

The Life Cycle of Knowledge in Big Language Models: A Survey

Boxi Cao^{1,3}, Hongyu Lin¹, Xianpei Han^{1,2✉}, Le Sun^{1,2}

¹Chinese Information Processing Laboratory ²State Key Laboratory of Computer Science
Institute of Software, Chinese Academy of Sciences, Beijing, China

³University of Chinese Academy of Sciences, Beijing, China
{boxi2020, hongyu, xianpei, sunle}@iscas.ac.cn

Abstract

Knowledge plays a critical role in artificial intelligence. Recently, the extensive success of pre-trained language models (PLMs) has raised significant attention about how knowledge can be acquired, maintained, updated and used by language models. Despite the enormous amount of related studies, there still lacks a unified view of how knowledge circulates within language models throughout the learning, tuning, and application processes, which may prevent us from further understanding the connections between current progress or realizing existing limitations. In this survey, we revisit PLMs as knowledge-based systems by dividing the life circle of knowledge in PLMs into five critical periods, and investigating how knowledge circulates when it is built, maintained and used. To this end, we systematically review existing studies of each period of the knowledge life cycle, summarize the main challenges and current limitations, and discuss future directions¹.

1 Introduction

Fundamentally, AI is the science of knowledge – how to represent knowledge and how to obtain and use knowledge.

Nilson (1974)

Knowledge is the key to high-level intelligence. How a model obtains, stores, understands and applies knowledge has long been a critical research topic in machine intelligence. Recent years have witnessed the rapid development of pre-trained language models (PLMs). Through self-supervised pre-training on large-scale unlabeled corpora, PLMs show strong generalization and transferring abilities across different tasks/datasets/settings over previous methods, and

¹We openly released a corresponding paper list which will be regularly updated on <https://github.com/c-box/KnowledgeLifecycle>.

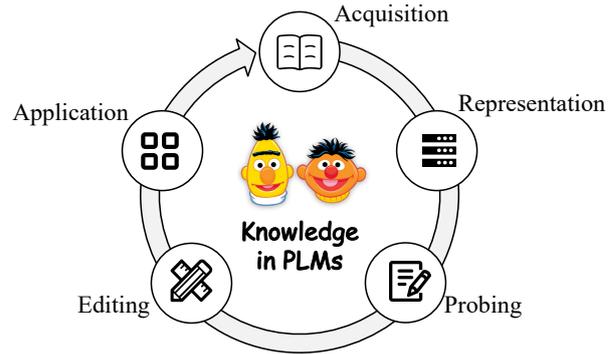


Figure 1: Five critical periods in life circle of knowledge in language models.

therefore have achieved remarkable success in natural language processing (Devlin et al., 2019; Liu et al., 2019c; Raffel et al., 2020; Radford et al., 2019b; Brown et al., 2020; Lewis et al., 2020a).

The success of pre-trained language models has raised great attention about the nature of their entailed knowledge. There have been numerous studies focusing on how knowledge can be acquired, maintained, and used by pre-trained language models. Along these lines, many novel research directions have been explored. For example, knowledge infusing devotes to injecting explicit structured knowledge into PLMs (Sun et al., 2019; Zhang et al., 2019; Sachan et al., 2021). Knowledge probing aims to evaluate the type and amount of knowledge stored in PLMs’ parameters (Petroni et al., 2019; Lin et al., 2019; Hewitt and Manning, 2019). And knowledge editing is dedicated to modifying the incorrect or undesirable knowledge acquired by PLMs (Zhu et al., 2020; De Cao et al., 2021; Mitchell et al., 2021).

Despite the large amount of related studies, current studies primarily focus on one specific stage of knowledge process in PLMs, thereby lacking a unified perspective on how knowledge circulates throughout the entire model learning, tuning, and

application phases. The absence of such comprehensive studies makes it hard to better understand the connections between different knowledge-based tasks, discover the correlations between different periods during the knowledge life cycle in PLMs, exploit the missing links and tasks for investigating knowledge in PLMs, or explore the shortcomings and limitations of existing studies. For example, while numerous studies attempt to assess the knowledge in language models that are already pre-trained, there are few studies dedicated to investigating why PLMs can learn from pure text without any supervision about knowledge, as well as how PLMs represent or store these knowledge. Meanwhile, many researchers have tried to explicitly inject various kinds of structural knowledge into PLMs, but few studies propose to help PLMs better acquire specific kinds of knowledge from pure text by exploiting the knowledge acquisition mechanisms behind. As a result, related research may be overly focused on several directions but fail to comprehensively understand, maintain and control knowledge in PLMs, and therefore limits the improvements and further application.

In this survey, we propose to systematically review the knowledge-related studies in pre-trained language models from a knowledge engineering perspective. Inspired by research in cognitive science (Zimbardo and Ruch, 1975; Churchland and Sejnowski, 1988) and knowledge engineering (Studer et al., 1998; Schreiber et al., 2000), we regard pre-trained language models as knowledge-based systems, and investigate the life cycle of how knowledge circulates when it is acquired, maintained and used in pre-trained models (Studer et al., 1998; Schreiber et al., 2000). Specifically, we divide the life cycle of knowledge in pre-trained language models into the following five critical periods as shown in Fig. 1:

- **Knowledge Acquisition**, which focuses on the procedure of language models learning various knowledge from text or other knowledge sources.
- **Knowledge Representation**, which focuses on the underlying mechanism of how different kinds of knowledge are transformed, encoded, and distributed in PLMs’ parameters.
- **Knowledge Probing**, which aims to evaluate how well current PLMs entailing different types of knowledge.

- **Knowledge Editing**, which tries to edit or delete knowledge containing in language models.
- **Knowledge Application**, which tries to distill or leverage knowledge in pre-trained language models for practical application.

For each of these periods, we sort out the existing studies, summarize the main challenges and limitations, and discuss future directions. Based on the unified perspective, we are able to understand and utilize the close connections between different periods instead of consider them as independent tasks. For instance, understanding the knowledge representation mechanism of PLMs is valuable for researchers to design better knowledge acquisition objectives and knowledge editing strategies. Proposing reliable knowledge probing methods could help us find the suitable applications for PLMs, and gain insight into their limitations, thereby facilitating improvement. Through this survey, we are willing to comprehensively conclude the progress, challenges and limitations of current studies, help researchers better understand the whole field from a novel perspective, and shed light on the future directions about how to better regulate, represent and apply the knowledge in language models from a unified perspective.

We summarize our contributions as follows:

- We propose to revisit pre-trained language models as knowledge-based systems, and divide the life cycle of knowledge in PLMs into five critical periods.
- For each period, we review existing studies, summarize the main challenges and shortcomings for each direction.
- Based on this review, we discuss about the limitations of the current research, and shed light to potential future directions.

2 Overview

In this section, we present the overall structure of this survey, describe our taxonomy shown in Fig. 2 in detail, and discuss the topics in each critical period.

Knowledge Acquisition is the knowledge learning procedure of language models. Currently, there are two main sources for knowledge acquisition:

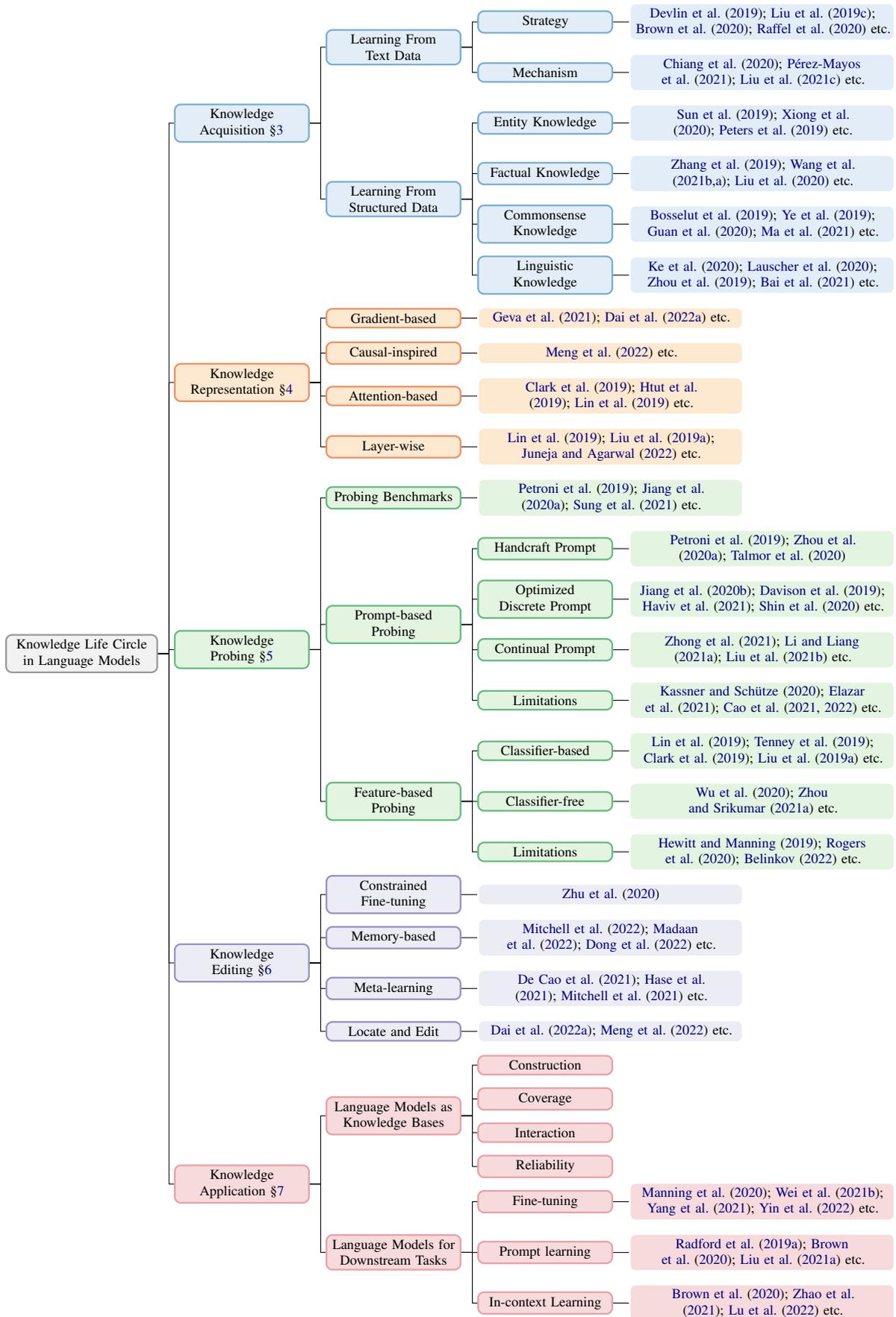


Figure 2: Typology of knowledge life cycle in big language models.

the plain text data and the structured data. For acquiring knowledge from text data, LMs typically conduct self-supervised learning on large-scale text corpora (Devlin et al., 2019; Liu et al., 2019c; Brown et al., 2020; Raffel et al., 2020). This survey will focus on the methods and mechanisms of how pre-trained language models obtaining knowledge from pure texts (Chiang et al., 2020; Pérez-Mayos et al., 2021; Liu et al., 2021c). For acquiring knowledge from structured data, current research focus on knowledge injection from different kinds of structured data into PLMs. The primary categories of structured data contains entity knowledge (Sun et al., 2019; Xiong et al., 2020; Peters et al., 2019), factual knowledge (Zhang et al., 2019; Wang et al., 2021b,a; Liu et al., 2020), commonsense knowledge (Bosselut et al., 2019; Ye et al., 2019; Guan et al., 2020; Ma et al., 2021) and linguistic knowledge (Ke et al., 2020; Lauscher et al., 2020; Zhou et al., 2019; Bai et al., 2021). We will discuss all of them in Section 3.

Knowledge Representation aims to investigate how language models encode, store and represent knowledge in their dense parameters. The investigation about the knowledge representation mechanisms will aid in a better understanding and control of knowledge in PLMs, and may also inspire researchers for better understanding the knowledge representation in human brains. Currently, the strategies for knowledge representation analysis in PLMs include gradient-based (Geva et al., 2021; Dai et al., 2022a), causal-inspired (Meng et al., 2022), attention-based (Clark et al., 2019; Htut et al., 2019; Lin et al., 2019), and layer-wise (Lin et al., 2019; Liu et al., 2019a; Juneja and Agarwal, 2022) methods. We will discuss them in Section 4.

Knowledge Probing aims to evaluate how well current PLMs entailing specific types of knowledge. Currently, two primary strategies are used to probe the knowledge in PLMs: 1) Prompt-based probing, which usually constructs knowledge-instructed prompt, then query PLMs using these natural language expressions (Petroni et al., 2019; Jiang et al., 2020a; Sung et al., 2021; Forbes et al., 2019; Zhou et al., 2020a). For example, querying PLMs with “The capital of France is ___.” to evaluate whether PLMs have stored the corresponding knowledge <France, capital, Paris>. Meanwhile, to improve PLMs’ performance, a series of studies devote to optimizing prompts in both discrete (Jiang et al., 2020b; Davison et al., 2019; Haviv et al., 2021;

Shin et al., 2020) and continual space (Zhong et al., 2021; Li and Liang, 2021a; Liu et al., 2021b). Despite the widely application of prompt-based probing, lots of studies also point out that there still exist some pending issues such as inconsistent (Elazar et al., 2021; Kassner and Schütze, 2020; Jang et al., 2022; Cao et al., 2022), inaccurate (Poerner et al., 2020; Zhong et al., 2021; Cao et al., 2021) and unreliable (Cao et al., 2021; Li et al., 2022a), and question the quantity results of prompt-based probing. 2) Feature-based probing, which normally freezes the parameters of original PLMs, and evaluates PLMs on probing tasks based on their internal representation or attention weights. We categorize existing feature-based probing studies into classifier-based probing (Lin et al., 2019; Tenney et al., 2019; Clark et al., 2019; Liu et al., 2019a) and classifier-free probing (Wu et al., 2020; Zhou and Srikumar, 2021a) according to whether an additional classifier is introduced. Since most methods introduce additional parameters or training data, the main shortcoming of feature-based probing is whether the results should attribute to knowledge in PLMs or probing task learned by additional probes. We will discuss them in Section 5.

Knowledge Editing aims to modify the incorrect knowledge or delete the undesirable information in PLMs. Because of inevitable mistakes learned by PLMs and the update of knowledge, reliable and effective knowledge editing approaches are essential for the sustainable application of PLMs. Current approaches include constrained fine-tuning (Zhu et al., 2020), memory-based (Mitchell et al., 2022; Madaan et al., 2022; Dong et al., 2022), meta-learning inspired (De Cao et al., 2021; Hase et al., 2021; Mitchell et al., 2021) and location-based methods (Dai et al., 2022a; Meng et al., 2022). We will discuss them in Section 6.

Knowledge Application aims to distill or leverage specific knowledge from PLMs to benefit further applications. Currently, there are two main kinds of application paradigms for knowledge in PLMs: 1) Language models as knowledge bases (LMs-as-KBs), which regards language models as dense knowledge bases that can be directly queried with natural language to obtain specific types of knowledge (Petroni et al., 2019; Heinzlerling and Inui, 2021; Jiang et al., 2020b; Wang et al., 2020; Cao et al., 2021; Razniewski et al., 2021; AlKhamissi et al., 2022). And we provide

a comprehensive comparison between structured knowledge bases and LMs-as-KBs (Razniewski et al., 2021) from four aspects, including construction, coverage, interaction and reliability; 2) Language models for downstream task, which directly uses PLMs entailing specific kinds of knowledge in downstream NLP tasks via fine-tuning (Manning et al., 2020; Wei et al., 2021b; Yang et al., 2021; Yin et al., 2022), prompt-learning (Radford et al., 2019a; Brown et al., 2020; Liu et al., 2021a) and in-context learning (Brown et al., 2020; Zhao et al., 2021; Lu et al., 2022). We will discuss them in Section 7.

3 Knowledge Acquisition

During the knowledge acquisition period, pre-trained language models learn knowledge from different knowledge sources. In this section, we categorize and describe knowledge acquisition strategies according to knowledge sources, and then discuss the future directions.

3.1 Learning from Text Data

Currently, pre-trained language models usually acquire various knowledge from pure text through self-supervised learning on a large-scale text corpus. In this section, we will first introduce several widely used learning objectives (Qiu et al., 2020), and then discuss the learning mechanisms behind them.

Causal Language modeling aims to autoregressively predict the next token in the input sequence, which is the most popular pre-training tasks (Radford et al., 2019b; Brown et al., 2020; Ouyang et al., 2022; Scao et al., 2022) and has demonstrated excellent effectiveness in capturing context dependency and text generation paradigms. One limitation of causal language modeling is unidirectional, which can only capture contextual information from left to right.

Masked Language Modeling aims to mask some tokens in the input randomly, and then predict the masked token conditioned on the rest of sequence (Devlin et al., 2019; Liu et al., 2019c). Unlike causal language modeling, which can only obtain information in a unidirectional manner, masked language modeling can capture contextual information from both left-to-right and right-to-left directions.

Seq2seq Masked Language Modeling uses an encoder-decoder architecture for pre-training,

which first feeds the encoder with masked sequence, and the decoder is supposed to predict the masked tokens autoregressively (Raffel et al., 2020; Song et al., 2019).

Denosing Autoencoder first corrupts the input sequence with randomly mask symbols, then feed the input into a bidirectional encoder, and the likelihood of the whole original input is calculated with an auto-regressive decoder (Lewis et al., 2020a).

Although PLMs are pre-trained without any supervision from external knowledge sources, they have been shown to capture a diverse range of knowledge within their parameters, such as linguistic knowledge (Lin et al., 2019; Tenney et al., 2019; Liu et al., 2019b; Htut et al., 2019; Hewitt and Manning, 2019; Goldberg, 2019; Warstadt et al., 2019), semantic knowledge (Tenney et al., 2019; Wallace et al., 2019; Ettinger, 2020) and world knowledge (Davison et al., 2019; Bouraoui et al., 2020; Forbes et al., 2019; Zhou et al., 2020b; Roberts et al., 2020; Lin et al., 2020a; Tamborrino et al., 2020).

Intuitively, PLMs learn such knowledge because they can abstract, generalize and store the implicit knowledge in the text through self-supervised learning. Unfortunately, the underlying mechanism of how and why PLMs acquire or forget knowledge still remains to be explored. And it will be valuable to understand the behaviors of PLMs and inspire better knowledge acquisition strategies.

To understand the underlying mechanisms, some studies dive into the dynamics of LMs' pre-training procedure. Many researchers study the training dynamics of neural networks. For example, Achille et al. (2019) try to figure out whether there exist critical periods in the learning process of neural networks. Liu et al. (2021c) devote to finding a mathematical solution for the semantic development in deep linear networks. Other studies (Saphra and Lopez, 2019, 2020) analyze the training dynamics of LSTM (Hochreiter and Schmidhuber, 1997) with techniques such as SVCCA (Raghu et al., 2017).

While most existing studies focus on neural networks with relatively simple architectures. Only a few studies consider knowledge in large-scale pre-trained language models. Chiang et al. (2020) first systematically investigate the knowledge acquisition process during the training of ALBERT (Lan et al., 2020). Specifically, they study the syntactic knowledge, semantic knowledge, and world knowledge development during pre-training, and find that

the learning process varies across knowledge, and having more pre-trained steps could not necessarily increase the knowledge in PLMs. Pérez-Mayos et al. (2021) investigate the effect of the size of the pre-trained corpus on the syntactic ability of the RoBERTa (Liu et al., 2019c) model, and find that models pre-trained on more data typically contain more syntactic knowledge and perform better in related downstream tasks. Liu et al. (2021c) also investigate the knowledge acquisition process of RoBERTa (Liu et al., 2019c) on various knowledge. And find that compared with linguistic knowledge which can be learned quickly and robustly, world knowledge is learned slowly and domain-sensitive.

3.2 Learning from Structured Data

Apart from acquiring knowledge from pure text, PLMs can also acquire knowledge by injecting explicit structured knowledge into them. In this section, we review these studies according to the category of structured knowledge sources.

Entity Knowledge To learn entity knowledge explicitly, lots of studies propose entity-guided tasks for language model pretraining. For example, Sun et al. (2019) and Shen et al. (2020) use entity-level masking to enhance language models, which first recognize named entities in a sentence, then all the tokens corresponding to these entities as masked and predicted at once. Xiong et al. (2020) present replaced entity detection, which randomly replaces the named entities in a sentence with another mention of the same entity or other entities of the same type, and LMs are supposed to determine which entities are replaced. Yamada et al. (2020) treat words and entities as independent tokens, and conduct mask language modeling separately to learn both contextualized word representation and entity representation. Févry et al. (2020) combine the mention detection and entity linking pre-training objectives with mask language modeling to match the entities in text with specific entity memories. In addition to the entity mentions themselves, researchers have also introduced other meta-information such as entity description to further assist the entity knowledge learning (Logeswaran et al., 2019; Gillick et al., 2019). Another efficient way to enrich PLMs' text representation with entity knowledge is utilizing word-to-entity attention (Peters et al., 2019; Yamada et al., 2020).

Factual Knowledge In structured knowledge bases, factual knowledge is generally represented

as triples (subject entity, relation, object entity). For a long time, researchers have been dedicated to aiding PLMs to acquire more factual knowledge to perform better on downstream tasks. On the one hand, introducing knowledge graph embedding into the pre-training procedure could be effective. Zhang et al. (2019) propose an aggregator to combine the corresponding knowledge embedding of the entities in text and token embedding. Wang et al. (2021b) co-train mask language modeling and knowledge graph embedding objectives, which could produce both informative text and knowledge embedding. On the other hand, some studies propose designing factual knowledge-guided auxiliary tasks. Wang et al. (2021a) add an adapter to infuse knowledge into PLMs without updating the original parameters. The adapter is trained with predication prediction to determine the relation type between tokens. Qin et al. (2021) propose the entity discrimination tasks to predict the object entity given subject entity and relation, as well as relation discrimination tasks to predict the semantic connection between relation pairs. Banerjee and Baral (2020) directly pre-train language model on the knowledge graph, the model is given two elements of a knowledge triple to predict the rest one. Liu et al. (2020) argue that incorporating a whole knowledge base into PLMs might induce the knowledge noise issue, and propose to learn from a specific sub-graph related to each input sentence. Moreover, Baldini Soares et al. (2019) propose to learn relational knowledge solely from entity-link text through "matching in the blank" objective, which first replaces the entities in text with blank symbols and then brings the relation representations closer when they have the same pair of entities.

Commonsense Knowledge One of the most common strategy for PLMs learning commonsense knowledge is converting the knowledge to natural language expressions before learning. Bosselut et al. (2019); Guan et al. (2020); Schwartz et al. (2020) first transfer the commonsense knowledge triples to natural language with prompt, then pre-train LMs on these knowledge-augmented data. Ye et al. (2019) post-training LMs on commonsense QA datasets created by AWS (align, mask, select). Ma et al. (2021) transform structured commonsense knowledge into natural language questions for model learning.

Linguistic Knowledge By designing the corresponding pre-training tasks, PLMs could also learn linguistic knowledge explicitly, such as sentiment knowledge (Ke et al., 2020; Tian et al., 2020), lexical knowledge (Lauscher et al., 2020; Levine et al., 2020; Zhou et al., 2019), syntax knowledge (Zhou et al., 2019; Sachan et al., 2021; Bai et al., 2021), etc. For example, to equip LMs with sentiment knowledge, Ke et al. (2020) first label each word with a POS tag and sentiment polarity, and then incorporate both the word-level and sentence-level sentiment label with the mask language modeling objective. Similarly, Tian et al. (2020) first mine sentiment knowledge from unlabeled data based on pointwise mutual information (PMI), and then conduct pre-training tasks such as sentiment masking, sentiment word prediction and word polarity prediction with these sentiment information. As for lexical knowledge, Lauscher et al. (2020) first acquire word similarity information from WordNet (Miller, 1992) and BabelNet (Navigli and Ponzetto, 2010), and then add word relation classification tasks in addition to BERT’s original pre-training tasks. Levine et al. (2020) also introduces the lexical information from WordNet and adds a masked-word prediction task. To incorporate dependency knowledge with PLMs, Song et al. (2022) construct a dependency matrix for attention alignment calibration and a fusion module to integrate dependency information. Explicitly learning syntax knowledge also raises the researchers’ attention, Sachan et al. (2021) investigate infusing syntax knowledge by either adding a syntax-GNN on the output of transformers or incorporating with text embedding using attention. To further capture the syntax knowledge, Bai et al. (2021) using multiple attention networks, with each one encoding one relation from the syntax tree.

3.3 Discussions and Future Work

As we mentioned above, there have been extensive studies for better knowledge acquisition of language models, and most of them focus on infusing existing structured knowledge sources into PLMs. The learning from text data methods can be easily scaled, and the knowledge sources is easily obtained. But the underlying mechanism is still mostly unclear, the knowledge acquisition process is implicitly and thus is hard to control, and may lead to inconsistent prediction, undesirable bias and unforeseen risks. The learning from structured

data methods can explicitly inject knowledge into PLMs, but are limited by the cost, domain, scale and quality of knowledge sources. Furthermore, since the knowledge injection methods are often specialized to specific kinds of knowledge, it is often difficult to extend or produce new knowledge.

Furthermore, because all knowledge in PLMs are implicitly encoded as parameters, it is often very difficult to control and validate the knowledge acquisition process. There are also several studies such as retrieval-based PLMs, focus on retrieving related knowledge or context to enhance original PLMs (Guu et al., 2020; Lewis et al., 2020b; Yasunaga et al., 2022), rather than injecting knowledge into PLMs’ parameters.

Several future directions of knowledge acquisition in PLMs may lie in: 1) For the knowledge acquisition from existing structured knowledge sources, it is critical to develop universal knowledge injection methods which can uniformly injecting different types of knowledge from different knowledge sources, and ensures continuous learning and avoid catastrophic forgetting in the meantime. 2) For the knowledge acquisition from pure text data, it is helpful to fully understand the underlying mechanism of knowledge learning in PLMs, and develop effective knowledge learning algorithms which can learn specific knowledge from text data in a controllable and predicable way. 3) Furthermore, it is also important to build comprehensive benchmarks for investigating and assessing the knowledge acquisition process of PLMs.

4 Knowledge Representation

Knowledge representation studies investigate how pre-trained language models encode, transform and store the acquired knowledge. In PLMs, knowledge is encoded to dense vector representations and held in their distributed parameters, but how each kind of knowledge is encoded, transformed, and stored into the parameters is still unclear and needs further investigation. Currently, a few studies have investigated the knowledge representation in language models, and we will first review these studies according to their analysis techniques.

4.1 Analyzing Knowledge Representations in PLMs

Currently, the analyzing approaches for knowledge representation in PLMs can be classified into four categories: gradient-based, causal-inspired,

attention-based and layer-wise methods. The first three methods aim to locate specific knowledge in PLMs' corresponding neurons or attention heads, and the layer-wise methods hypothesize that knowledge is represented in different layers of PLMs.

Gradient-based Dai et al. (2022a) first introduce the concept of knowledge neurons, which are neurons in transformer (Vaswani et al., 2017) related to certain factual knowledge. Specifically, they hypothesize the knowledge neurons are located in feed-forward networks, which are considered as key-value memories (Geva et al., 2021). Then by feeding the LM with knowledge-expressing prompts such as "Michael Jordan was born in [MASK]", the corresponding knowledge neuron is identified as the neurons in the feed-forward networks with higher attribution scores, which are calculated based on integrated gradients.

Causal-inspired Meng et al. (2022) identify knowledge neurons as the neuron activations in transformers that have the strongest causal effect on predicting certain factual knowledge. Such neurons are located through a causal mediation analysis. Specifically, they calculate the causal effect on factual prediction by comparing probability variation of object prediction between the clean and corrupted token embedding. Their experiments also demonstrate that the mid-layer feed-forward modules play a decisive role in factual knowledge representation.

Attention-based In addition to the feed-forward layers, the attention heads are also be considered as representations which may encode the knowledge-related information. Clark et al. (2019); Htut et al. (2019) investigate the linguistic knowledge encoded in attention heads, and find that while some individual attention heads are associated with specific aspects of syntax, the linguistic knowledge is distributed and represented by multiple attention heads. Lin et al. (2019) find that PLMs' attention weights could encode syntactic properties such as subject-verb agreement and reflexive dependencies, and higher layers represent these syntactic properties more accurately.

Layer-wise Lin et al. (2019) conduct a layer-wise probing for linguistic knowledge, which trains a specific classifier for each layer, and find that the lower layers encode the positional information of tokens, and higher layers encode more composi-

tional information. Liu et al. (2019a) analyze the layerwise transferability of PLMs on a wide range of tasks and find that the middle layers usually have better performance and transferability. Wallat et al. (2020) proposes to probe the captured factual knowledge with LAMA (Petroni et al., 2019) of each layer in PLMs, and finds that a significant amount of knowledge is stored in the intermediate layers. Juneja and Agarwal (2022) also conduct a layer-wised factual knowledge analysis based on knowledge neuron (Dai et al., 2022a), and demonstrate that most relational knowledge (e.g., Paris is the capital of "some nation".) can be attributed to the middle layers, which would be refined into facts (e.g., Paris is the capital of France.) in the last few layers.

4.2 Discussions and Future Works

The above studies reach some consensus about knowledge representation in PLMs, including: 1) Factual knowledge can be associated with feed-forward modules in middle or higher layers. 2) Linguistic knowledge is distributed and represented in multiple attention heads, while a single attention head can only associate with a specific aspect of linguistics. 3) The lower layers of PLMs often encode the coarse-grained and general information of knowledge, while the fine-grained and task-specific knowledge are mostly stored in higher layers. These findings are valuable for us to understand knowledge representation in language models but are also limited to specific knowledge types or model architectures. Therefore the knowledge representation in PLMs is still an open problem which needs further exploration.

In the future, several directions of knowledge representation in PLMs may lie in the following: 1) Because knowledge representation is a long-standing concern in cognitive science, neuroscience, psychology, and artificial intelligence, it is helpful to borrow ideas from other related areas and design cognitively-inspired analysis methods. 2) Current knowledge representation studies in PLMs mostly focus on a specific type of knowledge and often result in local and specific conclusions. It is important to comprehensively investigate different types of knowledge together, e.g., compare the differences and commonalities of knowledge representations of different knowledge types, pre-training tasks, or model architectures, and come up with more universal and insightful conclusions.

5 Knowledge Probing

Knowledge probing aims to assess how well pre-trained language models entail different kinds of knowledge. A comprehensive and accurate assessment of PLMs’ knowledge can help us identify and understand language models’ capabilities and deficiencies, allow a fair comparison between LMs with different architectures and pre-training tasks, guide the improvement of a specific model, and select suitable models for different real-world scenarios. In this section, we will first introduce existing benchmarks for knowledge probing, then introduce the representative prompt-based and feature-based probing methods and analyze their corresponding limitations, and discuss future directions.

5.1 Benchmarks for Knowledge Probing

To assess the knowledge in PLMs, lots of benchmarks have been proposed to probe various knowledge contained in PLMs, for example, linguistic knowledge (Ettinger, 2020; Warstadt et al., 2020; Lin et al., 2019; Warstadt et al., 2019; Tenney et al., 2019), syntactic knowledge (Clark et al., 2019; Hewitt and Manning, 2019), factual knowledge (Petroni et al., 2019; Jiang et al., 2020a; Kassner et al., 2021; Sung et al., 2021), commonsense knowledge (Forbes et al., 2019; Zhou et al., 2020a), etc. Table 1 summarizes several representative knowledge probing benchmarks.

5.2 Prompt-based Knowledge Probing

Prompt-based probing is one of the most popular approaches for knowledge probing. To evaluate whether LMs know a specific knowledge such as the birthplace of Michael Jordan, we could query LMs with knowledge queries such as “Michael Jordan was born in ___.”, where “was born in” is a prompt for a specific type of knowledge. As shown in Table 1, prompt-based probing has been widely used in benchmarks such as LAMA (Petroni et al., 2019), oLMpics (Talmor et al., 2020), LM diagnostics (Ettinger, 2020), BIG-bench (Srivastava et al., 2022), etc.

For prompt-based probing, the main challenge is how to design effective prompts which are suitable for different kinds of knowledge and different PLMs. In the following we will introduce the typical prompt types for knowledge probing and discuss their limitations.

5.2.1 Prompt Development

Handcraft Prompt Early methods often manually write prompts for different kinds of knowledge. There are two primary advantages of manually created prompts: the readability and without the need of any other resources or training. For example, LAMA (Petroni et al., 2019) manually creates one cloze-style prompt for each relation, which is used to probe the factual knowledge in language models. CAT (Zhou et al., 2020a) reframes the instances in existing commonsense datasets into paired sentences with task-specific prompts, and determine whether PLMs contain specific commonsense knowledge by comparing the sentence scores, e.g., “money can be used to buy cars” v.s. “money can be used to buy stars”. oLMpics (Talmor et al., 2020) convert the probing tasks for reasoning ability into multi-choice questions with manually created prompts, and compare the LMs’ probability of candidate choices.

Optimized Discrete Prompt Despite the mentioned advantages, Jiang et al. (2020b) argues that handcraft prompts could be sub-optimal. Therefore, a series of studies have been proposed to optimize the prompts in a discrete space so that PLMs could achieve better performance. Jiang et al. (2020b) propose a mining-based method in order to find prompts with higher performance from text corpus. They first retrieve potential prompts which contain both the subject and object entity, then select prompts using a validation dataset. Davison et al. (2019) select prompt from a handcrafted candidate set according to the log-likelihood calculated by LMs. Haviv et al. (2021) propose a paraphrasing-based method, where each query is first reframed by a trained rewriter and then fed into PLMs. Shin et al. (2020) propose an automatic prompt generation method based on gradient-guided search, where a prompt is iteratively updated from “[MASK]” token by maximizing the label likelihood of training instances.

Continual Prompt Although the prompts generated by Shin et al. (2020) are discrete text, they are very difficult to be understood by humans. Therefore several studies directly search better-performed prompts on continual space rather than confining to discrete space, i.e., representing prompts as dense vectors. Continual prompts have shown good performance for knowledge probing, and further extensions including handcraft prompts

Method	Benchmarks	Knowledge Type	Formulation
Prompt-based	LM diagnostics (Ettinger, 2020)	linguistic	text filling
	BLiMP (Warstadt et al., 2020)	linguistic	sentence scores comparison
	LAMA (Petroni et al., 2019)	factual, commonsense	text filling
	X-FACTR (Jiang et al., 2020a)	factual, multilingual	
	Multilingual LAMA (Kassner et al., 2021)	factual, multilingual	sentence scores comparison
	Bio LAMA (Sung et al., 2021)	factual, biological	
	CAT (Zhou et al., 2020a)	commonsense	text filling
	NumerSense (Lin et al., 2020b)	commonsense, numerical	multiple choices
oLMPICS (Talmor et al., 2020)	reasoning		
Feature-based	Open Sesame (Lin et al., 2019)	linguistic	diagnostic classifier and attention
	LKT (Liu et al., 2019b)	linguistic	token or token pair labeling
	NPI probe (Warstadt et al., 2019)	linguistic	probing classifier
	Edge probe (Tenney et al., 2019)	linguistic, semantic	edge probing
	MDL probe (Voita and Titov, 2020)	linguistic	minimum description length
	Structural probe (Hewitt and Manning, 2019)	syntactic	structural probing
	Physical Commonsense (Forbes et al., 2019)	commonsense, physical	probing classifier

Table 1: Summary about some representative knowledge probing benchmarks.

initialization (Zhong et al., 2021), adding continual prompts on both input and transformer blocks (Li and Liang, 2021a) or adding LSTM layers above the input embeddings (Liu et al., 2021b).

5.2.2 Limitations of Prompt-based Probing

Although prompts have been widely used to probe the knowledge in PLMs, there are still lots of pending issues unresolved, which make the probing results unstable and the the assessment of knowledge in PLMs unreliable.

Inconsistent Prompt-based probing have been shown often result in inconsistent results due to prompt selection, instance verbalization, negation, etc. Firstly, Elazar et al. (2021) find semantically equivalent prompts may result in different predictions, Cao et al. (2022) further find that PLMs would prefer specific prompts with the same linguistic regularity with the pre-training corpus, such a prompt preference will significantly affect the probing results, and result in inconsistent comparisons between PLMs. Besides prompts, the instance verbalization process also leads to inconsistent predictions. For example, when we ask BERT “The capital of the U.S. is [MASK]”, the answer is Washington, but when we replace the U.S. with its alias America, the prediction will change to Chicago. In addition, PLMs also exhibit inconsistency when facing negation (Kassner and Schütze, 2020; Jang et al., 2022). For instance, PLMs would generate highly similar predictions between a fact (“Birds can fly”) and its incorrect negation (“Birds cannot fly”) (Kassner and Schütze, 2020). Jang et al. (2022) conduct the negation experiments on

PLMs of varying sizes and various downstream tasks, and find that not only PLMs cannot well understand negation prompts, but also show an inverse scaling law.

Inaccurate The performance of PLMs under prompt-based probing may also be overestimated. Poerner et al. (2020) find that many samples in the probing datasets could be easily “guessed” by only relying on the surface form association. For example, the object entity is a substring of the subject entity (e.g., “Apple Watch is produced by Apple”). Furthermore, the training dataset for prompt optimization may correlate with probing dataset, which results in spurious correlations (Zhong et al., 2021) and the performance improvements may come from these spurious correlations. Cao et al. (2021) also find that many prompts with better performance are prompts which over-fit to answer distributions, rather than a better semantic description of the target relation.

Unreliable To reach a faithful probing results, it is essential to understand why PLMs make a specific prediction. However, studies find that PLMs do not always make predictions based on specific knowledge. In that case, the knowledge probing results could be unreliable. Cao et al. (2021) find that the prompts but not the answers dominate the prediction distribution of PLMs, resulting in severely prompt-biased probing conclusions. Li et al. (2022a) conduct a causal-inspired analysis and find that PLMs’ predictions rely more on words that are close in position and frequently co-occur, rather than those related to knowledge.

Bias Analysis While lots of studies conduct empirical experiments on the biases in prompt-based probing, few have investigated the source and interpretation of these biases. Several studies employ causal analysis for bias analysis, which has been widely used to identify undesirable biases and fairness concerns (Hardt et al., 2016; Kilbertus et al., 2017; Kusner et al., 2017; Vig et al., 2020; Feder et al., 2021). Cao et al. (2022) propose a causal analysis framework to identify, interpret and eliminate biases that exist in prompt-based probing with a theoretical guarantee. Similarly, Elazar et al. (2022) also propose a causal framework to estimate the causal effects of the data statistics in training corpus on the factual predictions of PLMs. Finlayson et al. (2021) apply causal mediation analysis to investigate the syntactic agreement mechanisms in PLMs.

5.3 Feature-based Knowledge Probing

Feature-based knowledge probing is also widely used to probe the knowledge in PLMs, where the parameters of original PLMs are frozen, and the probing tasks are accomplished based on the internal representation or attention weights produced by PLMs. In this section, we introduce and discuss the feature-based probing approaches.

5.3.1 Classifier-based Probing

Classifier-based probing trains a classifier to predict specific knowledge properties on the top of the fixed PLMs, and assesses the effectiveness of PLMs using the classifier’s performance (Belinkov, 2022). Such approaches are first propose to evaluate the linguistic properties (e.g., morphological, syntactic) associated with static embeddings (Köhn, 2015; Gupta et al., 2015), and have been widely used to probe the linguistic knowledge (Lin et al., 2019; Tenney et al., 2019; Clark et al., 2019; Liu et al., 2019a; Hewitt and Manning, 2019) and semantic knowledge (Tenney et al., 2019; Wallace et al., 2019; Yaghoobzadeh et al., 2019; Liu et al., 2019a) in PLMs. Popular classifiers include linear classifier, logistics regression, multi-layer perceptron, etc.

5.3.2 Classifier-free Probing

Since the results and conclusions of classifier-based methods are dependent on the training quality and selection of the classifier, some studies have developed feature-based probing approaches without an additional classifier. For example, Wu et al. (2020)

propose perturbed masking, which calculates an impact matrix through a two-stage perturbation, where the matrix captures the impacts a token has on the prediction of another token, and is further used for the syntactic probe. Zhou and Sriku-mar (2021b) introduce DirectProbe, which directly probes the geometric properties of PLMs’ representation without an additional classifier. Clark et al. (2019) probe syntactic knowledge in language models by investigating the attention weights without a classifier, e.g., analyze the most attended word of the given token.

5.3.3 Limitations of Feature-based Probing

There are two main limitations of current feature-based probing approaches (Rogers et al., 2020; Belinkov, 2022). The first limitation concerns the attribution of results, which is originally pointed out by Hewitt and Manning (2019). While most probes introduce additional training data and parameters, it’s difficult to attribute evaluation results to the knowledge in PLMs, or the probe itself, which may learn to perform the probing task. The second limitation pertains to the inconsistency between different probe designs for the same type of knowledge. There are various probe selections for each kind of knowledge, but the probe results between simple probes like linear classifier or complex probes could be inconsistent.

5.4 Discussions and Future Works

With the growing scale and abilities of big language models, the comprehensive, accurate and reliable measurements of the actual knowledge and capabilities of LMs become increasingly important. However, the accurate, robust and reliable probing approach is still an open problem. Firstly, as we discussed above, both prompt-based probing and feature-based probing have their own limitations, which might result in unreliable or even contradicting conclusions. Secondly, most existing benchmarks are specialized to specific knowledge types and specific model architectures.

In the future, the main directions of knowledge probing may lie in: 1) Comprehensive benchmark construction. As we demonstrate in table 1, current knowledge probing benchmarks are mostly too specialized, which may lead to inconsistent, biased or unreliable results. Therefore it is critical to build a comprehensive and unbiased benchmark. 2) Debiased probing approaches. Currently prompt-based probing is the dominant knowledge probing meth-

ods due to its simplicity. However, there still exist lots of issues in prompt-based probing. Therefore, the design of unbiased datasets and better probing frameworks is another useful direction worth investigating.

6 Knowledge Editing

Knowledge editing is the process which modifies the stored knowledge in pre-trained language models, either by replacing it with new knowledge (e.g., changing the current prime minister of the UK to Rishi Sunak) or by removing it entirely (e.g., some personal privacy information). There are two primary motivations for editing knowledge in language models: 1) even the state-of-the-art language models (e.g. ChatGPT²) could learn lots of incorrect knowledge; 2) many facts are time-sensitive, requiring regular updates to their corresponding knowledge.

Unfortunately, editing knowledge in PLMs poses significant challenges. Firstly, naive solutions such as retraining are often impractical due to the massive size of large-scale language models. Secondly, due to the black box and non-linear nature of PLMs, any minor modification might result in a significant undesirable change in model predictions. As a result, it can be challenging to precisely edit the target knowledge.

To promote the development of relevant studies, De Cao et al. (2021) formulate three desiderata for knowledge editing methods: 1) **Generality**: the method is able to edit the language models already pre-trained without the need for specialized re-training. 2) **Reliability**: the method is supposed to successfully edit knowledge required modification while not influencing the rest of knowledge in LMs. 3) **Consistency**: the modification should be consistent across paraphrases with equivalent semantics (e.g., Michael Jordan was born in [MASK]. v.s. The birthplace of Michael Jordan is [MASK].) and relevant knowledge required modification accordingly (e.g., Rishi Sunak becomes the prime minister of the UK. v.s. Liz Truss is not the primer minister of the UK.).

In this section, we divide current strategies for knowledge editing into four categories and the summary of comparisons between these approaches is shown in Table 2. In the following we will describe and discuss these methods.

6.1 Constrained Fine-tuning

The naive solution to edit knowledge in a PLM is to re-train it using the updated training dataset, but such a naive solution is computationally expensive and may be impractical because PLMs are involved. Therefore, a better solution is to fine-tune PLMs only on a small subset which only contains the target samples. However, such a method may suffer from catastrophic forgetting, and affects the rest knowledge which is not intended to be edited. Therefore, Zhu et al. (2020) propose to modify the knowledge in PLMs with constrained fine-tuning, specifically, they use an \mathcal{L}_2 or \mathcal{L}_∞ normalization to constrain the parameters change of models. Furthermore, they find that only fine-tuning the initial and final layers while keeping the rest of the model frozen yields better performance than fine-tuning the whole model. However, in deep neural networks like PLMs, even a minor change of the parameters could change the model’s predictions on a lot of samples. Therefore, such methods could potentially affect other knowledge stored in PLMs which is not required modification.

6.2 Memory-Based Editing

Instead of directly modifying parameters of PLMs, another natural solution is to maintain a knowledge cache which stores all new knowledge, and replace the original predictions when a input hits the cache. However, a symbolic knowledge cache may suffer from robustness issues, i.e., the inputs with the same meaning can differ in natural language expressions, therefore they may result in different predictions. To address this problem, Mitchell et al. (2022) propose a memory-based approach for knowledge editing. Specifically, the model contains five modules: an edit memory that stores the modified knowledge, a classifier, a counterfactual model, and the frozen original language model. Given an input, the classifier determines whether it hits a sample in the edit memory, and the counterfactual model’s prediction will overrule the original language model’s prediction if it hits a memory cache. This method is effective but does not actually edit the knowledge encoded in the parameters of language models, thus cannot benefit downstream tasks. Meanwhile, Dong et al. (2022) add additional trainable parameters in the feed-forward module of PLMs, which are trained on a modified knowledge dataset while the original parameters are frozen. They also demonstrate that the modified

²<https://openai.com/blog/chatgpt/>

Approach	Knowledge Support	Training Required	Online Edit	Batch Edit	Downstream Benefit	Unforeseen Side Effects
Constrained Tuning						
FTM (Zhu et al., 2020)	Factual	YES	NO	YES	Potential	YES
Memory-based						
SERAC (Mitchell et al., 2022)	Factual, QA	YES	YES	YES	NO	NO
MEM-PROMPT (Madaan et al., 2022)	Linguistic, Ethics	NO	YES	YES	Potential	Unlikely
CALINET (Dong et al., 2022)	Factual	YES	NO	YES	Potential	YES
Meta-learning						
KNOWEDITOR (De Cao et al., 2021)	Factual	YES	Possible	YES	Potential	YES
MEMD (Mitchell et al., 2021)	Factual	YES	NO	YES	Potential	YES
Locate and Edit						
Knowledge Neuron (Dai et al., 2022a)	Factual	NO	YES	NO	NO	YES
ROME (Meng et al., 2022)	Factual	NO	YES	NO	NO	Possible

Table 2: Comparisons between existing knowledge editing approaches. “Online Edit” refers to quickly editing an individual target knowledge. “Batch Edit” refers to editing a set of target knowledge simultaneously. “Downstream Benefit” refers to the potential for the modified knowledge to be utilized by the edited language model for downstream tasks. “Unforeseen Side Effects” refers to the impact of knowledge editing on the language model beyond the modification of target knowledge.

knowledge could benefit related QA tasks. Moreover, Madaan et al. (2022) introduce the users’ feedback for PLMs’ error correction. Specifically, they maintain a memory of models’ mistake and users’ feedback, which enhance the model to produce updated prompt and avoid similar mistakes.

6.3 Meta-learning-based Editing

Sinitsin et al. (2020) first propose editable training to conduct model editing based on meta-learning, which aims to train the model parameters to suit model editing. By constraining the training objective, the editing procedure could be accomplished under k gradient step while ensuring reliability, locality, and efficiency. However, such a method is not practical for pre-trained language models since it requires expensive specialized retraining. A different strategy is to utilize a hyper network, which uses one network to generate the weights of another network (Ha et al., 2017). De Cao et al. (2021); Hase et al. (2021) train a hyper-network to predict the parameter changes for each data point, with the constraint of editing target knowledge without affecting others. Although computationally efficient, Mitchell et al. (2021) argues that this method fails to edit very large models, and proposes model editor networks with gradient decomposition (MEND). Specifically, by decomposing the gradient of standard fine-tuning into a low-rank form, they could train multiple MLPs to generate local model parameter changes, without damaging models’ predictions on unrelated knowledge. Experiments show that MEND can be applied to large

pre-trained models for fast model editing. One limitation of existing meta-learning-based methods is that their robustness and generalization are still questionable, as they ensure locality by constraining the parameter space change or the predictions on specific datasets. In that case, the knowledge that requires no modifications or the knowledge that is related to edited knowledge but not paraphrasing could also be incorrect.

6.4 Locate and Edit

Based on the assumption that “knowledge is locally stored in PLMs”, the “locate and edit” strategy first locates the parameters corresponding to specific knowledge, and edit them by directly replacing with updated ones. This approach is also introduced in Section 4.1. Dai et al. (2022a) present a case study of factual knowledge editing in PLMs with corresponding knowledge neurons. By directly modifying the value of knowledge neurons, they achieve knowledge editing with a relatively low but nontrivial success rate. Although the editing procedure is straightforward once the corresponding knowledge neuron is located, this method has not proved its effectiveness on large-scale editing or the effects of unrelated knowledge. Similarly, Meng et al. (2022) first connect the knowledge required modification with a key-value pair in one of the middle MLP layers, and modify the corresponding knowledge by directly updating the key-value pair. Since these methods are based on the locality hypothesis of factual knowledge, which has not been widely confirmed yet, the changes in certain

parameters may affect irrelevant knowledge and lead to unexpected results.

6.5 Discussions and Future Works

To utilize pre-trained language models as a sustainable knowledge resource, the precise, effective, reliable and consistent knowledge editing is essential. However, as discussed above, all current editing methods have their own limitations. Therefore it is worthwhile to enhance current methods and develop new knowledge editing strategies.

In the future, several useful directions of knowledge editing may lie in: 1) **Broader range of target knowledge.** As shown in Table 2, current studies mostly focus on the editing of factual knowledge, which is relatively easy to formalize and evaluate. In the future, researchers could explore the editing methods towards other kinds of knowledge, and develop universal approaches which can edit all kinds of knowledge in the same way. 2) **Comprehensive evaluation.** Current most knowledge editing studies are evaluated using metrics such as editing success rate on target knowledge, predictions invariance rate on unrelated knowledge for assessing generality, and accuracy on paraphrases of target knowledge for assessing consistency. However, we find that these metrics are limited to comprehensively evaluate the knowledge editing capability of different approaches. For instance, most evaluations only sample unrelated knowledge from the same distribution of target knowledge. However, the influence of a knowledge edit could be much broader, e.g., affecting the performance on downstream tasks or the knowledge from other distributions and categories. In addition, as mentioned in Mitchell et al. (2021), most studies measure the consistency of samples generated through back translation, which ignores the knowledge affected by knowledge editing except the paraphrases, e.g., the country with the largest population would be affected by the population modification of the countries. Therefore, it is important to design comprehensive benchmark which can better assess the capabilities of editing strategies. 3) **More effective editing approaches.** Ideally, a knowledge editing approach should satisfy the desiderata of generality, reliability and consistency, and can handle large-scale and individual knowledge editing tasks with high efficiency. To this end, we may borrow ideas from other fields, such as meta-learning, continual learning, and life-long learning. Furthermore, it is

useful to connect knowledge editing studies with knowledge representation studies (§ 4).

7 Knowledge Application

Knowledge application studies how to effectively distill and leverage the knowledge in PLMs for other applications. Specifically, we divide knowledge applications into two categories: language models as knowledge bases and language models for downstream tasks, and in following we describe them in detail.

7.1 Language Models as Knowledge Bases

The impressive performance of large-scale pre-trained language models, as well as the potentially enormous amount of implicitly stored knowledge, raises extensive attention about using language models as an alternative to conventional structured knowledge bases (LMs-as-KBs) (Petroni et al., 2019; Heinzerling and Inui, 2021; Jiang et al., 2020b; Wang et al., 2020; Cao et al., 2021; Razniewski et al., 2021; AlKhamissi et al., 2022).

Unfortunately, along with the promising advantages and potentials compared with structured knowledge bases, there also exist intrinsic flaws for language models as knowledge base (Razniewski et al., 2021), which are summarized in Table 3. In following we describe them in detail.

Construction procedure is one of the biggest advantages of LMs-as-KBs compared with structured KBs. Constructing large-scale structured KBs such as Freebase (Bollacker et al., 2008) and Wikidata (Vrandečić and Krötzsch, 2014) often requires extremely complex pipelines (Petroni et al., 2019), e.g., ontology construction, knowledge acquisition, knowledge verification, knowledge fusion, knowledge storage, and knowledge population. Such a complex pipeline involves lots of NLP techniques, including ontology engineering, entity linking, entity recognition, relation extraction, entity matching and so on. And each technique requires corresponding expert knowledge, supervised data and human efforts. Moreover, due to the pipeline nature, error propagation is always a critical issue.

In contrast, the knowledge of language models can be easily learned from pure text using self-supervised learning, without any explicit supervision signal (§3.1). Furthermore, the construction procedure is end-to-end, therefore no ontology engineering, expert knowledge, or human annotations are needed.

Perspectives	Structured KB	LMs-as-KBs
Construction		
Ontology/Schema	Pre-defined	Open-ended 😊
Process	Pipeline	End-to-End 😊
Human Effort	Data annotation	Self-supervised 😊
Expert Knowledge	Common	Not required 😊
Coverage		
Domain	Constrained	Open 😊
Amount	Limited	Potential
Knowledge Fusing	Complex	Easy 😊
Interaction		
Query	Structured	Natural Language 😊
Prediction	Deterministic 😊	Probabilistic
Rejection	Yes 😊	Hard
Editing	Easy 😊	Limited
Reliability		
Ambiguity	Low 😊	High
Correctness	Relatively High 😊	Questionable
Current Practicality	Extensive 😊	Limited yet

Table 3: The comparisons between conventional structured knowledge bases and using language models as knowledge bases (LMs-as-KBs). Part of this table is inspired by Razniewski et al. (2021). The advantages are marked in bold. From the table, we can easily find that although LMs-as-KBs are more advantageous on construction and coverage, the critical current limitations of interaction and reliability significantly hinder its real-word applications, and far from substitution of structured knowledge bases.

Coverage is another big advantage of LMs-as-KBs. Traditional structured KBs are often limited by its pre-defined schemas, and the difficulty of acquiring knowledge further limits their coverage. In comparison, by directly representing knowledge in parameters, there is no schema limitations for LMs-as-KBs. And all knowledge is learned from un-annotated text corpus, therefore the knowledge coverage is mostly only determined by the coverage of pre-training corpus.

The above advantages make LMs-as-KBs an extremely attractive and promising idea. However, there are also some intrinsic flaws which hinder LMs from fully substituting structured KBs.

Interaction with structured KB and LMs-as-KBs are quite different. Structured KBs often use structural querying methods such as SPARQL (Pérez et al., 2009), e.g., querying the birthplace of Michael Jordan using <Michael Jordan, Birthplace, ?>. In the case of language model-based KBs, the queries are mostly natural language expressions such as “The birthplace of Michael Jordan is [MASK]”.

Compared with structural queries, natural language-based queries are more natural and

friendly for users. However, structured KBs can return deterministic answers (e.g., Brooklyn), but LM-based KBs can only generate candidates with different probabilities (e.g., <Brooklyn, 0.8>). The probabilistic predictions may be incorrect, inconsistent and confusing. Furthermore, structured KBs can identify the queries they cannot answer, but current LM-based KBs can hardly reject the queries it cannot answer, thus resulting in the knowledge hallucination problem. Concretely, if we query some knowledge that is not stored in a structured KB, the answer could be blank when no tuples are matched. However, no matter what we ask, language models will always “guess” the answers, even such knowledge is never learned by LMs. Although there are some naive solutions to this problem such as rejecting answers with a low probability, this is still an open problem currently.

Finally, it is difficult to edit knowledge in LM-based KBs, as discussed in Section 6. In comparison, it is easy to add, modify and delete knowledge in structured KBs.

Reliability is another concern for LMs-as-KBs. The first problem is ambiguity. In structured KBs, all entities and facts have their own IDs (e.g., Q89

for Apple the fruit and Q312 for Apple Inc. in Wikidata), therefore there is no ambiguity problem. However, in LM-based KBs, all pieces of knowledge are represented as natural language expressions and will therefore suffer from the ambiguity problem of natural language. For example, do "U.S.A" and "America" represent the same entity in a language model? Previous studies have observed that such verbalization requirements will result in prompt preference bias and instance verbalization bias in LMs-as-KBs (Cao et al., 2022). The consistency of predictions is another drawback of LMs-as-KBs, i.e., a LM-based KB may return different answers to the semantically equivalent queries.

7.2 Language Models for Downstream Tasks

Besides using language models as knowledge bases, the knowledge in PLMs can also benefit many downstream tasks in different ways. Fig. 3 shows three main paradigms and we describe them in detail.

7.2.1 Fine-tuning

Fine-tuning is a common way to leverage knowledge in language models, which learns to distill and leverage knowledge by further tuning PLMs using task-specific datasets. Firstly, implicitly learned knowledge from text has been recognized as one of the main reasons for PLMs' remarkable performance and strong generalization ability across so many NLP tasks (Manning et al., 2020; Wei et al., 2021b; Yang et al., 2021; Yin et al., 2022). Secondly, many studies have shown that injecting knowledge into language models can lead to better performance on downstream tasks. For instance, integrating entity knowledge into PLMs could improve the performance of a wide range of language understanding tasks Sun et al. (2019); Shen et al. (2020), infusing factual knowledge into PLMs could benefit their performance on tasks such as relation extraction, entity typing, etc. (Zhang et al., 2019; Wang et al., 2021b,a; Liu et al., 2020), and incorporating linguistic knowledge with PLMs could increase their performance on benchmarks such as GLUE (Levine et al., 2020; Sachan et al., 2021; Bai et al., 2021).

7.2.2 Prompt Learning

Prompt-based learning is another way to leverage the knowledge in PLMs for downstream tasks. For example, to classify the sentiment polarity of the

sentence "Best movie ever.", we can add a prompt and transform the input into "Best movie ever. It is ___.". And the polarity can be determined by comparing the PLMs' prediction probability between candidate answers "good" and "bad". By selecting appropriate prompts, PLMs have been shown competitive zero-shot performance on some downstream tasks without any supervised training (Radford et al., 2019a; Brown et al., 2020; Liu et al., 2021a).

Because handcraft prompts often suffer from unstable performance across different prompts and cannot utilize the information from supervised data, many prompt optimization approaches have been proposed to acquire better-performing prompts (Liu et al., 2021a), such as paraphrasing (Jiang et al., 2020b; Haviv et al., 2021), gradient-based search (Shin et al., 2020), model generation (Gao et al., 2021), knowledge enhanced (Hu et al., 2022), etc. Furthermore, prompt-tuning, which adds some trainable vectors to the inputs as continuous prompts, while keeping the parameters of LMs freezing, has achieved competitive performance with fine-tuning (Li and Liang, 2021b; Liu et al., 2021b; Hambardzumyan et al., 2021; Lester et al., 2021). In addition to optimizing single prompts, ensembling (Jiang et al., 2020b; Qin and Eisner, 2021), compositing (Han et al., 2021), or decoupling (Ozturkler et al., 2022) multiple prompts could also improve model performance. Moreover, prompt has also been applied to data augmentation (Schick and Schütze, 2021), domain adaptation (Ben-David et al., 2021), debiasing (Schick et al., 2021) and so on.

More recently, instruction-tuning, which pre-trains LMs on a wide range of datasets given the natural language description of tasks as instructions, has achieved significant performance and generalization ability improvements of language models (Wei et al., 2021a; Sanh et al., 2021; Ouyang et al., 2022; Chung et al., 2022).

7.2.3 In-context Learning

Applications Currently the parameters of PLMs have been scaled to 175B (e.g., GPT-3 (Brown et al., 2020), OPT (Zhang et al., 2022b), BLOOM³) or even larger (e.g., PaLM (Chowdhery et al., 2022)), making the computational expense of fine-tuning and prompt-tuning infeasible for most researchers. Therefore, tuning-free in-context learn-

³<https://huggingface.co/bigscience/bloom>

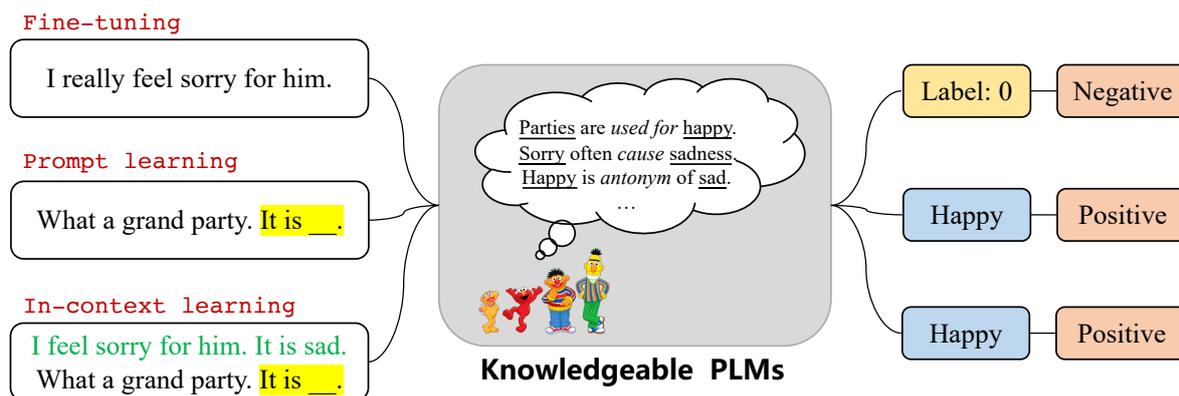


Figure 3: The primary paradigms that apply the knowledge in PLMs to downstream tasks.

ing has become one of the most popular approaches to apply the knowledge in large-scale PLMs in downstream tasks (Dong et al., 2023). For instance, for the sentiment classification task, in-context learning will first sample several demonstrations, such as (what a horrible meal, negative), and combine them with the original query. In this way, the input becomes “What a horrible meal. It is bad. [SEP] Best movie ever. It is ___.” The provided demonstrations offer extra information about the task and enable PLMs to utilize the analogy ability to predict the correct answer. In-context learning has achieved good performance on lots of downstream tasks such as language understanding (Brown et al., 2020; Zhao et al., 2021; Lee et al., 2022; Eisenstein et al., 2022; Zhang et al., 2022a), data generation (Li et al., 2022b; Dai et al., 2022b; Yu et al., 2022), or reasoning (Wei et al., 2022; Lampinen et al., 2022; Zhou et al., 2022).

Bias Problem One drawback of in-context learning is the bias problem, i.e., the performance is sensitive to demonstration selections, demonstration orders, label distribution of demonstrations and prompt selection, etc. (Zhao et al., 2021; Lu et al., 2022; Liu et al., 2022). Therefore, to achieve better performance of in-context learning, Zhao et al. (2021) first propose to estimate the biases by feeding the model with an uninformative input (e.g., [MASK] or N/A), and then calibrate the prediction probabilities uniformly distributed for eliminating the models’ bias towards specific answers. For demonstration selection, Gao et al. (2021); Liu et al. (2022) propose to select demonstrations that are semantically close to the input query. Rubin et al. (2022) train a dense retriever on LM-scored datasets to select demonstrations. Su et al. (2022)

introduce a graph-based selection method to ensure the demonstration’s diversity and representativeness. For demonstration sort, Lu et al. (2022) first construct a development dataset by sampling from language models, and then use entropy-based metrics to determine the optimal demonstration permutation. For prompt selection, Gao et al. (2021) use a language model to generate candidate prompts and select ones with better performance on the development set.

Mechanism Although in-context learning has been widely applied on various downstream tasks, its underlying mechanism is still unclear. Reynolds and McDonnell (2021) find that zero-shot prompting sometimes can significantly outperform in-context learning, and argue that the additional demonstrations do not help PLMs to learn a new task, but rather locate the task they have already learned. Cao et al. (2021) investigate the in-context learning for knowledge probing, and find that the demonstrations can only provide type-level guidance but factual information. Min et al. (2022) find that randomly replacing the demonstrations’ labels hardly affects the performance, and show that the effectiveness of in-context learning relies more on the label space and input distribution restriction provided by demonstrations rather than the precise input-label mapping. Chan et al. (2022) find that only when the data includes both burstiness and large-scale of rarely occurring classes, in-context learning capability can emerge in transformer model. von Oswald et al. (2022) investigate the connections between in-context learning and gradient descent, and demonstrate the similarity between in-context learning and the gradient-based few-shot learning.

7.3 Discussions and Future Works

Leveraging knowledge in PLMs is both promising and challenging. On the one hand, it is obvious that the large amount of implicit knowledge stored in PLMs will benefit different downstream tasks. On the other hand, all current application paradigms have their own limitations. For instance, the consistency and reliability of LMs-as-KBs hinder PLMs to replace structured KBs. Moreover, fine-tuning, prompt learning and in-context learning methods often suffer from catastrophic forgetting, computational cost, inconsistent and unstable predictions, social bias, etc.

To address these challenges, several main future directions of knowledge application may lie in the following: 1) For LMs-as-KBs, we need to propose specific pre-training approaches to address current shortcomings in consistency and reliability. 2) For LMs for downstream tasks, we suggest explore more application strategies, such as new tuning-free methods to address the computational cost issue and black-box tuning (Sun et al., 2022) methods to tune pre-trained language models without access to their parameters.

8 Conclusions

In this survey, we conduct a comprehensive review about the life circle of knowledge in pre-trained language models, including knowledge acquisition, knowledge representation, knowledge probing, knowledge editing and knowledge application. We systematically review related studies for each period, discuss the advantages and limitations of different methods, summarize the main challenge, and present some future directions. We believe this survey will benefit researchers in many areas such as language models, knowledge graph, knowledge base, etc.

Acknowledgments

We sincerely thank all anonymous reviewers for their insightful comments and valuable suggestions. This research work is supported by the National Natural Science Foundation of China under Grants no. 62122077 and CAS Project for Young Scientists in Basic Research under Grant No. YSBR-040.

References

- Alessandro Achille, Matteo Rovere, and Stefano Soatto. 2019. [Critical learning periods in deep networks](#). In *Proc. of ICLR*. OpenReview.net.
- Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. 2022. [A review on language models as knowledge bases](#). *ArXiv preprint*, abs/2204.06031.
- Jiangang Bai, Yujing Wang, Yiren Chen, Yaming Yang, Jing Bai, Jing Yu, and Yunhai Tong. 2021. [Syntax-BERT: Improving pre-trained transformers with syntax trees](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3011–3020, Online. Association for Computational Linguistics.
- Livio Baldini Soares, Nicholas FitzGerald, Jeffrey Ling, and Tom Kwiatkowski. 2019. [Matching the blanks: Distributional similarity for relation learning](#). In *Proc. of ACL*, pages 2895–2905, Florence, Italy. Association for Computational Linguistics.
- Pratyay Banerjee and Chitta Baral. 2020. [Self-supervised knowledge triplet learning for zero-shot question answering](#). In *Proc. of EMNLP*, pages 151–162, Online. Association for Computational Linguistics.
- Yonatan Belinkov. 2022. [Probing classifiers: Promises, shortcomings, and advances](#). *Computational Linguistics*, 48(1):207–219.
- Eyal Ben-David, Nadav Oved, and Roi Reichart. 2021. [Pada: A prompt-based autoregressive approach for adaptation to unseen domains](#). *ArXiv preprint*, abs/2102.12206.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. [Freebase: A collaboratively created graph database for structuring human knowledge](#). In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, SIGMOD '08, page 1247–1250, New York, NY, USA. Association for Computing Machinery.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. [COMET: Commonsense transformers for automatic knowledge graph construction](#). In *Proc. of ACL*, pages 4762–4779, Florence, Italy. Association for Computational Linguistics.
- Zied Bouraoui, José Camacho-Collados, and Steven Schockaert. 2020. [Inducing relational knowledge from BERT](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7456–7463. AAAI Press.

- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Boxi Cao, Hongyu Lin, Xianpei Han, Fangchao Liu, and Le Sun. 2022. [Can prompt probe pretrained language models? understanding the invisible risks from a causal view](#). In *Proc. of ACL*, pages 5796–5808, Dublin, Ireland. Association for Computational Linguistics.
- Boxi Cao, Hongyu Lin, Xianpei Han, Le Sun, Lingyong Yan, Meng Liao, Tong Xue, and Jin Xu. 2021. [Knowledgeable or educated guess? revisiting language models as knowledge bases](#). In *Proc. of ACL*, pages 1860–1874, Online. Association for Computational Linguistics.
- Stephanie C. Y. Chan, Adam Santoro, Andrew K. Lampinen, Jane X. Wang, Aaditya Singh, Pierre H. Richemond, Jay McClelland, and Felix Hill. 2022. [Data distributional properties drive emergent in-context learning in transformers](#). *CoRR*, abs/2205.05055.
- Cheng-Han Chiang, Sung-Feng Huang, and Hung-yi Lee. 2020. [Pretrained language model embryology: The birth of ALBERT](#). In *Proc. of EMNLP*, pages 6813–6828, Online. Association for Computational Linguistics.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. [Palm: Scaling language modeling with pathways](#). *ArXiv preprint*, abs/2204.02311.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. [Scaling instruction-finetuned language models](#). *ArXiv preprint*, abs/2210.11416.
- Patricia S Churchland and Terrence J Sejnowski. 1988. Perspectives on cognitive neuroscience. *Science*, 242(4879):741–745.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022a. [Knowledge neurons in pretrained transformers](#). In *Proc. of ACL*, pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.
- Zhuyun Dai, Vincent Y Zhao, Ji Ma, Yi Luan, Jianmo Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith B Hall, and Ming-Wei Chang. 2022b. [Promptagator: Few-shot dense retrieval from 8 examples](#). *ArXiv preprint*, abs/2209.11755.
- Joe Davison, Joshua Feldman, and Alexander Rush. 2019. [Commonsense knowledge mining from pretrained models](#). In *Proc. of EMNLP*, pages 1173–1178, Hong Kong, China. Association for Computational Linguistics.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. [Editing factual knowledge in language models](#). In *Proc. of EMNLP*, pages 6491–6506, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proc. of NAACL-HLT*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu, Zhifang Sui, and Lei Li. 2022. [Calibrating factual knowledge in pretrained language models](#). *ArXiv preprint*, abs/2210.03329.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2023. [A survey for in-context learning](#). *ArXiv preprint*, abs/2301.00234.
- Jacob Eisenstein, Daniel Andor, Bernd Bohnet, Michael Collins, and David Mimno. 2022. [Honest students from untrusted teachers: Learning an interpretable question-answering pipeline from a pretrained language model](#). *ArXiv preprint*, abs/2210.02498.
- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Amir Feder, Abhilasha Ravichander, Marius Mosbach, Yonatan Belinkov, Hinrich Schütze, and Yoav Goldberg. 2022. [Measuring causal effects of data statistics on language model’sfactual’predictions](#). *ArXiv preprint*, abs/2207.14251.
- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. [Measuring and improving consistency in pretrained language models](#). *Transactions of the Association for Computational Linguistics*, 9:1012–1031.

- Allyson Ettinger. 2020. [What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models](#). *Transactions of the Association for Computational Linguistics*, 8:34–48.
- Amir Feder, Katherine A Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E Roberts, et al. 2021. [Causal inference in natural language processing: Estimation, prediction, interpretation and beyond](#). *ArXiv preprint*, abs/2109.00725.
- Thibault Févry, Livio Baldini Soares, Nicholas FitzGerald, Eunsol Choi, and Tom Kwiatkowski. 2020. [Entities as experts: Sparse memory access with entity supervision](#). In *Proc. of EMNLP*, pages 4937–4951, Online. Association for Computational Linguistics.
- Matthew Finlayson, Aaron Mueller, Sebastian Gehrmann, Stuart Shieber, Tal Linzen, and Yonatan Belinkov. 2021. [Causal analysis of syntactic agreement mechanisms in neural language models](#). In *Proc. of ACL*, pages 1828–1843, Online. Association for Computational Linguistics.
- Maxwell Forbes, Ari Holtzman, and Yejin Choi. 2019. [Do Neural Language Representations Learn Physical Commonsense?](#) *ArXiv preprint*, abs/1908.02899.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. [Making pre-trained language models better few-shot learners](#). In *Proc. of ACL*, pages 3816–3830, Online. Association for Computational Linguistics.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. [Transformer feed-forward layers are key-value memories](#). In *Proc. of EMNLP*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Daniel Gillick, Sayali Kulkarni, Larry Lansing, Alessandro Presta, Jason Baldridge, Eugene Ie, and Diego Garcia-Olano. 2019. [Learning dense representations for entity retrieval](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 528–537, Hong Kong, China. Association for Computational Linguistics.
- Yoav Goldberg. 2019. [Assessing BERT’s Syntactic Abilities](#). *ArXiv preprint*, abs/1901.05287.
- Jian Guan, Fei Huang, Zhihao Zhao, Xiaoyan Zhu, and Minlie Huang. 2020. [A knowledge-enhanced pre-training model for commonsense story generation](#). *Transactions of the Association for Computational Linguistics*, 8:93–108.
- Abhijeet Gupta, Gemma Boleda, Marco Baroni, and Sebastian Padó. 2015. [Distributional vectors encode referential attributes](#). In *Proc. of EMNLP*, pages 12–21, Lisbon, Portugal. Association for Computational Linguistics.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. [REALM: retrieval-augmented language model pre-training](#). *ArXiv preprint*, abs/2002.08909.
- David Ha, Andrew M. Dai, and Quoc V. Le. 2017. [Hypernetworks](#). In *Proc. of ICLR*. OpenReview.net.
- Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. 2021. [WARP: Word-level Adversarial ReProgramming](#). In *Proc. of ACL*, pages 4921–4933, Online. Association for Computational Linguistics.
- Xu Han, Weilin Zhao, Ning Ding, Zhiyuan Liu, and Maosong Sun. 2021. [Ptr: Prompt tuning with rules for text classification](#). *ArXiv preprint*, abs/2105.11259.
- Moritz Hardt, Eric Price, and Nati Srebro. 2016. [Equality of opportunity in supervised learning](#). In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 3315–3323.
- Peter Hase, Mona Diab, Asli Celikyilmaz, Xian Li, Zornitsa Kozareva, Veselin Stoyanov, Mohit Bansal, and Srinivasan Iyer. 2021. [Do language models have beliefs? methods for detecting, updating, and visualizing model beliefs](#). *ArXiv preprint*, abs/2111.13654.
- Adi Haviv, Jonathan Berant, and Amir Globerson. 2021. [BERTese: Learning to speak to BERT](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3618–3623, Online. Association for Computational Linguistics.
- Benjamin Heinzerling and Kentaro Inui. 2021. [Language models as knowledge bases: On entity representations, storage capacity, and paraphrased queries](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1772–1791, Online. Association for Computational Linguistics.
- John Hewitt and Christopher D. Manning. 2019. [A structural probe for finding syntax in word representations](#). In *Proc. of NAACL-HLT*, pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Phu Mon Htut, Jason Phang, Shikha Bordia, and Samuel R. Bowman. 2019. [Do Attention Heads in BERT Track Syntactic Dependencies?](#) *ArXiv preprint*, abs/1911.12246.
- Shengding Hu, Ning Ding, Huadong Wang, Zhiyuan Liu, Jingang Wang, Juanzi Li, Wei Wu, and

- Maosong Sun. 2022. [Knowledgeable prompting: Incorporating knowledge into prompt verbalizer for text classification](#). In *Proc. of ACL*, pages 2225–2240, Dublin, Ireland. Association for Computational Linguistics.
- Joel Jang, Seonghyeon Ye, and Minjoon Seo. 2022. [Can large language models truly understand prompts? a case study with negated prompts](#). *ArXiv preprint*, abs/2209.12711.
- Zhengbao Jiang, Antonios Anastasopoulos, Jun Araki, Haibo Ding, and Graham Neubig. 2020a. [X-FACTR: Multilingual factual knowledge retrieval from pretrained language models](#). In *Proc. of EMNLP*, pages 5943–5959, Online. Association for Computational Linguistics.
- Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020b. [How can we know what language models know?](#) *Transactions of the Association for Computational Linguistics*, 8:423–438.
- Jeevesh Juneja and Ritu Agarwal. 2022. [Finding patterns in knowledge attribution for transformers](#). *ArXiv preprint*, abs/2205.01366.
- Nora Kassner, Philipp Dufter, and Hinrich Schütze. 2021. [Multilingual LAMA: Investigating knowledge in multilingual pretrained language models](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3250–3258, Online. Association for Computational Linguistics.
- Nora Kassner and Hinrich Schütze. 2020. [Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly](#). In *Proc. of ACL*, pages 7811–7818, Online. Association for Computational Linguistics.
- Pei Ke, Haozhe Ji, Siyang Liu, Xiaoyan Zhu, and Minlie Huang. 2020. [SentiLARE: Sentiment-aware language representation learning with linguistic knowledge](#). In *Proc. of EMNLP*, pages 6975–6988, Online. Association for Computational Linguistics.
- Niki Kilbertus, Mateo Rojas-Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. 2017. [Avoiding discrimination through causal reasoning](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 656–666.
- Arne Köhn. 2015. [What’s in an embedding? analyzing word embeddings through multilingual evaluation](#). In *Proc. of EMNLP*, pages 2067–2073, Lisbon, Portugal. Association for Computational Linguistics.
- Matt J. Kusner, Joshua R. Loftus, Chris Russell, and Ricardo Silva. 2017. [Counterfactual fairness](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4066–4076.
- Andrew K Lampinen, Ishita Dasgupta, Stephanie CY Chan, Kory Matthewson, Michael Henry Tessler, Antonia Creswell, James L McClelland, Jane X Wang, and Felix Hill. 2022. [Can language models learn from explanations in context?](#) *ArXiv preprint*, abs/2204.02329.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. [ALBERT: A lite BERT for self-supervised learning of language representations](#). In *Proc. of ICLR*. OpenReview.net.
- Anne Lauscher, Ivan Vulić, Edoardo Maria Ponti, Anna Korhonen, and Goran Glavaš. 2020. [Specializing unsupervised pretraining models for word-level semantic similarity](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1371–1383, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Dong-Ho Lee, Akshen Kadakia, Kangmin Tan, Mahak Agarwal, Xinyu Feng, Takashi Shibuya, Ryosuke Mitani, Toshiyuki Sekiya, Jay Pujara, and Xiang Ren. 2022. [Good examples make a faster learner: Simple demonstration-based learning for low-resource NER](#). In *Proc. of ACL*, pages 2687–2700, Dublin, Ireland. Association for Computational Linguistics.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proc. of EMNLP*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yoav Levine, Barak Lenz, Or Dagan, Ori Ram, Dan Padnos, Or Sharir, Shai Shalev-Shwartz, Amnon Shashua, and Yoav Shoham. 2020. [SenseBERT: Driving some sense into BERT](#). In *Proc. of ACL*, pages 4656–4667, Online. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proc. of ACL*, pages 7871–7880, Online. Association for Computational Linguistics.
- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020b. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Shaobo Li, Xiaoguang Li, Lifeng Shang, Zhenhua Dong, Chengjie Sun, Bingquan Liu, Zhenzhou Ji,

- Xin Jiang, and Qun Liu. 2022a. [How pre-trained language models capture factual knowledge? a causal-inspired analysis.](#) In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1720–1732, Dublin, Ireland. Association for Computational Linguistics.
- Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, et al. 2022b. [Explanations from large language models make small reasoners better.](#) *ArXiv preprint*, abs/2210.06726.
- Xiang Lisa Li and Percy Liang. 2021a. [Prefix-tuning: Optimizing continuous prompts for generation.](#) In *Proc. of ACL*, pages 4582–4597, Online. Association for Computational Linguistics.
- Xiang Lisa Li and Percy Liang. 2021b. [Prefix-tuning: Optimizing continuous prompts for generation.](#) In *Proc. of ACL*, pages 4582–4597, Online. Association for Computational Linguistics.
- Bill Yuchen Lin, Seyeon Lee, Rahul Khanna, and Xiang Ren. 2020a. [Birds have four legs?! NumerSense: Probing Numerical Commonsense Knowledge of Pre-Trained Language Models.](#) In *Proc. of EMNLP*, pages 6862–6868, Online. Association for Computational Linguistics.
- Bill Yuchen Lin, Seyeon Lee, Rahul Khanna, and Xiang Ren. 2020b. [Birds have four legs?! NumerSense: Probing Numerical Commonsense Knowledge of Pre-Trained Language Models.](#) In *Proc. of EMNLP*, pages 6862–6868, Online. Association for Computational Linguistics.
- Yongjie Lin, Yi Chern Tan, and Robert Frank. 2019. [Open sesame: Getting inside BERT’s linguistic knowledge.](#) In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 241–253, Florence, Italy. Association for Computational Linguistics.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. [What makes good in-context examples for GPT-3?](#) In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.
- Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew E. Peters, and Noah A. Smith. 2019a. [Linguistic knowledge and transferability of contextual representations.](#) In *Proc. of NAACL-HLT*, pages 1073–1094, Minneapolis, Minnesota. Association for Computational Linguistics.
- Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew E. Peters, and Noah A. Smith. 2019b. [Linguistic knowledge and transferability of contextual representations.](#) In *Proc. of NAACL-HLT*, pages 1073–1094, Minneapolis, Minnesota. Association for Computational Linguistics.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021a. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing.](#) *ArXiv preprint*, abs/2107.13586.
- Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. 2020. [K-BERT: enabling language representation with knowledge graph.](#) In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 2901–2908. AAAI Press.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021b. [Gpt understands, too.](#) *ArXiv preprint*, abs/2103.10385.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019c. [Roberta: A robustly optimized bert pretraining approach.](#) *ArXiv preprint*, abs/1907.11692.
- Zeyu Liu, Yizhong Wang, Jungo Kasai, Hannaneh Hajishirzi, and Noah A. Smith. 2021c. [Probing across time: What does RoBERTa know and when?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 820–842, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Lajanugen Logeswaran, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, Jacob Devlin, and Honglak Lee. 2019. [Zero-shot entity linking by reading entity descriptions.](#) In *Proc. of ACL*, pages 3449–3460, Florence, Italy. Association for Computational Linguistics.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. [Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity.](#) In *Proc. of ACL*, pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.
- Kaixin Ma, Filip Ilievski, Jonathan Francis, Yonatan Bisk, Eric Nyberg, and Alessandro Oltramari. 2021. [Knowledge-driven data construction for zero-shot evaluation in commonsense question answering.](#) In *Proc. of AAAI*, volume 35, pages 13507–13515.
- Aman Madaan, Niket Tandon, Peter Clark, and Yiming Yang. 2022. [Memory-assisted prompt editing to improve gpt-3 after deployment.](#) *ArXiv preprint*, abs/2201.06009.
- Christopher D Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. 2020. [Emergent linguistic structure in artificial neural networks trained by self-supervision.](#) *Proceedings of the National Academy of Sciences*, 117(48):30046–30054.

- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. [Locating and editing factual knowledge in gpt](#). *ArXiv preprint*, abs/2202.05262.
- George A. Miller. 1992. [WordNet: A lexical database for English](#). In *Speech and Natural Language: Proceedings of a Workshop Held at Harriman, New York, February 23-26, 1992*.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. [Rethinking the role of demonstrations: What makes in-context learning work?](#) *ArXiv preprint*, abs/2202.12837.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2021. [Fast model editing at scale](#). *ArXiv preprint*, abs/2110.11309.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memory-based model editing at scale. In *International Conference on Machine Learning*, pages 15817–15831. PMLR.
- Roberto Navigli and Simone Paolo Ponzetto. 2010. [BabelNet: Building a very large multilingual semantic network](#). In *Proc. of ACL*, pages 216–225, Uppsala, Sweden. Association for Computational Linguistics.
- Nils J. Nilson. 1974. Artificial intelligence. In *Information Processing*, pages 778–801.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. [Training language models to follow instructions with human feedback](#). *ArXiv preprint*, abs/2203.02155.
- Batu Ozturkler, Nikolay Malkin, Zhen Wang, and Nebojsa Jojic. 2022. [Thinksum: Probabilistic reasoning over sets using large language models](#). *ArXiv preprint*, abs/2210.01293.
- Jorge Pérez, Marcelo Arenas, and Claudio Gutierrez. 2009. Semantics and complexity of sparql. *ACM Transactions on Database Systems (TODS)*, 34(3):1–45.
- Laura Pérez-Mayos, Miguel Ballesteros, and Leo Wanner. 2021. [How much pretraining data do language models need to learn syntax?](#) In *Proc. of EMNLP*, pages 1571–1582, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. 2019. [Knowledge enhanced contextual word representations](#). In *Proc. of EMNLP*, 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 *Proc. of EMNLP*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Nina Poerner, Ulli Waltinger, and Hinrich Schütze. 2020. [E-BERT: Efficient-yet-effective entity embeddings for BERT](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 803–818, Online. Association for Computational Linguistics.
- Guanghui Qin and Jason Eisner. 2021. [Learning how to ask: Querying LMs with mixtures of soft prompts](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5203–5212, Online. Association for Computational Linguistics.
- Yujia Qin, Yankai Lin, Ryuichi Takanobu, Zhiyuan Liu, Peng Li, Heng Ji, Minlie Huang, Maosong Sun, and Jie Zhou. 2021. [ERICA: Improving entity and relation understanding for pre-trained language models via contrastive learning](#). In *Proc. of ACL*, pages 3350–3363, Online. Association for Computational Linguistics.
- Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. 2020. [Pre-trained models for natural language processing: A survey](#). *ArXiv preprint*, abs/2003.08271.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019a. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019b. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67.
- Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. 2017. [SVCCA: singular vector canonical correlation analysis for deep learning dynamics and interpretability](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 6076–6085.
- Simon Razniewski, Andrew Yates, Nora Kassner, and Gerhard Weikum. 2021. [Language models as or for knowledge bases](#). *ArXiv preprint*, abs/2110.04888.

- Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–7.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In *Proc. of EMNLP*, pages 5418–5426, Online. Association for Computational Linguistics.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866.
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2671, Seattle, United States. Association for Computational Linguistics.
- Devendra Sachan, Yuhao Zhang, Peng Qi, and William L. Hamilton. 2021. Do syntax trees help pre-trained transformers extract information? In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2647–2661, Online. Association for Computational Linguistics.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. *ArXiv preprint*, abs/2110.08207.
- Naomi Saphra and Adam Lopez. 2019. Understanding learning dynamics of language models with SVCCA. In *Proc. of NAACL-HLT*, pages 3257–3267, Minneapolis, Minnesota. Association for Computational Linguistics.
- Naomi Saphra and Adam Lopez. 2020. LSTMs compose—and Learn—Bottom-up. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2797–2809, Online. Association for Computational Linguistics.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *ArXiv preprint*, abs/2211.05100.
- Timo Schick and Hinrich Schütze. 2021. Exploiting cloze-questions for few-shot text classification and natural language inference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, Online. Association for Computational Linguistics.
- Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021. Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in NLP. *Transactions of the Association for Computational Linguistics*, 9:1408–1424.
- August Th Schreiber, Guus Schreiber, Hans Akkermans, Anjo Anjewierden, Nigel Shadbolt, Robert de Hoog, Walter Van de Velde, and Bob Wielinga. 2000. *Knowledge engineering and management: the CommonKADS methodology*. MIT press.
- Tao Shen, Yi Mao, Pengcheng He, Guodong Long, Adam Trischler, and Weizhu Chen. 2020. Exploiting structured knowledge in text via graph-guided representation learning. In *Proc. of EMNLP*, pages 8980–8994, Online. Association for Computational Linguistics.
- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. Auto-Prompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proc. of EMNLP*, pages 4222–4235, Online. Association for Computational Linguistics.
- Vered Shwartz, Peter West, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Unsupervised commonsense question answering with self-talk. In *Proc. of EMNLP*, pages 4615–4629, Online. Association for Computational Linguistics.
- Anton Sinitin, Vsevolod Plokhotnyuk, Dmitriy Pyrkov, Sergei Popov, and Artem Babenko. 2020. Editable neural networks. In *Proc. of ICLR*. OpenReview.net.
- Jian Song, Di Liang, Rumei Li, Yuntao Li, Sirui Wang, Minlong Peng, Wei Wu, and Yongxin Yu. 2022. Improving semantic matching through dependency-enhanced pre-trained model with adaptive fusion. *ArXiv preprint*, abs/2210.08471.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. MASS: masked sequence to sequence pre-training for language generation. In *Proc. of ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages 5926–5936. PMLR.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *ArXiv preprint*, abs/2206.04615.
- Rudi Studer, V Richard Benjamins, and Dieter Fensel. 1998. Knowledge engineering: principles and methods. *Data & knowledge engineering*, 25(1-2):161–197.
- Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, et al. 2022. Selective annotation makes language models better few-shot learners. *ArXiv preprint*, abs/2209.01975.

- Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. 2022. Black-box tuning for language-model-as-a-service. In *International Conference on Machine Learning*, pages 20841–20855. PMLR.
- Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. 2019. [Ernie: Enhanced representation through knowledge integration](#). *ArXiv preprint*, abs/1904.09223.
- Mujeen Sung, Jinhyuk Lee, Sean Yi, Minji Jeon, Sungdong Kim, and Jaewoo Kang. 2021. [Can language models be biomedical knowledge bases?](#) In *Proc. of EMNLP*, pages 4723–4734, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. 2020. [oLMpics-on what language model pre-training captures](#). *Transactions of the Association for Computational Linguistics*, 8:743–758.
- Alexandre Tamborrino, Nicola Pellicanò, Baptiste Pannier, Pascal Voitot, and Louise Naudin. 2020. [Pre-training is \(almost\) all you need: An application to commonsense reasoning](#). In *Proc. of ACL*, pages 3878–3887, Online. Association for Computational Linguistics.
- 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 *Proc. of ICLR*. OpenReview.net.
- Hao Tian, Can Gao, Xinyan Xiao, Hao Liu, Bolei He, Hua Wu, Haifeng Wang, and Feng Wu. 2020. [SKEP: Sentiment knowledge enhanced pre-training for sentiment analysis](#). In *Proc. of ACL*, pages 4067–4076, Online. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4–9, 2017, Long Beach, CA, USA*, pages 5998–6008.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart M. Shieber. 2020. [Investigating gender bias in language models using causal mediation analysis](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*.
- Elena Voita and Ivan Titov. 2020. [Information-theoretic probing with minimum description length](#). In *Proc. of EMNLP*, pages 183–196, Online. Association for Computational Linguistics.
- Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. 2022. [Transformers learn in-context by gradient descent](#). *ArXiv preprint*, abs/2212.07677.
- Denny Vrandečić and Markus Krötzsch. 2014. [Wiki-data: A free collaborative knowledgebase](#). *Commun. ACM*, 57(10):78–85.
- Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh, and Matt Gardner. 2019. [Do NLP models know numbers? probing numeracy in embeddings](#). In *Proc. of EMNLP*, pages 5307–5315, Hong Kong, China. Association for Computational Linguistics.
- Jonas Wallat, Jaspreet Singh, and Avishek Anand. 2020. [BERTnesia: Investigating the capture and forgetting of knowledge in BERT](#). In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 174–183, Online. Association for Computational Linguistics.
- Chenguang Wang, Xiao Liu, and Dawn Song. 2020. [Language models are open knowledge graphs](#). *ArXiv preprint*, abs/2010.11967.
- Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Jianshu Ji, Guihong Cao, Daxin Jiang, and Ming Zhou. 2021a. [K-Adapter: Infusing Knowledge into Pre-Trained Models with Adapters](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1405–1418, Online. Association for Computational Linguistics.
- Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. 2021b. [KEPLER: A unified model for knowledge embedding and pre-trained language representation](#). *Transactions of the Association for Computational Linguistics*, 9:176–194.
- Alex Warstadt, Yu Cao, Ioana Grosu, Wei Peng, Hagen Blix, Yining Nie, Anna Alsop, Shikha Bordia, Haokun Liu, Alicia Parrish, Sheng-Fu Wang, Jason Phang, Anhad Mohananey, Phu Mon Htut, Paloma Jeretic, and Samuel R. Bowman. 2019. [Investigating BERT’s knowledge of language: Five analysis methods with NPIs](#). In *Proc. of EMNLP*, pages 2877–2887, Hong Kong, China. Association for Computational Linguistics.
- Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. [BLiMP: The benchmark of linguistic minimal pairs for English](#). *Transactions of the Association for Computational Linguistics*, 8:377–392.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021a. [Finetuned language models are zero-shot learners](#). *ArXiv preprint*, abs/2109.01652.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. [Chain of thought prompting elicits reasoning in large language models](#). *ArXiv preprint*, abs/2201.11903.
- Xiaokai Wei, Shen Wang, Dejiao Zhang, Parminder Bhatia, and Andrew Arnold. 2021b. [Knowledge enhanced pretrained language models: A comprehensive survey](#). *ArXiv preprint*, abs/2110.08455.
- Zhiyong Wu, Yun Chen, Ben Kao, and Qun Liu. 2020. [Perturbed masking: Parameter-free probing for analyzing and interpreting BERT](#). In *Proc. of ACL*, pages 4166–4176, Online. Association for Computational Linguistics.
- Wenhan Xiong, Jingfei Du, William Yang Wang, and Veselin Stoyanov. 2020. [Pretrained encyclopedia: Weakly supervised knowledge-pretrained language model](#). In *Proc. of ICLR*. OpenReview.net.
- Yadollah Yaghoobzadeh, Katharina Kann, T. J. Hazen, Eneko Agirre, and Hinrich Schütze. 2019. [Probing for semantic classes: Diagnosing the meaning content of word embeddings](#). In *Proc. of ACL*, pages 5740–5753, Florence, Italy. Association for Computational Linguistics.
- Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. [LUKE: Deep contextualized entity representations with entity-aware self-attention](#). In *Proc. of EMNLP*, pages 6442–6454, Online. Association for Computational Linguistics.
- Jian Yang, Gang Xiao, Yulong Shen, Wei Jiang, Xinyu Hu, Ying Zhang, and Jinghui Peng. 2021. [A survey of knowledge enhanced pre-trained models](#). *ArXiv preprint*, abs/2110.00269.
- Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Rich James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2022. [Retrieval-augmented multimodal language modeling](#). *CoRR*, abs/2211.12561.
- Zhi-Xiu Ye, Qian Chen, Wen Wang, and Zhen-Hua Ling. 2019. [Align, mask and select: A simple method for incorporating commonsense knowledge into language representation models](#). *ArXiv preprint*, abs/1908.06725.
- Da Yin, Li Dong, Hao Cheng, Xiaodong Liu, Kai-Wei Chang, Furu Wei, and Jianfeng Gao. 2022. [A survey of knowledge-intensive nlp with pre-trained language models](#). *ArXiv preprint*, abs/2202.08772.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2022. [Generate rather than retrieve: Large language models are strong context generators](#). *ArXiv preprint*, abs/2209.10063.
- Hongxin Zhang, Yanzhe Zhang, Ruiyi Zhang, and Diyi Yang. 2022a. [Robustness of demonstration-based learning under limited data scenario](#). *ArXiv preprint*, abs/2210.10693.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022b. [Opt: Open pre-trained transformer language models](#).
- Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. [ERNIE: Enhanced language representation with informative entities](#). In *Proc. of ACL*, pages 1441–1451, Florence, Italy. Association for Computational Linguistics.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. [Calibrate before use: Improving few-shot performance of language models](#). In *Proc. of ICML*, volume 139 of *Proceedings of Machine Learning Research*, pages 12697–12706. PMLR.
- Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021. [Factual probing is \[MASK\]: Learning vs. learning to recall](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5017–5033, Online. Association for Computational Linguistics.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Olivier Bousquet, Quoc Le, and Ed Chi. 2022. [Least-to-most prompting enables complex reasoning in large language models](#). *ArXiv preprint*, abs/2205.10625.
- Junru Zhou, Zhuosheng Zhang, Hai Zhao, and Shuailiang Zhang. 2019. [Limit-bert: Linguistic informed multi-task bert](#). *ArXiv preprint*, abs/1910.14296.
- Xuhui Zhou, Yue Zhang, Leyang Cui, and Dandan Huang. 2020a. [Evaluating commonsense in pre-trained language models](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 9733–9740. AAAI Press.
- Xuhui Zhou, Yue Zhang, Leyang Cui, and Dandan Huang. 2020b. [Evaluating commonsense in pre-trained language models](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 9733–9740. AAAI Press.

Yichu Zhou and Vivek Srikumar. 2021a. [DirectProbe: Studying representations without classifiers](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5070–5083, Online. Association for Computational Linguistics.

Yichu Zhou and Vivek Srikumar. 2021b. [DirectProbe: Studying representations without classifiers](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5070–5083, Online. Association for Computational Linguistics.

Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv Kumar. 2020. [Modifying memories in transformer models](#). *ArXiv preprint*, abs/2012.00363.

Philip G Zimbardo and Floyd L Ruch. 1975. Psychology and life.