

# REPLUG: Retrieval-Augmented Black-Box Language Models

Anonymous ACL submission

## Abstract

We introduce REPLUG, a retrieval-augmented language modeling framework that treats the language model (LM) as a black box and augments it with a tuneable retrieval model. Unlike prior retrieval-augmented LMs that train language models with special cross attention mechanisms to encode the retrieved text, REPLUG simply prepends retrieved documents to the input for the frozen black-box LM. This simple design can be easily applied to any existing language models. Furthermore, we show that the LM can be used to supervise the retrieval model, which can then find documents that help the LM make better predictions. Our experiments demonstrate that REPLUG with the tuned retriever significantly improves the performance of GPT-3 (175B) on language modeling by 6.3%, as well as the performance of Codex on five-shot MMLU by 5.1%.

## 1 Introduction

Large language models (LMs) such as GPT-3 (Brown et al., 2020a) and Codex (Chen et al., 2021a), have demonstrated impressive performance on a wide range of language tasks. These models are typically trained on very large datasets and store a substantial amount of world or domain knowledge implicitly in their parameters. However, they are also prone to hallucination and cannot represent the full long tail of knowledge from the training corpus. Retrieval-augmented language models (Khandelwal et al., 2020; Borgeaud et al., 2022; Izacard et al., 2022b; Yasunaga et al., 2022), in contrast, can retrieve knowledge from an external datastore when needed, potentially reducing hallucination and increasing coverage. Previous approaches of retrieval-augmented language models require access to the internal LM representations (e.g., to train the model (Borgeaud et al., 2022; Izacard et al., 2022b) or to index the datastore (Khandelwal et al., 2020)), and are thus difficult to be applied to very large LMs. In addition, many best-in-class

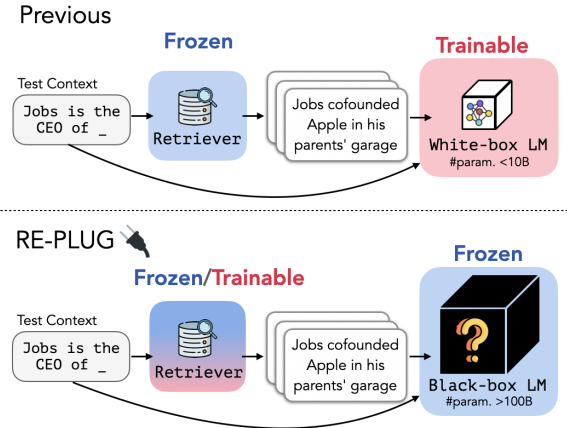


Figure 1: Different from previous retrieval-augmented approaches (Borgeaud et al., 2022) that enhance a language model with retrieval by updating the LM’s parameters, REPLUG treats the LM as a black box and augments it with a frozen or tunable retriever. This black-box assumption makes REPLUG applicable to large LMs, which are often served via APIs.

LLMs can only be accessed through APIs. Internal representations of such models are not exposed and fine-tuning is not supported.

In this work, we introduce REPLUG (**R**etrieve and **P**lug), a new retrieval-augmented LM framework where the language model is viewed as a black box and the retrieval component is added as a tuneable plug-and-play module. Given an input context, REPLUG first retrieves relevant documents from an external corpus using an *off-the-shelf* retrieval model. The retrieved documents are prepended to the input context and fed into the black-box LM to make the final prediction. Because the LM context length limits the number of documents that can be prepended, we also adopt an ensemble scheme that encodes the retrieved documents in parallel with the same black-box LM, allowing us to easily trade compute for accuracy. As shown in Figure 1, REPLUG is extremely flexible and can be used with any existing black-box

LM and retrieval model.

We also introduce REPLUG LSR (REPLUG with LM-Supervised Retrieval), a training scheme that can further improve the initial retrieval model in REPLUG with supervision signals from a black-box language model. The key idea is to adapt the retriever to the LM, which is in contrast to prior work (Borgeaud et al., 2022) that adapts language models to the retriever. We use a training objective which prefers retrieving documents that improve language model perplexity, while treating the LM as a frozen, black-box scoring function.

Our experiments show that REPLUG can improve the performance of diverse black-box LMs on both language modeling and downstream tasks, including MMLU (Hendrycks et al., 2021) and open-domain QA (Kwiatkowski et al., 2019; Joshi et al., 2017). For instance, REPLUG can improve Codex (175B) performance on MMLU by 4.5%, achieving comparable results to the 540B, instruction-finetuned Flan-PaLM. Furthermore, tuning the retriever with our training scheme (i.e., REPLUG LSR) outperforms various off-the-shelf retrievers and leads to additional improvements, including up to 6.3% increase in GPT-3 175B language modeling. To the best of our knowledge, our work is the first to show the benefits of retrieval to large LMs (>100B model parameters), for both reducing LM perplexity and improving in-context learning performance. We summarize our contributions as follows:

- We introduce REPLUG (§3), the first retrieval-augmented language modeling framework for enhancing black-box LMs with retrieval. Unlike previous methods that require updating the LM’s parameters, REPLUG could be easily plugged into any existing LM without additional finetuning.
- We propose a training scheme (§4) to further adapt an off-the-shelf retrieval model to the LM, using the language modeling scores as supervision signals, resulting in improved retrieval quality.
- We are the first to demonstrate that retrieval can benefit large-scale, state-of-the-art LMs on language modeling (§6) and in-context learning tasks. Evaluations show that REPLUG can improve the performance of various language models such as GPT, OPT and

BLOOM, including very large models with up to 175B parameters.

## 2 Background and Related Work

**Black-box Language Models** Large language models, such as GPT-3 (Brown et al., 2020a), Codex (Chen et al., 2021a), are not open-sourced due to commercial considerations and are only available as black-box APIs, through which users can send queries and receive responses. On the other hand, even open sourced language models such as BLOOM-176B (Scao et al., 2022) require significant computational resources to run and fine-tune locally. For example, finetuning BLOOM-176B requires 72 A100 GPUs (Younes Belkda, 2022), making them inaccessible to researchers and developers with limited resources. Traditionally, retrieval-augmented model frameworks (Khandelwal et al., 2020; Borgeaud et al., 2022; Yu, 2022; Izacard et al., 2022b; Goyal et al., 2022) have focused on the white-box setting, where language models are fine-tuned to incorporate retrieved documents. However, the increasing scale and black-box nature of LLMs makes this approach infeasible. To address these challenges, we investigate retrieval-augmentation in the **black-box setting**, where users only have access to the model predictions and cannot access or modify its parameters.

**Retrieval-augmented Models** Augmenting language models with relevant information retrieved from knowledge stores has shown to be effective in improving performance on various NLP tasks, including language modeling (Min et al., 2022; Borgeaud et al., 2022; Khandelwal et al., 2020) and open-domain question answering (Lewis et al., 2020; Izacard et al., 2022b; Hu et al., 2022). Specifically, using the input as query, (1) a retriever first retrieves a set of documents from a corpus and then (2) a language model incorporates the retrieved documents as additional information to make a final prediction. Previous retrieval-augmented LMs require updating the model parameters, which cannot be applied to black-box LMs, which cannot be applied to black-box LMs. For example, Atlas (Izacard et al., 2022b) finetunes an *encoder-decoder* model jointly with the retriever by modeling documents as latent variables, while RETRO (Borgeaud et al., 2022) changes the *decoder-only* architecture to incorporate retrieved texts and pretrains the language model from scratch. Another line of retrieval-augmented LMs such as kNN-LM (Khan-

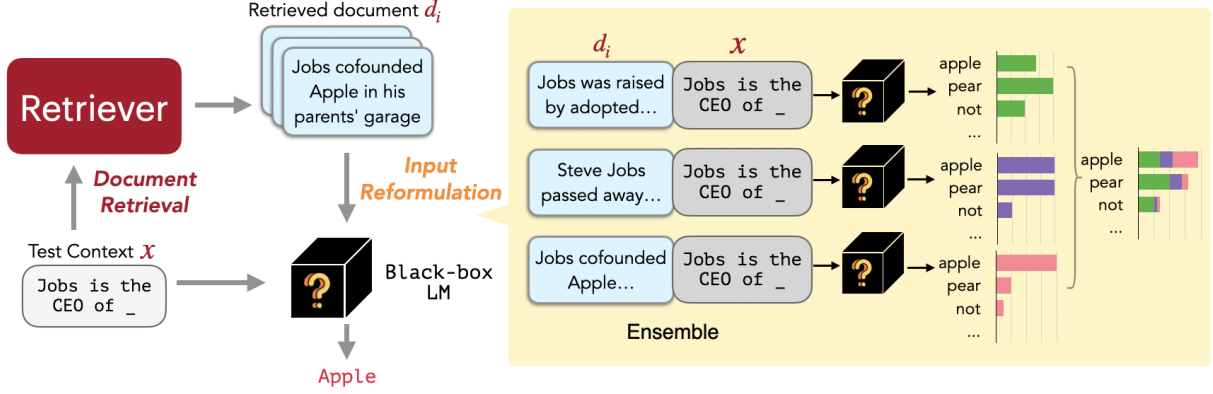


Figure 2: **REPLUG at inference** (§3). Given an input context, REPLUG first retrieves a small set of relevant documents from an external corpus using a retriever (§3.1 *Document Retrieval*). Then it prepends each document separately to the input context and ensembles output probabilities from different passes (§3.2 *Input Reformulation*).

delwal et al., 2020; Zhong et al., 2022) retrieves a set of tokens and interpolates between the LM’s next token distribution and kNN distributions computed from the retrieved tokens at inference. kNN-LM requires access to internal LM representations to compute the kNN distribution, which are not available for black-box LMs such as GPT-3. In this work, we investigate ways to improve large black-box language models with retrieval. While concurrent work (Mallen et al., 2022; Si et al., 2023) has demonstrated that using a frozen retriever can improve GPT-3 performance on open-domain question answering, we approach the problem in a more general setting, including language modeling and understanding tasks. We additionally adopt an ensemble method to incorporate more documents and a training scheme to further adapt the retriever to large LMs.

### 3 REPLUG

We introduce REPLUG (**R**etrieve and **P**lug), a new retrieval-augmented LM paradigm where the LM is treated as black box and the retrieval component is added as a potentially tuneable module.

As shown in Figure 2, given an input context, REPLUG first retrieves a small set of relevant documents from an external corpus using a retriever (§3.1). Then we pass the concatenation of each retrieved document with the input context through the LM in parallel, and ensemble the predicted probabilities (§3.2).

#### 3.1 Document Retrieval

Given an input context  $x$ , the retriever aims to retrieve a small set of documents from a corpus

$\mathcal{D} = \{d_1 \dots d_m\}$  that are relevant to  $x$ . Following prior work (Qu et al., 2021; Izacard and Grave, 2021a; Ni et al., 2021), we use a dense retriever based on the dual encoder architecture, where an encoder is used to encode both the input context  $x$  and the document  $d$ . Specifically, the encoder maps each document  $d \in \mathcal{D}$  to an embedding  $\mathbf{E}(d)$  by taking the mean pooling of the last hidden representation over the tokens in  $d$ . At query time, the same encoder is applied to the input context  $x$  to obtain a query embedding  $\mathbf{E}(x)$ . The similarity between the query embedding and the document embedding is computed by their cosine similarity:

$$s(d, x) = \cos(\mathbf{E}(d), \mathbf{E}(x)) \quad (1)$$

The top- $k$  documents that have the highest similarity scores when compared with the input  $x$  are retrieved in this step. For efficient retrieval, we precompute the embedding of each document  $d \in \mathcal{D}$  and construct FAISS index (Johnson et al., 2019) over these embeddings.

#### 3.2 Input Reformulation

The retrieved top- $k$  documents provide rich information about the original input context  $x$  and can potentially help the LM to make a better prediction. One simple way to incorporate the retrieved documents as part of the input to the LM is to prepend  $x$  with all  $k$  documents. However, this simple scheme is fundamentally restricted by the number of documents (i.e.,  $k$ ) we can include, given the language model’s context window size. To address this limitation, we adopt an ensemble strategy described as follows. Assume  $\mathcal{D}' \subset \mathcal{D}$  consists of  $k$  most relevant documents to  $x$ , according to the scoring

function in Eq. (1). We prepend each document  $d \in \mathcal{D}'$  to  $x$ , pass this concatenation to the LM separately, and then ensemble output probabilities from all  $k$  passes. Formally, given the input context  $x$  and its top- $k$  relevant documents  $\mathcal{D}'$ , the output probability of the next token  $y$  is computed as a weighted average ensemble:

$$p(y | x, \mathcal{D}') = \sum_{d \in \mathcal{D}'} p(y | d \circ x) \cdot \lambda(d, x),$$

where  $\circ$  denotes the concatenation of two sequences and the weight  $\lambda(d, x)$  is based on the similarity score between the document  $d$  and the input context  $x$ :

$$\lambda(d, x) = \frac{e^{s(d,x)}}{\sum_{d \in \mathcal{D}'} e^{s(d,x)}}$$

## 4 REPLUG LSR: Training the Dense Retriever

Instead of relying only on existing neural dense retrieval models (Karpukhin et al., 2020a; Izacard et al., 2022a; Su et al., 2022), we further propose REPLUG LSR (REPLUG with LM-Supervised Retrieval), which adapts the retriever in REPLUG by using the LM itself to provide supervision about which documents should be retrieved.

Inspired by Sachan et al. (2022), our approach can be seen as adjusting the probabilities of the retrieved documents to match the probabilities of the output sequence perplexities of the language model. In other words, we would like the retriever to find documents that result in lower perplexity scores. As shown in Figure 3, our training algorithm consists of the four steps: (1) retrieving documents and computing the retrieval likelihood (§4.1), (2) scoring the retrieved documents by the language model (§4.2), (3) updating the retrieval model parameters by minimizing the KL divergence between the retrieval likelihood and the LM’s score distribution (§4.3), and (4) asynchronous update of the datastore index (§4.4).

### 4.1 Computing Retrieval Likelihood

We retrieve  $k$  documents  $\mathcal{D}' \subset \mathcal{D}$  with the highest similarity scores from a corpus  $\mathcal{D}$  given an input context  $x$ , as described in §3.1. We then compute the retrieval likelihood of each retrieved document  $d$ :

$$P_R(d | x) = \frac{e^{s(d,x)/\gamma}}{\sum_{d \in \mathcal{D}'} e^{s(d,x)/\gamma}}$$

where  $\gamma$  is a hyperparameter that controls the temperature of the softmax. Ideally, the retrieval likelihood is computed by marginalizing over all the documents in the corpus  $\mathcal{D}$ , which is intractable in practice. Therefore, we approximate the retrieval likelihood by only marginalizing over the retrieved documents  $\mathcal{D}'$ .

### 4.2 Computing LM likelihood

We use the LM as a scoring function to measure how much each document could improve the LM perplexity. Specifically, we first compute  $P_{LM}(y | d, x)$ , the LM probability of the ground truth output  $y$  given the input context  $x$  and a document  $d$ . The higher the probability, the better the document  $d_i$  is at improving the LM’s perplexity. We then compute the LM likelihood of each document  $d$  as follows:

$$Q(d | x, y) = \frac{e^{P_{LM}(y|d,x)/\beta}}{\sum_{d \in \mathcal{D}'} e^{P_{LM}(y|d,x)/\beta}}$$

where  $\beta$  is another hyperparameter.

### 4.3 Loss Function

Given the input context  $x$  and the corresponding ground truth continuation  $y$ , we compute the retrieval likelihood and the language model likelihood. The dense retriever is trained by minimizing the KL divergence between these two distributions:

$$\mathcal{L} = \frac{1}{|\mathcal{B}|} \sum_{x \in \mathcal{B}} KL(Q_{LM}(d | x, y) \| P_R(d | x)),$$

where  $\mathcal{B}$  is a set of input contexts. When minimizing the loss, we can only update the retrieval model parameters. The LM parameters are fixed due to our black-box assumption.

### 4.4 Asynchronous Update of the Datastore Index

Because the parameters in the retriever are updated during the training process, the previously computed document embeddings are no longer up to date. Therefore, following Guu et al. (2020), we recompute the document embeddings and rebuild the efficient search index using the new embeddings every  $T$  training steps. Then we use the new document embeddings and index for retrieval, and repeat the training procedure.

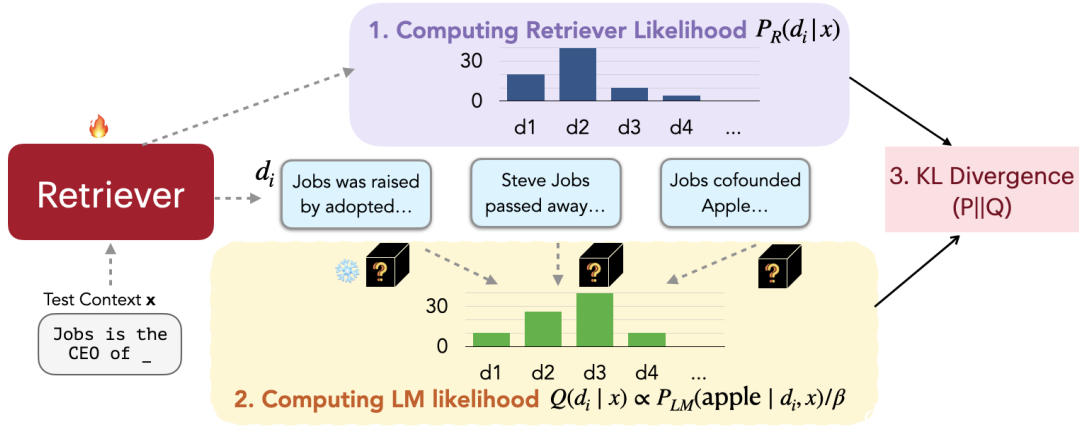


Figure 3: **REPLUG LSR training process (§4)**. The retriever is trained using the output of a frozen language model as supervision signals.

## 5 Training Setup

In this section, we describe the details of our training procedure. We first describe the model setting in REPLUG (§5.1) and then describe the procedure for training the retriever in REPLUG LSR (§5.2).

### 5.1 REPLUG

In theory, any type of retriever, either dense (Karpukhin et al., 2020b; Ni et al., 2021) or sparse (Robertson et al., 2009), could be used for REPLUG. Following prior work (Izacard et al., 2022b), we use the Contriever (Izacard et al., 2022a) as the retrieval model for REPLUG, as it has demonstrated strong performance.

### 5.2 REPLUG LSR

For REPLUG LSR, we initialize the retriever with the Contriever model (Izacard et al., 2022a). We use GPT-3 Curie (Brown et al., 2020b) as the supervision LM to compute the LM likelihood.

**Training data** We use 800K sequences of 256 tokens each, sampled from the Pile training data (Gao et al., 2020), as our training queries. Each query is split into two parts: the first 128 tokens are used as the input context  $x$ , and the last 128 tokens are used as the ground truth continuation  $y$ . For the external corpus  $D$ , we sample 36M documents of 128 tokens from the Pile training data. To avoid trivial retrieval, we ensure that the external corpus documents do not overlap with the documents from which the training queries are sampled.

**Training details** To make the training process more efficient, we pre-compute the document embeddings of the external corpus  $D$  and create a

FAISS index (Johnson et al., 2019) for fast similarity search. Given a query  $x$ , we retrieve the top 20 documents from the FAISS index and compute the retrieval likelihood and the LM likelihood with a temperature of 0.1. We train the retriever using the Adam optimizer (Kingma and Ba, 2015) with a learning rate of  $2e-5$ , a batch size of 64, and a warmup ratio of 0.1. We re-compute the document embeddings every 3k steps and fine-tune the retriever for a total of 25k steps.

## 6 Experiments

We perform evaluations on both language modeling (§6.1) and downstream tasks such as MMLU (§6.2) and open-domain QA (§6.3). In all settings, REPLUG improve the performance of various black-box language models, showing the effectiveness and generality of our approach.

### 6.1 Language Modeling

**Datasets** The Pile (Gao et al., 2020) is a language modeling benchmark that consists of text sources from diverse domains such as web pages, code and academic papers. Following prior work, we report bits per UTF-8 encoded byte (BPB) as the metric on each subset domain.

**Baselines** We consider GPT-3 and GPT-2 family LMs as the baselines. The four models from GPT-3 (Davinci, Curie, Babbage and Ada) are black-box models that are only accessible through API.

**Our model** We add REPLUG and REPLUG LSR to the baselines. We randomly subsampled Pile training data (36M documents of 128 tokens) and use them as the retrieval corpus for all models. As

Model		# Parameters	Original	+ REPLUG	Gain %	+ REPLUG LSR	Gain %
GPT-2	Small	117M	1.33	1.26	5.3	1.21	9.0
	Medium	345M	1.20	1.14	5.0	1.11	7.5
	Large	774M	1.19	1.15	3.4	1.09	8.4
	XL	1.5B	1.16	1.09	6.0	1.07	7.8
GPT-3 (black-box)	Ada	350M	1.05	0.98	6.7	0.96	8.6
	Babbage	1.3B	0.95	0.90	5.3	0.88	7.4
	Curie	6.7B	0.88	0.85	3.4	0.82	6.8
	Davinci	175B	0.80	0.77	3.8	0.75	6.3

Table 1: **Both REPLUG and REPLUG LSR consistently enhanced the performance of different language models.** Bits per byte (BPB) of the Pile using GPT-3 and GPT-2 family models (Original) and their retrieval-augmented versions (+REPLUG and +REPLUG LSR. The gain % shows the relative improvement of our models compared to the original language model.

the Pile dataset has made efforts to deduplicate documents across train, validation and test splits (Gao et al., 2020), we did not do additional filtering. For both REPLUG and REPLUG LSR, we use a length of 128-token context to do retrieval and adopt the ensemble method (Section 3.2) to incorporate top 10 retrieved documents during inference.

**Results** Table 1 reports the results of the original baselines, baselines augmented with the REPLUG, and baselines augmented with the REPLUG LSR. We observe that both REPLUG and REPLUG LSR significantly outperform the baselines. This demonstrates that simply adding a retrieval module to a frozen language model (i.e., the black-box setting) is effective at improving the performance of different sized language models on language modeling tasks. Furthermore, REPLUG LSR consistently performs better than REPLUG by a large margin. Specifically, REPLUG LSR results in 7.7% improvement over baselines compared to 4.7% improvement of REPLUG averaged over the 8 models. This indicates that further adapting the retriever to the target LM is beneficial.

## 6.2 MMLU

**Datasets** MMLU (Hendrycks et al., 2021) is a multiple choice QA dataset that covers exam questions from 57 tasks including mathematics, US history and etc. The 57 tasks are grouped into 4 categories: humanities, STEM, social sciences and other. Following Chung et al. (2022a), we evaluate REPLUG in the 5-shot in-context learning setting.

**Baselines** We consider two groups of strong previous models as baselines for comparisons. The first group of baselines is the state-of-the-art LLMs including Codex<sup>1</sup> (Chen et al.,

<sup>1</sup>Code-Davinci-002

2021b), PaLM (Chowdhery et al., 2022), and Flan-PaLM (Chung et al., 2022b). According to Chung et al. (2022b), these three models rank top-3 in the leaderboard of MMLU. Additionally, we include strong open-source LMs such as LLaMA (Touvron et al., 2023). The second group of baselines consists of retrieval-augmented language models. We only include Atlas (Izacard et al., 2022b) in this group, as no other retrieval-augmented LMs have been evaluated on the MMLU dataset. Atlas trains both the retriever and the language model, which we consider a white-box retrieval LM setting.

**Our model** We add REPLUG and REPLUG LSR to Codex and LLaMA because other models such as PaLM and Flan-PaLM are not accessible to the public. We use the test question as the query to retrieve 10 relevant documents from Wikipedia (2018, December) and prepend each retrieved document to the test question, resulting in 10 separate inputs. These inputs are then separately fed into the language models, and the output probabilities are ensemble together. The retriever interacts with Codex and LLaMA through black-box access.

**Results** Table 2 presents the results from the baselines, REPLUG, and REPLUG LSR on the MMLU dataset. We observe that both the REPLUG and REPLUG LSR improve the original Codex model by 4.5% and 5.1%, respectively. In addition, REPLUG LSR largely outperforms the previous retrieval-augmented language model, Atlas, demonstrating the effectiveness of our black-box retrieval language model setting. Although our models slightly underperform Flan-PaLM, this is still a strong result because Flan-PaLM has three times more parameters. We would expect that the REPLUG LSR could further improve Flan-PaLM, if we had access to the model.

Model	# Parameters	Humanities	Social.	STEM	Other	All
Codex	175B	74.2	76.9	57.8	70.1	68.3
PaLM	540B	77.0	81.0	55.6	69.6	69.3
Flan-PaLM	540B	-	-	-	-	72.2
LLaMA	13B	-	-	-	-	55.6
Atlas	11B	46.1	54.6	38.8	52.8	47.9
Codex + REPLUG	175B	76.0	79.7	58.8	72.1	71.4
Codex + REPLUG LSR	175B	76.5	79.9	58.9	73.2	71.8
LLaMA + REPLUG	13B	-	-	-	-	58.8
LLaMA + REPLUG LSR	13B	-	-	-	-	59.3

Table 2: **REPLUG and REPLUG LSR improves Codex by 4.5% and 5.1% respectively.** Performance on MMLU broken down into 4 categories. The last column averages the performance over these categories. All models are evaluated based on 5-shot in-context learning with direct prompting.

Another interesting observation is that the REPLUG LSR outperforms the original model by 1.9% even in the STEM category. This suggests that retrieval may improve a language model’s problem-solving abilities.

### 6.3 Open Domain QA

Lastly, we conduct evaluation on two open-domain QA datasets: Natural Questions (NQ) (Kwiatkowski et al., 2019) and TriviaQA (Joshi et al., 2017).

Model	NQ		TQA	
	k-shot	Full	k-shot	Full
Chinchilla	35.5	-	64.6	-
PaLM	39.6	-	-	-
Codex	40.6	-	73.6	-
LLaMA	29.0	-	69.6	-
RETRO <sup>†</sup>	-	45.5	-	-
R2-D2 <sup>†</sup>	-	55.9	-	69.9
Atlas <sup>†</sup>	30.9	<b>60.4</b>	74.5	<b>79.8</b>
Codex + REPLUG	44.7	-	76.8	-
Codex + REPLUG LSR	<b>45.5</b>	-	<b>77.3</b>	-
LLaMA + REPLUG	36.1	-	73.3	-
LLaMA + REPLUG LSR	37.2	-	74.1	-

Table 3: Performance on NQ and TQA. We report results for both k-shot (64 shots for Chinchilla, PaLM, and Atlas; 16 shots for Codex-based models) and full data settings. Note that models with <sup>†</sup> are finetuned using training examples, while others use in-context learning.

**Datasets** NQ and TriviaQA are two open-domain QA datasets. Following prior work (Izcard and Grave, 2021b; Si et al., 2023), we report Exact Match for the filtered set of TriviaQA. We consider the k-shot setting where the model is only given a few training examples and full data setting where the model is given all the training examples.

**Baselines** We compare our model with several state-of-the-art baselines, both in a few-shot set-

ting and with full training data. The first group of models consists of powerful large language models, including Chinchilla (Hoffmann et al., 2022), PaLM (Chowdhery et al., 2022), Codex and LLaMA 13B (Touvron et al., 2023). These models are all evaluated using in-context learning under the few-shot setting, with Chinchilla and PaLM evaluated using 64 shots, and Codex using 16 shots. The second group of models for comparison includes retrieval-augmented language models such as RETRO (Borgeaud et al., 2021), R2-D2 (Fajcik et al., 2021), and Atlas (Izcard et al., 2022b). All of these retrieval-augmented models are finetuned on the training data, either in a few-shot setting or with full training data. Specifically, Atlas is finetuned on 64 examples in the few-shot setting.

**Our model** We add REPLUG and REPLUG LSR to Codex and LLaMA 13B with Wikipedia as the retrieval corpus and evaluate them in a 16-shot in context learning. We incorporate top-10 retrieved documents using our proposed ensemble method.

**Results** As shown in Table 3, REPLUG LSR significantly improves the performance of the original Codex by 12.0% on NQ and 5.0% on TQA. It outperforms the previous best model, Atlas, which was fine-tuned with 64 training examples, achieving a new state-of-the-art in the few-shot setting. However, this result still lags behind the performance of retrieval-augmented language models fine-tuned on the full training data. This is likely due to the presence of near-duplicate test questions in the training set (e.g., Lewis et al. (2021) found that 32.5% of test questions overlap with the training sets in NQ).

## 7 Analysis

### 7.1 REPLUG is applicable to diverse models

Here we further study whether REPLUG could enhance *diverse* language model families that have

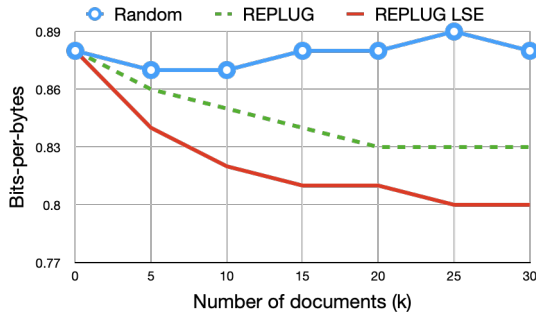


Figure 4: **Ensembling random documents does not result in improved performance.** BPB of Curie augmented with different methods (random, REPLUG and REPLUG LSR) when varying the number of documents.

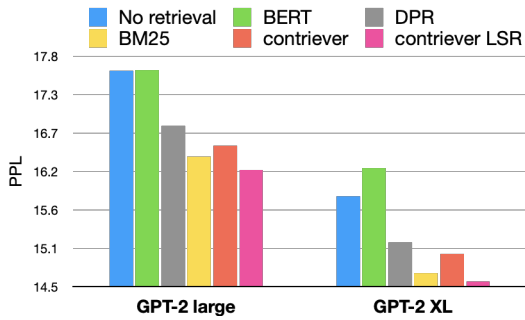


Figure 5: **LM-supervised retriever (Contriever LSR) outperforms other off-the-shelf retrievers.**

504 been pre-trained using different data and methods.  
 505 Specifically, we focus on three groups of language  
 506 models with varying sizes: GPT-2 (117M, 345M,  
 507 774M, 1.5B parameters) (Brown et al., 2020a),  
 508 OPT (125M, 350M, 1.3B, 2.7B, 6.7B, 13B, 30B,  
 509 66B) (Zhang et al., 2022) and BLOOM (560M,  
 510 1.1B, 1.7B, 3B and 7B) (Scao et al., 2022). We  
 511 evaluate each model on Wikitext-103 (Stephen  
 512 et al., 2017) test data and report its perplexity.  
 513 For comparison, we augment each model with RE-  
 514 PLUG that adopts the ensemble method to incorpo-  
 515 rate top 10 retrieved documents. Following prior  
 516 work (Khandelwal et al., 2020), we use Wikitext-  
 517 103 training data as the retrieval corpus.

518 Figure 6 in Appendix A shows the performance  
 519 of different-sized LMs with and without REPLUG.  
 520 We observe that the performance gain brought by  
 521 REPLUG stays consistent with model size. For  
 522 example, OPT-125M achieves 6.9% perplexity im-  
 523 provement, while OPT-66B achieves 5.6% perplex-  
 524 ity improvement. Additionally, REPLUG improves  
 525 the perplexity of all the model families, which in-  
 526 dicates that REPLUG is applicable to diverse lan-  
 527 guage models with different sizes.

## 7.2 REPLUG performance gain does not simply come from the ensembling effect

528  
 529  
 530 The core of our method design is the use of an en-  
 531 semble method that combines output probabilities  
 532 of different passes, in which each retrieved docu-  
 533 ment is prepended separately to the input and fed  
 534 into a language model. To study whether the gains  
 535 come solely from the ensemble method, we com-  
 536 pare our method to ensembling random documents.  
 537 For this, we randomly sample several documents,  
 538 concatenated each random document with the input,  
 539 and ensemble the outputs of different runs (referred  
 540 to as "random"). As shown in Figure 4, we evalu-  
 541 ated the performance of GPT-3 Curie on Pile when  
 542 augmented with random documents, documents  
 543 retrieved by REPLUG, and documents retrieved  
 544 by REPLUG LSR. We observed that ensembling  
 545 random documents leads to worse performance, in-  
 546 dicating that the performance gains of REPLUG  
 547 do not come from the ensembling effect. Instead,  
 548 ensembling the **relevant** documents is crucial for  
 549 the success of REPLUG. Additionally, as more docu-  
 550 ments were ensembled, the performance of RE-  
 551 PLUG and REPLUG LSR improved monotonically.  
 552 However, a small number of documents (e.g., 10)  
 553 was sufficient to achieve large performance gains.

## 7.3 LSR retriever outperforms other off-the-shelf retrievers

554  
 555  
 556 We investigate the effectiveness of tunable retriever  
 557 (LSR) compared with off-the-shelf retrievers.  
 558 Specifically, we compare LM-supervised contriever  
 559 (LSR) with other dense retrievers such as BERT-  
 560 base (Borgeaud et al., 2022), DPR (Karpukhin  
 561 et al., 2020b) and a sparse retriever BM25 (Robert-  
 562 son et al., 2009). Figure 5 shows Wikitext-  
 563 103 perplexity of GPT-2 XL (1.5B) and GPT-2  
 564 Large (774M) augmented with different retrievers.  
 565 Among all off-the-shelf retrievers, the sparse re-  
 566 triever BM25 performs best. However, it still lags  
 567 behind our LM supervised retriever (Contriever  
 568 LSR), demonstrating the effectiveness of our train-  
 569 ing scheme that adapts the retriever to LMs.

## 8 Conclusion

570  
 571 We introduce REPLUG, a retrieval-augmented LM  
 572 paradigm that augments black-box LMs with a  
 573 tuneable retriever. This work opens up new possi-  
 574 bilities for integrating retrieval into large black-box  
 575 LMs and is the first to demonstrate even the state-  
 576 of-the-art LLMs could benefit from retrieval.



577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627

## 9 Limitations

**Interpretability** REPLUG exhibits limitations in interpretability. It’s unclear when the model relies on retrieved knowledge or on knowledge encoded within its own parameters. Future research could work towards the development of more interpretable retrieval-augmented language models. Such models could trace the source of the generated answers, whether it’s from retrieved data or internal parameters, thus providing a clear knowledge provenance.

**On-demand retrieval** REPLUG always perform retrieval no matter if the external information is needed. This approach runs the risk of presenting irrelevant documents, which can potentially distract the models, while also incurring additional computational overheads. Future studies could explore methods that allow the language model to determine when external knowledge is required.

**Database size** In line with prior research, REPLUG uses Wikipedia and Pile as the targeted search databases. However, these resources might only encompass a minor fraction of the external knowledge needed by LMs. Future research should explore methods to efficiently expand these databases and examine how an LM’s performance scales with the size of the database.

## References

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. 2021. Improving language models by retrieving from trillions of tokens. *arXiv preprint arXiv:2112.04426*.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International Conference on Machine Learning*, pages 2206–2240. PMLR.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020a.

Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc. 628  
629  
630  
631

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020b. [Language models are few-shot learners](#). In *Proc. of NeurIPS*. 632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebguss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021a. [Evaluating large language models trained on code](#). *CoRR*, abs/2107.03374. 644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebguss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021b. [Evaluating large language models trained on code](#). *CoRR*, abs/2107.03374. 665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts,

688	Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. <i>arXiv preprint arXiv:2204.02311</i> .	743
689		744
690		745
691		
692	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022a. Scaling instruction-finetuned language models. <i>arXiv preprint arXiv:2210.11416</i> .	746
693		747
694		748
695		
696		
697	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022b. Scaling instruction-finetuned language models. <i>arXiv preprint arXiv:2210.11416</i> .	749
698		750
699		751
700		752
701		753
702	Martin Fajcik, Martin Docekal, Karel Ondrej, and Pavel Smrz. 2021. R2-D2: A modular baseline for open-domain question answering. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 854–870, Punta Cana, Dominican Republic. Association for Computational Linguistics.	754
703		755
704		
705		
706		
707		
708	Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The Pile: An 800gb dataset of diverse text for language modeling. <i>arXiv preprint arXiv:2101.00027</i> .	756
709		757
710		758
711		759
712		760
713		761
714	Anirudh Goyal, Abram Friesen, Andrea Banino, Theophane Weber, Nan Rosemary Ke, Adria Puigdomenech Badia, Arthur Guez, Mehdi Mirza, Peter C Humphreys, Ksenia Konyushova, et al. 2022. Retrieval-augmented reinforcement learning. In <i>International Conference on Machine Learning</i> , pages 7740–7765. PMLR.	762
715		763
716		764
717		
718		
719		
720		
721	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In <i>International Conference on Machine Learning</i> , pages 3929–3938. PMLR.	765
722		766
723		767
724		768
725		769
726	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In <i>International Conference on Learning Representations</i> .	770
727		771
728		
729		
730		
731	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. 2022. Training compute-optimal large language models. <i>arXiv preprint arXiv:2203.15556</i> .	772
732		773
733		774
734		775
735		776
736		777
737	Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A Smith, and Jiebo Luo. 2022. Promptcap: Prompt-guided task-aware image captioning. <i>arXiv preprint arXiv:2211.09699</i> .	778
738		779
739		780
740		781
741	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and	782
742		783
	Edouard Grave. 2022a. Unsupervised dense information retrieval with contrastive learning. <i>Transactions on Machine Learning Research</i> .	784
		785
		786
	Gautier Izacard and Edouard Grave. 2021a. Leveraging passage retrieval with generative models for open domain question answering. In <i>Proc. of EACL</i> .	787
		788
		789
		790
		791
	Gautier Izacard and Edouard Grave. 2021b. Leveraging passage retrieval with generative models for open domain question answering. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 874–880, Online. Association for Computational Linguistics.	792
		793
		794
		795
		796
		797
		798
	Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022b. Few-shot learning with retrieval augmented language models. <i>arXiv preprint arXiv:2208.03299</i> .	799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900

799	Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natural questions: A benchmark for question answering research</a> . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	853
800		854
801		855
802		856
803	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	857
804	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	858
805	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	859
806	täschel, Sebastian Riedel, and Douwe Kiela. 2020.	860
807	<a href="#">Retrieval-augmented generation for knowledge-</a>	861
808	<a href="#">intensive nlp tasks</a> .	
809	Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel.	
810	2021. Question and answer test-train overlap in open-	
811	domain question answering datasets. In <i>Proceedings</i>	
812	<i>of the 16th Conference of the European Chapter of</i>	
813	<i>the Association for Computational Linguistics: Main</i>	
814	<i>Volume</i> , pages 1000–1008.	
815	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi	
816	Das, Hannaneh Hajishirzi, and Daniel Khashabi.	
817	2022. When not to trust language models: Investi-	
818	gating effectiveness and limitations of paramet-	
819	ric and non-parametric memories. <i>arXiv preprint</i>	
820	<i>arXiv:2212.10511</i> .	
821	Sewon Min, Weijia Shi, Mike Lewis, Xilun Chen, Wen-	
822	tau Yih, Hannaneh Hajishirzi, and Luke Zettlemoyer.	
823	2022. Nonparametric masked language modeling.	
824	<i>arXiv preprint arXiv:2212.01349</i> .	
825	Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gus-	
826	tavo Hernández Ábrego, Ji Ma, Vincent Y. Zhao,	
827	Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei	
828	Yang. 2021. <a href="#">Large dual encoders are generalizable</a>	
829	<a href="#">retrievers</a> .	
830	Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang	
831	Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and	
832	Haifeng Wang. 2021. <a href="#">RocketQA: An optimized training</a>	
833	<a href="#">approach to dense passage retrieval for open-</a>	
834	<a href="#">domain question answering</a> . In <i>Proceedings of the</i>	
835	<i>2021 Conference of the North American Chapter of</i>	
836	<i>the Association for Computational Linguistics: Hu-</i>	
837	<i>man Language Technologies</i> , pages 5835–5847, On-	
838	line. Association for Computational Linguistics.	
839	Stephen Robertson, Hugo Zaragoza, et al. 2009. The	
840	probabilistic relevance framework: Bm25 and be-	
841	yond. <i>Foundations and Trends® in Information Re-</i>	
842	<i>trieval</i> , 3(4):333–389.	
843	Devendra Singh Sachan, Mike Lewis, Dani Yogatama,	
844	Luke Zettlemoyer, Joelle Pineau, and Manzil Zaheer.	
845	2022. Questions are all you need to train a dense	
846	passage retriever. <i>arXiv preprint arXiv:2206.10658</i> .	
847	Teven Le Scao, Angela Fan, Christopher Akiki, El-	
848	lie Pavlick, Suzana Ilić, Daniel Hesslow, Roman	
849	Castagné, Alexandra Sasha Luccioni, François Yvon,	
850	Matthias Gallé, et al. 2022. Bloom: A 176b-	
851	parameter open-access multilingual language model.	
852	<i>arXiv preprint arXiv:2211.05100</i> .	
	Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang	853
	Wang, Jianfeng Wang, Jordan Boyd-Graber, and Li-	854
	juan Wang. 2023. <a href="#">Prompting gpt-3 to be reliable</a> . In	855
	<i>Proc. of ICLR</i> .	856
	Merity Stephen, Xiong Caiming, Bradbury James, and	857
	Richard Socher. 2017. Pointer sentinel mixture mod-	858
	els. In <i>5th International Conference on Learning</i>	859
	<i>Representations, ICLR 2017, Toulon, France, April</i>	860
	<i>24-26, 2017, Conference Track Proceedings</i> .	861
	Hongjin Su, Jungo Kasai, Yizhong Wang, Yushi Hu,	862
	Mari Ostendorf, Wen-tau Yih, Noah A Smith, Luke	863
	Zettlemoyer, Tao Yu, et al. 2022. One embedder, any	864
	task: Instruction-finetuned text embeddings. <i>arXiv</i>	865
	<i>preprint arXiv:2212.09741</i> .	866
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	867
	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	868
	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	869
	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	870
	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	871
	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	872
	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	873
	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	874
	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	875
	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	876
	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	877
	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	878
	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	879
	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	880
	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	881
	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	882
	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	883
	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	884
	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	885
	Melanie Kambadur, Sharan Narang, Aurelien Rod-	886
	riguez, Robert Stojnic, Sergey Edunov, and Thomas	887
	Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-</a>	888
	<a href="#">tuned chat models</a> .	889
	Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi,	890
	Rich James, Jure Leskovec, Percy Liang, Mike Lewis,	891
	Luke Zettlemoyer, and Wen-tau Yih. 2022. Retrieval-	892
	augmented multimodal language modeling. <i>arXiv</i>	893
	<i>preprint arXiv:2211.12561</i> .	894
	Tim Dettmers Younes Belkda. 2022. <a href="#">A gentle introduc-</a>	895
	<a href="#">tion to 8-bit matrix multiplication</a> .	896
	Wenhao Yu. 2022. <a href="#">Retrieval-augmented generation</a>	897
	<a href="#">across heterogeneous knowledge</a> . In <i>Proceedings</i>	898
	<i>of the 2022 Conference of the North American Chap-</i>	899
	<i>ter of the Association for Computational Linguistics:</i>	900
	<i>Human Language Technologies: Student Research</i>	901
	<i>Workshop</i> , pages 52–58, Hybrid: Seattle, Washington	902
	+ Online. Association for Computational Linguistics.	903
	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	904
	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	905
	wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.	906
	Opt: Open pre-trained transformer language models.	907
	<i>arXiv preprint arXiv:2205.01068</i> .	908
	Zexuan Zhong, Tao Lei, and Danqi Chen. 2022. Train-	909
	ing language models with memory augmentation. In	910

**A REPLUG is applicable to diverse models** 913  
914

**B Qualitative Analysis: Rare Entities Benefit from Retrieval** 915  
916

To understand why the REPLUG improves language modeling performance, we conducted manual analysis of examples in which the REPLUG results in a decrease in perplexity. We find that REPLUG is more helpful when texts contain rare entities. [Figure 7](#) shows a test context and its continuation from the Wikitext-103 test set. For REPLUG, we use the test context as a query to retrieve a relevant document from Wikitext-103 training data. We then compute the perplexity of the continuation using the original GPT-2 1.5B and its REPLUG enhanced version. After incorporating the retrieved document, the perplexity of the continuation improves by 11%. Among all tokens in the continuation, we found that REPLUG is most helpful for the rare entity name "Li Bai". This is likely because the original LM does not have sufficient information about this rare entity name. However, by incorporating the retrieved document, REPLUG was able to match the name with the relevant information in the retrieved document, resulting in better performance. 917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938

**C Prompts used for MMLU and open-domain QA** 939  
940

Please see [Table 4](#) and [Table 5](#). 941

**D Dense Retriever vs. Sparse Retriever** 942

The proposed model uses Contriever, a dense retriever, as its retriever backbone. Additionally, we investigate the performance of a sparse retriever in comparison to the dense retriever. For our sparse model, we employ BM25. As depicted in [Figure 8](#), we observe that BM25 consistently outperforms Contriever but falls short when compared to LM-supervised Contriever, thus highlighting the effectiveness of our proposed training scheme. 943  
944  
945  
946  
947  
948  
949  
950  
951

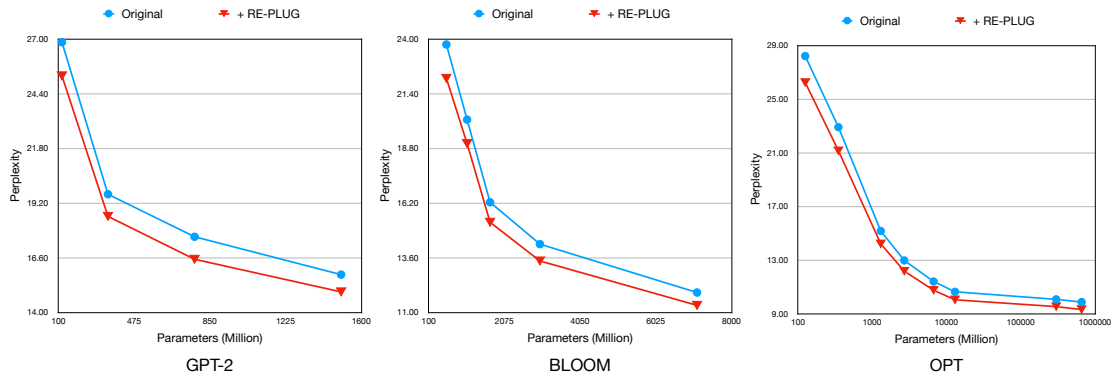


Figure 6: **GPT-2, BLOOM and OPT models of varying sizes consistently benefit from REPLUG.** The x-axis indicates the size of the language model and the y-axis is its perplexity on Wikitext-103.

---

**Knowledge:** Arctic Ocean. Although over half of Europe’s original forests disappeared through the centuries of deforestation, Europe still has over one quarter of its land area as forest, such as the broadleaf and mixed forests, taiga of Scandinavia and Russia, mixed rainforests of the Caucasus and the Cork oak forests in the western Mediterranean. During recent times, deforestation has been slowed and many trees have been planted. However, in many cases monoculture plantations of conifers have replaced the original mixed natural forest, because these grow quicker. The plantations now cover vast areas of land, but offer poorer habitats for many European

**Question:** As of 2015, since 1990 forests have \_\_\_\_\_in Europe and have \_\_\_\_\_in Africa and the Americas.

A. "increased, increased" B. "increased, decreased" C. "decreased, increased" D. "decreased, decreased"

**Answer:** B

**Knowledge:** Over the past decades, the political outlook of Americans has become more progressive, with those below the age of thirty being considerably more liberal than the overall population. According to recent polls, 56% of those age 18 to 29 favor gay marriage, 68% state environmental protection to be as important as job creation, 52% "think immigrants strengthen the country with their hard work and talents," 62% favor a "tax financed, government-administrated universal health care" program and 74% "say people’s willshould have more influence on U.S. laws than the Bible, compared to 37%, 49%, 38%, 47% and 58% among the

**Question:** As of 2019, about what percentage of Americans agree that the state is run for the benefit of all the people?

A. 31% B. 46% C. 61% D. 76%

**Answer:** B

...

**Knowledge:** last week at a United Nations climate meeting in Germany, China and India should easily exceed the targets they set for themselves in the 2015 Paris Agreement... India is now expected to obtain 40 percent of its electricity from non-fossil fuel sources by 2022, eight years ahead of schedule." Solar power in Japan has been expanding since the late 1990s. By the end of 2017, cumulative installed PV capacity reached over 50 GW with nearly 8 GW installed in the year 2017. The country is a leading manufacturer of solar panels and is in the top 4 ranking for countries

**Question:** Which of the following countries generated the most total energy from solar sources in 2019?

A. China B. United States C. Germany D. Japan

---

Table 4: Prompt for MMLU

---

**Knowledge:** received 122,000 buys (excluding WWE Network views), down from the previous year's 199,000 buys. The event is named after the Money In The Bank ladder match, in which multiple wrestlers use ladders to retrieve a briefcase hanging above the ring. The winner is guaranteed a match for the WWE World Heavyweight Championship at a time of their choosing within the next year. On the June 2 episode of "Raw", Alberto Del Rio qualified for the match by defeating Dolph Ziggler. The following week, following Daniel Bryan being stripped of his WWE World Championship due to injury, Stephanie McMahon changed the

**Question:** Who won the mens money in the bank match?

**Answer:** Braun Strowman

**Knowledge:** in 3D on March 17, 2017. The first official presentation of the film took place at Disney's three-day D23 Expo in August 2015. The world premiere of "Beauty and the Beast" took place at Spencer House in London, England on February 23, 2017; and the film later premiered at the El Capitan Theatre in Hollywood, California, on March 2, 2017. The stream was broadcast onto YouTube. A sing along version of the film released in over 1,200 US theaters nationwide on April 7, 2017. The United Kingdom received the same version on April 21, 2017. The film was re-released in

**Question:** When does beaty and the beast take place

**Answer:** Rococo-era

...

**Knowledge:** Love Yourself "Love Yourself" is a song recorded by Canadian singer Justin Bieber for his fourth studio album "Purpose" (2015). The song was released first as a promotional single on November 8, 2015, and later was released as the album's third single. It was written by Ed Sheeran, Benny Blanco and Bieber, and produced by Blanco. An acoustic pop song, "Love Yourself" features an electric guitar and a brief flurry of trumpets as its main instrumentation. During the song, Bieber uses a husky tone in the lower registers. Lyrically, the song is a kiss-off to a narcissistic ex-lover who did

**Question:** love yourself by justin bieber is about who

---

Table 5: Prompt for open-domain QA

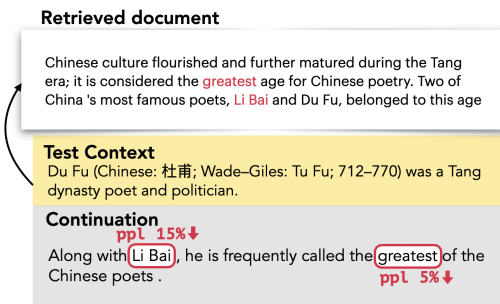


Figure 7: **Rare entities benefit from retrieval.** After incorporating the retrieved document during inference, the entity "Li Bai" and the token "greatest" in the continuation show the most improvement in perplexity (15% for "Li Bai" and 5% for "greatest"). Other tokens' perplexity changes are within 5%.



Figure 8: PPL of GPT-2 models on Witext-103 with no retrieval (Origin), Contriever (REPLUG), LM-supervised Contriever (REPLUG LSR) and BM25.