

Exploring and Exploiting Latent Commonsense Knowledge in Pretrained Masked Language Models

Anonymous EMNLP submission

Abstract

Pretrained masked language models (PLMs) were shown to be inheriting a considerable amount of relational knowledge from the source corpora. In this paper, we present an in-depth analysis concerning eliciting relational commonsense knowledge already present in PLMs from the perspective of network pruning. We show that it is possible to find sub-networks capable of representing grounded commonsense relations at non-trivial sparsity meanwhile being generalizable to downstream commonsense knowledge base completion and commonsense reasoning tasks.

1 Introduction

The past few years have witnessed the revolution of NLP methods with the advent of pretrained language models (PLMs) such as BERT (Devlin et al., 2019) and ROBERTA (Liu et al., 2019a). These bi-directional Transformer (Vaswani et al., 2017) encoders are first pretrained on vast amount of unlabeled text corpora and then fine-tuned on task-specific data, offering a surge of improvements on a wealth of downstream NLP tasks. Although this transfer learning paradigm has become the de-facto standard, we know very little about *what* and *how much* knowledge embedded in PLMs actually contributes to the success. Notable endeavors toward this understanding focus on probing linguistic knowledge therein. They demonstrated that pretraining did impart useful linguistic abstraction about syntax and semantics into PLMs (Peters et al., 2018; Goldberg, 2019; Tenney et al., 2019).

More recently, several works are presenting intriguing results examining the relational knowledge within PLMs. Relational knowledge (Speer and Havasi, 2012; Vrandečić and Krötzsch, 2014) is typically defined as describing the abstract relationship between a pair of concepts or entities, which is crucial for facilitating language understanding. Petroni et al. (2020) first posed the LAMA probe,

an English benchmark comprising multiple sets of prompts. Each prompt is a cloze-like sentence transformed from a relational knowledge triple:

Knowledge Triple: $\langle bus, HasA, ? \rangle$

Object Label: seats.

Sentence: you are likely to find _ in a bus.

By substituting _ with a special [MASK] token and reusing the masked language modeling (MLM) head, prompt-based relational knowledge probing provides an estimated lower bound of what PLMs know without training additional layer as in the previous linguistic probe. They showed that, even without grounded supervision, PLMs capture such relational knowledge at a level competitive to supervised alternatives. Subsequent works further showed that some specific prompts, acquired either through heuristical mining (Jiang et al., 2020) or gradient-guided search (Shin et al., 2020), can better trigger the models to correctly predict the missing object.

Despite the mounting evidence for the existence of relational knowledge in PLMs, it remains unclear how such knowledge are represented internally. It hinders the utility of PLMs being extended to more structured tasks, such as knowledge base completion. In light of this, we raise the core question in this paper: *Can we disentangle PLMs into relation-specific knowledge models and take advantage of them in a more flexible way?*

We approach this question by first drawing inspiration from recent findings (Saunshi et al., 2020; Lee et al., 2020; Zhang and Hashimoto, 2021): *the more cloze-like MLM pretraining simulates downstream task, the more successful the transfer will be.* For example, filling in *like* or *hate* into a cloze like *I [MASK] this film, it's great.* provide a clear way in which the model can implicitly learn to perform sentiment classification. Similarly, we hypothesis that MLM pretraining on clozes expressing certain relation r between masked word and remaining context would lead to good performance on

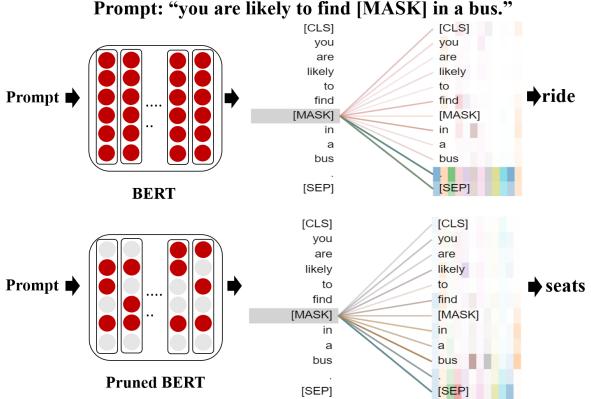


Figure 1: Querying original/pruned BERT-BASE with prompts of relation *HasA*. The color spectrum indicates the 12 attention heads in the last layer (Vig, 2019).

knowledge probing that targets relation r . Instead of designing cloze instances resembling specific downstream tasks, we exploit such correlation conversely: *the more successful the transfer is, the more cloze-like MLM pretraining simulates downstream tasks*. Specifically, we propose **weakly supervised weights pruning**, an end-to-end differentiable procedure to search for different subnetworks within PLMs that target zero-shot knowledge probing of different relations.

We show in experiment that it is possible to find subnetworks capable of representing grounded commonsense relations at non-trivial sparsity. Figure 1 exemplifies a cloze prompt where the identified subnetwork produces the valid answer *seats* by attending to relevant context, i.e., bus, while the full-scale BERT fails. We further investigate the possibility of repurposing these subnetworks on various downstream tasks. Experimental results on commonsense knowledge base completion show that the identified subnetworks perform on par with or even better than strong supervised knowledge base completion methods. We also explore the usage of the subnetworks on multiple commonsense reasoning tasks and empirically find that, when combined properly, the subnetworks can outperform the original PLMs in both many-shot and zero-shot settings.

In summary, our **contributions** include: (i) We present a novel way of eliciting relational knowledge hidden in PLMs from the perspective of network pruning (Section 2.3). (ii) Grounding on a concept-centric relation schema, we show that the proposed pruning procedure successfully identified sparse subnetworks specializing in miscella-

neous commonsense knowledge remarkably better than their full-scale counterparts (Section 3.1). (iii) We showcase the effectiveness of these subnetworks on commonsense knowledge base completion tasks (Section 3.2) as well as a heuristic application on multiple commonsense reasoning tasks, gleaning insight on the transformation from language representation to knowledge representation.

We release code and all versions of our pruned PLMs at <https://anonymous.4open.science>.

2 Methodology

We first provide background on pretrained masked language models and the formulation of cloze prompt for querying these models, then we proceed to elaborate our proposed pruning procedure.

2.1 Pretrained Masked Language Models

Formally, given a sequence of tokens $w = [w_1, w_2, \dots, w_n]$, where n is its length, the model outputs a sequence of fixed-size hidden representations $h = [h_1, h_2, \dots, h_n]$ for each token. In standard MLM pretraining, the corresponding representation h_i is fed into a designated MLM head for computing the reconstruction probability $P(w_i | w_{-i})$ of the masked i -th token w_i , where w_{-i} are the remaining unmasked tokens. We denote the original pretrained model \mathcal{LM} with unpruned parameter θ as \mathcal{LM}_θ in following sections.

2.2 Knowledge Probing with Cloze Prompts

The natural language cloze prompts, such as “*you are likely to find a basement in below your [MASK]*”, offer a straightforward mean of querying pretrained masked language models that conform to their interfaces.

We follow the formulation of Petroni et al. (2020), where relational knowledge is in the form of triplets $\langle subj, r, obj \rangle$. Here $subj$ refers to the subject, obj refers to the object, and r indicates their corresponding relation. To query a model \mathcal{LM}_θ , each relation r is associated with a set of cloze template prompts T_r , each of which consists of a sequence of tokens, two of which are place-holders for $subj$ and obj (e.g., “*you are likely to find [subj] in [obj]*”). We can check the existence of the knowledge in \mathcal{LM}_θ by substituting the $[subj]$ place-holder with the surface form of real subject and asking the model to predict the missing

164 object:

$$165 \hat{obj} = \arg \max_{w \in \mathcal{V}} P_{\mathcal{LM}_\theta}([obj] = w | subj, T_r)$$

166 where \mathcal{V} is the vocabulary of \mathcal{LM}_θ . We say that
167 \mathcal{LM}_θ grasps the knowledge if $\hat{obj} = obj$.

168 2.3 Weakly Supervised Weights Pruning

169 Given a pretrained masked language model \mathcal{LM}
170 and the associated set of pretrained parameters
171 $\theta \in R^d$, where d is the dimensionality, we are
172 interested in finding the subnetwork \mathcal{LM}_{θ_r} that
173 is maximally predictive of prompts of relation r .
174 The intuition is that if a subnetwork specializes ex-
175clusively on relation r , the parameters it reserves
176 should inherit the corresponding knowledge from
177 MLM pretraining on cloze instances of r .

178 Similar to Zhao et al. (2020), for each weight
179 matrix W^l from the set of all weight matrices \mathbf{W}^l
180 in the l -th transformer layer, we assign learnable
181 pruning mask generator G_r^l that is element-wise
182 initialized from a prior distribution $\phi(\cdot)$. Each
183 entry $g_{i,j}^l \in G_r^l$ is a real-valued scalar that deter-
184mines whether its corresponding weight $w_{i,j}^l \in W^l$
185 should be pruned. To investigate if $w_{i,j}^l$ should be
186 softly scaled or entirely removed to effectively re-
187 cover \mathcal{LM}_{θ_r} , we explore two different schemes of
188 converting G_r^l into a masking matrix M_r^l from a
189 probabilistic view.

190 2.3.1 Stochastic Pruning

191 The first variant is to establish a probabilistic for-
192 mulation for determining the importance of indi-
193 vidual weights. Formally, $g_{i,j}^l$ is taken as input to a
194 sigmoid function for parametrizing a Bernoulli dis-
195 tribution $B(\sigma(g_{i,j}^l))$, from which a binary masking
196 random variable $m_{i,j}^l$ is sampled:

$$197 m_{i,j}^l \sim B(\sigma(g_{i,j}^l)) \quad (1)$$

198 where $m_{i,j}^l \in M_r^l$. The resulting masking matrix
199 M_r^l can then be used to select weights within origi-
200 nal linear layer W^l by Hadamard product:

$$201 W_r^l = W^l \odot M_r^l \quad (2)$$

202 Due to the non-differentiability introduced by sam-
203 pling, the gradient w.r.t. loss function (described
204 in Section 2.3.3) cannot be back-propagated to $g_{i,j}^l$.
205 As a remedy, we use the re-parametrization tech-
206 nique (Li et al., 2018) to approximate $m_{i,j}^l$ with
207 another differentiable variable $\tilde{m}_{i,j}^l$:

$$208 \tilde{m}_{i,j}^l = \sigma\left(\frac{g_{i,j}^l + \log U - \log(1-U)}{\tau}\right) \quad (3)$$

209 where $U \sim Uniform(0, 1)$ and τ is a small posi-
210 tive temperature parameter. As τ approaches zero,
211 $\tilde{m}_{i,j}^l$ will match sampled $m_{i,j}^l$ more accurately (de-
212tailed proof can be found in Appendix A).

213 In this way, Eq. (2) becomes:

$$214 W_r^l = W^l \odot \tilde{M}_r^l \quad (4)$$

215 2.3.2 Deterministic Pruning

216 While our first probabilistic pruning formulation
217 considers flexible weights combination, the sec-
218 ond proposed variant utilizes a hard thresholding
219 function to directly generate the masking matrix.

220 Let t denotes the predefined thresholding hyper-
221 parameter ranging from 0 to 1, then we have:

$$222 \hat{m}_{i,j}^l = \begin{cases} 1, & \sigma(g_{i,j}^l) \geq t, \\ 0, & otherwise. \end{cases} \quad (5)$$

223 where σ is the sigmoid function. Similar to Section
224 2.3.1, the resulting binary masking matrix \hat{M}_r^l is
225 then used to select weights relevant to relation r by
226 Hadamard product:

$$227 W_r^l = W^l \odot \hat{M}_r^l \quad (6)$$

228 Note that the hard thresholding operation in Eq. (5)
229 also blocks the gradient propagation to $g_{i,j}^l$. Here
230 we employ the Straight-Through gradient estimator
231 (Bengio et al., 2013; Hubara et al., 2016) and
232 use $\frac{\partial \mathcal{L}_r}{\partial \hat{m}_{i,j}^l}$ as a proxy of $\frac{\partial \mathcal{L}_r}{\partial g_{i,j}^l}$. We elaborate on the
233 loss function \mathcal{L}_r w.r.t relation r in the next section.

234 2.3.3 Training and Inference

235 The resultant pruned model (i.e., subnetwork)
236 \mathcal{LM}_{θ_r} is expected to behave like a specialized neu-
237 ral knowledge base. That is, given a prompt requir-
238 ing knowledge about relation r , \mathcal{LM}_{θ_r} should be
239 able to fill in the missing object more accurately
240 than its full-scale counterpart \mathcal{LM}_θ . To this end,
241 the learning objective for pruning mask generator
242 $\{G_r^l\}_{l_b \leq l \leq l_t}$, where l_b and l_t indicate the range of
243 transformer layers, is to find the subnetwork \mathcal{LM}_{θ_r}
244 that minimizes the following objective:

$$245 \mathcal{L}_r = -\mathbb{E}_{(subj, T_r, obj) \sim D_r} [\log P_{\mathcal{LM}_{\theta_r}}(obj | subj, T_r)]$$

246 where D_r is the collection of prompts under re-
247 lation r . The training procedure is conducted for
248 each relation $r \in \mathcal{R}$ of interest and finally, we ac-
249 quire a set of trained $\{G_r\}_{r \in \mathcal{R}}$ for the designated
250 pretrained model \mathcal{LM} .

251 During inference, for deterministic pruning, M_r
252 is obtained from G_r by Eq. (5). For stochastic
253 pruning, M_r is obtained by taking the expectation
254 value (i.e., $\sigma(G_r)$) of Bernoulli variables.

Model	P@1 (%)	P@2 (%)	P@3 (%)	Sparsity	$l_b - l_t$	# Param.
DISTILBERT-BASE w/o pruning	11.4	16.6	19.9	0%	-	66M
DISTILBERT-BASE w/ stochastic pruning	14.8	21.5	26.3	~30%	4-6	66M
DISTILBERT-BASE w/ deterministic pruning	44.1	52.9	57.6	~50%	4-6	56M
BERT-BASE w/o pruning	12.9	18.4	21.8	0%	-	110M
BERT-BASE w/ stochastic pruning	17.2	25.1	29.6	~30%	7-12	110M
BERT-BASE w/ deterministic pruning	57.6	63.8	67.2	~50%	7-12	88M
ROBERTA-BASE w/o pruning	15.4	21.2	24.6	0%	-	125M
ROBERTA-BASE w/ stochastic pruning	16.6	22.2	25.8	~30%	7-12	125M
ROBERTA-BASE w/ deterministic pruning	38.3	42.8	44.6	~50%	7-12	100M
MPNET-BASE w/o pruning	14.8	20.7	24.0	0%	-	110M
MPNET-BASE w/ stochastic pruning	19.8	27.9	33.2	~30%	7-12	110M
MPNET-BASE w/ deterministic pruning	62.7	68.7	71.4	~50%	7-12	88M

Table 1: Relational knowledge probing results on C-LAMA. We show one representative pruning configuration for each type of model and relegate the complete results to Appendix due to space limits.

3 Experiments

We first expound our pruning setting and provide evidences of its ability to identify relation-specific subnetworks in PLM. Then we experiment on several commonsense-intensive scenarios to seek good practices for using these subnetworks.

3.1 Disentangling PLMs into Relation-specific Subnetworks

Corpus	#Facts	#Rels	#Sentences
ConceptNet	11,458	16	29,774

Table 2: Statistics of C-LAMA.

Dataset. We use the ConceptNet (Speer and Havasi, 2012) subset of LAMA benchmark as weak supervision, denoted as C-LAMA. C-LAMA contains facts from the English part of ConceptNet that has single-token objects covering 16 relations. Prompts that will be used for pruning and probing are extracted from Open Mind Common Sense (OMCS). Statistics is listed in Table 2. Since our goal is to explore the existence of specialized subnetworks within general-purpose PLMs, we utilize the whole dataset to fully exploit the high-precision knowledge therein. This naturally deviates from the conventional supervised experimental train/valid/test setting.

Models. For the choices of \mathcal{LM} , we consider the 6-layer DISTILBERT-BASE (Sanh, 2019), 12-layer BERT-BASE, 12-layer ROBERTA-BASE (Liu et al., 2019a). We also include the more recent powerful 12-layer MPNET-BASE (Song et al., 2020) model. All models are implemented with HuggingFace’s transformers (Wolf et al., 2019) library.

Setup. The prior distribution $\phi(\cdot)$ is a Gaussian $\mathcal{N}(\mu, 1)$ where μ is the mean controlling initial sparsity of pruned model (e.g., $\mu = 0$ indicates 50% initial sparsity). We set l_t to be the top layer of a given model and set l_b to be within 3 – 4 for DISTILBERT, 5 – 9 for BERT, ROBERTA, and MPNET. The temperature τ is fixed as 0.1. The threshold t is fixed as 0.5. We use Adam (Kingma and Ba, 2015) with a batch size of 32 and a linear warm-up scheduler with 0.1 warm-up ratio for training the mask up to 6 epochs. The learning rate is fixed as 3×10^{-4} . All experiments are conducted on a GTX 1080 Ti GPU with 11G RAM.

To what extent can we achieve disentanglement? As stated in Section 2.3, for each model, our pruning procedure will identify 16 subnetworks for 16 commonsense relations. We record the P@K (K=1,2,3) scores of each subnetwork \mathcal{LM}_{θ_r} evaluated on its corresponding subset D_r of C-LAMA and shows the macro-averaged results in Table 1.

Among all models without pruning, ROBERTA achieves the highest P@1 score of 15.4 while DISTILBERT gets the lowest 11.4. It indicates that while PLMs are shown to be helpful for downstream learning, they cannot accurately complete cloze-like prompts that require commonsense relation knowledge. This observation also coincide with previous finding (Zhang and Hashimoto, 2021) that the uniform masking adopted by PLMs is biased towards extracting statistical and syntactic dependencies. Comparing the results for each pair of original and subnetworks, we consistently observe a surprisingly significant increase (37.0 on average), especially for deterministically pruned ones. This large performance gap provides unique

new evidence of sparse latent relational knowledge structures in PLMs, which are weakened by pre-trained weights that are *reserved* for more general-purpose use.

We also observe that the deterministic pruning excels by a huge margin across all models, which suggests the importance of removing noisy transformation of input representations to regain the expressivity of specific commonsense knowledge. Another advantage of deterministic pruning in memory footprint is that only sets of 1-bit masks rather than 32-bits float parameters need to be saved for solving multiple tasks. For the above reasons, we focus our analysis on and use *pruned* to denote deterministically pruned PLMs throughout this paper henceforth.

How specialized are these subnetworks? We next investigate alternative ways of building the subnetworks. We attempted to instantiate the pruning mask upon BERT-BASE via either creating a set of mismatched masks or assigning a randomly initialized mask with comparable sparsity. For *mismatched* one, we corrupt the correspondence of relation between masks and prompts by shifting the order of masks 15 times, as there are 16 relations in total. Then we calculate the micro-averaged P@K for each shift and average the results. For *random* one, we initialize the mask element-wise with a Bernoulli $B(0.5)$ and average the results from 5 different random seeds. The results are shown in

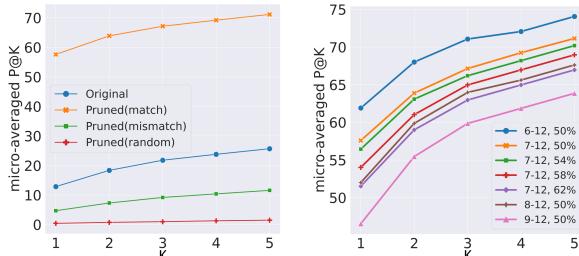


Figure 2: Ablation on the pruning masks (left) and effect of initial sparsity and pruned layers (right).

Figure 2 (left). If we apply the random masks with sparsity comparable to learned ones, the P@1 drops drastically to 0.4. This notable gap proves that the effective subnetworks cannot be trivially identified through random weights sampling. If we apply the mismatched masks from other relations, the P@1 also significantly drops to 4.8, even inferior to the unpruned one. It shows that the latent structure for different types of commonsense knowledge exhibits remarkably distinct geometry.

Factors impacting performance. We also investigate how initial masking sparsity and choice of layers to prune influence the probing performance. We experiment on BERT-BASE with l_b in {6, 7, 8, 9} and initial sparsity in {50%, 54%, 58%, 62%}. Figure 2 (right) shows that (i) increasing the number of pruned layers helps distill more knowledge. (ii) larger initial sparsity is more likely to prune away weights important to certain knowledge and cannot be recovered in the later training process. In general, we find an initial sparsity around 50% yields decent performance both in probing and downstream applications (see Section 3.3 and Section 3.2).

Visualization of attention weights and representations. To explain how the subnetworks accommodate more accurate commonsense knowledge despite having far fewer weights than the full-scale models, we randomly sample several prompts that the subnetworks correctly answered but the full-scale model (BERT-BASE) failed to and visualize the attention patterns in the last layer. Specifically,

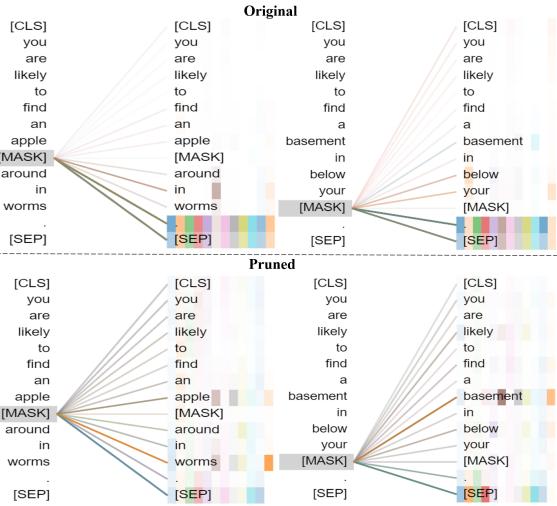


Figure 3: Attention weight visualization. *AtLocation* is required for prompt in the left column and *PartOf* is required for prompt in the right column.

we focus on the attention weights between [MASK] token and other tokens in the prompt. A first glance of change of attention pattern is given in Figure 1 and we show more examples of other ConceptNet relations in Figure 3. We observe that while the original pretrained model tends to attend to special tokens like period and [SEP], the subnetwork successfully grasps the relevant concepts (i.e., apple, worms, and basement) in the prompt hence produces the right object. We also use t-SNE (van der Maaten and Hinton, 2008) to visualize the last

Model	Development Set				Test Set			
	MRR (%)	P@1 (%)	P@2 (%)	P@3 (%)	MRR (%)	P@1 (%)	P@2 (%)	P@3 (%)
Supervised								
DISTMULT (Yang et al., 2015)	8.5	4.2	6.6	8.3	10.5	5.4	8.4	10.9
COMPLEX (Trouillon et al., 2016)	10.7	6.5	9.0	11.0	13.6	8.2	12.4	15.7
CONVE (Dettmers et al., 2018)	18.9	11.5	16.6	19.0	21.9	13.5	18.9	24.0
TUCKER (Balazevic et al., 2019)	17.3	10.9	14.8	18.8	21.6	14.0	20.4	24.0
CONTRANSE (Shang et al., 2019)	19.8	13.2	17.8	21.3	24.0	15.6	21.9	<u>26.5</u>
SACN (Shang et al., 2019)	21.2	13.2	19.8	23.2	24.2	14.4	<u>22.1</u>	28.0
Unsupervised								
DISTILBERT-BASE	9.0	3.1	6.9	10.3	10.8	5.8	9.6	11.2
BERT-BASE	12.4	7.2	10.0	13.7	14.3	8.3	13.7	16.6
ROBERTA-BASE	8.3	4.2	6.0	7.1	9.4	5.1	7.1	9.3
MPNET-BASE	11.7	7.2	9.4	11.1	11.1	6.0	9.9	11.7
DISTILBERT-BASE (pruned)	24.1	15.8	24.1	<u>26.4</u>	<u>23.4</u>	<u>14.8</u>	22.2	<u>26.5</u>
BERT-BASE (pruned)	<u>23.7</u>	<u>15.5</u>	<u>22.1</u>	27.0	22.8	14.3	20.9	26.0
ROBERTA-BASE (pruned)	9.0	4.9	7.1	8.9	9.5	6.1	7.6	11.4
MPNET-BASE (pruned)	22.1	12.9	21.2	25.5	20.0	11.4	18.8	22.9

Table 3: Link prediction results. Best results are marked with **bold** font and second best with underline.

layer’s representation of [CLS] for each prompt. From Figure 4, the representations computed by original pretrained model are hardly separable as different types of knowledge are intermingled. In contrast, the pruned subnetwork can extract meaningful and disentangled representations for different commonsense relations.

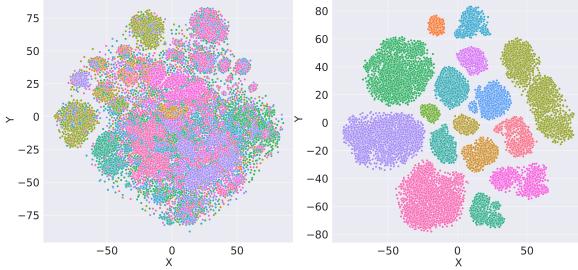


Figure 4: t-SNE visualization of [CLS]’s representation from original (left) and pruned (right) BERT-BASE.

3.2 Commonsense Knowledge Base Completion (CKBC)

We evaluate the utility of identified relation-specific subnetworks on CKBC in an unsupervised manner. Specifically, we use the ConceptNet-100K benchmark provided by Li et al. (2016). To allow a fair evaluation, we manually create a subset of ConceptNet-100K consisting of triples with single-token subject/object, of which the dev/test set have **no overlap** with C-LAMA. Each relation is associated with a sentence template (provided in Appendix) (Kwon et al., 2019) of which the wording is distinct from those in C-LAMA. We acknowledge that these sentence templates might be suboptimal for certain relations, but prompt optimization is out of the scope of this paper. The resulting dataset contains 17,891 training instances, 349 development

instances, and 446 test instances.

Link prediction. We first formulate CKBC as a link prediction task and compare subnetworks (i.e., \mathcal{LM}_{θ_r} is queried to predict missing link for instance of relation r) as well as original PLMs against strong supervised KB completion methods.

Table 3 shows the results. Most of the supervised models outperform full-scale PLMs by a large margin, which suggests the inefficacy of directly using PLMs to perform link prediction. However, the subnetworks identified by our pruning procedure can acquire performance on par with or better than state-of-the-art supervised models. Surprisingly, the pruned DISTILBERT get the highest MRR, outperforming other larger and more advanced PLMs. ROBERTA struggles to predict correct objects, perhaps due to its larger vocabulary size compared to WordPiece (50,265 vs 30,522) and less lexicon overlap (53% vs 59%) with the dataset.

Triple classification We can also formulate CKBC as a triple classification task. Following Davison et al. (2019), we use estimated point-wise mutual information (PMI) computed by pretrained language model as a surrogate of a triple’s validity. An expectation-maximization-based Gaussian mixture clustering method is used and instances in the cluster with higher mean PMI are labeled as valid. In our preliminary experiments, we found that the model pruned by the mask of a single relation might not be robust for PMI estimation and generally performed inferior to the intact model. In the same spirit as model ensembling, we then perform grid search over combinations of multiple knowledge, which is similar to what we did in zero-shot commonsense reasoning. For all four PLMs considered in Table 4, we observe that there exists

393
394
395
396
397
398
399

417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452

Model	F1 Score
DISTILBERT-BASE	74.1
DISTILBERT-BASE (pruned)	76.3
BERT-BASE	73.7
BERT-BASE (pruned)	76.7
ROBERTA-BASE	74.8
ROBERTA-BASE (pruned)	76.9
MPNET-BASE	76.5
MPNET-BASE (pruned)	78.0

Table 4: Triple classification on ConceptNet-100K.

one or multiple knowledge combinations delivering F1 score higher than the original models.

Triple extraction. We then investigate the ability of specialized subnetworks to extract novel commonsense knowledge triples absent from the dataset. We randomly sample 100 triples from the test set of ConceptNet-100K and for each sample use top- K predictions from pruned DISTILBERT-BASE as candidate objects. Three human annotators are asked to first determine the correctness of each candidate object and further determine their novelty (i.e., not present in any of train/validation/test set) if deemed to be correct. The Fleiss Kappa inter-annotator agreement κ is 0.66/0.65 for precision and novelty, respectively. Figure 5 shows the change of precision-novelty

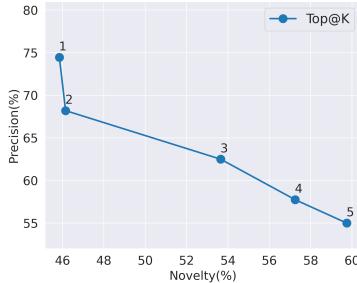


Figure 5: Precision-novelty curve with varied K .

with varied K . We observe a clear trade-off between the validity and novelty of triples extracted by the pruned model. As expected, a large K inevitably makes noisy predictions but is more likely to extract unseen knowledge. For the purpose of knowledge enrichment, one might choose a large K to ensure a desirable recall. We list the obtained novel triples in the Appendix D due to space limits.

3.3 Commonsense Reasoning (CSR)

After identifying sparse subnetworks within PLMs that specialize in different commonsense knowledge, we now evaluate their generalization ability

Task	Original	Pruned	p -value
RTE	69.2 ± 2.3	69.8 ± 2.0	0.12
COPA	62.4 ± 5.0	63.0 ± 4.7	0.33
CommonsenseQA	53.1 ± 0.6	54.1 ± 0.7	0.08
SWAG	73.9 ± 0.3	74.2 ± 0.1	0.09
HellaSWAG	38.9 ± 0.4	39.1 ± 0.5	0.32
aNLI	63.7 ± 0.4	64.0 ± 0.4	0.19
CosmosQA	61.3 ± 1.0	61.8 ± 0.2	0.26

Table 5: Fine-tuning results of BERT-BASE for CSR.

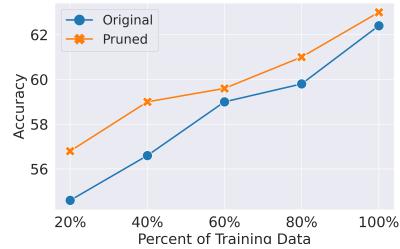


Figure 6: BERT-BASE results on COPA.

in the context of commonsense reasoning.

Many-shot learning. We experiment with BERT-BASE and its deterministically pruned version using supervised fine-tuning on 7 datasets: RTE (Dagan et al., 2009), COPA (Roemmele et al., 2011), CommonsenseQA (Talmor et al., 2019), SWAG (Zellers et al., 2018), HellaSWAG (Zellers et al., 2019), aNLI (Bhagavatula et al., 2020) and CosmosQA (Huang et al., 2019). For each task, we heuristically identify the commonsense knowledge it might requires. If multiple types of knowledge are required, we simply take the union of all masks and apply the resultant mask to the pre-trained model as initialization for finetuning. We repeat training three times with different random seeds for each task. The choice of mask combination for each task can be found in Appendix B.

The results in Table 5 shows that, when initialized with proper weights, the model can be better fine-tuned on downstream commonsense reasoning tasks via more useful *prior* knowledge. We further analyze the change of performance under the low-resource regime on COPA dataset. Figure 6 shows that the pruned BERT exhibits a notable advantage when training data is extremely scarce. As more training data is seen, the benefit of the pruned model becomes less prominent, i.e., $p > 0.05$.

Zero-shot learning. We next assess the ability of specialized subnetworks to perform zero-shot commonsense reasoning, a setting where the knowledge relied on to complete the task is solely determined by the model parameters. Here we focus on: COPA, CommonsenseQA, Conjunction Accept-

Model	COPA (Tra.)	COPA (Dev.)	CSQA	CA	WSC	SM	ARCT1	ARCT2	Average
DISTILBERT-BASE	58.3	60.0	29.6	84.6	53.3	71.6	48.6	50.4	57.0
DISTILBERT-BASE (pruned)	61.5	69.0	31.5	89.6	56.9	72.1	53.4	51.6	60.7
BERT-BASE	60.2	54.0	26.5	89.0	57.3	69.7	46.8	50.3	56.7
BERT-BASE (pruned)	63.0	64.0	28.5	91.8	59.0	71.7	50.0	52.0	60.0
ROBERTA-BASE	60.7	59.0	39.9	90.1	61.8	73.1	48.6	53.1	60.7
ROBERTA-BASE (pruned)	65.3	72.0	40.4	93.4	62.9	74.4	53.2	55.1	64.6
MPNET-BASE	66.5	69.0	40.0	94.5	64.3	75.8	52.9	56.7	64.9
MPNET-BASE (pruned)	71.0	74.0	41.7	97.3	66.4	77.5	56.1	57.7	67.7

Table 6: Zero-shot results of accuracy (%) on commonsense reasoning tasks. Better results of each pair is in **bold**.

ability (CA) (Zhou et al., 2019), Winograd Schema Challenge (WSC) (Levesque et al., 2012), Sense-Making (SM) (Wang et al., 2019), ARCT1 (Habernal et al., 2018) and ARCT2 (Niven and Kao, 2019). Each sample in the above datasets can be formulated as *[CLS] premise [SEP] hypothesis_i [SEP]*, where *i* is the subscript. Hypothesis with the highest plausibility scored by PLMs is the predicted answer.

Since multiple types of knowledge are typically required for effectively reasoning over concepts, for each task, we perform grid search over combinations of 3-4 different commonsense knowledge out of the 16 total types and reported the best accuracy in Table 6. We put the best combination for each model on each task in Appendix B for space constraints. By combining multiple commonsense knowledge useful for the task, we show that the pruned models can actually surpass their full-scale version in all tasks considered in our experiments. The most likely explanation is that knowledge irrelevant to the specific task in the original models hurt the in-domain zero-shot reasoning capability. It also manifests that the most important reasoning skills vary from task to task.

4 Related Work

Since the emergence of large pretrained language models, much work has focused on understanding internal contextual representations produced by such models. Most prior work (Shi et al., 2016; Belinkov et al., 2017) pays special attention to either using extraneous probing tasks to examine whether certain linguistic properties can be identified from those representations or ablating the models to observe how behavior changes. More recently, some studies (Peters et al., 2018; Goldberg, 2019; Tenney et al., 2019) have shown the existence of linguistic knowledge (e.g., syntax) in various but generally lower layers of pretrained transformers.

It is worthy of exploring and shedding more light

on how PLMs memorize abstract knowledge rather than trivial statistical co-occurrence patterns. We extend previous probe (Petroni et al., 2020) on relational knowledge. Specifically, we are concerned with commonsense knowledge that is grounded on ConceptNet relations. Our work differs in that we focus on not only probing but also bringing potentially more implicit commonsense knowledge to the surface and unleashing more potential in knowledge-intensive applications.

Another line of researches relevant to our work is network pruning (Liu et al., 2019b; Lin et al., 2020) and the lottery ticket hypothesis (Frankle and Carbin, 2019; Prasanna et al., 2020; Chen et al., 2020). The former aims at reducing the size of model parameters without compromising test accuracy and the latter reveals subnetworks whose initializations made them capable of being trained effectively comparable to the original model. In contrast, we seek to uncover subnetworks in over-parametrized PLMs that specialize on commonsense knowledge tailored for downstream tasks rather than focusing on good global initialization, and achieve good results.

5 Conclusion

We apply network pruning, a novel approach to explore the latent relational knowledge hidden in PLMs. With a preliminary focus on commonsense knowledge, we find evidence of latent sparse subnetworks capable of representing grounded commonsense relations in a plethora of PLMs. Further experiments on downstream tasks showed that such subnetworks can be effectively utilized as auspicious neural knowledge bases, fine-tuning starting points, and robust zero-shot reasoners. Our work raises a new viewpoint about the inner storage scheme as well as practical utilization of relational knowledge in PLMs, opening up avenues to future work on better understanding and adapting pretrained language representations.

References

- 594
 595 Ivana Balazevic, Carl Allen, and Timothy Hospedales.
 596 2019. **TuckER: Tensor factorization for knowledge**
 597 **graph completion.** In *Proceedings of the 2019 Con-*
 598 *ference on Empirical Methods in Natural Language*
 599 *Processing and the 9th International Joint Conference*
 600 *on Natural Language Processing (EMNLP-IJCNLP)*,
 601 pages 5185–5194, Hong Kong, China. As-
 602 *sociation for Computational Linguistics.*
- 603 Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Has-
 604 san Sajjad, and James Glass. 2017. **What do neu-**
 605 **ral machine translation models learn about morphol-**
 606 **ogy?** In *Proceedings of the 55th Annual Meeting of*
 607 *the Association for Computational Linguistics (Vol-*
 608 *ume 1: Long Papers)*, pages 861–872, Vancouver,
 609 Canada. Association for Computational Linguistics.
- 610 Yoshua Bengio, Nicholas Léonard, and Aaron C.
 611 Courville. 2013. **Estimating or propagating gra-**
 612 **dients through stochastic neurons for conditional com-**
 613 **putation.** *CoRR*, abs/1308.3432.
- 614 Chandra Bhagavatula, Ronan Le Bras, Chaitanya
 615 Malaviya, Keisuke Sakaguchi, Ari Holtzman, Han-
 616 nah Rashkin, Doug Downey, Wen-tau Yih, and Yejin
 617 Choi. 2020. **Abductive commonsense reasoning.** In
 618 *8th International Conference on Learning Repre-*
 619 *sentations, ICLR 2020, Addis Ababa, Ethiopia, April*
 620 *26-30, 2020.* OpenReview.net.
- 621 Tianlong Chen, Jonathan Frankle, Shiyu Chang, Sijia
 622 Liu, Yang Zhang, Zhangyang Wang, and Michael
 623 Carbin. 2020. **The lottery ticket hypothesis for pre-**
 624 **trained BERT networks.** In *Advances in Neural*
 625 *Information Processing Systems 33: Annual Con-*
 626 *ference on Neural Information Processing Systems*
 627 *2020, NeurIPS 2020, December 6-12, 2020, virtual.*
- 628 Ido Dagan, Bill Dolan, Bernardo Magnini, and Dan
 629 Roth. 2009. **Recognizing textual entailment: Ratio-**
 630 **nal, evaluation and approaches.** *Natural Language*
 631 *Engineering*, 15(Special Issue 04):i–xvii.
- 632 Joe Davison, Joshua Feldman, and Alexander Rush.
 633 2019. **Commonsense knowledge mining from pre-**
 634 **trained models.** In *Proceedings of the 2019 Con-*
 635 *ference on Empirical Methods in Natural Language*
 636 *Processing and the 9th International Joint Conference*
 637 *on Natural Language Processing (EMNLP-IJCNLP)*,
 638 pages 1173–1178, Hong Kong, China. As-
 639 *sociation for Computational Linguistics.*
- 640 Tim Dettmers, Pasquale Minervini, Pontus Stenetorp,
 641 and Sebastian Riedel. 2018. **Convolutional 2d**
 642 **knowledge graph embeddings.** In *Proceedings of*
 643 *the Thirty-Second AAAI Conference on Artificial*
 644 *Intelligence, (AAAI-18), the 30th innovative Ap-*
 645 *plications of Artificial Intelligence (IAAI-18), and*
 646 *the 8th AAAI Symposium on Educational Advances*
 647 *in Artificial Intelligence (EAAI-18), New Orleans,*
 648 *Louisiana, USA, February 2-7, 2018*, pages 1811–
 649 1818. AAAI Press.
- 650 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
 651 Kristina Toutanova. 2019. **BERT: Pre-training of**
 652 **deep bidirectional transformers for language under-**
 653 **standing.** In *Proceedings of the 2019 Conference*
 654 *of the North American Chapter of the Association*
 655 *for Computational Linguistics: Human Language*
 656 *Technologies, Volume 1 (Long and Short Papers)*,
 657 pages 4171–4186, Minneapolis, Minnesota. Asso-
 658 *ciation for Computational Linguistics.*
- 659 Jonathan Frankle and Michael Carbin. 2019. **The lot-**
 660 **tery ticket hypothesis: Finding sparse, trainable neu-**
 661 **ral networks.** In *7th International Conference on*
 662 *Learning Representations, ICLR 2019, New Orleans,*
 663 *LA, USA, May 6-9, 2019.* OpenReview.net.
- 664 Yoav Goldberg. 2019. **Assessing bert’s syntactic abili-**
 665 **ties.** *CoRR*, abs/1901.05287.
- 666 Ivan Habernal, Henning Wachsmuth, Iryna Gurevych,
 667 and Benno Stein. 2018. **The argument reasoning**
 668 **comprehension task: Identification and reconstruc-**
 669 **tion of implicit warrants.** In *Proceedings of the 2018*
 670 *Conference of the North American Chapter of the*
 671 *Association for Computational Linguistics: Human*
 672 *Language Technologies, Volume 1 (Long Papers)*,
 673 pages 1930–1940, New Orleans, Louisiana. Asso-
 674 *ciation for Computational Linguistics.*
- 675 Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and
 676 Yejin Choi. 2019. **Cosmos QA: Machine reading**
 677 **comprehension with contextual commonsense rea-**
 678 **soning.** In *Proceedings of the 2019 Conference on*
 679 *Empirical Methods in Natural Language Processing*
 680 *and the 9th International Joint Conference on Natu-*
 681 *ral Language Processing (EMNLP-IJCNLP)*, pages
 682 2391–2401, Hong Kong, China. Association for
 683 *Computational Linguistics.*
- 684 Itay Hubara, Matthieu Courbariaux, Daniel Soudry,
 685 Ran El-Yaniv, and Yoshua Bengio. 2016. **Binarized**
 686 **neural networks.** In *Advances in Neural Information*
 687 *Processing Systems 29: Annual Conference on Neu-*
 688 *ral Information Processing Systems 2016, Decem-*
 689 *ber 5-10, 2016, Barcelona, Spain*, pages 4107–4115.
- 690 Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham
 691 Neubig. 2020. **How can we know what language**
 692 **models know?** *Transactions of the Association for*
 693 *Computational Linguistics*, 8:423–438.
- 694 Diederik P. Kingma and Jimmy Ba. 2015. **Adam: A**
 695 **method for stochastic optimization.** In *3rd Inter-*
 696 *national Conference on Learning Representations,*
 697 *ICLR 2015, San Diego, CA, USA, May 7-9, 2015,*
 698 *Conference Track Proceedings.*
- 699 Sunjae Kwon, Cheongwoong Kang, Jiyeon Han, and
 700 Jaesik Choi. 2019. **Why do masked neural language**
 701 **models still need common sense knowledge?** *arXiv.*
- 702 Jason D. Lee, Qi Lei, Nikunj Saunshi, and Jiacheng
 703 Zhuo. 2020. **Predicting what you already know**
 704 **helps: Provable self-supervised learning.** *CoRR*,
 705 abs/2008.01064.

706	Hector J. Levesque, Ernest Davis, and Leora Morgenstern. 2012. The Winograd Schema Challenge . In <i>Proceedings of the Thirteenth International Conference on Principles of Knowledge Representation and Reasoning</i> , KR’12, pages 552–561. AAAI Press, Rome, Italy.	763
707		764
708		765
709		766
710		767
711		768
712	Xiang Li, Aynaz Taheri, Lifu Tu, and Kevin Gimpel. 2016. Commonsense knowledge base completion . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1445–1455, Berlin, Germany. Association for Computational Linguistics.	769
713		770
714		771
715		772
716		773
717		774
718	Zhuohan Li, Di He, Fei Tian, Wei Chen, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2018. Towards binary-valued gates for robust LSTM training . In <i>Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018</i> , volume 80 of <i>Proceedings of Machine Learning Research</i> , pages 3001–3010. PMLR.	775
719		776
720		777
721		778
722		779
723		780
724		781
725		782
726	Tao Lin, Sebastian U. Stich, Luis Barba, Daniil Dmitriev, and Martin Jaggi. 2020. Dynamic model pruning with feedback . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	783
727		784
728		785
729		786
730		787
731		788
732	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019a. Roberta: A robustly optimized BERT pretraining approach . <i>CoRR</i> , abs/1907.11692.	789
733		790
734		791
735		792
736		793
737	Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. 2019b. Rethinking the value of network pruning . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	794
738		795
739		796
740		797
741		798
742	Timothy Niven and Hung-Yu Kao. 2019. Probing neural network comprehension of natural language arguments . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4658–4664, Florence, Italy. Association for Computational Linguistics.	799
743		800
744		801
745		802
746		803
747		804
748	Matthew Peters, Mark Neumann, Luke Zettlemoyer, and Wen-tau Yih. 2018. Dissecting contextual word embeddings: Architecture and representation . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1499–1509, Brussels, Belgium. Association for Computational Linguistics.	805
749		806
750		807
751		808
752		809
753		810
754		811
755	Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2020. Language models as knowledge bases? <i>EMNLP-IJCNLP 2019 - 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, Proceedings of the Conference</i> , pages 2463–2473.	812
756		813
757		814
758		815
759		816
760		817
761		818
762		819
763	Sai Prasanna, Anna Rogers, and Anna Rumshisky. 2020. When BERT Plays the Lottery, All Tickets Are Winning . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3208–3229, Online. Association for Computational Linguistics.	819
764		820
765		821
766		822
767		823
768		824
769	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S. Gordon. 2011. Choice of Plausible Alternatives: An Evaluation of Commonsense Causal Reasoning . In <i>AAAI Spring Symposium on Logical Formalizations of Commonsense Reasoning</i> , Stanford University.	825
770		826
771		827
772		828
773		829
774		830
775	Victor Sanh. 2019. Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter . <i>CoRR</i> , abs/1910.01108.	831
776		832
777		833
778	Nikunj Saunshi, Sadhika Malladi, and Sanjeev Arora. 2020. A mathematical exploration of why language models help solve downstream tasks . <i>CoRR</i> , abs/2010.03648.	834
779		835
780		836
781		837
782	Chao Shang, Yun Tang, Jing Huang, Jinbo Bi, Xiaodong He, and Bowen Zhou. 2019. End-to-end structure-aware convolutional networks for knowledge base completion . In <i>The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019</i> , pages 3060–3067. AAAI Press.	838
783		839
784		840
785		841
786		842
787		843
788		844
789		845
790		846
791		847
792		848
793	Xing Shi, Inkit Padhi, and Kevin Knight. 2016. Does string-based neural MT learn source syntax? In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1526–1534, Austin, Texas. Association for Computational Linguistics.	849
794		850
795		851
796		852
797		853
798		854
799	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 <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4222–4235, Online. Association for Computational Linguistics.	855
800		856
801		857
802		858
803		859
804		860
805		861
806		862
807	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. Mpnet: Masked and permuted pre-training for language understanding . <i>arXiv preprint arXiv:2004.09297</i> .	863
808		864
809		865
810		866
811	Robyn Speer and Catherine Havasi. 2012. Representing general relational knowledge in ConceptNet 5 . In <i>Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)</i> , pages 3679–3686, Istanbul, Turkey. European Language Resources Association (ELRA).	867
812		868
813		869
814		870
815		871
816		872
817	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense	873
818		874
819		875

820	knowledge. In <i>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)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	874
821		875
822		876
823		877
824		878
825		879
826		880
827		
828		
829		
830		
831		
832		
833		
834		
835	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. 2019. What do you learn from context? probing for sentence structure in contextualized word representations . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	881
836		882
837		883
838		884
839		885
840		886
841		887
842	Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. 2016. Complex embeddings for simple link prediction . <i>CoRR</i> , abs/1606.06357.	888
843		889
844		890
845		891
846		
847		
848		
849	Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE . <i>Journal of Machine Learning Research</i> , 9:2579–2605.	892
850		893
851		894
852	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . In <i>Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA</i> , pages 5998–6008.	895
853		896
854		897
855		898
856		
857		
858		
859		
860		
861		
862	Jesse Vig. 2019. Visualizing attention in transformer-based language representation models . <i>CoRR</i> , abs/1904.02679.	899
863		900
864		901
865		
866		
867		
868	Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: A free collaborative knowledgebase . <i>Commun. ACM</i> , 57(10):78–85.	902
869		
870		
871		
872		
873		
874	Cunxiang Wang, Shuaileong Liang, Yue Zhang, Xiaonan Li, and Tian Gao. 2019. Does it make sense? and why? a pilot study for sense making and explanation . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4020–4026, Florence, Italy. Association for Computational Linguistics.	903
875		
876		
877		
878		
879		
880	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierrette Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface’s transformers: State-of-the-art natural language processing . <i>CoRR</i> , abs/1910.03771.	904
881		
882		
883		
884		
885		
886		
887		
888	Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. Embedding entities and relations for learning and inference in knowledge bases . In <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> .	905
889		
890		
891		

Review:

META-REVIEW

Comments: This paper presents a set of experiments in which pruning methods are used to identify (relatively sparse) subnetworks within pretrained MLMs that encode "commonsense" knowledge. While the practical use of this approach is not immediately clear, reviewers agreed that this is nonetheless an interesting direction and analysis, and the authors clarified key remaining questions about the work in their response. These clarifications should be integrated into future versions of the manuscript.

REVIEWER #1

The core review

This paper presents an in-depth analysis of the nature of relational knowledge elicited by pretrained language models. Along these lines, the authors posit that it is possible to find subnetworks at non-trivial sparsity that can generalize to downstream commonsense reasoning tasks.

Strengths:

- Understanding the nature of relational knowledge present in pretrained language models is an important challenge and this paper tackles that by identifying a smaller sparse subnetwork
- They validate this by deploying the pruned model to downstream tasks

Weaknesses:

- It is hard to understand the motivation of the paper. What does pruning really show ? Presence of relation-only subnetwork ? If yes, the authors have not made an attempt to show that such networks do not capture any additional information. From an end-task perspective, are those the only "subnetworks" that can be used to enhance performance ? If yes, is the relational knowledge distributed differently than the authors claim ?
- It is unclear how their network specifically targets "relational commonsense" knowledge. In particular, the current "knowledge probing with cloze prompts" is not sufficiently convincing that it is meant to capture relational knowledge in the first place
- The paper is fraught with terms that are either vague or not defined
 - Line 97 - what does "embodies" mean ?
 - Line 89 - what does "discrepancy" mean ?
 - Line 415 - what does "special" mean ?

- The authors don't adequately motivate why the visualization of attention weights is necessary to establish any analysis.
-

Reasons to Accept

The paper is not ready to be accepted for scholarly publication

Reasons to Reject

The paper is not motivated well, and currently lacks some foundational basis albeit having a good idea to understand the relationship between pruning and knowledge.

Reviewer's Scores

Overall Recommendation: 3.5

REVIEWER #2

The core review

This paper addresses the question of disentangling relational knowledge from pretrained language embeddings. They propose an end-to-end weakly supervised weights pruning method to search for subnetworks within pretrained language models in which relational knowledge is elicited. They show that the proposed method identifies sparse subnetworks which performs well on several downstream tasks. For pruning they propose two different methods - (1) stochastic pruning - where the first method a binary masking variable is sampled from a Bernoulli distribution (2) deterministic pruning - where a hard thresholding function is used.

Reasons to Accept

The paper provides new insights into how relational knowledge is captured in pretrained language models. Experiments are thorough. The authors show a huge performance difference between finetuning the original language model vs finetuning the pruned model.

Reasons to Reject

There is very little discussion about the practical implications of identifying these subnetworks. Apart from the understanding that they exist and can be finetuned for downstream tasks, there is no discussion on why there is such a huge performance difference between finetuning original BERT vs finetuning pruned BERT.

Reviewer's Scores

Overall Recommendation: 3.5

Questions for the Author(s)

In Table 1, the precision is zero for models finetuned on the original bert. Any reason why pure Bert finetuning would do so poorly? Also, how does finetuning RoBERTA (without pruning) perform for this task?

REVIEWER #3

The core review

In this paper, the authors probe into the common sense knowledge already present in the pretrained masked language models by using weakly supervised weight pruning techniques. The idea is to find subnetworks within PLMs in which relational knowledge is better represented. The authors show that grounding on external relation schema successfully identifies sparse subnetworks. They specifically investigate two questions (a) if they can disentangle the pretrained general-purpose knowledge representation into a relation-specific knowledge representation, (b) if the relation-specific knowledge can be used for downstream knowledge-intensive tasks.

Strengths:

- The paper proposes novel pruning strategies to find subnetworks that take into account relational knowledge.
- Empirical evaluation is conducted on multiple PLMs and tested on 7 common sense reasoning tasks.

Weakness:

- As Petroni et. al 2019 and authors mention PLMs are parameterized to have a knowledge representation that is entangled in the shared parameter space. However, when fine-tuning the general knowledge representation is tuned towards a specific task. I wonder if the results presented by the pruning methods proposed would be significant after fine-tuning to the task.
 - Some of the experiments and results are unclear. Why BERT-BASE-FINETUNED-SQuAD does not perform on the LAMA task (0.0 for P@1) as the SQuAD dataset is also a subset of LAMA?
-

Reasons to Reject

- It is understandable that subnetworks of PLMs are useful for representing knowledge. However, it is unclear their usefulness after fine-tuning PLMs for a task. A fine-tuned model would ideally capture the required common sense relational representation for the specific task.
-
-

Reviewer's Scores

Overall Recommendation: 3.5

Rebuttal:

Response to Review #1:

Thanks for your valuable review. Regarding your questions:

Q1: "...What does pruning really show..."

A1: The pruning technique learns parts of a PLM that are specialized for a certain type of commonsense relation by keeping the relevant parameters and resetting the rest of the parameters. The resulting subnetwork can accurately model $P(\text{object}|\text{subject, relation})$ for a particular relation. This is shown in Section 3.1 line 420-425, 444-447, as well as in Figure 2.

The main motivation of the paper is to show that PLM can be partitioned into these relation-specific subnetworks, and not to improve the performance of any particular down-stream tasks, although we have successfully shown that some of the down-stream tasks can indeed benefit from

the PLM pruning, as a by-product. In contrast, Chen et al.[1] prunes BERT to obtain a subnetwork that targets a specific application and not a relation like in our work, which is more general.

References:

[1]. The Lottery Ticket Hypothesis for Pretrained BERT Networks. NIPS 2020.

Q2: "knowledge probing with cloze prompts" is not convincing...

A2: The cloze prompts are written by human annotators in OMCS. Each prompt describes a commonsense relation that can be formulated as a subject-predicate-object triple. The blank is carefully chosen by human that is either the subject or the object of the triple.

We use cloze prompts here because they well unify the human language model and structured knowledge.

Q3: "The paper is fraught with...":

A3:

1. "embody" means some MLM pretraining instances require commonsense knowledge of certain relations to correctly predict masked word. Done

2. "discrepancy" means the difference between MLM pretraining formulation $P(\text{word_masked}|\text{word_unmasked})$ and relational knowledge evaluation formulation $P(\text{object}|\text{subject, relation})$. Done

3. "special" means the uniqueness of a particular subnetwork for representing a specific commonsense relation. Done

We will correct some of the language uses according to your suggestion.

Q4: "no motivation of visualization of attention weights...":

A4: Before pruning, the attention between words in the example cloze prompt shows no particular interesting pattern; after pruning, there is clearly strong attention between the subject and the object words for the same example. This is a good evidence that our technique strengthens the PLM representation of the relation depicted in the prompt.

Response to Review #2:

Thanks for your valuable review. Regarding your questions:

Q1: "the precision is zero for models finetuned on ..." Done

A1: First of all, all the P@1 scores of Table 1 are results of LAMA, which tests the model's understanding of structured knowledge, and it's not a downstream task. Second, fine-tuning BERT on downstream tasks is essentially adapting the pretrained representations to capture task-specific features, usually by heavily modifying the outer-most layers of BERT. It has been shown that large language models, when fine-tuned, are prone to exploiting superficial statistical cues (Naik et al.[1], Sanchez et al.[2], McCoy et al.[3]) in the dataset, e.g., trigger word or n-gram indicative for certain label, to achieve high performance without truly utilizing "knowledge". That is why BERT fine-tuned by CONLL and SQuAD suffers the most on the LAMA test, because it "forgets" some of its previously acquired knowledge. The same goes for fine-tuned RoBERTa. The precisions for RoBERTa fine-tuned on SST-2 and SQuADv1.0 are both near zero.

When we apply pruning of fine-tuned BERT, Table 1 shows that we can recover the representation of

some knowledge and the P@1 is improved to 27.1 and 22.5. But some of the knowledge maybe permanently

lost, so the precision is still much lower than BERT without finetuning at 57.6.

References:

- [1]. Stress Test Evaluation for Natural Language Inference. COLING 2019.
- [2]. Behavior Analysis of NLI models: Uncovering the Influence of Three Factors on Robustness. NAACL 2018.
- [3]. Right for the Wrong Reasons: Diagnosing Syntactic Heuristics in Natural Language Inference. ACL 2019.

Q2: "huge performance difference between finetuning original BERT vs finetuning pruned BERT"

Done

A2: Please refer to the previous answer. Notice in Table 1, "BERT-BASE-FINETUNED-CONLL03 w/ deterministic pruning" means finetune BERT first and then do deterministic pruning, not the other way around.

This also applies to other fine-tune + pruning experiments.

Response to Review #3:

Thanks for your valuable review. Regarding your questions:

Q1: "...the SQuAD dataset is also a subset of LAMA?" Done

A1: In this paper we only used the subset of LAMA exclusively targeting relational commonsense knowledge (line 331-335). It has nothing to do with SQuAD.

Q2: "...results of the pruning will hold even after comparing with fine-tuned models?"

A2: We think you are referring to Table 1, which only includes comparisons with BERT fine-tuned on CONLL03 and SQuAD. We agree that it is conceivable that BERT fine-tuned on commonsense related tasks

in Table 3 might strengthen the representation of commonsense knowledge and improve the

its performance on LAMA. But in reality, that's the case. Our experiments on BERT fine-tuned on SWAG and aNLI shows that the P@1 on LAMA is only slightly better than zero, significantly lower than 12.9. The reason is given in A1 to Review #2. With your permission we can add this discussion into the revised version.

Q3: "What is the value of t in the deterministic pruning?" Done

A3: The pruning threshold value t is 0.5 for all models.

Q4: "unclear how a specific combination of these relations is useful for a task?"

A4: This framework provides a mask for each relation type. One can try to union the masks of a set of relation types to prune the model. Table 2 of Appendix shows the best combination of relation masks

for different PLMs and different tasks.

When all relations are used, the results on the downstream commonsense tasks deteriorate. This is because most tasks use only a few relations. A model masked by all relations behaves toward the original network without pruning.