

# Nearest Neighbor Zero-Shot Inference

Weijia Shi Julian Michael Suchin Gururangan Luke Zettlemoyer

Paul G. Allen School of Computer Science & Engineering,

University of Washington, Seattle, WA

{swj0419, julianjm, sg01, lsz}@cs.washington.edu

## Abstract

We introduce kNN-Prompt, a simple and effective technique to use k-nearest neighbor (kNN) retrieval augmentation (Khandelwal et al., 2021) for zero-shot inference with language models (LMs). Key to our approach is the introduction of *fuzzy verbalizers* which leverage the sparse kNN distribution for downstream tasks by automatically associating each classification label with a set of natural language tokens. Across eleven diverse end-tasks (spanning text classification, fact retrieval and question answering), using kNN-Prompt with GPT-2 Large yields significant performance boosts over zero-shot baselines (14% absolute improvement over the base LM on average). Extensive experiments show that kNN-Prompt is effective for domain adaptation with no further training, and that the benefits of retrieval increase with the size of the model used for kNN retrieval. Overall, we show that augmenting a language model with retrieval can bring significant gains for zero-shot inference, with the possibility that larger retrieval models may yield even greater benefits.

## 1 Introduction

Retrieval-augmented language models (LMs) have access to a non-parametric memory, allowing them to directly access a large external text collection during inference. Previous work has shown that these models substantially outperform their non-retrieval-based counterparts on language modeling tasks (Khandelwal et al., 2020; He et al., 2021; Borgeaud et al., 2021), but it is an open question whether they also achieve similar gains in few-shot and zero-shot end task evaluations (Radford et al., 2019; Brown et al., 2020a). In this paper, we demonstrate that, with some extensions to improve coverage, the performance gains of retrieval-augmented LMs generalize well to a wide range of downstream tasks.

We study the k-nearest neighbors language model (Khandelwal et al., 2020, kNN-LM), which

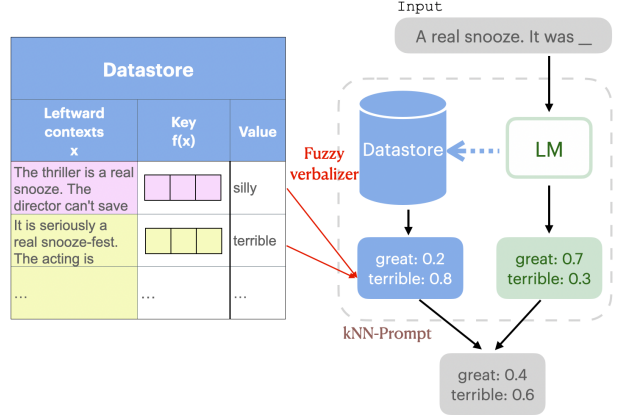


Figure 1: kNN-Prompt incorporates information from a large, heterogeneous corpus to facilitate zero-shot inference. The datastore contains key-value pairs where the key is an encoding of a leftward context and the value is the next token following the context.

interpolates a neural LM’s output distribution with a nearest-neighbor distribution constructed by retrieving tokens from a corpus using the LM’s output embeddings. We are the first to study kNN-LM’s zero-shot application to end tasks, and what we find is that applying the technique naïvely only produces marginal improvements (Section 4). The main challenge is that the support of the kNN distribution is sparse (covering at most k tokens, often less), as it only assigns probability mass to nearest neighbors. This means it often entirely misses the tokens that are used to verbalize the output label in the standard application of LMs to zero-shot classification: across the datasets we test, an output label receives nonzero probability under the kNN distribution only 45.8% of the time (see Section 6).

To address this challenge, we introduce kNN-Prompt, a simple and effective method built on kNN-LM for improving zero-shot inference with no further training. Key to our approach are *fuzzy verbalizers*, which automatically expand the set of tokens corresponding to each output label. For

example, in Figure 1, the verbalized label of the negative sentiment is “terrible”. Our fuzzy verbalizer also maps “silly” to negative sentiment, allowing the model to better leverage the information available in the kNN distribution.

Extensive experiments (Section 3) show that using kNN-Prompt with a heterogeneous datastore consistently improves an LM’s zero-shot abilities on eleven tasks, including sentiment analysis, topic classification, entailment, fact retrieval and question answering. These improvements hold for every model in the GPT-2 family. Furthermore, kNN-Prompt can be adapted to new domains and tasks with no further training (Section 5). With a domain-specific datastore corpus, we achieve comparable or better performance to prompting the LM after domain-adaptive pretraining (Gururangan et al., 2020) on that corpus. To better understand these gains, we conduct a thorough analysis (Section 6), showing that fuzzy verbalizers are essential for leveraging the kNN distribution, the benefits of retrieval increase with retrieval model size, and even relatively small datastores can yield sizeable performance gains if they are tailored to the domain or task.

Overall, our results show how retrieval can benefit zero-shot inference with LMs on a wide variety of tasks, and suggest that applying retrieval with larger models may yield even greater benefits. Code and data are available at [github.com/swj0419/kNN\\_prompt](https://github.com/swj0419/kNN_prompt).

## 2 Method

To perform zero-shot prediction on a downstream task using a pretrained language model, we recast the task as language modeling (Radford et al., 2019) by converting each input instance into a natural language prompt (Section 2.1). We then augment the pretrained model with the k-nearest-neighbors language modeling technique from Khandelwal et al. (2020). To better benefit from the sparse kNN distribution, we introduce *fuzzy verbalizers* for mapping from the LM’s outputs to a distribution over task-specific labels (Section 2.3). Finally, we decode the output from this label distribution using the domain-conditional PMI scoring method of Holtzman et al. (2021).

### 2.1 Prompting and Verbalizers

We address classification problems where an instance consists of an input sequence of tokens

$\mathbf{x} = (x_0, x_1, \dots, x_{|\mathbf{x}|})$  from a vocabulary  $\mathcal{V}$  and an output label  $y \in Y$ . The output label set  $Y$  may be fixed for the task (*text classification*) or provided for each instance as a set of expressions in  $\mathcal{V}^*$  (*multiple-choice*). For example, in the sentiment analysis example in Figure 2, the input is  $\mathbf{x}$  = “Mr. Tsai is one of world cinema’s most gifted artists.” The output labels are  $Y = \{y^+, y^-\}$ , referring to positive and negative sentiment.

To cast the task as language modeling, we deterministically transform each input example  $\mathbf{x}$  into a **prompt**  $p(\mathbf{x})$ . Providing this prompt to an LM yields a probability distribution

$$P_{\text{LM}}(\mathbf{z} \mid p(\mathbf{x})) = \prod_{z_i} P_{\text{LM}}(z_i \mid p(\mathbf{x}), \mathbf{z}_{<i})$$

over continuation sequences  $\mathbf{z} \in \mathcal{V}^*$ . To extract an output label from these continuations, we apply *verbalizers*  $V : y \rightarrow \mathcal{V}^*$  (Schick and Schütze, 2021) which map each output label  $y \in Y$  to a natural language expression  $V(y) = \mathbf{z}$ . We can then compute a probability for each label:

$$P(y \mid \mathbf{x}) \propto P_{\text{LM}}(V(y) \mid p(\mathbf{x})), \quad (1)$$

normalizing over all  $y \in Y$ .

For example, our prompt transformation for sentiment analysis adds *It was* after the input, and uses the verbalizer  $V(y^+) = \textit{great}$ ,  $V(y^-) = \textit{terrible}$ , which classifies sentiment according to the relative probabilities of *It was great* and *It was terrible* after the input sequence (see Figure 2, bottom left). In the case of multiple-choice problems, our verbalizer is just the identity function.

### 2.2 k-Nearest Neighbors Language Modeling

Following Khandelwal et al. (2020), we augment the LM with a *datastore* from which it can retrieve tokens that inform its predictions, improving performance without further training.

The datastore is a key-value store generated by running the LM over a corpus of text. Each value is a token  $w \in \mathcal{V}$  from the corpus, and its key is the vector hidden representation at the output layer of the LM running forward on the left context  $\mathbf{c} \in \mathcal{V}^*$  (call this  $f(\mathbf{c})$ ). At inference time, when predicting the next token for an input sequence  $\mathbf{c}$ , the kNN-LM retrieves the  $k$  nearest neighbors of  $\mathbf{c}$  from the datastore according to the distance  $d(\cdot, f(\mathbf{c}))$  of their key vectors.<sup>1</sup>

<sup>1</sup>In this work, we use Euclidean distance; Khandelwal et al. use its square.

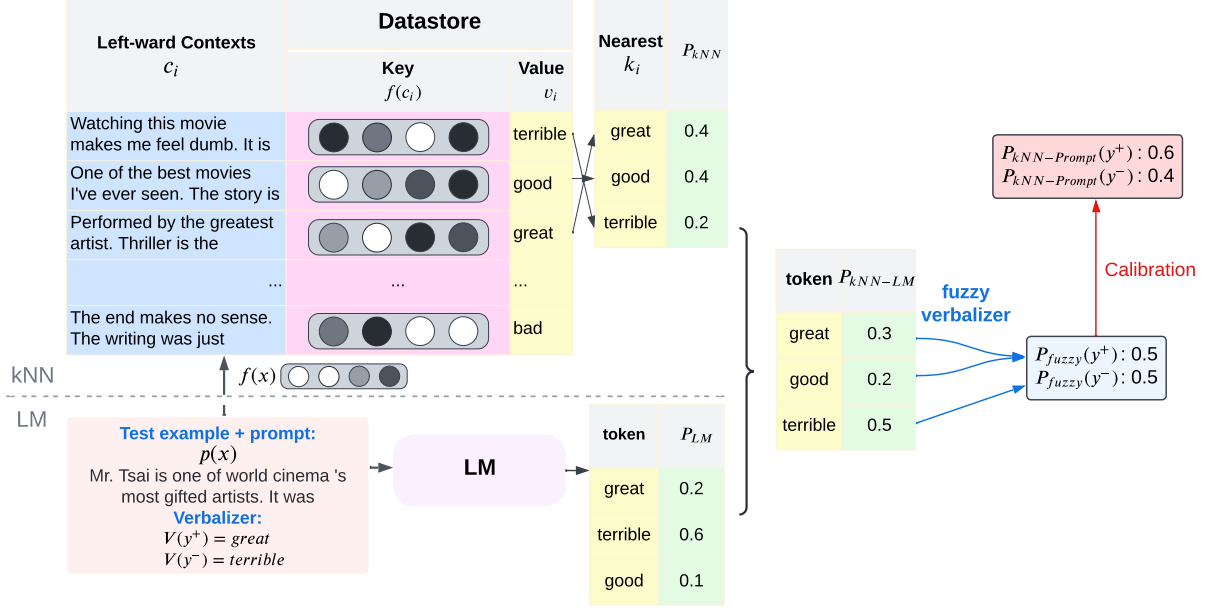


Figure 2: An illustration of kNN-Prompt applying to sentiment analysis tasks. Texts are encoded in the datastore, where each entry consists of a representation of a leftward context and its next token. During inference, a test example is mapped to a prompt form and used to retrieve the  $k$  most similar contexts and their next tokens from the datastore. kNN distribution is a multinomial computed on the distance of the text example and similar contexts. The final prediction is formed by combining the kNN distribution with the language model’s output distribution.

A softmax over the (negative) distances induces a distribution over the the tokens  $w_i$  in the nearest neighbor set:

$$P_{kNN}(v \mid c) \propto \sum_{(f(c_i), w_i)} \mathbb{1}_{v=w_i} e^{\frac{-d(f(c_i), f(c))}{t}}$$

where  $t$  is a temperature parameter.<sup>2</sup> We can then interpolate this with the original LM as follows:

$$P_{kNN-LM}(v \mid c) = (1 - \lambda)P_{LM}(v|c) + \lambda P_{kNN}(v|c).$$

The hyperparameters for the kNN-LM approach are the number  $k$  of nearest neighbors, the interpolation constant  $\lambda$ , the temperature  $t$ , and the choice of datastore.

### 2.3 Fuzzy verbalizers

One challenge in performing zero-shot inference with LMs on downstream tasks is the choice of verbalizer. On one hand, LMs may be highly sensitive to the particular surface form in ways that are irrelevant to the classification task (Holtzman et al., 2021). On the other hand, for a kNN model, the  $k$  nearest neighbor set is sparse and may fail

to cover any of the tokens in the set of verbalizers (i.e.,  $P_{kNN}(V(y)) = 0$  for all  $y \in Y$ ), limiting its utility in those cases. To address these issues, we introduce *fuzzy verbalizers*, which associate each label  $y$  with a neighborhood of token sequences around a specific verbalization  $V(y) \in \mathcal{V}^*$ .

To do this, we first associate each token  $v \in \mathcal{V}$  with a neighborhood  $\mathcal{N}(v) \subseteq \mathcal{V}$  of similar tokens. In particular, we use  $v$ ’s top-5 most similar words according to the cosine similarity of their GloVe embeddings (Pennington et al., 2014), as well as any of  $v$ ’s synonyms in ConceptNet (Speer et al., 2017).<sup>3</sup> Then, for the purposes of calculating the probability of a verbalized label  $z = V(y)$ , we treat a prediction of any token in each  $z_i$ ’s neighborhood as a viable substitute for it, marginalizing over  $\mathcal{N}(z_i)$  at each timestep:

$$P_{FV}(y \mid x) \propto \prod_{z_i \in V(y)} \sum_{v \in \mathcal{N}(z_i)} P(v \mid p(x), z_{<i}) \quad (2)$$

This incorporates more information from the LM to inform the induced distribution over labels  $P_{FV}(y \mid x)$ , and in the case of a kNN-based model, helps mitigate the effect the sparsity of the kNN distribution has on zero-shot prediction (see Section 6).

<sup>2</sup>We have added the temperature adjustment in the softmax on top of Khandelwal et al.’s kNN-LM formulation.

<sup>3</sup><https://conceptnet.io>

Corpus	Size	# Tokens
Wikitext-103	181MB	114M
Amazon Reviews	89MB	19M
CC-NEWS	457MB	324M
IMDB	45MB	8M
Total	722MB	465M

Table 1: Statistics of our heterogeneous datastore corpora.

## 2.4 Full model

To make a zero-shot prediction for an input  $\mathbf{x}$ , we first transform it into a prompt  $p(\mathbf{x})$  and obtain a distribution over continuations ( $z$ ) with a kNN-LM:  $P_{\text{kNN-LM}}(z \mid p(\mathbf{x}))$ . We then convert this to a probability distribution over output labels  $P(y \mid p(\mathbf{x}))$  using a fuzzy verbalizer (Section 2.3, Equation 2). Finally, we output the best label according to the *domain-conditional PMI* scoring rule (Holtzman et al., 2021):

$$\text{PMI}_{\text{DC}}(y, p(\mathbf{x})) = \log \frac{P(y \mid p(\mathbf{x}))}{P(y \mid \mathbf{p})},$$

where  $\mathbf{p}$  is a task-dependent string which is independent of the particular input (generally the local context at the end of the prompt, e.g., we use  $\mathbf{p} = \text{"It was"}$  for sentiment analysis, as shown in Figure 2).

## 3 Experimental Setup

### 3.1 Tasks

We experiment with 11 tasks, including fact retrieval, question answering, topic classification, sentiment analysis, entailment and partisanship classification.

**Topic Classification** We use the AG News (AGN) and Yahoo! Answers (Yahoo) corpora from Zhang et al. (2015) for topic classification.

**Sentiment and Partisanship** We study sentiment analysis using the Rotten Tomatoes (RT) and SST-2 corpora of Socher et al. (2013), movie reviews from Pang and Lee (2005, MR), the customer review dataset from Hu and Liu (2004, CR) consisting of Amazon and Yelp reviews, and the hyperpartisan news detection dataset from Kiesel et al. (2019, HYP), which focuses on classifying whether a text exhibits extreme political views.

**Entailment** Entailment datasets focus on classifying whether one sentence plausibly implies an-

other to be true or false. We evaluate on the CommitmentBank (de Marneffe et al., 2019, CB) and the Recognizing Textual Entailment (Dagan et al., 2010, RTE) dataset provided in GLUE (Wang et al., 2018).

### Fact Retrieval and Question Answering

We evaluate fact retrieval with the LAMA probe (Petroni et al., 2019), which tests an LM’s recovery of factual subject-relation-object triples using a cloze format (e.g., Dante was born in [Mask]). We use test examples where the missing token is at the end of the sentence (suitable for left-to-right LMs) and we report the mean accuracy across all triples. For question answering, we consider CommonsenseQA (Talmor et al., 2019, CQA), consisting of multiple-choice common-sense questions authored by crowd-workers on the basis of knowledge encoded in ConceptNet (Speer et al., 2017). Since the the answers for LAMA and CommonsenseQA can be any string, we perform beam decoding from our LM to produce a set of possible outputs and proceed as in the multiple-choice case.

### 3.2 kNN-Prompt Model Details

**Inference Model** For our main experiments, we directly use GPT-2 large from Huggingface<sup>4</sup> as our base LM. We consider other model sizes in Section 6.

**Retriever Model** Following the inference model, we use GPT-2 large to build the datastore. The keys are the 1280-dimensional hidden representations before the final MLP which predicts the token distribution at each timestep, produced using a single forward pass over the datastore corpus. For efficient similarity search, we create a FAISS (Johnson et al., 2019) index and search for nearest neighbors by Euclidean distance.

**Datastore Corpus** For our datastore, we aim to curate a large, heterogeneous corpus of data broadly relevant to the tasks we evaluate. To this end, we combine four sources of data including Wikitext-103 (Merity et al., 2016), the Amazon review corpus of He and McAuley (2016), and subsets of CC-NEWS<sup>5</sup> and IMDB<sup>6</sup> sampled uniformly from each. Table 1 lists the specifics of each data source.

<sup>4</sup><https://github.com/huggingface/transformers>

<sup>5</sup>[https://huggingface.co/datasets/cc\\_news](https://huggingface.co/datasets/cc_news)

<sup>6</sup><https://datasets.imdbws.com>



	RTE	CQA	CB	Yahoo	LAMA	RT	SST-2	CR	MR	HYP	AGN	Avg
LM	47.6	33.3	48.1	37.7	12.1	53.8	77.6	73.5	55.6	58.4	75.1	52.5
LM+PMI	54.6	46.5	52.0	45.6	18.3	78.9	74.1	67.3	78.7	51.2	66.9	57.6
kNN-LM	52.3	42.9	49.6	36.2	16.8	55.3	78.2	76.8	53.0	59.2	70.2	53.8
kNN-Prompt	<b>54.9</b>	<b>49.2</b>	<b>57.2</b>	<b>53.4</b>	<b>29.5</b>	<b>84.1</b>	<b>87.8</b>	<b>84.7</b>	<b>82.3</b>	<b>60.5</b>	<b>77.1</b>	<b>66.6</b>

Table 2: Zero-shot results on a variety of tasks. Our model, kNN-Prompt, handily outperforms Holtzman et al. (2021)’s PMI scoring method alone (LM+PMI) as well as the base kNN-LM method of Khandelwal et al. (2020).

	CR	HYP	LAMA
LM	82.6 <sub>4.1</sub>	59.0 <sub>0.5</sub>	10.5 <sub>2.4</sub>
LM+PMI	73.3 <sub>5.5</sub>	58.8 <sub>2.6</sub>	18.2 <sub>3.7</sub>
kNN-LM	82.3 <sub>4.2</sub>	58.9 <sub>1.5</sub>	18.9 <sub>1.6</sub>
kNN-prompt	<b>84.8<sub>1.7</sub></b>	<b>63.4<sub>1.1</sub></b>	<b>29.2<sub>1.2</sub></b>

Table 3: The mean and standard deviation for 4 uniformly sampled sets of 4 demonstration examples used for few-shot inference.

**Inference Procedure** We retrieve  $k=512$  neighbors, soften the kNN distribution with a temperature value of 3 and use an interpolation factor of  $\lambda = 0.3$ . Our primary evaluation is zero-shot. All hyperparameters were chosen on the basis of development experiments (see Section 6 for more details).

### 3.3 Baselines

**LM** is the result of prompting the base language model (GPT-2 Large), choosing the output label whose verbalizer has the highest probability under the language model  $P_{LM}(V(y) | p(\mathbf{x}))$ .

**LM+PMI** is the approach of Holtzman et al. (2021), calibrating **LM** with domain-conditional PMI scoring (Section 2.4).

**kNN-LM** directly applies the kNN-LM of Khandelwal et al. (2020) in the same way as **LM**, choosing the highest-probability output label.

## 4 Experimental Results

Results for zero-shot prediction are in Table 2. kNN-Prompt outperforms all baselines in all tasks, improving over the base LM by 14.1% on average. The gains are particularly pronounced for MR and RT (sentiment analysis on movie reviews), Yahoo (topic classification), and LAMA (fact recovery). For MR and RT, the gains seem to come mostly from PMI calibration. On the other hand, large performance boosts on LAMA only come with the full kNN-Prompt model, which indicates the importance of combining retrieval, fuzzy verbalization,

and PMI calibration for this task.

Interestingly, the kNN-LM alone yields a fairly small improvement over the base LM (about 1–2% on average). It is not strong enough to outperform LM+PMI even on LAMA, which intuitively should benefit from retrieval. This suggests that the fuzzy verbalizer and PMI calibration methods may help kNN-Prompt better leverage the information in the  $k$ -nearest neighbors distribution. We carefully examine possible sources of kNN-Prompt’s performance gains in Section 6.

**Few-shot inference** For a subset of tasks, we additionally compare to a few-shot setting where we prepend four examples uniformly sampled from the training data to the input (Table 3). The demonstration examples are converted to prompt and verbalizer format. We report the mean accuracy and standard deviation with 4 different random seeds. We find that kNN-Prompt consistently outperform baselines, demonstrating that kNN-Prompt is applicable to the few-shot setting as well. We leave further exploration of this phenomenon to future work.

## 5 kNN-Prompt for Domain Adaptation

One of the advantages of retrieval-based LMs is that they can be adapted to new domains with no further training.

To test this capability, we replace our heterogeneous datastore (Section 3.2) with domain-specific ones for several tasks. To build these domain-specific datastores, we select Amazon Reviews for CR, CC-NEWS for HYP and Wikitext-103 for LAMA, and encode them separately.

For comparison, we consider domain-adaptive pretraining (Gururangan et al., 2020, DAPT), which further trains the LM on the domain-specific corpus. We train GPT-2 Large on each domain corpus for a single pass, then apply it to downstream tasks using our prompting and verbalizer setup and domain-conditional PMI scoring.

As shown in Table 4, kNN-Prompt performs

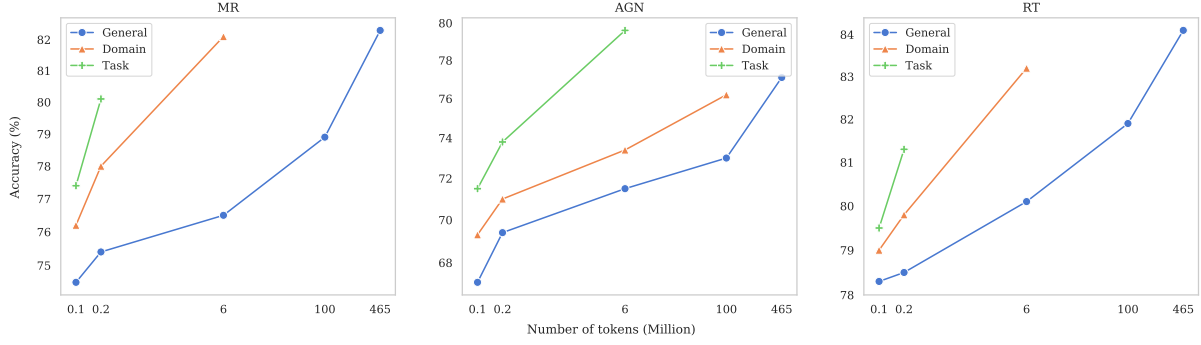


Figure 3: Effect of the number of tokens in the datastore on MR, AGN and RT. Each line represents the kNN-Prompt model with a different datastore and the line ends when the entire available datastore is used. General, Domain, and Task refer to the heterogeneous corpus (Table 1), domain-specific corpus, and task-specific corpus, respectively. We use IMDB as the domain-specific corpus for MR and RT, and CC-NEWS for AGN. The task-specific corpus is the unlabeled training data of each task. GPT-2 Large is used for both retriever and inference models.

	CR	HYP	LAMA
LM + PMI	67.3	51.2	18.3
kNN-prompt	<b>85.1</b>	60.8	<b>28.3</b>
DAPT (LM + PMI)	84.9	<b>61.3</b>	24.6

Table 4: Domain adaptation experiments using domain-specific datastores. DAPT requires training the LM on the corresponding datastore, while kNN-Prompt can use it as the datastore with no further training.

comparably with DAPT. Specifically, while DAPT slightly outperforms kNN-Prompt on CR and HYP, kNN-Prompt shows a clear advantage over DAPT on LAMA. These results indicate that kNN-Prompt is an effective method for domain adaptation. Critically, unlike DAPT, kNN-Prompt does not require further training, which is more practical and efficient for adapting very large LMs.

**Effect of datastore distribution and size** To better understand kNN-Prompt’s potential for domain adaptation, we experiment with varying sizes and distributions of the datastore. For each task, we consider three options for the datastore corpus: our heterogeneous corpus (Section 3.2), a domain-specific corpus, and a task-specific corpus constructed from the task’s (unlabeled) training data. Each of these data sources exhibits increasing levels of relevance to the task.

Figure 3 shows how model performance varies with the choice of datastore across different datastore sizes. For a fixed number of tokens, retrieving from a task-specific datastore is best. Furthermore, token-for-token, adding task-specific data leads to more gains than domain-specific data, which in turn is better than our heterogeneous corpus.

Model	Acc.	$\Delta$ Acc.
LM	52.5	0
LM+kNN (kNN-LM)	53.8	+1.3
LM+Fuzzy	58.0	+5.5
LM+PMI	57.8	+5.3
LM+Fuzzy+PMI	60.1	+7.6
LM+kNN+Fuzzy	62.9	+10.4
LM+kNN+PMI	58.2	+5.7
LM+kNN+Fuzzy+PMI (kNN-Prompt)	<b>66.6</b>	<b>+14.1</b>

Table 5: Effect of different components on the average zero-shot accuracy across the eleven tasks.

When a sufficient amount of task-specific data is available, using it for the datastore can outperform a vastly larger corpus. For example, for AGN, using 6M tokens of unlabeled training data outperforms using our 465M token heterogeneous corpus. However, in cases where the task training data is smaller, using large heterogeneous data can be more effective, as increasing the number of tokens in the datastore always improves task performance for all of the tasks and data distributions we test. These results suggest that while having a large datastore is beneficial, curating task-specific data can also be an effective way of improving model performance, especially if datastore size is limited (e.g., due to memory constraints).

## 6 Analysis

We perform several experiments to understand the contribution of each component of kNN-Prompt and inform our choice of hyperparameters.

**Model ablations** kNN-Prompt incorporates three features on top of the base LM: kNN retrieval and

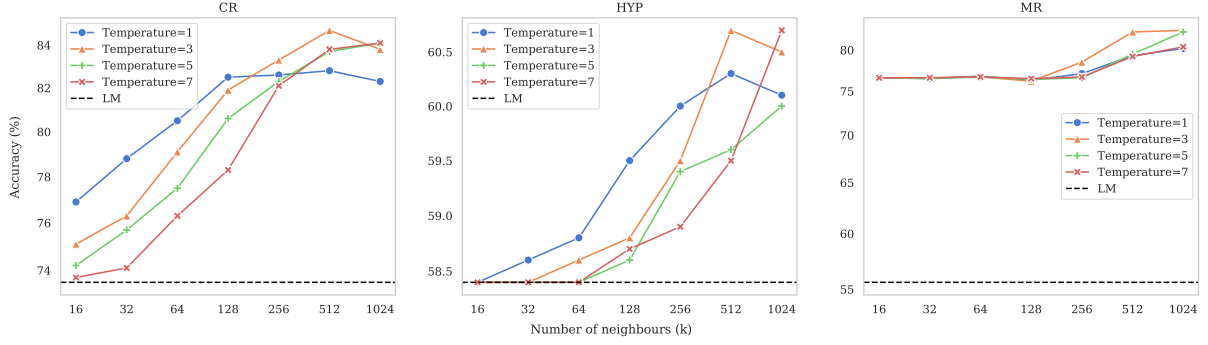


Figure 4: Effect of the number of retrieved neighbors and softmax temperature on kNN-Prompt performance for three tasks: CR, HYP and MR. Task performance monotonically improves with the number of neighbors as  $k$  is increased to 512.

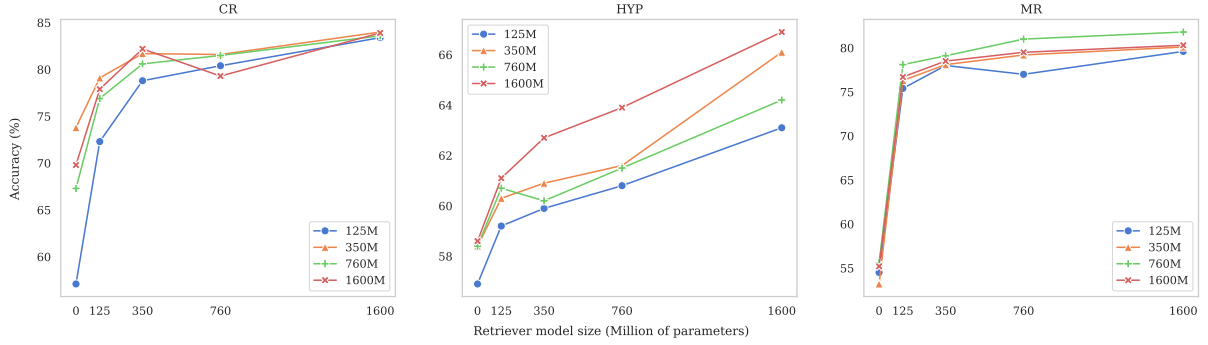


Figure 5: Effect of the retriever model size (GPT-2) on three tasks: CR, HYP and RTE. A size of 0 indicates that no retriever is used. Different lines represent different-sized inference models (GPT-2). The benefits of kNN-Prompt scales with the retriever model size.

interpolation (Section 2.2), fuzzy verbalizers (Section 2.3), and PMI scoring (Section 2.4). Table 5 shows the results of ablation experiments analyzing the contribution of each component.

First, we find that the gains from kNN retrieval alone are modest (+1.3%), but much greater once we add fuzzy verbalizers on top of them (+10.4%), exceeding the contribution of the two components independently (with fuzzy verbalizers alone at +5.5%). This supports the argument that fuzzy verbalizers allow the model to make better use of the sparse support of the kNN distribution. Indeed, we find that across all tasks, an output label receives nonzero probability under the kNN distribution in kNN-LM only 45.8% of the time. With fuzzy verbalizers, this increases to 78%.

Second, we find that for the base LM, fuzzy verbalizers bring gains (+5.5%) similar to PMI scoring (+5.3%), but the gains are only partially additive when combining the two techniques (+7.6%). This suggests that by incorporating more varied surface forms into the score for each label, fuzzy verbalizers may partially — but not completely — mitigate

the surface form competition problem which PMI scoring was designed to tackle (Holtzman et al., 2021). The effect of PMI scoring is increased, however, when we use fuzzy verbalizers and kNN retrieval together (+14.1% for the full model versus +10.4% for kNN with fuzzy verbalizers), suggesting that the kNN distribution might suffer from more surface form competition problems than the base LM distribution.

**kNN retrieval hyperparameters** Figure 4 shows how the number of retrieved nearest neighbors ( $k$ ) and softmax temperature affect model performance on three datasets. In most cases, performance monotonically improves with the number of neighbors when  $k$  is smaller than 512 and deteriorates after that. When  $k < 256$ , a temperature of 1 performs best, while flattening the distribution is useful when retrieving more neighbors. Overall, using 512 neighbors and a temperature value of 3 performs consistently well across the tasks we test.

**Retrieval model size and inference model size** Figure 5 shows how performance varies with the

size of the retriever and inference models on three tasks. We observe substantial gains as the size of the retriever increases, which hold regardless of inference model size.

It should be noted that a larger retriever leads to a larger datastore and slower retrieval: increasing the retriever size from 125M to 1600M parameters doubles the memory footprint of the datastore, which scales with the size of the retriever model’s output embeddings. These computational tradeoffs may inform the retriever size best suited for a particular application.

## 7 Related Work

**Retrieval-augmented LMs** Several studies propose the use of retrieval mechanisms with external datastores to improve language modeling performance (Khandelwal et al., 2020) and open-domain question answering (Izacard and Grave, 2020; Lewis et al., 2020). Other work incorporates retrieval directly into the LM pretraining process (Guu et al., 2020; Borgeaud et al., 2021). Khandelwal et al. (2021) applies nearest neighbor retrieval to conditional sequence generation to improve the quality of machine translation systems. Our work is the first to show that retrieval augmentation, introduced at test time, improves the zero-shot inference of language models on a variety of end tasks.

**Zero-shot and few-shot inference** Brown et al. (2020b) demonstrate that large LMs can perform zero-shot (given only a prompt) and few-shot learning (using a concatenation of training examples as demonstrations) without any finetuning. Subsequent work further improves their zero-shot and few-shot abilities with calibration (Holtzman et al., 2021; Zhao et al., 2021; Min et al., 2021a), prompt engineering (Lu et al., 2021; Shin et al., 2020) and meta-tuning (Min et al., 2021b; Wei et al., 2022; Zhong et al., 2021). Rubin et al. (2021) and Liu et al. (2021) apply retrieval methods to select in-context learning examples that are semantically similar to a test example for few-shot inference. However, these retrieval methods require access to a large set of labeled data. In contrast, kNN-Prompt only assumes the availability of a heterogeneous unlabeled corpus.

## 8 Conclusions and Future Work

We present kNN-Prompt, a new technique to augment LMs with nearest neighbor retrieval for zero-

shot inference on end tasks. kNN-Prompt substantially improves zero-shot performance on a wide range of multiple-choice and classification tasks. With a domain- or task-relevant datastore, kNN-Prompt enables efficient domain adaptation with no additional training, and its benefits scale with the size of the retrieval model.

Future work may study the usefulness of kNN-Prompt with larger inference models, which, combined with larger retrieval models, may result in better zero-shot performance. Careful analysis could explore datastore curation methods to balance task-relevancy, domain generality, and size; datastores could also be compressed for efficient retrieval. Finally, future work may explore the possibility of retrieving and interpolating contexts at different levels of granularity, from tokens and spans to documents.

## 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*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020a. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher 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*. *CoRR*, abs/2005.14165.
- Ido Dagan, Bill Dolan, Bernardo Magnini, and Dan Roth. 2010. *Recognizing textual entailment: Rational, evaluation and approaches – erratum*. *Natural Language Engineering*, 16(1):105–105.
- Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. 2019. *The commitmentbank: Investigating projection in naturally occurring discourse*. *Proceedings of Sinn und Bedeutung*, 23(2):107–124.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,



- and Noah A. Smith. 2020. [Don't stop pretraining: Adapt language models to domains and tasks](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International Conference on Machine Learning*, pages 3929–3938. PMLR.
- Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick. 2021. Efficient nearest neighbor language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5703–5714.
- Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*, pages 507–517.
- Ari Holtzman, Peter West, Vered Shwartz, Yejin Choi, and Luke Zettlemoyer. 2021. [Surface form competition: Why the highest probability answer isn't always right](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7038–7051, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Minqing Hu and Bing Liu. 2004. [Mining and summarizing customer reviews](#). In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '04, page 168–177, New York, NY, USA. Association for Computing Machinery.
- Gautier Izacard and Edouard Grave. 2020. [Leveraging passage retrieval with generative models for open domain question answering](#).
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2021. [Nearest neighbor machine translation](#). In *International Conference on Learning Representations*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. [Generalization through memorization: Nearest neighbor language models](#). In *International Conference on Learning Representations*.
- Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, David Corney, Benno Stein, and Martin Potthast. 2019. [SemEval-2019 task 4: Hyperpartisan news detection](#). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 829–839, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#).
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2021a. Noisy channel language model prompting for few-shot text classification. *arXiv preprint*.
- Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2021b. Metaicl: Learning to learn in context. *arXiv preprint arXiv:2110.15943*.
- Bo Pang and Lillian Lee. 2005. [Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales](#). In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)*, pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. [Language models as knowledge bases?](#) In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2021. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*.

- Timo Schick and Hinrich Schütze. 2021. [Exploiting cloze-questions for few-shot text classification and natural language inference](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, Online. Association for Computational Linguistics.
- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. [AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4222–4235, Online. Association for Computational Linguistics.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-first AAAI conference on artificial intelligence*.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. [Finetuned language models are zero-shot learners](#). In *International Conference on Learning Representations*.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.
- Tony Z Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models.
- Ruiqi Zhong, Kristy Lee, Zheng Zhang, and Dan Klein. 2021. [Adapting language models for zero-shot learning by meta-tuning on dataset and prompt collections](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2856–2878, Punta Cana, Dominican Republic. Association for Computational Linguistics.