

How Can We Know What Language Models Know?

Zhengbao Jiang^{1*} Frank F. Xu^{1*} Jun Araki² Graham Neubig¹

¹Language Technologies Institute, Carnegie Mellon University

²Bosch Research North America

{zhengbaj, fangzhex, gneubig}@cs.cmu.edu jun.araki@us.bosch.com

Abstract

Recent work has presented intriguing results examining the knowledge contained in language models (LMs) by having the LM fill in the blanks of prompts such as “*Obama is a _ by profession*”. These prompts are usually manually created, and quite possibly sub-optimal; another prompt such as “*Obama worked as a _*” may result in more accurately predicting the correct profession. **Because of this, given an inappropriate prompt, we might fail to retrieve facts that the LM *does* know, and thus any given prompt only provides a lower bound estimate of the knowledge contained in an LM.** In this paper, we **attempt to more accurately estimate the knowledge contained in LMs by automatically discovering better prompts to use in this querying process.** Specifically, we propose mining-based and paraphrasing-based methods to automatically generate high-quality and diverse prompts, as well as ensemble methods to combine answers from different prompts. Extensive experiments on the LAMA benchmark for extracting relational knowledge from LMs demonstrate that our methods can improve accuracy from 31.1% to 39.6%, providing a tighter lower bound on what LMs know. We have released the code and the resulting LM Prompt And Query Archive (LPAQA) at <https://github.com/jzbjyb/LPAQA>.

1 Introduction

Recent years have seen the primary role of language models (LMs) transition from generating or evaluating the fluency of natural text (Mikolov and Zweig, 2012; Merity et al., 2018; Melis et al., 2018; Gamon et al., 2005) to being a powerful tool for text understanding. This understanding has mainly been achieved through the use of language modeling as a pre-training task for *feature extractors*, where the hidden vectors learned through a language modeling objective are then used in

down-stream language understanding systems (Dai and Le, 2015; Melamud et al., 2016; Peters et al., 2018; Devlin et al., 2019).

Interestingly, it is also becoming apparent that LMs¹ *themselves* can be used as a tool for text understanding by formulating queries in natural language and either generating textual answers directly (McCann et al., 2018; Radford et al., 2019), or assessing multiple choices and picking the most likely one (Zweig and Burges, 2011; Rajani et al., 2019). For example, LMs have been used to answer factoid questions (Radford et al., 2019), answer common sense queries (Trinh and Le, 2018; Sap et al., 2019), or extract factual knowledge about relations between entities (Petroni et al., 2019; Baldini Soares et al., 2019). Regardless of the end task, the knowledge contained in LMs is probed by providing a prompt, and letting the LM either generate the continuation of a prefix (e.g., “*Barack Obama was born in _*”), or predict missing words in a cloze-style template (e.g., “*Barack Obama is a _ by profession*”).

However, while this paradigm has been used to achieve a number of intriguing results regarding the knowledge expressed by LMs, they usually rely on prompts that were manually created based on the intuition of the experimenter. These manually created prompts (e.g., “*Barack Obama was born in _*”) might be sub-optimal because LMs might have learned target knowledge from substantially different contexts (e.g., “*The birth place of Barack Obama is Honolulu, Hawaii.*”) during their training. Thus it is quite possible that a fact that the LM *does* know cannot be retrieved due to the prompts not being effective queries for the fact. Thus, existing results are simply a *lower bound* on the extent of knowledge contained

¹Some models we use in this paper, e.g., BERT (Devlin et al., 2019), are bi-directional, and do not directly define probability distribution over text, which is the underlying definition of an LM. Nonetheless, we call them LMs for simplicity.

* The first two authors contributed equally.

| Prompts | | | | | |
|---|------------------|--|----------------------|-------|-------------------|
| | manual | DirectX is developed by y_{man} | | | |
| | mined | y_{mine} | released the DirectX | | |
| | paraphrased | DirectX is created by y_{para} | | | |
| Top 5 predictions and log probabilities | | | | | |
| | y_{man} | | y_{mine} | | y_{para} |
| 1 | Intel | -1.06 | Microsoft | -1.77 | Microsoft -2.23 |
| 2 | Microsoft | -2.21 | They | -2.43 | Intel -2.30 |
| 3 | IBM | -2.76 | It | -2.80 | default -2.96 |
| 4 | Google | -3.40 | Sega | -3.01 | Apple -3.44 |
| 5 | Nokia | -3.58 | Sony | -3.19 | Google -3.45 |

Figure 1: Top-5 predictions and their log probabilities using different prompts (manual, mined, and paraphrased) to query BERT. Correct answer is underlined.

in LMs, and in fact, LMs may be even more knowledgeable than these initial results indicate. In this paper we ask the question: “How can we tighten this lower bound and get a more accurate estimate of the knowledge contained in state-of-the-art LMs?” This is interesting both scientifically, as a probe of the knowledge that LMs contain, and from an engineering perspective, as it will result in higher recall when using LMs as part of a knowledge extraction system.

In particular, we focus on the setting of Petroni et al. (2019) who examine extracting knowledge regarding the relations between entities (definitions in § 2). We propose two automatic methods to systematically improve the breadth and quality of the prompts used to query the existence of a relation (§ 3). Specifically, as shown in Figure 1, these are *mining-based* methods inspired by previous relation extraction methods (Ravichandran and Hovy, 2002), and *paraphrasing-based* methods that take a seed prompt (either manually created or automatically mined), and paraphrase it into several other semantically similar expressions. Further, because different prompts may work better when querying for different subject-object pairs, we also investigate lightweight ensemble methods to combine the answers from different prompts together (§ 4).

We experiment on the LAMA benchmark (Petroni et al., 2019), which is an English-language benchmark devised to test the ability of LMs to retrieve relations between entities (§ 5). We first demonstrate that improved prompts significantly improve accuracy on this task, with the one-best prompt extracted by our method raising accuracy from 31.1% to 34.1% on BERT-base (Devlin et al., 2019), with similar gains being obtained with

BERT-large as well. We further demonstrate that using a diversity of prompts through ensembling further improves accuracy to 39.6%. We perform extensive analysis and ablations, glean insights both about how to best query the knowledge stored in LMs and about potential directions for incorporating knowledge into LMs themselves. Finally, we have released the resulting LM Prompt And Query Archive (LPAQA) to facilitate future experiments on probing knowledge contained in LMs.

2 Knowledge Retrieval from LMs

Retrieving factual knowledge from LMs is quite different from querying standard declarative knowledge bases (KBs). In standard KBs, users formulate their information needs as a structured query defined by the KB schema and query language. For example, `SELECT ?y WHERE {wd:Q76 wdt:P19 ?y}` is a SPARQL query to search the birth place of `Barack_Obama`. In contrast, LMs must be queried by natural language prompts, such as “*Barack Obama was born in _*”, and the word assigned the highest probability in the blank will be returned as the answer. Unlike deterministic queries on KBs, this provides no guarantees of correctness or success.

While the idea of prompts is common to methods for extracting many varieties of knowledge from LMs, in this paper we specifically follow the formulation of Petroni et al. (2019), where factual knowledge is in the form of triples $\langle x, r, y \rangle$. Here x indicates the subject, y indicates the object, and r is their corresponding relation. To query the LM, r is associated with a cloze-style prompt t_r consisting of a sequence of tokens, two of which are placeholders for subjects and objects (e.g., “ *x plays at y position*”). The existence of the fact in the LM is assessed by replacing x with the surface form of the subject, and letting the model predict the missing object (e.g., “*LeBron James plays at _ position*”):²

$$\hat{y} = \arg \max_{y' \in \mathcal{V}} P_{\text{LM}}(y' | x, t_r),$$

²We can also go the other way around by filling in the objects and predicting the missing subjects. Since our focus is on improving prompts, we choose to be consistent with Petroni et al. (2019) to make a fair comparison, and leave exploring other settings to future work. Also notably, Petroni et al. (2019) only use objects consisting of a single token, so we only need to predict one word for the missing slot.

where \mathcal{V} is the vocabulary, and $P_{\text{LM}}(y'|x, t_r)$ is the LM probability of predicting y' in the blank conditioned on the other tokens (i.e., the subject and the prompt).³ We say that an LM has knowledge of a fact if \hat{y} is the same as the ground-truth y . **Because we would like our prompts to most effectively elicit any knowledge contained in the LM itself, a “good” prompt should trigger the LM to predict the ground-truth objects as often as possible.**

In previous work (McCann et al., 2018; Radford et al., 2019; Petroni et al., 2019), t_r has been a single manually defined prompt based on the intuition of the experimenter. As noted in the introduction, this method has no guarantee of being optimal, and thus we propose methods that *learn* effective prompts from a small set of training data consisting of gold subject-object pairs for each relation.

3 Prompt Generation

First, we tackle prompt generation: the task of generating a set of prompts $\{t_{r,i}\}_{i=1}^T$ for each relation r , where at least some of the prompts effectively trigger LMs to predict ground-truth objects. We employ two practical methods to either mine prompt candidates from a large corpus (§ 3.1) or diversify a seed prompt through paraphrasing (§ 3.2).

3.1 Mining-based Generation

Our first method is inspired by template-based relation extraction methods (Agichtein and Gravano, 2000; Ravichandran and Hovy, 2002), which are based on the observation that words in the vicinity of the subject x and object y in a large corpus often describe the relation r . Based on this intuition, we first identify all the Wikipedia sentences that contain both subjects and objects of a specific relation r using the assumption of distant supervision, then propose two methods to extract prompts.

Middle-word Prompts Following the observation that words in the middle of the subject and object are often indicative of the relation, we

³We restrict to masked LMs in this paper because the missing slot might not be the last token in the sentence and computing this probability in traditional left-to-right LMs using Bayes’ theorem is not tractable.

directly use those words as prompts. For example, “*Barack Obama was born in Hawaii*” is converted into a prompt “ *x was born in y* ” by replacing the subject and the object with placeholders.

Dependency-based Prompts Toutanova et al. (2015) note that in cases of templates where words do not appear in the middle (e.g., “*The capital of France is Paris*”), templates based on syntactic analysis of the sentence can be more effective for relation extraction. We follow this insight in our second strategy for prompt creation, which parses sentences with a dependency parser to identify the shortest dependency path between the subject and object, then uses the phrase spanning from the leftmost word to the rightmost word in the dependency path as a prompt. For instance, the dependency path in the above example is “*France \xleftarrow{pobj} of \xleftarrow{prep} capital \xleftarrow{nsbj} is \xrightarrow{attr} Paris*”, where the leftmost and rightmost words are “*capital*” and “*Paris*”, giving a prompt of “*capital of x is y* ”.

Notably, these mining-based methods do not rely on any manually created prompts, and can thus be flexibly applied to any relation where we can obtain a set of subject-object pairs. This will result in diverse prompts, covering a wide variety of ways that the relation may be expressed in text. However, it may also be prone to noise, as many prompts acquired in this way may not be very indicative of the relation (e.g., “ *x, y* ”), even if they are frequent.

3.2 Paraphrasing-based Generation

Our second method for generating prompts is more targeted—it aims to improve lexical diversity while remaining relatively faithful to the original prompt. Specifically, we do so by performing paraphrasing over the original prompt into other semantically similar or identical expressions. For example, if our original prompt is “ *x shares a border with y* ”, it may be paraphrased into “ *x has a common border with y* ” and “ *x adjoins y* ”. This is conceptually similar to query expansion techniques used in information retrieval that reformulate a given query to improve retrieval performance (Carpineto and Romano, 2012).

Although many methods could be used for paraphrasing (Romano et al., 2006; Bhagat and Ravichandran, 2008), we follow the simple

method of using **back-translation** (Sennrich et al., 2016; Mallinson et al., 2017) to **first translate the initial prompt into B candidates in another language, each of which is then back-translated into B candidates in the original language**. We then rank B^2 candidates based on their round-trip probability (i.e., $P_{\text{forward}}(\bar{t}|\hat{t}) \cdot P_{\text{backward}}(t|\bar{t})$, where \hat{t} is the initial prompt, \bar{t} is the translated prompt in the other language, and t is the final prompt), and keep the top T prompts.

4 Prompt Selection and Ensembling

In the previous section, we described methods to generate a set of candidate prompts $\{t_{r,i}\}_{i=1}^T$ for a particular relation r . Each of these prompts may be more or less effective at eliciting knowledge from the LM, and thus it is necessary to decide how to use these generated prompts at test time. In this section, we describe three methods to do so.

4.1 Top-1 Prompt Selection

For each prompt, we can measure its accuracy of predicting the ground-truth objects (on a training dataset) using:

$$A(t_{r,i}) = \frac{\sum_{(x,y) \in \mathcal{R}} \delta(y = \arg \max_{y'} P_{\text{LM}}(y'|x, t_{r,i}))}{|\mathcal{R}|},$$

where \mathcal{R} is a set of subject-object pairs with relation r , and $\delta(\cdot)$ is Kronecker’s delta function, returning 1 if the internal condition is true and 0 otherwise. In the simplest method for querying the LM, we choose the prompt with the highest accuracy and query using only this prompt.

4.2 Rank-based Ensemble

Next we examine methods that use not only the top-1 prompt, but combine together multiple prompts. The advantage to this is that the LM may have observed different entity pairs in different contexts within its training data, and having a variety of prompts may allow for elicitation of knowledge that appeared in these different contexts.

Our first method for ensembling is a parameter-free method that averages the predictions of the top-ranked prompts. We rank all the prompts based on their accuracy of predicting the objects on the training set, and use the average log

probabilities⁴ from the top K prompts to calculate the probability of the object:

$$s(y|x, r) = \sum_{i=1}^K \frac{1}{K} \log P_{\text{LM}}(y|x, t_{r,i}), \quad (1)$$

$$P(y|x, r) = \text{softmax}(s(\cdot|x, r))_y, \quad (2)$$

where $t_{r,i}$ is the prompt ranked at the i -th position. Here, K is a hyper-parameter, where a small K focuses on the few most accurate prompts, and a large K increases diversity of the prompts.

4.3 Optimized Ensemble

The above method treats the top K prompts equally, which is sub-optimal given some prompts are more reliable than others. Thus, we also propose a method that directly optimizes prompt weights. Formally, we re-define the score in Equation 1 as:

$$s(y|x, r) = \sum_{i=1}^T P_{\theta_r}(t_{r,i}|r) \log P_{\text{LM}}(y|x, t_{r,i}), \quad (3)$$

where $P_{\theta_r}(t_{r,i}|r) = \text{softmax}(\theta_r)$ is a distribution over prompts parameterized by θ_r , a T -sized real-value vector. For every relation, we learn to score a different set of T candidate prompts, so the total number of parameters is T times the number of relations. The parameter θ_r is optimized to maximize the probability of the gold-standard objects $P(y|x, r)$ over training data.

5 Main Experiments

5.1 Experimental Settings

In this section, we assess the extent to which our prompts can improve fact prediction performance, raising the lower bound on the knowledge we discern is contained in LMs.

Dataset As data, we use the T-REx subset (ElSahar et al., 2018) of the LAMA benchmark (Petroni et al., 2019), which has a broader set of 41 relations (compared with the Google-RE subset, which only covers 3). Each relation is associated with at most 1000 subject-object pairs from Wikidata, and a single manually designed

⁴Intuitively, because we are combining together scores in the log space, this has the effect of penalizing objects that are very unlikely given any certain prompt in the collection. We also compare with linear combination in ablations in § 5.3.

prompt. To learn to mine prompts (§ 3.1), rank prompts (§ 4.2), or learn ensemble weights (§ 4.3), we create a separate training set of subject-object pairs also from Wikidata for each relation that has no overlap with the T-REx dataset. We denote the training set as T-REx-train. For consistency with the T-REx dataset in LAMA, T-REx-train also is chosen to contain only single-token objects. To investigate the generality of our method, we also report the performance of our methods on the Google-RE subset,⁵ which takes a similar form to T-REx but is relatively small and only covers three relations.

Pörner et al. (2019) note that some facts in LAMA can be recalled solely based on surface forms of entities, without memorizing facts. They filter out those easy-to-guess facts and create a more difficult benchmark, denoted as LAMA-UHN. We also conduct experiments on the T-REx subset of LAMA-UHN (i.e., T-REx-UHN) to investigate whether our methods can still obtain improvements on this harder benchmark. Dataset statistics are summarized in Table 1.

Models As for the models to probe, in our main experiments we use the standard BERT-base and BERT-large models (Devlin et al., 2019). We also perform some experiments with other pre-trained models enhanced with external entity representations, namely, ERNIE (Zhang et al., 2019) and KnowBert (Peters et al., 2019), which we believe may do better on recall of entities.

Evaluation Metrics We use two metrics to evaluate the success of prompts in probing LMs. The first evaluation metric, *micro-averaged accuracy*, follows the LAMA benchmark⁶ in calculating the accuracy of all subject-object pairs for relation r :

$$\frac{1}{|\mathcal{R}|} \sum_{\langle x, y \rangle \in \mathcal{R}} \delta(\hat{y} = y),$$

where \hat{y} is the prediction and y is the ground truth. Then we average across all relations. However, we found that the object distributions

⁵<https://code.google.com/archive/p/relation-extraction-corpus/>.

⁶In LAMA, it is called ‘‘P@1.’’ There might be multiple correct answers for some cases, e.g., a person speaking multiple languages, but we only use one ground truth. We will leave exploring more advanced evaluation methods to future work.

| Properties | T-REx | T-REx-UHN | T-REx-train |
|-----------------|-------|-----------|-------------|
| #sub-obj pairs | 830.2 | 661.1 | 948.7 |
| #unique subject | 767.8 | 600.8 | 880.1 |
| #unique objects | 150.9 | 120.5 | 354.6 |
| object entropy | 3.6 | 3.4 | 4.4 |

Table 1: Dataset statistics. All the values are averaged across 41 relations.

of some relations are extremely skewed (e.g., more than half of the objects in relation `native_language` are French). This can lead to deceptively high scores, even for a majority-class baseline that picks the most common object for each relation, which achieves a score of 22.0%. To mitigate this problem, we also report *macro-averaged accuracy*, which computes accuracy for each unique object separately, then averages them together to get the relation-level accuracy:

$$\frac{1}{|\text{uni_obj}(\mathcal{R})|} \sum_{y' \in \text{uni_obj}(\mathcal{R})} \frac{\sum_{\langle x, y \rangle \in \mathcal{R}, y = y'} \delta(\hat{y} = y)}{|\{y | \langle x, y \rangle \in \mathcal{R}, y = y'\}|},$$

where `uni_obj(\mathcal{R})` returns a set of *unique* objects from relation r . This is a much stricter metric, with the majority-class baseline only achieving a score of 2.2%.

Methods We attempted different methods for prompt generation and selection/ensembling, and compare them with the manually designed prompts used in Petroni et al. (2019). **Majority** refers to predicting the majority object for each relation, as mentioned above. **Man** is the baseline from Petroni et al. (2019) that only uses the manually designed prompts for retrieval. **Mine** (§ 3.1) uses the prompts mined from Wikipedia through both middle words and dependency paths, and **Mine+Man** combines them with the manual prompts. **Mine+Para** (§ 3.2) paraphrases the highest-ranked mined prompt for each relation, while **Man+Para** uses the manual one instead.

The prompts are combined either by averaging the log probabilities from the **TopK** highest-ranked prompts (§ 4.2) or the weights after optimization (§ 4.3; **Opti.**). **Oracle** represents the upper bound of the performance of the generated prompts, where a fact is judged as correct if *any* one of the prompts allows the LM to successfully predict the object.

Implementation Details We use $T = 40$ most frequent prompts either generated through mining

or paraphrasing in all experiments, and the number of candidates in back-translation is set to $B = 7$. We remove prompts only containing stopwords/punctuations or longer than 10 words to reduce noise. We use the round-trip English-German neural machine translation models pre-trained on WMT’19 (Ng et al., 2019) for back-translation, as English-German is one of the most highly resourced language pairs.⁷ When optimizing ensemble parameters, we use Adam (Kingma and Ba, 2015) with default parameters and batch size of 32.

5.2 Evaluation Results

Micro- and macro-averaged accuracy of different methods are reported in Tables 2 and 3, respectively.

Single Prompt Experiments When only one prompt is used (in the first **Top1** column in both tables), the best of the proposed prompt generation methods increases micro-averaged accuracy from 31.1% to 34.1% on BERT-base, and from 32.3% to 39.4% on BERT-large. **This demonstrates that the manually created prompts are a somewhat weak lower bound; there are other prompts that further improve the ability to query knowledge from LMs.** Table 4 shows some of the mined prompts that resulted in a large performance gain compared with the manual ones. For the relation `religion`, “*x who converted to y*” improved 60.0% over the manually defined prompt of “*x is affiliated with the y religion*”, and for the relation `subclass_of`, “*x is a type of y*” raised the accuracy by 22.7% over “*x is a subclass of y*”. It can be seen that the largest gains from using mined prompts seem to occur in cases where the manually defined prompt is more complicated syntactically (e.g., the former), or when it uses less common wording (e.g., the latter) than the mined prompt.

Prompt Ensembling Next we turn to experiments that use multiple prompts to query the LM. Comparing the single-prompt results in column 1 to the ensembled results in the following three columns, we can see that ensembling multiple prompts almost always leads to better performance. The simple average used in **Top3** and

⁷<https://github.com/pytorch/fairseq/tree/master/examples/wmt19>.

| Prompts | Top1 | Top3 | Top5 | Opti. | Oracle |
|------------------------------|------|------|------|-------------|--------|
| <i>BERT-base (Man=31.1)</i> | | | | | |
| Mine | 31.4 | 34.2 | 34.7 | 38.9 | 50.7 |
| Mine+Man | 31.6 | 35.9 | 35.1 | 39.6 | 52.6 |
| Mine+Para | 32.7 | 34.0 | 34.5 | 36.2 | 48.1 |
| Man+Para | 34.1 | 35.8 | 36.6 | 37.3 | 47.9 |
| <i>BERT-large (Man=32.3)</i> | | | | | |
| Mine | 37.0 | 37.0 | 36.4 | 43.7 | 54.4 |
| Mine+Man | 39.4 | 40.6 | 38.4 | 43.9 | 56.1 |
| Mine+Para | 37.8 | 38.6 | 38.6 | 40.1 | 51.8 |
| Man+Para | 35.9 | 37.3 | 38.0 | 38.8 | 50.0 |

Table 2: Micro-averaged accuracy of different methods (%). **Majority** gives us 22.0%. Italic indicates best single-prompt accuracy, and bold indicates the best non-oracle accuracy overall.

| Prompts | Top1 | Top3 | Top5 | Opti. | Oracle |
|------------------------------|------|------|------|-------------|--------|
| <i>BERT-base (Man=22.8)</i> | | | | | |
| Mine | 20.7 | 22.7 | 23.9 | 25.7 | 36.2 |
| Mine+Man | 21.3 | 23.8 | 24.8 | 26.6 | 38.0 |
| Mine+Para | 21.2 | 22.4 | 23.0 | 23.6 | 34.1 |
| Man+Para | 22.8 | 23.8 | 24.6 | 25.0 | 34.9 |
| <i>BERT-large (Man=25.7)</i> | | | | | |
| Mine | 26.4 | 26.3 | 25.9 | 30.1 | 40.7 |
| Mine+Man | 28.1 | 28.3 | 27.3 | 30.7 | 42.2 |
| Mine+Para | 26.2 | 27.1 | 27.0 | 27.1 | 38.3 |
| Man+Para | 25.9 | 27.8 | 28.3 | 28.0 | 39.3 |

Table 3: Macro-averaged accuracy of different methods (%). **Majority** gives us 2.2%. Italic indicates best single-prompt accuracy, and bold indicates the best non-oracle accuracy overall.

Top5 outperforms **Top1** across different prompt generation methods. The optimized ensemble further raises micro-averaged accuracy to 38.9% and 43.7% on BERT-base and BERT-large respectively, outperforming the rank-based ensemble by a large margin. These two sets of results demonstrate that diverse prompts can indeed query the LM in different ways, and that the optimization-based method is able to find weights that effectively combine different prompts together.

We list the learned weights of top-3 mined prompts and accuracy gain over only using the top-1 prompt in Table 5. Weights tend to concentrate on one particular prompt, and the other prompts serve as complements. We also depict the performance of the rank-based ensemble method

| ID | Relations | Manual Prompts | Mined Prompts | Acc. Gain |
|------|-----------------------|---|-----------------------------|-----------|
| P140 | religion | x is affiliated with the y religion | x who converted to y | +60.0 |
| P159 | headquarters location | The headquarter of x is in y | x is based in y | +4.9 |
| P20 | place of death | x died in y | x died at his home in y | +4.6 |
| P264 | record label | x is represented by music label y | x recorded for y | +17.2 |
| P279 | subclass of | x is a subclass of y | x is a type of y | +22.7 |
| P39 | position held | x has the position of y | x is elected y | +7.9 |

Table 4: Micro-averaged accuracy gain (%) of the mined prompts over the manual prompts.

| ID | Relations | Prompts and Weights | Acc. Gain |
|------|--------------|---|-----------|
| P127 | owned by | x is owned by y .485 x was acquired by y .151 x division of y .151 | +7.0 |
| P140 | religion | x who converted to y .615 y tirthankara x .190 y dedicated to x .110 | +12.2 |
| P176 | manufacturer | y introduced the x .594 y announced the x .286 x attributed to the y .111 | +7.0 |

Table 5: Weights of top-3 mined prompts, and the micro-averaged accuracy gain (%) over using the top-1 prompt.

with respect to the number of prompts in Figure 2. For mined prompts, top-2 or top-3 usually gives us the best results, while for paraphrased prompts, top-5 is the best. Incorporating more prompts does not always improve accuracy, a finding consistent with the rapidly decreasing weights learned by the optimization-based method. The gap between **Oracle** and **Opti.** indicates that there is still space for improvement using better ensemble methods.

Mining vs. Paraphrasing For the rank-based ensembles (**Top1**, **3**, **5**), prompts generated by paraphrasing usually perform better than mined prompts, while for the optimization-based ensemble (**Opti.**), mined prompts perform better. We conjecture this is because mined prompts exhibit more variation compared to paraphrases, and proper weighting is of central importance. This difference in the variation can be observed in the average edit distance between the prompts of each class, which is 3.27 and 2.73 for mined and paraphrased prompts respectively. However, the improvement led by ensembling paraphrases is still significant over just using one prompt (**Top1** vs. **Opti.**), raising micro-averaged accuracy from 32.7% to 36.2% on BERT-base, and from 37.8% to 40.1% on BERT-large. This indicates that even small modifications to prompts can result in relatively large changes in predictions. Table 6 demonstrates cases where modification of one word (either function or content word) leads to significant accuracy

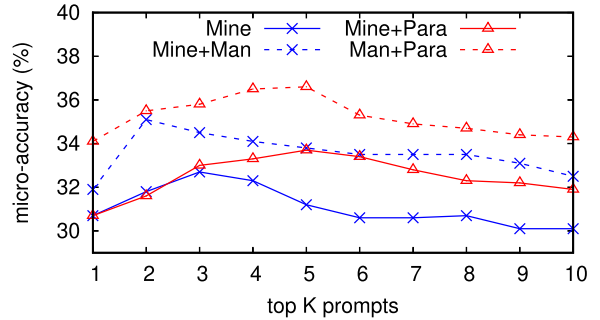


Figure 2: Performance for different top- K ensembles.

| ID | Modifications | Acc. Gain |
|------|------------------------------|-----------|
| P413 | x plays in→at y position | +23.2 |
| P495 | x was created→made in y | +10.8 |
| P495 | x was→is created in y | +10.0 |
| P361 | x is a part of y | +2.7 |
| P413 | x plays in→at y position | +2.2 |

Table 6: Small modifications (**update**, **insert**, and **delete**) in paraphrase lead to large accuracy gain (%).

improvements, indicating that large-scale LMs are still brittle to small changes in the ways they are queried.

Middle-word vs. Dependency-based We compare the performance of only using middle-word prompts and concatenating them with dependency-based prompts in Table 7. The

| Prompts | Top1 | Top3 | Top5 | Opti. | Oracle |
|---------|------|------|------|-------|--------|
| Mid | 30.7 | 32.7 | 31.2 | 36.9 | 45.1 |
| Mid+Dep | 31.4 | 34.2 | 34.7 | 38.9 | 50.7 |

Table 7: Ablation study of middle-word and dependency-based prompts on BERT-base.

| Model | Man | Mine | Mine +Man | Mine +Para | Man +Para |
|----------|------|------|--------------|---------------|--------------|
| BERT | 31.1 | 38.9 | 39.6 | 36.2 | 37.3 |
| ERNIE | 32.1 | 42.3 | 43.8 | 40.1 | 41.1 |
| KnowBert | 26.2 | 34.1 | 34.6 | 31.9 | 32.1 |

Table 8: Micro-averaged accuracy (%) of various LMs

improvements confirm our intuition that words belonging to the dependency path but not in the middle of the subject and object are also indicative of the relation.

Micro vs. Macro Comparing Tables 2 and 3, we can see that macro-averaged accuracy is much lower than micro-averaged accuracy, indicating that macro-averaged accuracy is a more challenging metric that evaluates how many unique objects LMs know. Our optimization-based method improves macro-averaged accuracy from 22.8% to 25.7% on BERT-base, and from 25.7% to 30.1% on BERT-base. This again confirms the effectiveness of ensembling multiple prompts, but the gains are somewhat smaller. Notably, in our optimization-based methods, the ensemble weights are optimized on each example in the training set, which is more conducive to optimizing micro-averaged accuracy. Optimization to improve macro-averaged accuracy is potentially an interesting direction for future work that may result in prompts more generally applicable to different types of objects.

Performance of Different LMs In Table 8, we compare BERT with ERNIE and KnowBert, which are enhanced with external knowledge by explicitly incorporating entity embeddings. ERNIE outperforms BERT by 1 point even with the manually defined prompts, but our prompt generation methods further emphasize the difference between the two methods, with the highest accuracy numbers differing by 4.2 points using the **Mine+Man** method. This

| Model | Man | Mine | Mine +Man | Mine +Para | Man +Para |
|------------|------|------|--------------|---------------|--------------|
| BERT-base | 21.3 | 28.7 | 29.4 | 26.8 | 27.0 |
| BERT-large | 24.2 | 34.5 | 34.5 | 31.6 | 29.8 |

Table 9: Micro-averaged accuracy (%) on LAMA-UHN.

| Model | Man | Mine | Mine +Man | Mine +Para | Man +Para |
|------------|------|------|--------------|---------------|--------------|
| BERT-base | 9.8 | 10.0 | 10.4 | 9.6 | 10.0 |
| BERT-large | 10.5 | 10.6 | 11.3 | 10.4 | 10.7 |

Table 10: Micro-averaged accuracy (%) on Google-RE.

indicates that if LMs are queried effectively, the differences between highly performant models may become more clear. KnowBert underperforms BERT on LAMA, which is opposite to the observation made in Peters et al. (2019). This is probably because that multi token subjects/objects are used to evaluate KnowBert in Peters et al. (2019), while LAMA contains only single-token objects.

LAMA-UHN Evaluation The performances on LAMA-UHN benchmark are reported in Table 9. Although the overall performances drop dramatically compared to the performances on the original LAMA benchmark (Table 2), optimized ensembles can still outperform manual prompts by a large margin, indicating that our methods are effective in retrieving knowledge that cannot be inferred based on surface forms.

5.3 Analysis

Next, we perform further analysis to better understand what type of prompts proved most suitable for facilitating retrieval of knowledge from LMs.

Prediction Consistency by Prompt We first analyze the conditions under which prompts will yield different predictions. We define the divergence between predictions of two prompts $t_{r,i}$ and $t_{r,j}$ using the following equation:

$$\text{Div}(t_{r,i}, t_{r,j}) = \frac{\sum_{(x,y) \in \mathcal{R}} \delta(C(x,y,t_{r,i}) \neq C(x,y,t_{r,j}))}{|\mathcal{R}|},$$

where $C(x,y,t_{r,i}) = 1$ if prompt $t_{r,i}$ can successfully predict y and 0 otherwise, and $\delta(\cdot)$ is

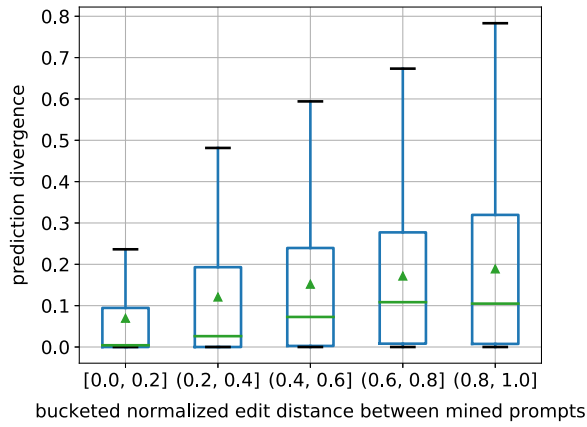


Figure 3: Correlation of edit distance between prompts and their prediction divergence.

| x/y V y/x | x/y V P y/x | x/y V W* P y/x |
|-------------------------------------|-----------------|--------------------|
| V = verb particle? adv? | | |
| W = (noun adj adv pron det) | | |
| P = (prep particle inf. marker) | | |

Table 11: Three part-of-speech-based regular expressions used in ReVerb to identify relational phrases.

Kronecker’s delta. For each relation, we normalize the edit distance of two prompts into $[0, 1]$ and bucket the normalized distance into five bins with intervals of 0.2. We plot a box chart for each bin to visualize the distribution of prediction divergence in Figure 3, with the green triangles representing mean values and the green bars in the box representing median values. As the edit distance becomes larger, the divergence increases, which confirms our intuition that very different prompts tend to cause different prediction results. The Pearson correlation coefficient is 0.25, which shows that there is a weak correlation between these two quantities.

Performance on Google-RE We also report the performance of optimized ensemble on the Google-RE subset in Table 10. Again, ensembling diverse prompts improves accuracies for both the BERT-base and BERT-large models. The gains are somewhat smaller than those on the T-REx subset, which might be caused by the fact that there are only three relations and one of them (predicting the `birth_date` of a person) is particularly hard to the extent that only one prompt yields non-zero accuracy.

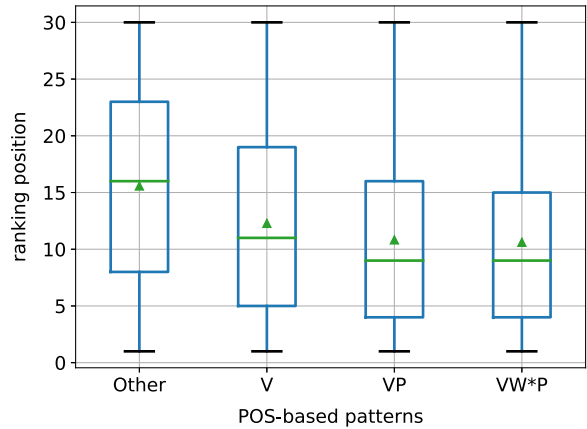


Figure 4: Ranking position distribution of prompts with different patterns. Lower is better.

POS-based Analysis Next, we try to examine which types of prompts tend to be effective in the abstract by examining the part-of-speech (POS) patterns of prompts that successfully extract knowledge from LMs. In open information extraction systems (Banko et al., 2007), manually defined patterns are often leveraged to filter out noisy relational phrases. For example, ReVerb (Fader et al., 2011) incorporates three syntactic constraints listed in Table 11 to improve the coherence and informativeness of the mined relational phrases. To test whether these patterns are also indicative of the ability of a prompt to retrieve knowledge from LMs, we use these three patterns to group prompts generated by our methods into four clusters, where the “other” cluster contains prompts that do not match any pattern. We then calculate the rank of each prompt within the extracted prompts, and plot the distribution of rank using box plots in Figure 4.⁸ We can see that the average rank of prompts matching these patterns is better than those in the “other” group, confirming our intuitions that good prompts should conform with those patterns. Some of the best performing prompts’ POS signatures are “ x VBD VBN IN y ” (e.g., “ x was born in y ”) and “ x VBZ DT NN IN y ” (e.g., “ x is the capital of y ”).

Cross-model Consistency Finally, it is of interest to know whether the prompts that we are extracting are highly tailored to a

⁸We use the ranking position of a prompt to represent its quality instead of its accuracy because accuracy distributions of different relations might span different ranges, making accuracy not directly comparable across relations.

| Test Train | BERT-base | | BERT-large | |
|------------|-----------|-------|------------|------|
| | base | large | large | base |
| Mine | 38.9 | 38.7 | 43.7 | 42.2 |
| Mine+Man | 39.6 | 40.1 | 43.9 | 42.2 |
| Mine+Para | 36.2 | 35.6 | 40.1 | 39.0 |
| Man+Para | 37.3 | 35.6 | 38.8 | 37.5 |

Table 12: Cross-model micro-averaged accuracy (%). The first row is the model to test, and the second row is the model on which prompt weights are learned.

specific model, or whether they can generalize across models. To do so, we use two settings: One compares BERT-base and BERT-large, the same model architecture with different sizes; the other compares BERT-base and ERNIE, different model architectures with a comparable size. In each setting, we compare when the optimization-based ensembles are trained on the same model, or when they are trained on one model and tested on the other. As shown in Tables 12 and 13, we found that in general there is usually some drop in performance in the cross-model scenario (third and fifth columns), but the losses tend to be small, and the highest performance when querying BERT-base is actually achieved by the weights optimized on BERT-large. Notably, the best accuracies of 40.1% and 42.2% (Table 12) and 39.5% and 40.5% (Table 13) with the weights optimized on the other model are still much higher than those obtained by the manual prompts, indicating that optimized prompts still afford large gains across models. Another interesting observation is that the drop in performance on ERNIE (last two columns in Table 13) is larger than that on BERT-large (last two columns in Table 12) using weights optimized on BERT-base, indicating that models sharing the same architecture benefit more from the same prompts.

Linear vs. Log-linear Combination As mentioned in § 4.2, we use log-linear combination of probabilities in our main experiments. However, it is also possible to calculate probabilities through regular linear interpolation:

$$P(y|x, r) = \sum_{i=1}^K \frac{1}{K} P_{\text{LM}}(y|x, t_{r,i}) \quad (4)$$

| Test Train | BERT | | ERNIE | |
|------------|------|-------|-------|------|
| | BERT | ERNIE | ERNIE | BERT |
| Mine | 38.9 | 38.0 | 42.3 | 38.7 |
| Mine+Man | 39.6 | 39.5 | 43.8 | 40.5 |
| Mine+Para | 36.2 | 34.2 | 40.1 | 39.0 |
| Man+Para | 37.3 | 35.2 | 41.1 | 40.3 |

Table 13: Cross-model micro-averaged accuracy (%). The first row is the model to test, and the second row is the model on which prompt weights are learned.

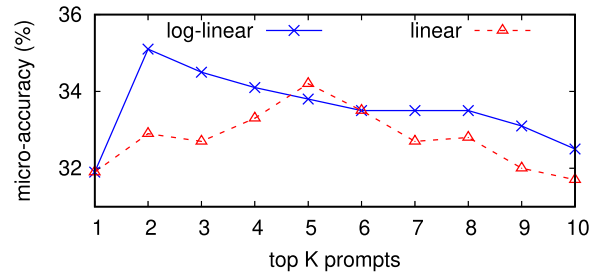


Figure 5: Performance of two interpolation methods.

We compare these two ways to combine predictions from multiple mined prompts in Figure 5 (§ 4.2). We assume that log-linear combination outperforms linear combination because log probabilities make it possible to penalize objects that are very unlikely given any certain prompt.

6 Omitted Design Elements

Finally, in addition to the elements of our main proposed methodology in § 3 and § 4, we experimented with a few additional methods that did not prove highly effective, and thus were omitted from our final design. We briefly describe these below, along with cursory experimental results.

6.1 LM-aware Prompt Generation

We examined methods to generate prompts by solving an optimization problem that maximizes the probability of producing the ground-truth objects with respect to the prompts:

$$t_r^* = \arg \max_{t_r} P_{\text{LM}}(y|x, t_r),$$

where $P_{\text{LM}}(y|x, t_r)$ is parameterized with a pre-trained LM. In other words, this method directly searches for a prompt that causes the LM to assign ground-truth objects the highest probability.

| Prompts | Top1 | Top3 | Top5 | Opti. | Oracle |
|---------|------|------|------|-------|--------|
| before | 31.9 | 34.5 | 33.8 | 38.1 | 47.9 |
| after | 30.2 | 32.5 | 34.7 | 37.5 | 50.8 |

Table 14: Micro-averaged accuracy (%) before and after LM-aware prompt fine-tuning.

Solving this problem of finding text sequences that optimize some continuous objective has been studied both in the context of end-to-end sequence generation (Hoang et al., 2017), and in the context of making small changes to an existing input for adversarial attacks (Ebrahimi et al., 2018; Wallace et al., 2019). However, we found that directly optimizing prompts guided by gradients was unstable and often yielded prompts in unnatural English in our preliminary experiments. Thus, we instead resorted to a more straightforward hill-climbing method that starts with an initial prompt, then masks out one token at a time and replaces it with the most probable token conditioned on the other tokens, inspired by the mask-predict decoding algorithm used in non-autoregressive machine translation (Ghazvininejad et al., 2019):⁹

$$P_{\text{LM}}(w_i | t_r \setminus i) = \frac{\sum_{(x,y) \in \mathcal{R}} P_{\text{LM}}(w_i | x, t_r \setminus i, y)}{|\mathcal{R}|},$$

where w_i is the i -th token in the prompt and $t_r \setminus i$ is the prompt with the i -th token masked out. We followed a simple rule that modifies a prompt from left to right, and this is repeated until convergence.

We used this method to refine all the mined and manual prompts on the T-REx-train dataset, and display their performance on the T-REx dataset in Table 14. After fine-tuning, the oracle performance increased significantly, while the ensemble performances (both rank-based and optimization-based) dropped slightly. This indicates that LM-aware fine-tuning has the potential to discover better prompts, but some portion of the refined prompts may have over-fit to the training set upon which they were optimized.

⁹In theory, this algorithm can be applied to both masked LMs like BERT and traditional left-to-right LMs, since the masked probability can be computed using Bayes’ theorem for traditional LMs. However, in practice, due to the large size of vocabulary, it can only be approximated with beam search, or computed with more complicated continuous optimization algorithms (Hoang et al., 2017).

| Features | Mine | | Paraphrase | |
|-----------|-------|-------|------------|-------|
| | macro | micro | macro | micro |
| forward | 38.1 | 25.2 | 37.3 | 25.0 |
| +backward | 38.2 | 25.5 | 37.4 | 25.2 |

Table 15: Performance (%) of using forward and backward features with BERT-base.

6.2 Forward and Backward Probabilities

Finally, given class imbalance and the propensity of the model to over-predict the majority object, we examine a method to encourage the model to predict subject-object pairs that are more strongly aligned. Inspired by the maximum mutual information objective used in Li et al. (2016a), we add the backward log probability $\log P_{\text{LM}}(x | y, t_{r,i})$ of each prompt to our optimization-based scoring function in Equation 3. Due to the large search space for objects, we turn to an approximation approach that only computes backward probability for the most probable B objects given by the forward probability at both training and test time. As shown in Table 15, the improvement resulting from backward probability is small, indicating that a diversity-promoting scoring function might not be necessary for knowledge retrieval from LMs.

7 Related Work

Much work has focused on understanding the internal representations in neural NLP models (Belinkov and Glass, 2019), either by using extrinsic probing tasks to examine whether certain linguistic properties can be predicted from those representations (Shi et al., 2016; Linzen et al., 2016; Belinkov et al., 2017), or by ablations to the models to investigate how behavior varies (Li et al., 2016b; Smith et al., 2017). For contextualized representations in particular, a broad suite of NLP tasks are used to analyze both syntactic and semantic properties, providing evidence that contextualized representations encode linguistic knowledge in different layers (Hewitt and Manning, 2019; Tenney et al., 2019a; Tenney et al., 2019b; Jawahar et al., 2019; Goldberg, 2019).

Different from analyses probing the representations themselves, our work follows Petroni et al. (2019); Pörner et al. (2019) in probing for factual

knowledge. They use manually defined prompts, which may be under-estimating the true performance obtainable by LMs. Concurrently to this work, Bouraoui et al. (2020) made a similar observation that using different prompts can help better extract relational knowledge from LMs, but they use models explicitly trained for relation extraction whereas our methods examine the knowledge included in LMs without any additional training.

Orthogonally, some previous works integrate external knowledge bases so that the language generation process is explicitly conditioned on symbolic knowledge (Ahn et al., 2016; Yang et al., 2017; Logan et al., 2019; Hayashi et al., 2020). Similar extensions have been applied to pre-trained LMs like BERT, where contextualized representations are enhanced with entity embeddings (Zhang et al., 2019; Peters et al., 2019; Pörner et al., 2019). In contrast, we focus on better knowledge retrieval through prompts from LMs as-is, without modifying them.

8 Conclusion

In this paper, we examined the importance of the prompts used in retrieving factual knowledge from language models. We **propose mining-based and paraphrasing-based methods to systematically generate diverse prompts to query specific pieces of relational knowledge**. Those prompts, when combined together, improve factual knowledge retrieval accuracy by 8%, outperforming manually designed prompts by a large margin. Our analysis indicates that LMs are indeed more knowledgeable than initially indicated by previous results, but they are also quite sensitive to how we query them. This indicates potential future directions such as (1) more robust LMs that can be queried in different ways but still return similar results, (2) methods to incorporate factual knowledge in LMs, and (3) further improvements in optimizing methods to query LMs for knowledge. Finally, we have released all our learned prompts to the community as the LM Prompt and Query Archive (LPAQA), available at: <https://github.com/jzbjyb/LPAQA>.

Acknowledgments

This work was supported by a gift from Bosch Research and NSF award no. 1815287. We would like to thank Paul Michel, Hiroaki Hayashi,

Pengcheng Yin, and Shuyan Zhou for their insightful comments and suggestions.

References

- Eugene Agichtein and Luis Gravano. 2000. *Snowball*: Extracting relations from large plain-text collections. In *Proceedings of the Fifth ACM Conference on Digital Libraries, June 2-7, 2000, San Antonio, TX, USA*, pages 85–94. ACM.
- Sungjin Ahn, Heeyoul Choi, Tanel Pärnamaa, and Yoshua Bengio. 2016. A neural knowledge language model. *CoRR*, abs/1608.00318v2.
- Livio Baldini Soares, Nicholas FitzGerald, Jeffrey Ling, and Tom Kwiatkowski. 2019. Matching the blanks: Distributional similarity for relation learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2895–2905, Florence, Italy. Association for Computational Linguistics.
- Michele Banko, Michael J. Cafarella, Stephen Soderland, Matthew Broadhead, and Oren Etzioni. 2007. Open information extraction from the web. In *IJCAI 2007, Proceedings of the 20th International Joint Conference on Artificial Intelligence, Hyderabad, India, January 6-12, 2007*, pages 2670–2676.
- Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. What do neural machine translation models learn about morphology? In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 861–872, Vancouver, Canada. Association for Computational Linguistics.
- Yonatan Belinkov and James R. Glass. 2019. Analysis methods in neural language processing: A survey. *Transactions of the Association for Computational Linguistics*, 7:49–72.
- Rahul Bhagat and Deepak Ravichandran. 2008. Large scale acquisition of paraphrases for learning surface patterns. In *Proceedings*

- of *ACL-08: HLT*, pages 674–682, Columbus, Ohio. Association for Computational Linguistics.
- Zied Bouraoui, Jose Camacho-Collados, and Steven Schockaert. 2020. Inducing relational knowledge from BERT. In *Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*, New York, USA.
- Claudio Carpineto and Giovanni Romano. 2012. A survey of automatic query expansion in information retrieval. *ACM, Computing Surveys*, 44(1):1:1–1:50.
- Andrew M. Dai and Quoc V. Le. 2015. Semi-supervised sequence learning. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 3079–3087.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. HotFlip: White-box adversarial examples for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 31–36, Melbourne, Australia, Association for Computational Linguistics.
- Hady ElSahar, Pavlos Vougiouklis, Arslan Remaci, Christophe Gravier, Jonathon S. Hare, Frédérique Laforest, and Elena Simperl. 2018. T-REx: A large scale alignment of natural language with knowledge base triples. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation, LREC 2018, Miyazaki, Japan, May 7-12, 2018*.
- Anthony Fader, Stephen Soderland, and Oren Etzioni. 2011. Identifying relations for open information extraction. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, EMNLP 2011, 27-31 July 2011, John McIntyre Conference Centre, Edinburgh, UK, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1535–1545.
- Michael Gamon, Anthony Aue, and Martine Smets. 2005. Sentence-level MT evaluation without reference translations: Beyond language modeling. In *Proceedings of EAMT*, pages 103–111.
- Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. 2019. Mask-predict: Parallel decoding of conditional masked language models. 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 6114–6123, Hong Kong, China. Association for Computational Linguistics.
- Yoav Goldberg. 2019. Assessing BERT’s syntactic abilities. *CoRR*, abs/1901.05287v1.
- Hiroaki Hayashi, Zecong Hu, Chenyan Xiong, and Graham Neubig. 2020. Latent relation language models. In *Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*, New York, USA.
- John Hewitt and Christopher D. Manning. 2019. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4129–4138.
- Cong Duy Vu Hoang, Gholamreza Haffari, and Trevor Cohn. 2017. Towards decoding as continuous optimisation in neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 146–156, Copenhagen, Denmark. Association for Computational Linguistics.
- Robert L. Logan IV, Nelson F. Liu, Matthew E. Peters, Matt Gardner, and Sameer Singh.

2019. Barack’s wife Hillary: Using knowledge graphs for fact-aware language modeling. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28-August 2, 2019, Volume 1: Long Papers*, pages 5962–5971.
- Ganesh Jawahar, Benoît Sagot, and Djamel Seddah. 2019. What does BERT learn about the structure of language? In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28-August 2, 2019, Volume 1: Long Papers*, pages 3651–3657.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016a. A diversity-promoting objective function for neural conversation models. In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 110–119.
- Jiwei Li, Will Monroe, and Dan Jurafsky. 2016b. Understanding neural networks through representation erasure. *CoRR*, abs/1612.08220v3.
- Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the ability of LSTMs to learn syntax-sensitive dependencies. *Transactions of the Association for Computational Linguistics*, 4:521–535.
- Jonathan Mallinson, Rico Sennrich, and Mirella Lapata. 2017. Paraphrasing revisited with neural machine translation. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 881–893, Valencia, Spain. Association for Computational Linguistics.
- Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. The natural language decathlon: Multitask learning as question answering. *CoRR*, abs/1806.08730v1.
- Oren Melamud, Jacob Goldberger, and Ido Dagan. 2016. context2vec: Learning generic context embedding with bidirectional LSTM. In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany, August 11-12, 2016*, pages 51–61.
- Gábor Melis, Chris Dyer, and Phil Blunsom. 2018. On the state of the art of evaluation in neural language models. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2018. Regularizing and optimizing LSTM language models. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*.
- Tomas Mikolov and Geoffrey Zweig. 2012. Context dependent recurrent neural network language model. In *2012 IEEE Spoken Language Technology Workshop (SLT)*, pages 234–239. IEEE.
- Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, and Sergey Edunov. 2019. Facebook FAIR’s WMT19 news translation task submission. In *Proceedings of the Fourth Conference on Machine Translation, WMT 2019, Florence, Italy, August 1-2, 2019 - Volume 2: Shared Task Papers, Day 1*, pages 314–319.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 2227–2237.
- Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. 2019. Knowledge enhanced contextual word representations. 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 43–54, Hong Kong, China. 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.
- Nina Pörner, Ulli Waltinger, and Hinrich Schütze. 2019. BERT is not a knowledge base (yet): Factual knowledge vs. Name-based reasoning in unsupervised QA. *CoRR*, abs/1911.03681v1.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8).
- Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! Leveraging language models for commonsense reasoning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4932–4942, Florence, Italy. Association for Computational Linguistics.
- Deepak Ravichandran and Eduard Hovy. 2002. Learning surface text patterns for a question answering system. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 41–47. Association for Computational Linguistics.
- Lorenza Romano, Milen Kouylekov, Idan Szpektor, Ido Dagan, and Alberto Lavelli. 2006. Investigating a generic paraphrase-based approach for relation extraction. In *11th Conference of the European Chapter of the Association for Computational Linguistics*, Trento, Italy. Association for Computational Linguistics.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*.
- Xing Shi, Inkit Padhi, and Kevin Knight. 2016. Does string-based neural MT learn source syntax? In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1526–1534, Austin, Texas. Association for Computational Linguistics.
- Noah A. Smith, Chris Dyer, Miguel Ballesteros, Graham Neubig, Lingpeng Kong, and Adhiguna Kuncoro. 2017. What do recurrent neural network grammars learn about syntax? In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2017, Valencia, Spain, April 3-7, 2017, Volume 1: Long Papers*, pages 1249–1258.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019a. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 4593–4601.
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. 2019b. What do you learn from context? Probing for sentence structure in contextualized word representations. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.
- Kristina Toutanova, Danqi Chen, Patrick Pantel, Hoifung Poon, Pallavi Choudhury, and Michael

- Gamon. 2015. Representing text for joint embedding of text and knowledge bases. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 1499–1509.
- Trieu H. Trinh and Quoc V. Le. 2018. A simple method for commonsense reasoning. *CoRR*, abs/1806.02847v2.
- Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. 2019. Universal adversarial triggers for attacking and analyzing NLP. 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 2153–2162, Hong Kong, China. Association for Computational Linguistics.
- Zichao Yang, Phil Blunsom, Chris Dyer, and Wang Ling. 2017. Reference-aware language models. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 1850–1859.
- Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced language representation with informative entities. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28-August 2, 2019, Volume 1: Long Papers*, pages 1441–1451.
- Geoffrey Zweig and Christopher J. C. Burges. 2011. The Microsoft Research sentence completion challenge. *Microsoft Research, Redmond, WA, USA, Technical Report MSR-TR-2011-129*.