to the Super Bowl this year. Therefore, it must be true that "they're going to do anything in the playoffs to make it to the Super Bowl this year"? Yes, no, or maybe?</p> <p>Target Text: Maybe</p> <p>#### Retrieved Documents ####</p> <p>Rank: 0</p> <p>Score: 41.988216</p> <p>Retrieved Text: Two-time super bowl champion and CNN Sport contributor Hines Ward shares his Week 9 takeaways with CNN's Jill Martin.</p> <p>We're at the halfway point, and you start to see teams separate the contenders from the pretenders. You really see what teams are made of. This is a crucial month for a lot of teams in the NFL. Let's start with the NFC South, where the Panthers and Saints need our attention.</p> <p>The Carolina Panthers -- I don't think anyone expected them to have the year that they're having.</p> <p>Cam Newton is looking like he's back to his MVP form, from back in 2015. What they're doing with running back Christian McCaffrey I just think is amazing. It's showing his versatility both running and catching the ball.</p> <p>They're only one game behind the New Orleans Saints, and they have key matchups at the end of the year. In the last three weeks of the season, they play each other twice. Right now, it looks like it should be for the division.</p> <p>Meanwhile, the Saints just knocked off the Rams. What, if anything does that performance show you?</p> <p>Well, it's a tough place to play. I think, right now, it's really a two-team race to try to get that home field advantage for the playoffs. I've played in New Orleans. I've been there. I know what their fans are like. It's one of the toughest places to play. It's loud. They get rowdy, and they love their Saints.</p> <p>Definitely having Drew Brees playing at home in the playoffs helps the Saints' chances of making it to the</p>	<p>===== Instance Index 7 =====</p> <p>Input Text: Given that It grew bigger with incredible speed, she was whizzing towards it. She must slow down or she 'd miss it. She took her foot off the accelerator and put it on the brake and as the car slowed she could see now that it was a child a toddler with a red woolly hat on. Therefore, it must be true that "it was a child"? Yes, no, or maybe?</p> <p>Target Text: Yes</p> <p>#### Retrieved Documents ####</p> <p>Rank: 0</p> <p>Score: 10.499066</p> <p>Retrieved Text: The Plan Has Been Executed – My Grace Is..</p> <p>This is my love letter to you son. Forever you will remain a child dear to me my daughter. I know you have read and heard that I have a plan to prosper you, to give you a hope and a future.</p> <p>Child oh my child, yes I had a plan for you back then, back, back then. It was all true. But here is a thing today for you grab hold of my child. To master and rejoice in. The plan has been executed. The plan is sealed and delivered. My child, yes I had a plan for you, a plan for you to live a happy life. To live a joyous life. My plan was great for you my child. My plan was great for you my precious child.</p> <p>Like every other parent, I had a plan for you my child. The plan was drawn down. Well designed, well traced and well set out. Just like a cartoonist would first draw before he brings the characters he has drawn to motion, I too, did that. I too my son had a plan in mind for you. I could not put you on earth and not have a plan at all. I did it and set it up my child. Worry not my son, the plan is executed. For long you heard the words that I had a plan for you, my precious child, please know this, the plan has been executed. The plan has come to life.</p> <p>My plan</p>

Figure 8: Example of retrieved documents on CB.

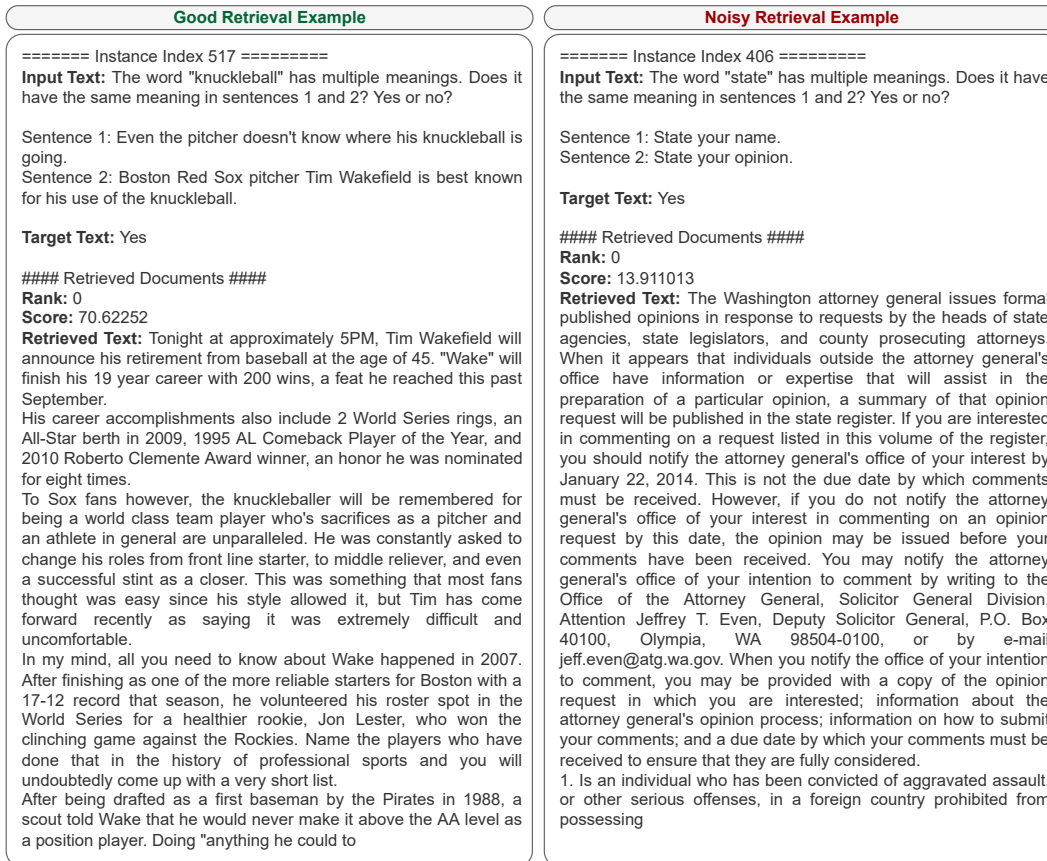


Figure 9: Example of retrieved documents on WiC.

Example of Inaccurate Ranking

=====
Instance Index 0
=====
Input Text: Sentence: How do I ready a guinea pig cage for it's new occupants?

Choice 1: Provide the guinea pig with a cage full of a few inches of bedding made of ripped paper strips, you will also need to supply it with a water bottle and a food dish.

Choice 2: Provide the guinea pig with a cage full of a few inches of bedding made of ripped jeans material, you will also need to supply it with a water bottle and a food dish.

What is the index of the correct choice for ending for the sentence?

Answer:

Target Text: 1

Retrieved Documents

Rank: 1
Score: 46.520477
Retrieved Text: how do I find neat names?
How to go about finding a vet?
guinea pig dali apparently on mend, again?
hamster cage Hammock pattern. . .
hamster cage Secure your cage doors!
ear infection, ear infections, inner ear infection degu sick am having rant!!!!
Do males hump each other?...

...

Rank: 10
Score: 41.217316
Retrieved Text: Contact Alittlebitiffy Animal Sanctuary at Alittlebitiffy Animal Rescue to express your interest.
Another successful adoption - amazing work Alittlebitiffy Animal Rescue!
More successful adoptions - amazing work Alittlebitiffy Animal Rescue! ...

...

Rank: 21
Score: 39.48997
Retrieved Text: Keeping your little furry friend healthy and happy should be a priority for any owner and, along with providing the right food for guinea pigs, finding an appropriate cage for them should be at the top of your priority list. Although, as you can see in this post here, there are numerous options on the market when it comes to commercially available guinea pig cages, some owners have opted towards a more do-it-yourself approach.
Many guinea pig parents complain that the regular pet store-sized cages are nothing but 'glorified litter boxes' and therefore are looking to improve the well-being of their covies by making them a healthy and large-enough living enclosure, rather than buying one.
If you are one of those owners, this article will guide you through what you need to know before you start making a DIY cage for your guinea pig and what options you have when it comes to materials, design and features....

Figure 10: Example of the inaccurate ranking of the retrieval. Here we show the ranked retrieved documents for instance 0 in Piqa. We can see that the 21th ranked document is more correlated than many of the higher ranked ones, such as rank 1 and rank 10.