# Give Us the Facts: Enhancing Large Language Models with Knowledge Graphs for Fact-aware Language Modeling

Linyao Yang, Hongyang Chen, *Senior Member, IEEE*, Zhao Li, Xiao Ding, Xindong Wu, *Fellow, IEEE*

*Abstract*—Recently, ChatGPT, a representative large language model (LLM), has gained considerable attention. Due to their powerful emergent abilities, recent LLMs are considered as a possible alternative to structured knowledge bases like knowledge graphs (KGs). However, while LLMs are proficient at learning probabilistic language patterns and engaging in conversations with humans, they, like previous smaller pre-trained language models (PLMs), still have difficulty in recalling facts while generating knowledge-grounded contents. To overcome these limitations, researchers have proposed enhancing data-driven PLMs with knowledge-based KGs to incorporate explicit factual knowledge into PLMs, thus improving their performance in generating texts requiring factual knowledge and providing more informed responses to user queries. This paper reviews the studies on enhancing PLMs with KGs, detailing existing knowledge graph enhanced pre-trained language models (KGPLMs) as well as their applications. Inspired by existing studies on KGPLM, this paper proposes enhancing LLMs with KGs by developing knowledge graph-enhanced large language models (KGLLMs). KGLLM provides a solution to enhance LLMs' factual reasoning ability, opening up new avenues for LLM research.

*Index Terms*—Large language model, Knowledge graph, Chat-GPT, Knowledge reasoning, Knowledge management.

## I. INTRODUCTION

IN recent years, the rapid development of big data [1]–[3] and high-speed computing has led to the emergence of pre-trained language models (PLMs). Plenty of PLMs, such as BERT [4], GPT [5], and T5 [6], have been proposed, which greatly improve the performance of various natural language processing (NLP) tasks. Recently, researchers have found that scaling model size or data size can improve model capacities on downstream tasks. Moreover, they found that when the parameter size exceeds a certain scale [7], these PLMs exhibit some surprising emergent abilities. Emergent abilities refer to the abilities that are not present in small models but arise in large models [7], which are utilized to distinguish large language models (LLMs) from PLMs.

On November 30, 2022, a chatbot program named ChatGPT was released by OpenAI, which is developed based on the

LLM GPT-3.5. By fine-tuning GPT with supervised learning and further optimizing the model using reinforcement learning from human feedback (RLHF), ChatGPT is capable of engaging in continuous conversation with humans based on chat context. It can even complete complex tasks such as coding and paper writing, showcasing its powerful emergent abilities [7]. Consequently, some researchers [8]–[11] explored whether LLMs can serve as parameterized knowledge bases to replace structured knowledge bases like knowledge graphs (KGs), as they also store a substantial amount of facts.

However, existing studies [12]–[15] have found that LLMs' ability to generate factually correct text is still limited. They are capable of remembering facts only during training. Consequently, these models often face challenges when attempting to recall relevant knowledge and apply the correct knowledge to generate knowledge grounded contents. On the other hand, as artificially constructed structured knowledge bases, KGs store a vast amount of knowledge closely related to real-world facts in a readable format. They explicitly express relationships between entities and intuitively display the overall structure of knowledge and reasoning chains, making them an ideal choice for knowledge modeling. As a result, there exists not only a competitive but also a complementary relationship between LLMs and KGs. LLMs have the ability to enhance knowledge extraction accuracy and improve the quality of KGs [16], while KGs can utilize explicit knowledge to guide the training of LLMs, improving their ability to recall and apply knowledge.

So far, numerous methods have been proposed for strengthening PLMs with KGs, which can be categorized into three types: before-training enhancement, during-training enhancement, and post-training enhancement. Although there exist a few surveys [17]–[19] of knowledge-enhanced PLMs, they focus on various forms of knowledge, lacking a systematic review of knowledge graph enhanced pre-trained language model (KGPLM) methods. For instance, Wei *et al.* [17] conducted a review of knowledge enhanced PLMs based on diverse knowledge sources but only covered a small set of KGPLMs. Similarly, Yang *et al.* [18] covered various forms of knowledge enhanced PLMs but provided only a partial review of KGPLMs without technical categorization. In another study, Zhen *et al.* [19] categorized knowledge enhanced PLMs into implicit incorporation and explicit incorporation methods, yet their review encompassed only a small subset of KGPLMs. Moreover, this field is rapidly evolving with numerous new technologies consistently being introduced. Therefore, to address questions of whether constructing KGs
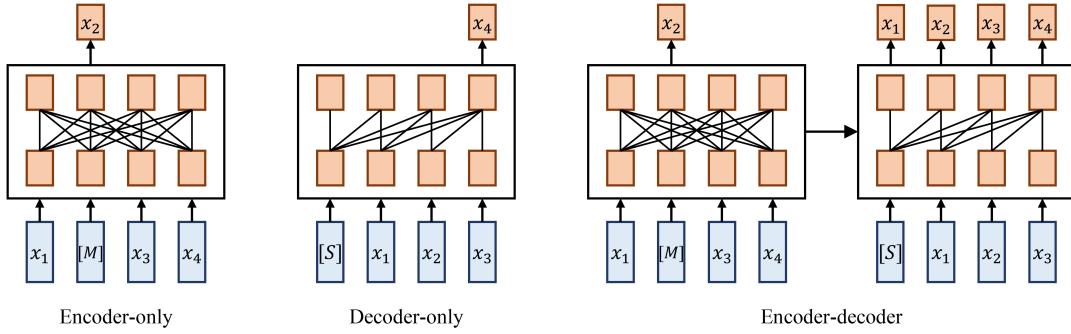
Fig. 1. Main frameworks of existing PLMs, in which $x_i$ is the $i$-th token of the input sentence, $[M]$ represents the masked token and $[S]$ is the start token.

is still necessary and how to improve the knowledge modeling ability of LLMs, we present a systematic review of relevant studies. We conducted a thorough search for papers related to the keywords "language model" and "knowledge graph". Subsequently, the papers that were most relevant to KGPLM were carefully refined and categorized. In comparison with existing surveys, this paper specifically concentrates on KGPLM and covers a broader range of up-to-date papers. Furthermore, we suggest the development of knowledge graph enhanced large language models (KGLLMs) to tackle the knowledge modeling challenge in LLMs. The main contributions of this paper are summarized as follows:

- We provide a comprehensive review for KGPLMs, which helps researchers to gain a deep insight of this field.
- We overview research on the evaluation of LLMs and draw comparisons between LLMs and KGs.
- We propose to enhance LLMs with KGs and suggest some possible future research directions, which may benefit researchers in the field of LLM.

The remainder of this paper is organized as follows. Section II overviews the background of LLMs. Section III categorizes the existing methods for KGPLMs and introduces representatives from each group. Section IV introduces the applications of KGPLMs. Section V discusses whether LLMs can replace KGs with the evidence from existing studies. Section VI proposes to enhance LLMs' ability to learn factual knowledge by developing KGLLMs and presents some future research directions. Section VII draws the conclusions.

## II. BACKGROUND

PLMs learn dense and continuous representations for words, addressing the issue of feature sparsity encountered in traditional encoding methods and significantly improving performance across various NLP tasks. Consequently, PLM-based methods have gained prominence, leading to the development of various types of PLMs. Recently, PLMs have been scaled to LLMs in order to achieve even better performance. In this section, we provide a comprehensive background of PLMs and offer an overview of their historical development.

### A. Background of PLMs

PLMs are a type of language model obtained through unsupervised learning [20] on a large corpus. They are capable of capturing the structure and characteristics of a language and generating universal representations for words. Following pre-training, PLMs can be fine-tuned for specific downstream tasks like text summarization, text classification, and text generation.

The model frameworks used by existing PLMs can be classified into three categories, as illustrated in Fig. 1: encoder-only, decoder-only, and encoder-decoder [21]. The encoder-only framework utilizes a bidirectional transformer to recover masked tokens based on the input sentences, which effectively utilizes contextual information to learn better text representations. More specifically, given an input token sequence $\mathcal{L} = (x_1, ..., x_T)$ with a few masked tokens $\mathcal{M}$, it models the likelihood of the masked tokens as $p(x) = \sum_{x_t \in \mathcal{M}} p(x_t | x_{\mathcal{Q}})$. However, due to the lack of a decoder, it cannot be directly applied to text generation tasks. BERT and its improved models mostly adopt the encoder-only framework. The decoder-only framework leverages a unidirectional transformer to predict tokens in an autoregressive fashion, making it suitable for text generation tasks. That is, given the text sequence $\mathcal{C} = (x_1, ..., x_T)$, this framework models the likelihood of the input token sequence as $p(x) = \prod_{t=1}^{T} p(x_t | x_{<t})$. GPT series and their improved models mostly adopt this framework. Nevertheless, compared with the other two frameworks, the decoder-only framework cannot make use of contextual information and cannot generalize well to other tasks. The encoder-decoder framework constructs a sequence-to-sequence model to predict the current token based on historical context with masked tokens. Its objective can be described as $\sum_{t=1}^{T} p(x_t | x_{<t, \mathcal{Q}})$. This framework excels at tasks that require generating output based on given inputs, yet its encoding and decoding speed is slow compared to the other two frameworks.

Multiple pre-training tasks for PLMs have been designed, which can be categorized into word-level, phrase-level, and sentence-level tasks. Typical word-level pre-training tasks include masked language modeling (MLM) [4] and replaced token detection (RTD) [22]. MLM randomly masks some tokens in the input sequence and trains PLMs to reconstruct the masked tokens based on context, whose loss function is:

$$\mathcal{L}_{\mathrm{MLM}} = - \sum_{x \in \mathcal{M}} \log p(x | x_{\mathcal{Q}}). \tag{1}$$

It can promote the learning of contextual information, thereby achieving better results in language understanding and language modeling tasks. RTD operates similarly to MLM but
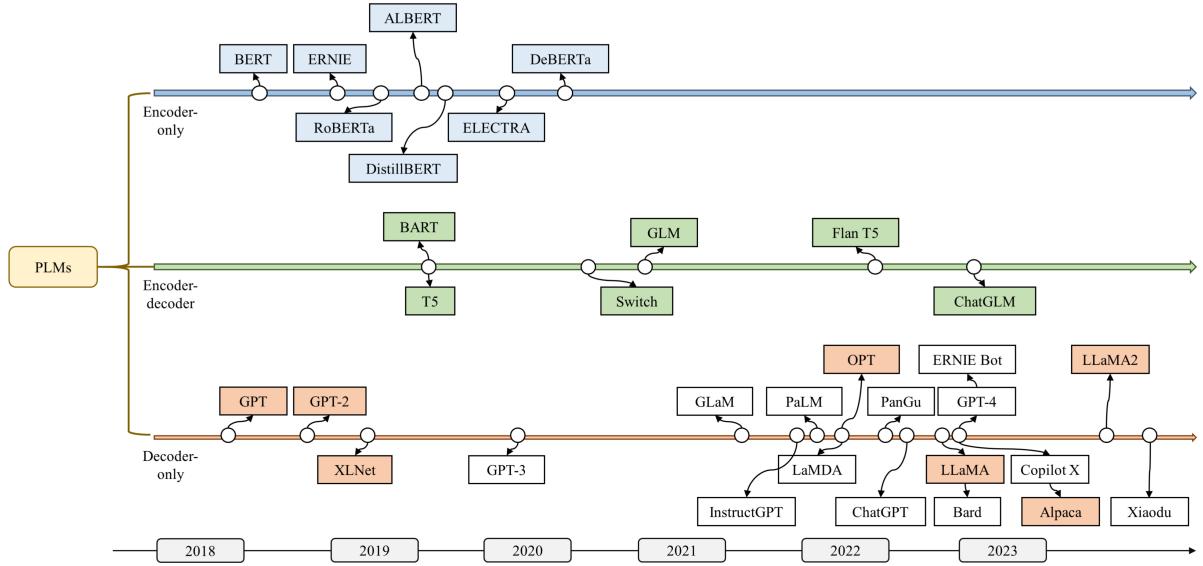
Fig. 2. Milestones of LLMs. Open-source models are represented by solid squares, while closed-source models are represented by hollow squares.

introduces greater randomness by substituting some tokens with alternative ones and training the model to predict the original tokens, whose loss function is defined as:

$$\mathcal{L}_{\text{RTD}} = -\sum_{t=1}^{T} \log p(y_t|\tilde{x}). \tag{2}$$

Here, $\tilde{x}$ is the corrupted token of $x$, while $y_t$ is 1 if $\tilde{x}_t = x_t$ and 0 otherwise. Compared with MLM, RTD can reflect changes in vocabulary in real texts more realistically and enable PLMs to handle unknown and misspelled words. The representative of phrase-level pre-training tasks is span boundary objective (SBO) [23], [24], which forces PLMs to predict each token of a masked span solely relying on the representations of the visible tokens at the boundaries, enhancing the syntactic structure analysis ability of PLMs and improving their performance in named entity recognition and sentiment analysis. The training objective of the SBO task can be expressed as:

$$\mathcal{L}_{\text{SBO}} = -\sum_{t=1}^{T} \log p(x_i|y_i), \tag{3}$$

where $y_i$ is token $x_i$'s representation in the span. Representatives of sentence-level pre-training tasks include next sentence prediction (NSP) [4] and sentence order prediction (SOP) [25]. NSP trains PLMs to distinguish whether two given sentences are continuous, thereby improving PLMs' performance in context-based tasks such as natural language inference and text classification. Similarly, SOP trains PLMs to determine the order of two randomly sampled and disrupted sentences, which improves their ability to capture sentence order information. The training objective of NSP and SOP is as follows:

$$\mathcal{L}_{\text{NSP/SOP}} = -\log p(y|s_1, s_2), \tag{4}$$

where $y = 1$ if $s_1$ and $s_2$ are two consecutive segments extracted from the corpus. Other tasks like deleted token detection (DTD), text infilling, sentence reordering (SR), and document reordering (DR) are also utilized by some PLMs [26], which improve their performance in some special tasks.

*B. Milestones*

As an early attempt, Elmo [27] employs a bidirectional long short term memory (LSTM) network to learn word representations capturing context. The model is trained with a bidirectional autoregressive language modeling objective, which involves maximizing the following log-likelihood:

$$\sum_{k=1}^{T} \Big( \log p \left( x_t \mid x_1, \ldots, x_{t-1}; \Theta_x, \overrightarrow{\Theta}_{LSTM} \right) \\ + \log p \left( x_t \mid x_{t+1}, \ldots, x_T; \Theta_x, \overleftarrow{\Theta}_{LSTM} \right) \Big), \tag{5}$$

where $p$ models the probability of token $x_t$ given the history context $(x_1, \ldots, x_{t-1})$ or the future context $(x_{t+1}, \ldots, x_T)$. $\Theta_x$ denotes the token representation. $\overrightarrow{\Theta}_{LSTM}$ and $\overleftarrow{\Theta}_{LSTM}$ denote the LSTM encoder in the forward direction and the backward direction, respectively. By learning context-aware word representations, Elmo largely raises the performance bar of NLP tasks. However, its feature extraction ability is limited since LSTM is difficult to handle long sequences. With the emergence of the highly parallelizable Transformer [28], more powerful contextualized PLMs have been developed. Notable PLMs with different frameworks are shown in Fig. 2.

Transformer employs a self-attention mechanism to capture the dependence among input sequences, allowing for parallel processing of tokens and improving efficiency. Specifically, the output from the self-attention mechanism is:

$$\mathbf{h} = softmax(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}})\mathbf{V}, \tag{6}$$

where $\mathbf{Q}$, $\mathbf{K}$, and $\mathbf{V}$ are the query matrix, key matrix, and value matrix. $d_k$ is the dimension of the key and query vectors.

Encoder-only PLMs utilize bidirectional Transformer as encoder and employ MLM and NSP tasks for self-supervised
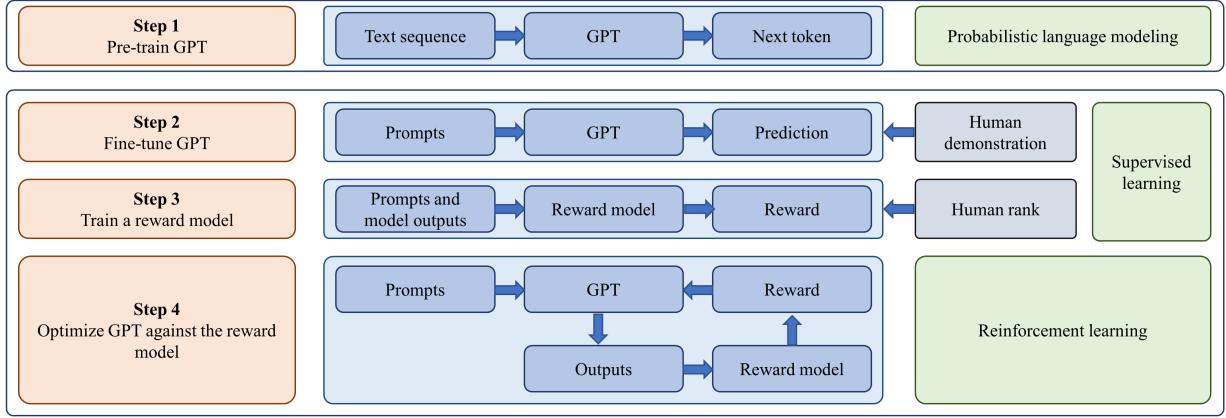
Fig. 3. The implementation process of ChatGPT.

training. RoBERTa [29] introduces a set of design choices and training strategies that lead to better performance, significantly enhancing BERT's performance on various benchmarks. DistilBERT [30] incorporates knowledge distillation into pre-training, which reduces the size of a BERT model by 40%. Other notable encoder-only PLMs include ERNIE [31], AL-BERT [25], ELECTRA [22], and DeBERTa [32].

In contrast, in decoder-only PLMs, a unidirectional Transformer is utilized as decoder, and the model is trained to predict the next token based on the preceding sequence. This training approach improves their language understanding and text generation abilities. Given an unsupervised corpus, GPT uses a unidirectional language modeling objective to optimize the model, maximizing the following log-likelihood:

$$\sum_i log p(x_t|x_{t-k}, ..., x_{t-1}; \Theta). \tag{7}$$

Here, $\Theta$ represents the parameters of the Transformer model. GPT-2 [33] improves upon GPT by increasing its model size and training corpus and enabling the model to automatically recognize task types for unsupervised training. XLNet [34] proposes a generalized autoregressive pretraining method, which enables learning bidirectional contexts.

In encoder-decoder PLMs, Transformer serves as both encoder and decoder. The encoder generates the latent representations for the input sequence, while the decoder generates the target output text. T5 [6] develops a unified framework that converts all NLP tasks into a text-to-text format, leading to exceptional performance on numerous benchmarks. In order to efficiently pre-train sequence-to-sequence models, BART [26] adopts a standard neural machine translation architecture and develops a denoising autoencoder.

### C. Scaling PLMs to LLMs

With the emergence of more and more PLMs, it has been revealed that model scaling can lead to improved performance. By increasing the parameter scale and data scale to a large enough size, it was found that these enlarged models exhibit some special abilities that do not possess by small-scale PLMs. Therefore, recent efforts have been devoted to scaling PLMs to LLMs to empower them with emergent abilities.

Typically, LLMs refer to PLMs that consist of hundreds of billions of parameters, such as GLM [35], Switch [36], Flan T5 [37], and ChatGLM [38] of the encoder-decoder framework. Besides, most existing LLMs adopt the decoder-only framework. Notable examples of decoder-only LLMs include GPT-3 [39], GLaM [40], InstructGPT [41], PaLM [42], LaMDA [43], OPT [44], LLaMA [45], Alpaca [46], GPT-4 [47], and LLaMA2 [48]. GPT-3 [39] further increases GPT-2's parameters and its training data size, and adopts zero-shot learning and diversity generation technologies, making it possible to learn and execute new tasks without annotated data and generate texts with diverse styles. GPT-3.5 not only increases the model size but also applies novel pre-training methods such as prompt-based extraction of templates (PET), which further improves the accuracy and fluency of generated texts. LLMs have stronger abilities to understand natural language and solve complex NLP tasks than smaller PLMs. GPT-3, for instance, exhibits a remarkable in-context learning ability. It can generate expected outputs for test cases by filling in the word sequence of the input text, relying solely on natural language instructions or demonstrations, without the need for additional training. Conversely, GPT-2 lacks this ability [49].

The most remarkable application of LLMs is ChatGPT, which adapts GPT-3.5 for dialogue and demonstrates an amazing conversation ability. The implementation process of ChatGPT is shown in Fig. 3 [50]. It first trains GPT on a large-scale corpus and then fine-tunes it on a dataset of labeler demonstrations. After that, it optimizes the model using RLHF [51], which trains a reward model to learn from direct feedback provided by human evaluators and optimizes the GPT model by formulating it as a reinforcement learning problem. In this setting, the pre-trained GPT model serves as the policy model that takes small pieces of prompts [52] as inputs and returns output texts. The GPT policy model is then optimized using the proximal policy optimization (PPO) algorithm [53] against the reward model. Based on the RLHF method, ChatGPT enables GPT to follow the expected instructions of humans and reduces the generation of toxic, biased, and harmful content. Besides, ChatGPT adopts the chain-of-thought strategy [54] and is additionally trained on code data, enabling it to solve tasks that require intermediate logical steps.

TABLE I
COMPARISON OF DIFFERENT PLMs

| Model framework | PLM | Year | Base model | Pre-training tasks | Pre-training data size | model size |
|---|---|---|---|---|---|---|
| Encoder-only | BERT | 2018 | Transformer | MLM, NSP | 3300M words | 340M |
| | ERNIE | 2019 | Transformer | MLM, NSP | 4500M subwords | 114M |
| | RoBERTa | 2019 | BERT | MLM | 160GB of text | 335M |
| | ALBERT | 2019 | BERT | SOP | 16GB of text | 233M |
| | DistillBERT | 2019 | BERT | MLM | 3300M words | 66M |
| | ELECTRA | 2020 | Transformer | RTD | 126GB of text | 110M |
| | DeBERTa | 2020 | Transformer | MLM | 78GB of text | 1.5B |
| Encoder-decoder | BART | 2019 | Transformer | MLM, DTD, text infilling, SR, DR | 160GB of text | 406M |
| | T5 | 2019 | Transformer | MLM | 20TB of text | 11B |
| | Switch | 2021 | Transformer | MLM | 180B tokens | 1.6T |
| | GLM | 2021 | Transformer | Blank infilling | 400B tokens | 130B |
| | Flan T5 | 2022 | T5 | 1800 fine-tuning tasks | - | 11B |
| | ChatGLM | 2023 | GLM | Blank infilling | 1T tokens | 6B |
| Decoder-only | GPT | 2018 | Transformer | Autoregressive language modeling | 800M words | 117M |
| | GPT-2 | 2019 | Transformer | Autoregressive language modeling | 40GB of text | 1.5B |
| | XLNet | 2019 | Transformer | Autoregressive language modeling | 33B tokens | 340M |
| | GPT-3 | 2020 | Transformer | Autoregressive language modeling | 45TB of text | 175B |
| | GLaM | 2021 | Transformer | Autoregressive language modeling | 1.6T tokens | 1.2T |
| | InstructGPT | 2022 | GPT-3 | Autoregressive language modeling | - | 175B |
| | PaLM | 2022 | Transformer | Autoregressive language modeling | 780B tokens | 540B |
| | LaMDA | 2022 | Transformer | Autoregressive language modeling | 768B tokens | 137B |
| | OPT | 2022 | Transformer | Autoregressive language modeling | 180B tokens | 175B |
| | ChatGPT | 2022 | GPT-3.5 | Autoregressive language modeling | - | - |
| | LLaMA | 2023 | Transformer | Autoregressive language modeling | 1.4T tokens | 65B |
| | GPT-4 | 2023 | Transformer | Autoregressive language modeling | 13T tokens | 1.8T |
| | Alpaca | 2023 | LLaMA | Autoregressive language modeling | 52K data | 7B |
| | LLaMA2 | 2023 | Transformer | Autoregressive language modeling | 2T tokens | 70B |

Another notable advancement is GPT-4 [47], a model that extends text input to multimodal signals and exhibits greater proficiency at solving tasks [55]. Furthermore, GPT-4 has undergone six months of iterative alignment, adding an additional safety reward in the RLHF training, which has made it more adept at generating helpful, honest, and harmless content. Additionally, GPT-4 implements some enhanced optimization methods, such as predictable scaling that accurately predicts GPT-4's final performance from smaller models trained with less computation.

Table I summarizes the characteristics of the above context-based PLMs and LLMs. As observed, the parameter size of the largest model has increased year by year.

### D. Pros and Cons of LLMs

A proliferation of benchmarks and tasks has been leveraged to evaluate the effectiveness and superiority of LLMs. Results from corresponding experiments demonstrate that LLMs achieve much better performance than previous deep learning models and smaller PLMs on a variety of NLP tasks. Besides, LLMs exhibit some emergent abilities and are capable of solving some complex tasks that traditional models and smaller PLMs cannot address. In summary, LLMs have the following superior characteristics.

**Zero-shot Learning.** LLMs outperform other models with zero-shot learning on most tasks and even perform better than fine-tuned models on some tasks. An empirical study [15] has shown that ChatGPT outperforms previous models with zero-shot learning on 9 of 13 datasets and even outperforms fully fine-tuned task-specific models on 4 tasks. This superior performance is attributed to the rich and diverse input data as well as the large parameter scale of LLMs, which allow them to capture the underlying patterns of natural language with high fidelity, leading to more robust and accurate inferences.

**In-context Learning.** In-context learning (ICL) is a paradigm that allows LLMs to learn tasks from only a few instances in the form of demonstration [56]. ICL was exhibited for the first time by GPT-3, which has become a common approach to use LLMs. ICL employs a formatted natural language prompt, which includes a description of the task and a handful of examples to illustrate the way to accomplish it. The ICL ability also benefits from the strong sequence processing ability and the rich knowledge reserve of LLMs.

**Step-by-step Reasoning.** By utilizing the chain-of-thought prompting strategy, LLMs can successfully complete some complex tasks, including arithmetic reasoning, commonsense reasoning, and symbolic reasoning. Such tasks are typically beyond the capability of smaller PLMs. The chain-of-thought is an improved prompting strategy, which integrates intermediate reasoning steps into the prompts to boost the performance of LLMs on complex reasoning tasks. Besides, the step-by-step reasoning ability is believed to be potentially acquired through training LLMs on well-structured code data [54].

**Instruction Following.** Instruction tuning is a unique fine-tuning approach that fine-tunes LLMs on a collection of natural language formatted instances. With this approach, LLMs are enabled to perform well on previously unseen tasks described through natural language instructions without relying on explicit examples [49]. For example, Wei et al. [57] fine-tuned a 137B parameter LLM on over 60 datasets based on instruction tuning and tested it on unseen task types. The experimental results demonstrated that the instruction-tuned model significantly outperformed its unmodified counterpart and zero-shot GPT-3.

**Human Alignment.** LLMs can be trained to generate high-quality, harmless responses that align with human values through the technique of RLHF, which involves incorporating humans into the training loop using carefully designed labeling strategies. RLHF comprises three steps: 1) collecting a labeled dataset consisting of input prompts and target outputs to fine-tune LLMs in a supervised way; 2) training a reward model on the assembled data, and 3) optimizing LLMs by formulating its optimization as a reinforcement learning problem. With this approach, LLMs are enabled to generate appropriate outputs that adhere to human expectations.

**Tools Manipulation.** Traditional PLMs are trained on plain text data, which limits their ability to solve non-textual tasks. Besides, their abilities are limited by the pre-training corpus, and cannot effectively solve tasks requiring real-time knowledge. In response to these limitations, recent LLMs are developed with the ability to manipulate external tools such as search engine, calculator, and compiler to enhance their performance in specialized domains [58]. More recently, the plugin mechanism has been supported in LLMs, providing an avenue for implementing novel functions. This mechanism has significantly broadened the range of capacities for LLMs, making them more flexible and adaptable to diverse tasks.

Although LLMs have made significant progress in natural language understanding and human-like content generation, they still have the following limitations and challenges [49].

**Unstructured Generation.** LLMs commonly rely on natural language prompts or instructions to generate text under specific conditions. This mechanism presents challenges for precisely constraining the generated outputs according to fine-grained or structural criteria. Ensuring specific text structures, such as the logical order of concepts throughout the entire text, can be difficult. This difficulty is amplified for tasks requiring formal rules or grammar. This is because LLMs mainly focus on the local context information of words and sentences during pre-training, while ignoring global syntactic and structural knowledge. A proposal for addressing this problem is to adopt an iterative prompting approach in generating text [59], mimicking the process of human writing. In contrast, KGs offer a structured summary and emphasize the correlation of relevant concepts when complex events involving the same entity extend across multiple sentences [60], thus enhancing the process of structured text generation.

**Hallucination.** When generating factual or knowledge-grounded texts, LLMs may produce content that contradicts existing sources or lack supporting evidence. This challenge widely occurs in existing LLMs and is known as the problem of hallucination, which results in a drop in their performance and poses risks when deploying them for real-world applications. The cause of this issue is related to LLMs' limited ability to utilize correct internal and external knowledge during task-solving. To alleviate this problem, existing studies have resorted to alignment tuning strategies, which incorporate human feedback to fine-tune LLMs. KGs provide structured and explicit representations of knowledge, which can be dynamically incorporated to augment LLMs, resulting in more factual rationales and reduced hallucination in generation [61].

**Inconsistency.** With the help of the chain-of-thought strategy, LLMs are capable of solving some complex reasoning tasks based on step-by-step reasoning. Despite their superior performance, LLMs may at times arrive at the desired answer based on an invalid reasoning path or produce an incorrect answer despite following a correct reasoning process. As a result, inconsistency arises between the derived answer and the underlying reasoning process. Additionally, research [62] has revealed that LLMs' abilities to forecast facts and answer queries are highly influenced by specific prompt templates and related entities. This is because that LLMs rely largely on simple heuristics to make predictions, their generations are correlated with co-occurrence frequencies between the target word and words in the prompt. Moreover, although LLMs' pre-training process helps them memorize facts, it fails to imbue them with the ability to generalize observed facts, leading to poor inferences. This issue can be partially addressed by introducing external KGs in LLM reasoning. By interactively exploring related entities and relations on KGs and performing reasoning based on the retrieved knowledge, LLMs can have better ability of knowledge traceability and knowledge correctability [63].

**Limited Reasoning Ability.** LLMs have demonstrated decent performance on some basic logical reasoning tasks when provided with question-answer examples. However, they exhibit poor performance on tasks that require the ability to comprehend and utilize supporting evidence for deriving conclusions. While LLMs typically generate valid reasoning steps, they face challenges when multiple candidate steps are deemed valid [64]. This results from LLMs being primed to solely choose the answer with the highest word overlapping with the input question. Additionally, LLMs struggle with predicting entity relationships due to their emphasis on shallow co-occurrence and sequence patterns of words. Moreover, despite exhibiting some basic numerical and symbolic reasoning abilities [65], LLMs face difficulties in numerical computation, especially for symbols infrequently encountered during pre-training. KGs explicitly capture the relations among concepts, which are essential for reasoning and can be utilized to enhance LLMs with structural reasoning capabilities. Previous studies have demonstrated that the integration of textual semantics and structural reasoning yields significant enhancement in the reasoning ability of LLMs [66], [67].

**Insufficient Domain Knowledge.** Because of the limited availability of domain-specific corpus, LLMs may not perform as well on domain-specific tasks as on general ones. For instance, while such models generally capture frequent patterns from general texts, generating medical reports, which involve numerous technical terms, may pose a great challenge for LLMs. This limitation suggests that during pre-training, it is difficult for LLMs to acquire sufficient domain knowledge, and injecting additional specialized knowledge may come at the cost of losing previously learned information, given the issue of catastrophic forgetting. Therefore, developing effective techniques for knowledge injection is of critical importance to enhance the performance of LLMs on specialized domains. Domain KGs are effective and standardized knowledge bases for specific domains, offering a feasible source for unified domain knowledge. For example, Ding *et al.* [68] proposed
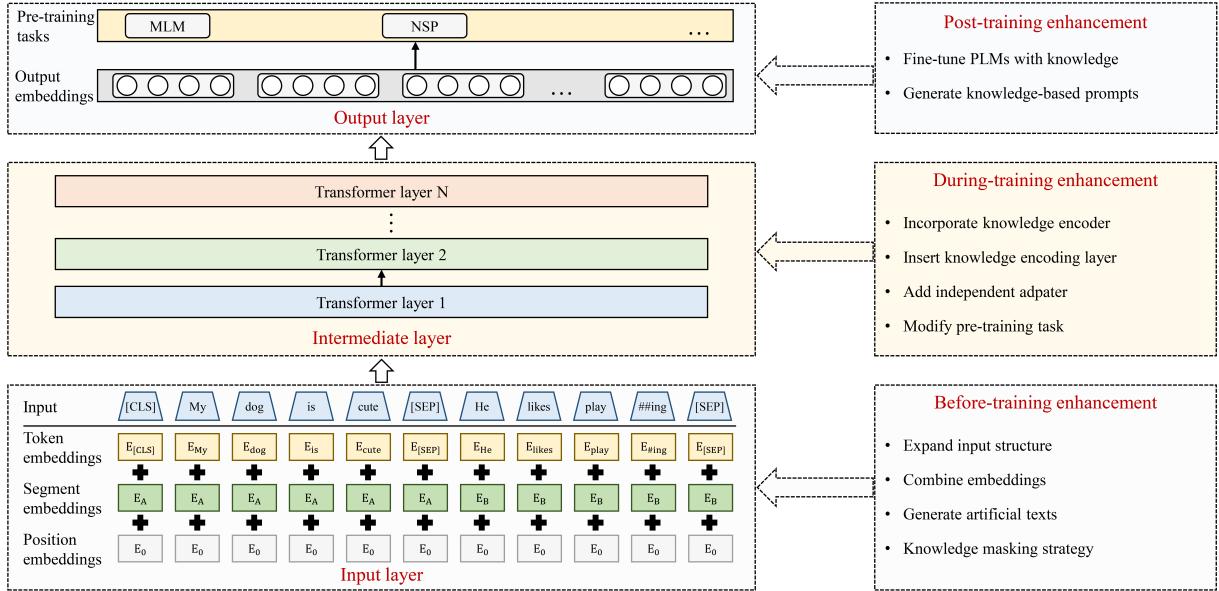
Fig. 4. Three types of KGPLMs according to the stage of knowledge graph participating in pre-training.

a unified domain LLM development service that leverages domain KGs to enhance the training process, which effectively improves LLMs' performance on domain-specific tasks.

**Knowledge Obsolescence.** LLMs are pre-trained on prior texts, thus limiting their ability to learn beyond the training corpus. This often results in poor performance when handling tasks that require most-recent knowledge. A simple solution to address this limitation is periodic retraining of LLMs on new data. However, the cost of such retraining is generally high. Hence, it is crucial to devise effective and efficient methods of incorporating current knowledge into LLMs. Prior studies have suggested using plugins as search engines for accessing up-to-date information. Nevertheless, these methods seem inadequate due to the difficulty of directly integrating specific knowledge into LLMs. Compared to LLMs, KGs offer a more straightforward update process that does not necessitate additional training. Updated knowledge can be incorporated into the input in the form of prompts, which are subsequently utilized by LLMs to generate accurate responses [69].

**Bias, Privacy, and Toxicity.** Although LLMs are trained to align with human expectations, they sometimes generate harmful, fully biased, offensive, and private content. When users interact with LLMs, models can be induced to generate such text, even without prior prompting or prompted with safe text. In fact, it has been observed that LLMs tend to degenerate into generating toxic text within just 25 generations [70]. Furthermore, despite their seemingly convincing text, LLMs generally tend to offer unhelpful and sometimes unsafe advice. For example, it has been revealed that GPT-3 produces worse advice than humans do in over 95% of the situations described on Reddit [71]. The reasons are that such biased, private, and toxic texts widely exist in the pre-training corpora and LLMs tend to generate memorized text or new text that is similar to the input text. KGs are commonly built from authoritative and reliable data sources, enabling the generation of high-quality

training data that align with human values, which is expected to enhance the security and reliability of LLMs.

**Computation-Intensive.** Training LLMs is computationally expensive, making it difficult to investigate their effectiveness with different techniques. The training process often requires thousands of GPUs and several weeks to complete. Moreover, LLMs are very computationally intensive and data hungry, making them difficult to deploy, especially in real-world applications where data and computing resources are limited. Through the integration of KGs, smaller LLMs have the potential to outperform larger ones, thereby reducing the cost associated with LLM deployment and application [63].

**Insufficient Interpretability.** Interpretability refers to how easily humans can comprehend a model's predictions, which is an essential gauge of the model's trustworthiness. LLMs are widely acknowledged as black boxes with opaque decision-making processes, making them challenging to interpret. KGs can be used to understand the knowledge learned by LLMs and interpret the reasoning process of LLMs, consequently enhancing the interpretability of LLMs [72].

Overall, LLMs have made noteworthy advancements and are considered a prototype of an artificial general intelligence system at its early stages. However, despite their ability to produce fluent and coherent text, they still encounter many obstacles. Among these obstacles, their struggle in recalling and accurately applying factual knowledge presents the primary challenge, and diminishes their ability to reason and accomplish knowledge-grounded tasks proficiently.

## III. KGPLMs

In light of the limitations posed by poor factual knowledge modeling ability, researchers have proposed incorporating knowledge into PLMs to improve their performance. In recent years, various KGPLMs have been proposed, which can be categorized into before-training enhancement, during-training

TABLE II
SUMMARY OF KGPLMS

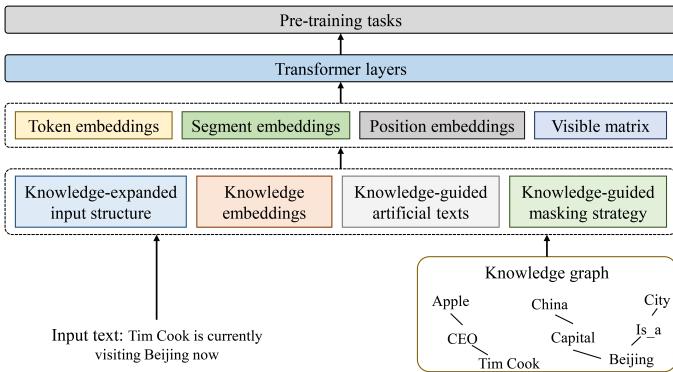| | Method | KGPLM |
|---|---|---|
| Before-training enhancement | Expand input structures | K-BERT [73], CoLAKE [74], Zhang *et al.* [75] |
| | Enrich input information | LUKE [76], E-BERT [77], KALM [78], OAG-BERT [79], DKPLM [80] |
| | Generate new data | AMS [81], KGPT [82], KGLM [83], ATOMIC [84], KEPLER [85] |
| | Optimize word masks | ERNIE [31], WKLM [86], GLM [35] |
| During-training enhancement | Incorporate knowledge encoders | ERNIE [31], ERNIE 3.0 [87], BERT-MK [88], CokeBERT [89], JointLK [90], KET [91], Liu et al. [92], QA-GNN [93], GreaseLM [67], KLMo [94] |
| | Insert knowledge encoding layers | KnowBERT [95], K-BERT [73], CoLAKE [74], JAKET [96], KGBART [97] |
| | Add independent adapters | K-Adapter [98], OM-ADAPT [99], DAKI-ALBERT [100], CKGA [101] |
| | Modify the pre-training task | ERNIE [31], LUKE [76], OAG-BERT [79], WKLM [86], SenseBERT [102], ERICA [103], SentiLARE [104], GLM [35], KEPLER [85], JAKET [96], ERNIE 2.0 [105], ERNIE 3.0 [87], DRAGON [106], LRLM [107] |
| Post-training enhancement | Fine-tune PLMs with knowledge | KALA [108], KeBioSum [109], KagNet [110], BioKGLM [111], Chang *et al.* [112] |
| | Generate knowledge-based prompts | Chang *et al.* [112], Andrus *et al.* [113], KP-PLM [114] |



Fig. 5. Main framework of before-training enhancement KGPLMs.

enhancement, and post-training enhancement methods according to the stage at which KGs participate in pre-training, as illustrated in Fig. 4.

### A. Before-training Enhancement KGPLMs

There are two challenges when integrating the knowledge from KGs into PLMs: heterogeneous embedding space and knowledge noise. The first challenge arises from the heterogeneity between text and KG. The second challenge occurs when unrelated knowledge diverts the sentence from its correct meaning. Before-training enhancement methods resolve these issues by unifying text and KG triples into the same input format, the framework of which is shown in Fig. 5. Existing studies propose diverse approaches to achieve this goal, including expanding input structures, enriching input information, generating new data, and optimizing word masks.

**Expand Input Structures.** Some methods expand the input text into graph structure to merge the structured knowledge of KGs and then convert the merged graph into text for PLM training. For example, K-BERT [73] converts texts to sentence trees to inject related triples by fusing them with KG subgraphs and introduces soft-position and visible matrix to overcome the problem of knowledge noise. Moreover, it proposes mask-self-attention, an extension of self-attention, to prevent erroneous semantic alterations by taking advantage of

the sentence structure information. Formally, the output from mask-self-attention is computed as:

$$h = softmax(\frac{\mathbf{QK}^T + \mathbf{M}}{\sqrt{d_k}})\mathbf{V}, \tag{8}$$

where $\mathbf{M}$ is the visible matrix. CoLAKE [74] addresses the heterogeneous embedding space challenge by combining knowledge context and language context into a unified word-knowledge graph. Zhang *et al.* [75] employed ConceptNet as the knowledge source and improved the visible matrix to control the information flow, which further improved the performance of K-BERT.

**Enrich Input Information.** Instead of merging data from texts and KGs, some studies incorporate entities as auxiliary information by combining their embeddings with text embeddings. LUKE [76] introduces entity type embedding to indicate that the corresponding token in a sentence is an entity, and trains the model with the masked entity prediction task in addition to the MLM task. Further, it extends the Transformer encoder using an entity-aware self-attention mechanism to simultaneously handle both types of tokens. E-BERT [77] aligns entity embeddings with wordpiece vectors through an unconstrained linear mapping matrix and feeds the aligned representations into BERT as if they were wordpiece vectors. KALM [78] signals the existence of entities to the input of the encoder in pre-training using an entity-extended tokenizer and adds an entity prediction task to train the model. Liu *et al.* [79] proposed OAG-BERT, a unified backbone language model for academic knowledge services, which integrates heterogeneous entity knowledge and scientific corpora in an open academic graph. They designed an entity type embedding to differentiate various entity types and used a span-aware entity masking strategy for MLM over entity names with different lengths. Besides, they designed the entity-aware 2D positional encoding to incorporate the entity span and sequence order information. Zhang *et al.* [80] decomposed the knowledge injection process of PLMs into pre-training, fine-tuning, and inference stages, and proposed DKPLM, which injects knowledge only during pre-training. Specifically, DKPLM detects long-tail entities according to their semantic importance in both texts and KGs and replaces the representations of detected long-tail entities with

the representations of the corresponding knowledge triples generated by shared PLM encoders. The most-commonly used knowledge embedding model is TransE [115], which learns entity and relation representations by minimizing the following loss function:

$$\mathcal{L}_{\text{KE}} = - \left\| \mathbf{e}_h + \mathbf{r} - \mathbf{e}_t \right\|_2^2, \tag{9}$$

where $\mathbf{e}_h$ and $\mathbf{e}_t$ are the embeddings of the head and tail entities, while $\mathbf{r}$ is the representation of the relation.

**Generate New Data.** There are also some studies that inject knowledge into PLMs by generating artificial text based on KGs. For example, AMS [81] constructs a commonsense-related question answering dataset for training PLMs based on an align-mask-select method. Specifically, it aligns sentences with commonsense knowledge triples, masks the aligned entities in the sentences and treats the masked sentences as questions. In the end, it selects several entities from KGs as distractor choices and trains the model to determine the correct answer. KGPT [82] crawls sentences with hyperlinks from Wikipedia and aligns the hyperlinked entities to the KG Wikidata to construct the knowledge-grounded corpus KGText. KGLM [83] constructs the Linked WikiText-2 dataset by aligning texts in WikiText-2 and entities in Wikidata. ATOMIC [84] organizes the inference knowledge in 877K textual descriptions into a KG and trains a PLM with a conditional sequence generation problem that encourages the model to generate the target sequence given an event phrase and an inference dimension. KEPLER [85] constructs a large-scale KG dataset with aligned entity descriptions from its corresponding Wikipedia pages for training KGPLMs.

**Optimize Word Masks.** MLM is the most commonly used pre-training task in PLMs, and the number and distribution of masks have a substantial influence on the performance of PLMs [116]. However, the random masking method may break the correlation between consecutive words, making it difficult for PLMs to learn semantic information. To address this issue, a few studies have proposed replacing the random masking strategy with a knowledge masking strategy that selects mask targets based on the knowledge from KGs, forcing models to learn enough knowledge to accurately predict the masked contents. For instance, ERNIE [31] recognizes named entities in texts and aligns them with their corresponding entities in KGs. It then randomly masks entities in the input text and trains the model to select their counterparts in KGs. In WKLM [86], entity mentions in the original texts are substituted with entities of identical types, and the model is trained to differentiate accurate entity mentions from those that are corrupted, which effectively improves its fact completion performance. GLM [35] reformulates the MLM objective to an entity-level masking strategy that identifies entities and selects informative ones by considering both document frequency and mutual reachability of the entities detected in the text.

Before-training enhancement methods can improve the semantic standardization and structural level of the corpus, which is helpful for improving the reasoning ability of PLMs [117] without improving the model size and training time. Besides, the training data enhanced by KGs can better describe commonsense knowledge, which helps to improve LLMs'

commonsense knowledge modeling ability. These methods are more suitable for those domains without sufficient training corpus and can effectively improve LLMs' performance and generalization ability in such domains. However, before-training enhancement processing requires additional computational resources and time, making the pre-training process more complex and cumbersome. Besides, it may introduce noise, which can have a negative impact on LLMs' training.

### B. During-training Enhancement KGPLMs

During-training enhancement methods enable PLMs to learn knowledge directly during training by improving their encoder and training task. Since plain PLMs cannot process text sequences and structured KG simultaneously, some studies have proposed incorporating knowledge encoders or external knowledge modules to enable learning from both text and KGs concurrently. Existing during-training enhancement KGPLMs can be divided into incorporating knowledge encoders, inserting knowledge encoding layers, adding independent adapters, and modifying the pre-training task, as shown in Fig. 6.

**Incorporate Knowledge Encoders.** ERNIE [31] integrates a knowledge encoder to incorporate KG information, which takes two types of input: the token embedding and the concatenation of the token and entity embeddings. Building on ERNIE, ERNIE 3.0 [87] builds a few task-specific modules upon the universal representation module to enable easy customization of the model for natural language understanding and generation tasks. BERT-MK [88] utilizes a graph contextualized knowledge embedding module to learn knowledge in subgraphs and incorporates the learned knowledge into the language model for knowledge generalization. CokeBERT [89] utilizes three modules to select contextual knowledge and embed knowledge context, where the text encoder computes embeddings for the input text, the knowledge context encoder dynamically selects knowledge context based on textual context and computes knowledge embeddings, while the knowledge fusion encoder fuses textual context and knowledge context embeddings for better language understanding. JointLK [90] performs joint reasoning between PLM and a graph neural network (GNN) through a dense bidirectional attention module to effectively fuse and reason over question and KG representations. KET [91] interprets contextual utterances using hierarchical self-attention and dynamically leverages external commonsense knowledge using a context-aware affective graph attention mechanism to detect emotions in textual conversations. Liu *et al.* [92] proposed a memory-augmented approach to condition a PLM on a KG, which represents the KG as a set of relation triples and retrieves pertinent relations for a given context to enhance text generation. QA-GNN [93] uses a PLM to estimate the importance of nodes to identify relevant knowledge from large KGs, and combines the QA context and KG to form a joint graph. Then, it mutually updates the representations of QA context and KG through graph-based message passing to perform joint reasoning. GreaseLM [67] integrates embeddings from a PLM and a GNN through several layers of modality interaction operations. KLMo [94] explicitly models the interaction between entity
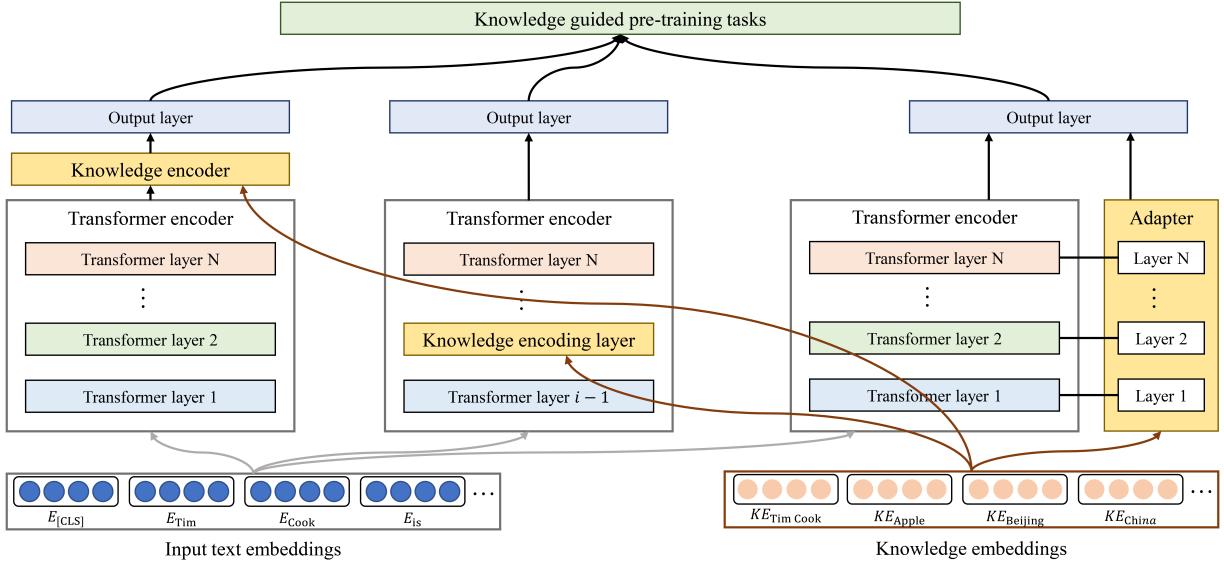
Fig. 6. Main framework of during-training enhancement KGPLMs.

spans in texts and all entities and relations in a contextual KG using a novel knowledge aggregator.

**Insert Knowledge Encoding Layers.** Some methods insert additional knowledge encoding layers in the middle of PLMs or adjust the encoding mechanism to enable PLMs to process knowledge. For instance, KnowBERT [95] incorporates a knowledge attention recontextualization module to integrate multiple KGs into a PLM. It explicitly models entity spans within the input text and uses an entity linker to retrieve relevant entity embeddings from the KG. These retrieved embeddings are then utilized to create knowledge-enhanced entity-span embeddings. K-BERT [73] changes the Transformer encoder to a mask-Transformer, which takes soft-position and visible matrix as input to control the influence of knowledge and avoid the knowledge noise issue. CoLAKE [74] slightly modifies the embedding layer and encoder layers of Transformer to adapt to input in the form of word-knowledge graph. This graph combines the knowledge context and language context into a unified data structure. JAKET [96] decomposes the encoder of a PLM into two modules, with the first providing embeddings for both the second and KG, while the second module takes text and entity embeddings to produce the final representation. KGBART [97] follows the BART architecture but replaces the traditional Transformer with an effective knowledge graph-augmented Transformer to capture relations between concept sets, where KGs serve as additional inputs to the graph attention mechanism.

**Add Independent Adapters.** Some methods add independent adapters to process knowledge, which are easy to train and whose training process does not affect the parameters of the original PLM. For instance, K-Adapter [98] enables the injection of various types of knowledge by training adapters independently on different tasks. This approach facilitates the continual fusion of knowledge. OM-ADAPT [99] complements BERT's distributional knowledge by incorporating conceptual knowledge from ConceptNet and the corresponding

Open Mind Common Sense corpus through adapter training. This approach avoids the expensive computational overhead of joint pre-training, as well as the problem of catastrophic forgetting associated with post-hoc fine-tuning. DAKI-ALBERT [100] proposes pre-training knowledge adapters for specific domain knowledge sources and integrating them through an attention-based knowledge controller to enhance PLMs with enriched knowledge. CKGA [101] introduces a novel commonsense KG-based adapter for sentiment classification tasks, which utilizes a PLM to encode commonsense knowledge and extracts corresponding knowledge with a GNN.

**Modify the Pre-training Task.** Several studies attempt to incorporate knowledge into PLMs by modifying the pre-training tasks. The most commonly used method is to change MLM to masked entity modeling (MEM) based on entities marked in texts. Examples of such methods include ERNIE [31], LUKE [76], OAG-BERT [79], WKLM [86], etc. SenseBERT [102] directly applies weak supervision at the word sense level, which trains a PLM to predict not only masked words but also their WordNet supersenses. ERICA [103] defines two novel pre-training tasks to explicitly model relational facts in texts through contrastive learning, in which the entity discrimination task trains the model to distinguish tail entities while the relation discrimination task is designed to train the model to distinguish the proximity between two relations. SentiLARE [104] introduces a context-aware sentiment attention mechanism to determine the sentiment polarity of each word based on its part-of-speech tag by querying SentiWordNet. It also proposes a novel pre-training task called label-aware masked language model to build knowledge-aware language representations. GLM [35] introduces a KG-guided masking scheme and then employs KGs to obtain distractors for masked entities and uses a novel distractor-suppressed ranking objective to optimize the model.

Other methods utilize the multi-task learning mechanism to integrate knowledge representation learning with the training
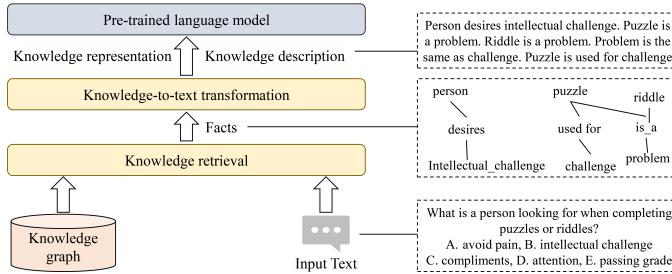
Fig. 7. Main framework of post-training enhancement KGPLMs.

of PLMs, simultaneously optimizing knowledge representation and model parameters. KEPLER [85] employs a shared encoder to encode texts and entities into a unified semantic space, while simultaneously optimizing knowledge embedding and MLM objectives. JAKET [96] jointly models KG and language using two modules, in which the language module and knowledge module mutually assist each other through embeddings. Building upon ERNIE, ERNIE 2.0 [105] proposes a continual multi-task learning framework that extracts valuable lexical, syntactic, and semantic information. ERNIE 3.0 [87] combines auto-regressive and auto-encoding networks to process multiple pre-training tasks at both language and knowledge levels. DRAGON [106] uses a cross-modal encoder that bidirectionally exchanges information between text tokens and KG nodes to produce fused representations and trains this encoder by unifying two self-supervised reasoning tasks: MLM and KG link prediction. LRLM [107] parameterizes the joint distribution over the words in a text and the entities therein, leveraging KGs through relations when modeling text.

During-training enhancement methods can adaptively incorporate external knowledge while learning parameters, often leading to improved performance on various downstream tasks. Moreover, they allow for customization to specific domains or tasks by introducing special information or modules. However, they may increase training time as they typically improve the parameter size and could be limited by the scope of knowledge included in the training data. Moreover, with more complex architecture and more parameters, LLMs are more susceptible to overfitting and require more training to maintain generalization. During-training enhancement methods are more suitable for those scenarios that require dealing with multiple complex tasks, and they often perform better on knowledge-grounded tasks than other methods.

### C. Post-training Enhancement KGPLMs

Post-training enhancement methods typically inject domain-specific knowledge into PLMs through fine-tuning them on additional data and tasks, which improves the model's performance on specific domain tasks. Additionally, with the rapid development of prompt learning [118], several recent investigations have proposed automatically generating prompts to improve the outputs of PLMs. The main framework of post-training enhancement KGPLMs is shown in Fig. 7.

**Fine-tune PLMs with Knowledge.** KALA [108] modulates PLMs' intermediate hidden representations with domain

knowledge, which largely outperforms adaptive pre-training models while still being computationally efficient. KeBioSum [109] investigates the integration of generative and discriminative training techniques to fuse knowledge into knowledge adapters. It applies adapter fusion to effectively incorporate these knowledge adapters into PLMs for the purpose of fine-tuning biomedical text summarization tasks. KagNet [110] proposes a textual inference framework for answering commonsense questions, which effectively utilizes KGs to provide human-readable results via intermediate attention scores. BioKGLM [111] presents a post-training procedure between pre-training and fine-tuning and uses diverse knowledge fusion strategies to facilitate the injection of KGs. Chang *et al.* [112] proposed attentively incorporating retrieved tuples from KGs to incorporate commonsense knowledge during fine-tuning.

**Generate Knowledge-based Prompts.** Bian *et al.* [119] presented a knowledge-to-text framework for knowledge-enhanced commonsense question-answering. It transforms structured knowledge into textual descriptions and utilizes machine reading comprehension models to predict answers by exploiting both original questions and textural knowledge descriptions. Andrus *et al.* [113] proposed using open information extraction models with rule-based post-processing to construct a custom dynamic KG. They further suggested utilizing few-shot learning with GPT-3 to verbalize extracted facts from the KG as natural language and incorporate them into prompts. KP-PLM [114] constructs a knowledge sub-graph from KGs for each context and adopts multiple continuous prompt rules to transform the knowledge sub-graph into natural language prompts. Furthermore, it leverages two novel knowledge-aware self-supervised tasks: prompt relevance inspection and masked prompt modeling, to optimize the model.

TABLE III
PERFORMANCE IMPROVEMENT OF SOME KGPLMs ON DIFFERENT
EVALUATION TASKS COMPARED WITH BERT

| KGPLM | Entity typing | Relation classification | Question answering |
|---|---|---|---|
| CoLAKE | 2.8 | 5.6 | — |
| LUKE | 4.6 | 6.7 | 19.2 |
| KEPLER | 2.6 | 6 | — |
| ERNIE | 2 | 3.4 | — |
| CokeBERT | 1.3 | 2.7 | — |
| K-Adapter | 4.1 | 1.9 | 5.4 |
| ERICA | 4.4 | 2.2 | 1.5 |
| KP-PLM | 4.6 | 3.8 | — |

Post-training enhancement methods are low-cost and easy to implement, which can effectively improve LLMs' performance on specific tasks. Besides, these methods can guide LLMs to generate text of specific styles and improve the quality and security of LLMs' output. Therefore, post-training enhancement methods are more suitable for domain-specific tasks and text generation scenarios that require sensitive information filtering and risk control. However, the labeling of fine-tuning data and the design of prompts rely on prior knowledge and external resources. If there is a lack of relevant prior knowledge, the optimization effect may be limited. Moreover, these methods may impose certain limitations on the flexibility of LLMs' generations. The generated text may be constrained by prompts and may not be able to be fully freely created.

## D. Effectiveness and Efficiency of KGPLMs

Most KGPLMs are designed for knowledge-grounded tasks. To evaluate their effectiveness in knowledge modeling, we report their performance on three knowledge-grounded tasks: entity typing, relation classification, and question answering. Table III provides a summary of KGPLMs and their respective improvements over the unenhanced BERT. The reported metric is F1-score. In Table III, the performances of these models on all tasks are higher than BERT, indicating that KGs enhance their knowledge modeling ability.

### TABLE IV
THE RUNNING TIME OF BERT AND DIFFERENT KGPLMS

| Model | Pre-training | Fine-tuning | Inference |
|---|---|---|---|
| BERT | 8.46 | 6.76 | 0.97 |
| RoBERTa | 9.60 | 7.09 | 1.55 |
| ERNIE | 14.71 | 8.19 | 1.95 |
| KEPLER | 18.12 | 7.53 | 1.86 |
| CoLAKE | 12.46 | 8.02 | 1.91 |
| DKPLM | 10.02 | 7.16 | 1.61 |

Typically, the incorporation of knowledge from KGs would lead to a larger parameter size compared with the base PLM. Consequently, the pre-training, fine-tuning and inference time of plain PLMs are consistently shorter than KGPLMs. As the statistical data shown in Table IV, due to these KGPLMs injecting the knowledge encoder module into PLMs, their running time of the three stages are consistently longer than BERT. However, with the incorporation of external knowledge, KGPLMs are easier to be trained with higher performance. For example, KALM with 775M parameters even performs better than GPT-2 on some downstream tasks [78], whose parameter size is 1.5B. This implies that we can obtain a satisfactory model with smaller parameter size and fewer training resources.

## IV. APPLICATIONS OF KGPLMS

KGPLMs outperform traditional PLMs in capturing factual and relational information, exhibiting stronger language understanding and generation abilities. These advantages lead to improved performance across a range of downstream applications. By employing diverse pre-training tasks and fine-tuning PLMs for specific applications, as illustrated in Fig. 8, KGPLMs have been successfully leveraged for multiple tasks.

**Named Entity Recognition.** Named entity recognition (NER) aims to identify entities with specific meanings from text, such as names of persons, places, and organizations. PLMs have successfully improved state-of-the-art word representations and demonstrated effectiveness on the NER task by modeling context information [120]. However, these models are trained to predict correlations between tokens, ignoring the underlying meanings behind them and the complete semantics of entities that consist of multiple tokens [121]. Previous work has already regarded NER as a knowledge intensive task and improved PLMs' NER performance by incorporating external knowledge into PLMs [122]. Therefore, researchers have developed KGPLMs for NER, which can leverage additional information beyond the training corpus for better performance, especially in domain-specific tasks where the training samples are often insufficient. For example, He *et al.* [123] incorporated prior knowledge of entities from an external knowledge base into word representations and introduced a KG augmented word representation framework for NER. Some other KGPLMs like K-BERT [73] and ERNIE [31] also demonstrate their superiority on diverse NER datasets.

**Relation Extraction.** Relation extraction involves distinguishing semantic relationships between entities and classifying them into predefined relation types. Although PLMs have improved the efficacy of relation extraction to some extent, when applied to small-scale and domain-specific texts, there is still a lack of information learning [124]. To address this limitation, several studies have suggested injecting prior knowledge from KGs into PLMs. KGPLMs have been demonstrated to be more effective than plain PLMs in relation extraction [125]. For example, Roy *et al.* [126] proposed merging KG embeddings with BERT to improve its performance on clinical relation extraction. BERT-MK [88] also demonstrates the effectiveness of KGPLMs on biomedical relation extraction. In addition to the biomedical field, KGPLMs such as KEPLER [85] and JAKET [96] are also commonly applied to public domain relation extraction tasks.

**Sentiment Analysis.** Sentiment analysis aims to analyze whether the emotions expressed in the text are positive, negative, or neutral. Recently, sentiment analysis has made remarkable advances with the help of PLMs, which achieve state-of-the-art performance on diverse benchmarks. However, current PLMs focus on acquiring semantic information through self-supervision techniques, disregarding sentiment-related knowledge throughout pre-training [127]. By integrating different types of sentiment knowledge into the pre-training process, the learned semantic representation would be more appropriate. For this reason, several KGPLMs have been applied to sentiment analysis, including SentiLARE [104], KCF-PLM [128], and KET [91], which have proven the effectiveness of injecting KGs into PLMs for sentiment analysis.

**Knowledge Graph Completion.** Due to the limitations in data quality and automatic extraction technology, KGs are often incomplete, and some relations between entities are missing [129]. Therefore, the knowledge graph completion task, aiming at inferring missing relations and improving the completeness of KGs, has been widely investigated. Given the triumph of PLMs, some PLM-based methods are proposed for the knowledge graph completion task. Nonetheless, most of these methods concentrate on modeling the textual representation of factual triples while neglecting the underlying topological contexts and logical rules that are essential for KG modeling [130], [131]. To address this challenge, some studies have suggested combining topology contexts and logical rules in KGs with textual semantics in PLMs to complete the KG. By integrating the structure information from KGs and the contextual information from texts, KGPLMs outperform those PLMs specifically designed for the KG completion task [35]. We can also extract the knowledge-enhanced embeddings to predict the rationality of given triples [85].

**Question Answering.** Question answering systems need to choose the correct answers for the given questions, which must be able to access relevant knowledge and reason over
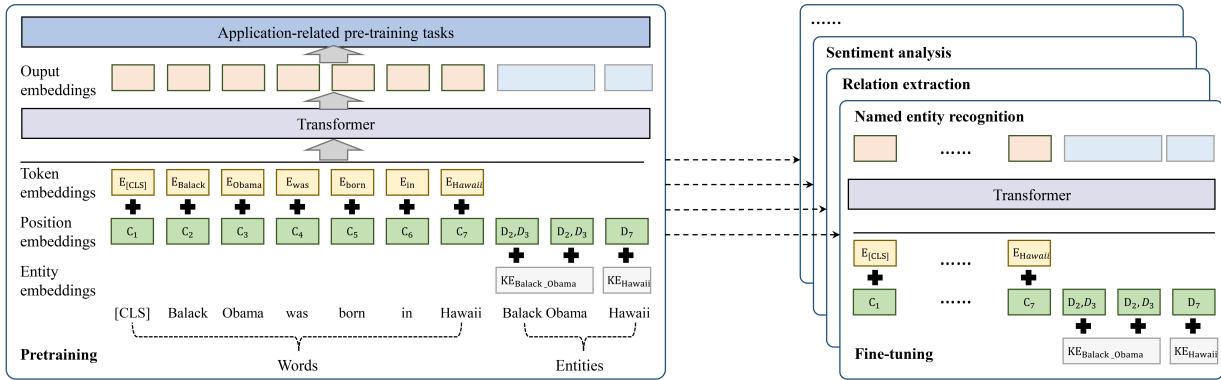
Fig. 8. The framework for KGPLMs to realize various applications.

it. Although PLMs have made remarkable achievements on many question answering tasks [132], they do not empirically perform well on structured reasoning. On the other hand, KGs are more suitable for structured reasoning and enable explainable predictions. Therefore, a few studies have proposed integrating PLMs with KGs to conduct structured reasoning and enable explainable predictions. Some methods incorporate KGs into PLMs while training them, such as QA-GNN [93] and WKLM [86]. Another line of research uses KGs to augment PLMs during answer inference. OreoLM [133], for example, incorporates a novel knowledge interaction layer into PLMs that interact with a differentiable knowledge graph reasoning module for collaborative reasoning. Here, PLMs guide KGs in walking towards desired answers while retrieved knowledge enhances PLMs. Experiments on common benchmarks illustrate that KGPLMs outperform traditional PLMs after KG incorporation.

**Natural Language Generation.** Natural language generation (NLG) serves as a fundamental building block for various applications in NLP, such as dialogue systems, neural machine translation, and story generation, and has been subject to numerous studies. Deep neural language models pre-trained on large corpus have caused remarkable improvements in multiple NLG benchmarks. However, even though they can memorize enough language patterns during pre-training, they merely capture average semantics of the data and most of them are not explicitly aware of domain-specific knowledge. Thus, when specific knowledge is required, contents generated by PLMs could be inappropriate. KGs, which store entity attributes and their relations, contain rich semantic contextual information. As a result, several studies have proposed incorporating KGs into PLMs to improve their NLG performance. For instance, Guan et al. [134] proposed improving GPT-2 with structured knowledge by post-training the model using knowledge examples sourced from KGs. They aimed to supply additional crucial information for story generation. Ji et al. [135] proposed GRF, a generation model that performs multi-hop reasoning on external KGs, enriching language generation with KG-derived data. Experimental results indicate that KGPLMs outperform PLMs in story ending generation [136], abductive reasoning [137], and question answering [93].

**Industrial Applications.** KGPLMs have been applied in many real-world applications. Typical applications include chatbots, such as ERNIE Bot[1] from Baidu, Qianwen[2] from Alibaba, and Bard[3] from Google, which incorporate KGs into PLMs to improve knowledge awareness while communicating with humans. Such applications have shown that KGPLMs can provide excellent language understanding and knowledge modeling abilities. PLMs have also been successfully applied in programming assistants, which can easily generate codes according to context or natural language prompts. However, there are still some issues encountered by PLM-based programming assistants, such as incorrect code recommendations and excessive reliance on code libraries. To tackle these challenges, GitHub and OpenAI released Copilot X[4], which incorporates KGs into the programming assistant to analyze the logical dependencies of the code and generate appropriate code recommendations. Aside from the above applications, KGPLMs are widely used in a variety of virtual assistants and search engines. Representatives of these applications include Xiaodu[5] from Baidu and PanGu[6] from Huawei, which can respond to a broad range of queries like weather forecasts, singing songs, and navigation.

## V. CAN LLMs REPLACE KGs?

Recent advancements in training PLMs on a large corpus have led to a surge of improvements for downstream NLP tasks. While primarily learning linguistic knowledge, PLMs may also store some relational knowledge present in the training data that enables them to answer complex queries. Although their knowledge cannot be directly queried like KGs, we can attempt to query them for factual knowledge by asking them to fill in masked tokens in sequences, as illustrated in Fig. 9. Consequently, some researchers believe that parametric PLMs can replace symbolic KGs as knowledge bases [138]. For example, Petroni et al. [8] proposed LAMA, a knowledge probe consisting of cloze-style queries, to measure relational knowledge contained in PLMs. Their results show

---

[1] https://yiyan.baidu.com/
[2] https://qianwen.aliyun.com/
[3] https://bard.google.com/
[4] https://github.com/features/preview/copilot-x
[5] https://dueros.baidu.com/en/index.html
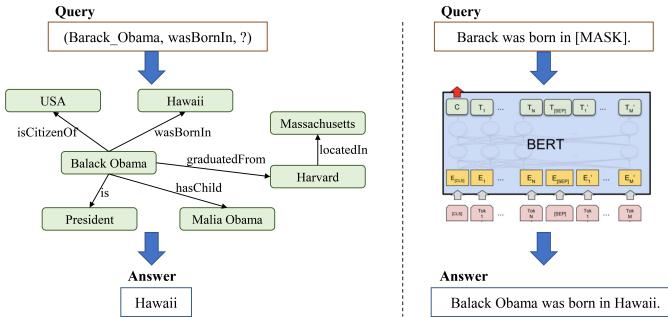[6] https://www.huaweicloud.com/product/pangu.html

Fig. 9. Querying KGs and PLMs for factual knowledge, in which the left part represents directly querying factual knowledge from KGs while the right part represents querying factual knowledge from PLMs by asking them to fill in masked tokens in sequences.

that PLMs contain relational knowledge and can recall stored facts without fine-tuning. Talmor *et al.* [139] developed eight cloze-style reasoning tasks to test the knowledge captured in BERT and RoBERTa. They found that different PLMs exhibit qualitatively different reasoning abilities and do not reason in an abstract manner but instead rely on context. Heinzerling and Inui [10] evaluated PLMs' ability to store millions of entity facts and query these facts via experimental tests with three entity representations. Their experimental results provide a proof-of-concept for PLMs as knowledge bases.

Nevertheless, after conducting extensive experimental analyses of PLMs, some studies have reported that PLMs struggle to accurately recall relational facts, raising doubts about their viability as knowledge bases. A surge of benchmark datasets and tasks have been proposed to examine the knowledge embodied within PLMs. For example, Wang *et al.* [140] released a benchmark to directly test a system's ability to differentiate natural language statements that make sense from those that do not. By comparing their performance with humans, they revealed that sense-making remains a technical challenge for PLMs. Sung *et al.* [141] created the BioLAMA benchmark that is comprised of 49K biomedical factual knowledge triples for probing biomedical PLMs. Their detailed analysis reveals that most PLMs' predictions are highly correlated with prompt templates without any subjects, hence producing similar results on each relation and hindering their capabilities to be used as biomedical knowledge bases. Wang *et al.* [12] constructed a new dataset of closed-book question answering and tested the BART's [26] ability to answer these questions. Experimental results show that it is challenging for BART to answer closed-book questions since it cannot remember training facts in high precision. Zhao *et al.* [142] introduced LAMA-TK, a dataset aimed at probing temporally-scoped knowledge. They investigated the capacity of PLMs for storing temporal knowledge that contains conflicting information and the ability to use stored knowledge for temporally-scoped knowledge queries. Their experimental results show that conflicting information poses great challenges to PLMs, which drops their storage accuracy and hinders their memorization of multiple answers. Kassner *et al.* [143] translated two established benchmarks into 53 languages to investigate the knowledge contained in the multilingual PLM mBERT [144]. They found that

mBERT yielded varying performance across languages. The above studies have proven that PLMs still face challenges in accurately storing knowledge, dealing with knowledge diversity, and retrieving correct knowledge to solve corresponding tasks. Additionally, Cao *et al.* [13] conducted a comprehensive investigation into the predictive mechanisms of PLMs across various extraction paradigms. They found that previous decent performance of PLMs mainly owes to the biased prompts which overfit dataset artifacts. AlKhamissi *et al.* [138] suggested five essential criteria that PLMs should meet in order to be considered proficient knowledge bases: access, edit, consistency, reasoning, and explainability and interpretability, and found that PLMs do not perform as well as KGs in terms of consistency, reasoning, and interpretability. They also reviewed the literature with respect to the five aspects and revealed that the community still has a long way to go to enable PLMs to serve as knowledge bases despite some recent breakthroughs. These studies raise doubts about PLMs' potential as knowledge bases and underscore the need for further research in this area.

Despite the fact that larger-sized LLMs seem to possess more fundamental knowledge of the world, their learned encyclopedic facts and common sense properties of objects are still unreliable. Furthermore, they have limited capabilities in inferring relationships between actions and events [65]. The ability of LLMs to predict facts is also significantly dependent on specific prompt templates and the included entities [145]. This owes to the fact that LLMs mainly rely on simple heuristics with most predictions correlated to co-occurrence frequencies of the target word and words in the prompt. Additionally, the accuracy of their predictions is highly reliant on the frequency of facts in the pre-training corpus [146].

To summarize, LLMs and KGs have their respective advantages and disadvantages. KGs lack the flexibility that LLMs offer, as KGs require substantial human effort to build and maintain, while LLMs provide more flexibility through unsupervised training on a large corpus. However, KGs are easier to access and edit, and have better consistency, reasoning ability, and interpretability. First, factual knowledge in KGs is often easily accessed through manual query instructions. In contrast, LLMs cannot be queried explicitly, as the knowledge is implicitly encoded in their parameters. Second, the triplets in KGs can be directly added, modified, and deleted. However, editing a specific fact in LLMs is not straightforward, since facts in LLMs cannot be directly accessed. To enable LLMs to learn up-to-date, correct, and unbiased knowledge, the whole model needs to be retrained on updated data, which is expensive and inflexible. Third, KGs are built with consistency in mind, and various algorithms have been proposed to eliminate conflicts that arise in KGs. On the other hand, LLMs may be inconsistent, as they may yield different answers to the same underlying factual questions. Fourth, it can be simple to follow the path of reasoning in KGs, while LLMs perform poorly on relational reasoning tasks. Finally, KGs have a clear reasoning path, so their outputs are easy to interpret. However, as typical black-box models, knowledge is hard to be identified by simply looking at LLMs' outputs.

Although current LLMs face limitations in directly serving

as knowledge bases, they contribute to constructing KGs that explicitly express their stored knowledge. One approach is to utilize LLMs as an information extraction tool to improve the accuracy of NER and relation extraction. Another way is to extract symbolic KGs from LLMs using prompts. For example, Hao *et al.* [147] proposed a novel framework to automatically construct KGs from LLMs that generates diverse prompts, searches for consistent outputs, and performs efficient knowledge search. Bosselut *et al.* [148] proposed a fine-tuned generative LLM for the automatic construction of commonsense KGs that generates tail entities based on given head entities and relations. These approaches demonstrate the potential of leveraging LLMs for effective KG construction.

To conclude, LLMs still face challenges in remembering large amounts of complex knowledge and retrieving the required information accurately. There are multiple aspects in which LLMs need to excel to qualify as comprehensive knowledge bases. On the other hand, KGs and LLMs complement each other, enhancing overall performance. Therefore, enhancing LLMs with KGs can significantly improve their performance on knowledge-grounded tasks.

## VI. ENHANCING LLMS WITH KGS

In the preceding sections, we have analyzed and compared existing KGPLMs. Despite demonstrating proficiency in a wide range of NLP tasks, the complexity of knowledge and language continues to pose unresolved challenges for KGPLMs. Furthermore, despite substantial improvements in generated text quality and learned facts with models scaling beyond 100B parameters, LLMs are still prone to unfactual responses and commonsense errors. Their predictions are highly dependent on input text, and minor variations in phrasing and word choice can lead to such errors. One potential solution is to enhance LLMs with KGs to improve their learning of factual knowledge, a topic that has not been thoroughly studied yet. Thus, we propose to enhance LLMs with KGs using techniques utilized by KGPLMs to achieve fact-aware language modeling.

### A. Overall Framework

The development framework for KGLLMs based on existing technologies is depicted in Fig. 10. Since LLMs primarily scale the size of parameters and training data from PLMs, their model architecture and training methods remain largely unchanged. Hence, all three types of KGPLM methods introduced before can be applied to developing KGLLMs. The before-training enhancement approaches can be utilized to construct KG-extended text, improving input quality and integrating factual information into the input. The during-training enhancement methods can be employed to adaptively fuse textual knowledge and structural knowledge to learn knowledge-enhanced word representations. Graph encoders, such as GNN, can serve as knowledge encoders, while attention mechanisms can be utilized to design the knowledge fusion module. Multi-task learning, including knowledge-guided pre-training tasks, helps improve LLMs' learning of factual knowledge. The post-training enhancement methods can be utilized to
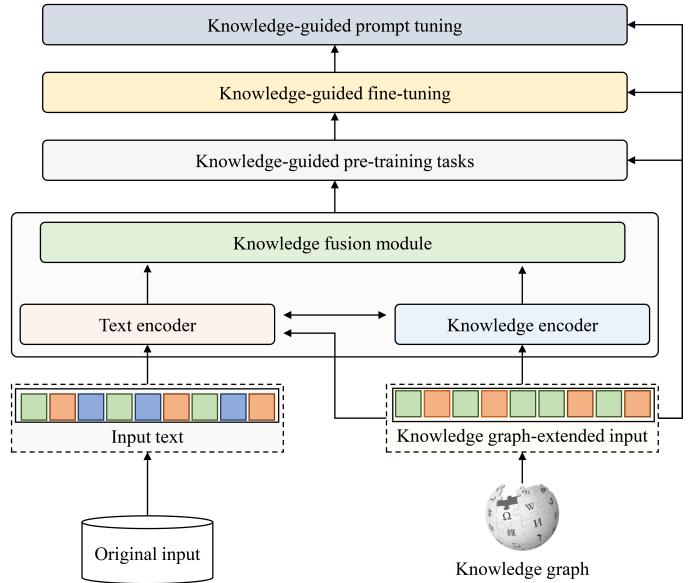


Fig. 10. Technical framework of developing KGLLMs.

further improve the performance of LLMs on some domain-specific tasks by fine-tuning them on knowledge-extended data or knowledge-grounded tasks. Moreover, one of the most important recent advancements of LLMs is prompt learning, which effectively improves the quality of generated text and enhance LLMs' generalization capability by inserting text pieces into the input. In prompt learning, selecting suitable prompt templates for specific tasks is crucial for enhancing model performance, requiring domain expertise. Therefore, KGs can be integrated into constructing prompt templates to make use of domain knowledge, which is expected to improve the model's understanding of domain factual knowledge by guiding LLMs with knowledge prompts.

### B. Discussion and Future Directions

In addition to knowledge graph enhancement methods, there are also other enhancement methods that can be used to improve LLMs' factual language modeling ability. Typically, these methods include data augmentation and retrieval augmentation. Data augmentation involves refining the training data during pretraining and emphasizing informative words, emphasizing the importance of the training corpus in equipping the model with factual knowledge. Compared with knowledge graph enhancement methods, these approaches utilize implicit knowledge to model factual knowledge in text and ignore the relationships between entities. Retrieval augmentation has emerged as a widely adopted approach, allowing LLMs to retrieve external data from databases [149] or tools and pass it to LLMs in the form of prompts or embeddings to improve LLMs' generations. These methods can address some challenges faced by plain LLMs, such as outdated information and the inability to memorize. However, they cannot fundamentally improve LLMs' knowledge modeling ability since they do not change LLMs' parameters.

Besides, some plugins have been developed to enhance the capabilities of LLMs in the context of a knowledge base. For

example, the Browsing plugin can call search engines to access real-time information on the website; the Retrieval plugin[7] uses OpenAI embeddings to index and search documents in vector databases; the Wolfram[8] plugin enables ChatGPT to provide more comprehensive and accurate answers by giving it access to the Wolfram Alpha knowledge base; the Expedia plugin[9] enables ChatGPT to provide personalized travel recommendations with the help of Expedia's entity graph.

Although KGLLMs have achieved some success, there are still many unresolved challenges. Here, we outline and discuss a few promising research directions for KGLLMs.

**Improving the efficiency of KGLLMs.** Due to the need for preprocessing and encoding knowledge from KGs, developing KGLLMs typically requires more computational resources and time compared to plain LLMs. However, the scaling law of KGLLMs may differ from that of plain LLMs. Previous studies on KGPLMs have demonstrated that smaller KGPLMs can even outperform larger PLMs. Therefore, a comprehensive investigation of the scaling law of KGLLMs is necessary to determine the optimal parameter size for their development. Based on this, we can potentially achieve a smaller model that satisfies performance requirements, resulting in reduced computational resources and time.

**Merging different knowledge in different ways.** Some common and well-defined knowledge could be stored within KGs for ease of access, while rarely used or implicit knowledge that cannot be expressed through triples should be incorporated into the parameters of LLMs. In particular, domain-specific knowledge, although infrequently accessed, may still require a significant amount of human effort to construct an associated KG due to the sparse nature of its related corpus.

**Incorporating more types of knowledge.** As introduced in Section III, the majority of existing KGPLMs only utilize a single modality and static KGs. However, there exist multimodal and temporal KGs that contain multimodal and temporal knowledge. These types of knowledge can complement textual and structural knowledge, enabling LLMs to learn the relationships between entities over time. Moreover, multimodal pre-trained models have gained popularity as they have been proven to improve the performance of pre-trained models on multimodal tasks [150] and enhance their cognitive ability. Therefore, incorporating multimodal and temporal KGs into LLMs has the potential to improve their performance, which is worth investigating. To achieve this goal, we need to align multimodal entities, design encoders capable of processing and fusing multimodal temporal data, and establish multimodal temporal learning tasks to extract useful information.

**Improving the effectiveness of knowledge incorporation.** By modifying inputs, model architecture, and the fine-tuning process, diverse methods have been proposed to incorporate relational triplets into PLMs. However, each method has its own set of advantages and disadvantages, with some performing well on particular tasks but underperforming on others. For example, LUKE [76] exhibits superior performance over

KEPLER [85] in most entity typing and relation classification tasks but performs worse in a few other tasks [89]. Besides, recent experimental analysis [151] reveals that existing KG-PLMs integrate only a small fraction of factual knowledge. Therefore, there is still a lot of room for research on effective knowledge integration methods. Further research is required on the selection of valuable knowledge and avoiding catastrophic forgetting when faced with vast and clashing knowledge.

**Enhancing the interpretability of KGLLMs.** Although it is widely believed that KGs can enhance the interpretability of LLMs, corresponding methods have not yet been thoroughly studied. Schuff *et al.* [152] investigated whether integrating external knowledge can improve natural language inference models' explainability by evaluating the scores of generated explanations on in-domain data and special transfer datasets. However, they found that the most commonly used metrics do not consistently align with human evaluations concerning the accuracy of explanations, incorporation of common knowledge, and grammatical and labeling correctness. To provide human-understandable explanations for LLMs, Chen *et al.* [153] proposed a knowledge-enhanced interpretation module that utilizes a KG and a GNN to extract key decision signals of LLMs. Despite a few studies attempting to improve the interpretability of PLMs, it remains unclear how to leverage KGs to improve the interpretability of KGPLMs. A feasible approach may involve searching for the relevant reasoning path in KGs based on the generated content and then generating an explanatory text based on the reasoning path.

**Exploring domain-specific KGLLMs.** Though there is already considerable research incorporating standard KGs with general PLMs, limited work has focused on domain-specific KGLLMs. However, the rise of artificial intelligence for science will lead to an increasing demand for domain-specific KGLLMs. In comparison to general LLMs, domain-specific LLMs require greater precision and specificity in incorporating domain knowledge. As a result, constructing accurate domain-specific KGs and integrating them with LLMs warrant further exploration. In order to develop domain-specific KGLLMs, it is essential to first construct a domain KG and gather relevant corpus data with the help of domain experts. Considering the generality of language patterns, it is advisable to blend common KGs with the domain-specific KG for enhancement.

## VII. CONCLUSION

The phenomenal success of ChatGPT has spurred the rapid advancement of LLMs. Given the impressive performance of LLMs on a variety of NLP tasks, some researchers wonder if they can be viewed as a type of parameterized knowledge base and replace KGs. However, LLMs still fall short in recalling and correctly using factual knowledge while generating knowledge-grounded text. In order to clarify the value of KGs in the era of LLMs, a comprehensive survey on KGPLMs was conducted in this paper. We began by examining the background of PLMs and the motivation for incorporating KGs into PLMs. Next, we categorized existing KGPLMs into three categories and provided details about each category. We

---

[7]https://github.com/openai/chatgpt-retrieval-plugin

[8]https://www.wolfram.com/wolfram-plugin-chatgpt/

[9]https://chatonai.org/expedia-chatgpt-plugin

also reviewed the applications of KGPLMs. After that, we analyzed whether PLMs and recent LLMs can replace KGs based on existing studies. In the end, we proposed enhancing LLMs with KGs to conduct fact-aware language modeling for improving their learning of factual knowledge. This paper addresses three questions: (1) What is the value of KGs in the era of LLMs? (2) How to incorporate KGs into LLMs to improve their performance? (3) What do we need to do for the future development of KGLLM? We hope this work will stimulate additional research advancements in LLM and KG.

## REFERENCES

[1] X. Zhou, C. Chai, G. Li, and J. Sun, "Database meets artificial intelligence: A survey," *IEEE Trans. Knowl. Data Eng.*, vol. 34, no. 3, pp. 1096–1116, 2022.

[2] Q. Wang, Y. Li, R. Zhang, K. Shu, Z. Zhang, and A. Zhou, "A scalable query-aware enormous database generator for database evaluation," *IEEE Trans. Knowl. Data Eng.*, vol. 35, no. 5, pp. 4395–4410, 2023.

[3] R. Lu, X. Jin, S. Zhang, M. Qiu, and X. Wu, "A study on big knowledge and its engineering issues," *IEEE Trans. Knowl. Data Eng.*, vol. 31, no. 9, pp. 1630–1644, 2019.

[4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proc. of the 17th Annu. Conf. of the North Amer. Chapter of the Assoc. for Comput. Linguistics: Hum. Lang. Technol.*, 2019, pp. 4171–4186.

[5] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever *et al.*, "Improving language understanding by generative pre-training," 2018.

[6] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. of Mach. Learn. Res.*, vol. 21, no. 140, pp. 1–67, 2020.

[7] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler *et al.*, "Emergent abilities of large language models," *arXiv:2206.07682*, 2022.

[8] F. Petroni, T. Rocktäschel, P. Lewis, A. Bakhtin, Y. Wu, A. H. Miller, and S. Riedel, "Language models as knowledge bases?" in *Proc. 2019 Conf. Empirical Methods Nat. Lang. Process. and 9th Int. Joint Conf. Nat. Lang. Process.*, 2019, pp. 2463–2473.

[9] C. Wang, X. Liu, and D. Song, "Language models are open knowledge graphs," *arXiv:2010.11967*, 2020.

[10] B. Heinzerling and K. Inui, "Language models as knowledge bases: On entity representations, storage capacity, and paraphrased queries," in *Proc. 16th Conf. Eur. Chapter Assoc. Comput. Linguist.*, 2021, pp. 1772–1791.

[11] N. Bian, X. Han, L. Sun, H. Lin, Y. Lu, and B. He, "Chatgpt is a knowledgeable but inexperienced solver: An investigation of commonsense problem in large language models," *arXiv:2303.16421*, 2023.

[12] C. Wang, P. Liu, and Y. Zhang, "Can generative pre-trained language models serve as knowledge bases for closed-book qa?" in *Proc. 59th Annu. Meet. Assoc. Comput. Linguist. and 11th Int. Joint Conf. Nat. Lang. Process.*, 2021, pp. 3241–3251.

[13] B. Cao, H. Lin, X. Han, L. Sun, L. Yan, M. Liao, T. Xue, and J. Xu, "Knowledgeable or educated guess? revisiting language models as knowledge bases," in *Proc. 59th Annu. Meet. Assoc. Comput. Linguist. and 11th Int. Joint Conf. Nat. Lang. Process.*, 2021, pp. 1860–1874.

[14] H. Liu, R. Ning, Z. Teng, J. Liu, Q. Zhou, and Y. Zhang, "Evaluating the logical reasoning ability of chatgpt and gpt-4," *arXiv:2304.03439*, 2023.

[15] Y. Bang, S. Cahyawijaya, N. Lee, W. Dai, D. Su, B. Wilie, H. Lovenia, Z. Ji, T. Yu, W. Chung *et al.*, "A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity," *arXiv:2302.04023*, 2023.

[16] H. M. Yohannes and T. Amagasa, "Named-entity recognition for a low-resource language using pre-trained language model," in *Proc. 37th ACM/SIGAPP Symp. Appl. Comput.*, 2022, p. 837–844.

[17] X. Wei, S. Wang, D. Zhang, P. Bhatia, and A. Arnold, "Knowledge enhanced pretrained language models: A comprehensive survey," *arXiv:2110.08455*, 2021.

[18] L. Hu, Z. Liu, Z. Zhao, L. Hou, L. Nie, and J. Li, "A survey of knowledge enhanced pre-trained language models," *IEEE Trans. Knowl. Data Eng.*, pp. 1–19, 2023.

[19] C. Zhen, Y. Shang, X. Liu, Y. Li, Y. Chen, and D. Zhang, "A survey on knowledge-enhanced pre-trained language models," *arXiv:2212.13428*, 2022.

[20] X. Liu, F. Zhang, Z. Hou, L. Mian, Z. Wang, J. Zhang, and J. Tang, "Self-supervised learning: Generative or contrastive," *IEEE Trans. on Knowl. Data Eng.*, vol. 35, no. 1, pp. 857–876, 2023.

[21] H. Wang, J. Li, H. Wu, E. Hovy, and Y. Sun, "Pre-trained language models and their applications," *Eng.*, 2022.

[22] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, "ELECTRA: Pre-training text encoders as discriminators rather than generators," in *Proc. 8th Int. Conf. Learn. Representations*, 2020.

[23] M. Joshi, D. Chen, Y. Liu, D. S. Weld, L. Zettlemoyer, and O. Levy, "SpanBERT: Improving pre-training by representing and predicting spans," *Trans. Assoc. Comput. Linguist.*, vol. 8, pp. 64–77, 2020.

[24] Y. Wang, C. Sun, Y. Wu, J. Yan, P. Gao, and G. Xie, "Pre-training entity relation encoder with intra-span and inter-span information," in *Proc. 2020 Conf. Empirical Methods Nat. Lang. Process.*, 2020, pp. 1692–1705.

[25] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "Albert: A lite bert for self-supervised learning of language representations," *arXiv:1909.11942*, 2019.

[26] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, "BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," in *Proc. 58th Ann. Meet. Assoc. Comput. Linguistics.*, 2020, pp. 7871–7880.

[27] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, "Deep contextualized word representations," *arXiv:1802.05365*, 2018.

[28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. 31st Conf. Neural Inform. Process. Syst.*, 2017.

[29] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert pretraining approach," *arXiv:1907.11692*, 2019.

[30] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," *arXiv:1910.01108*, 2019.

[31] Z. Zhang, X. Han, Z. Liu, X. Jiang, M. Sun, and Q. Liu, "ERNIE: Enhanced language representation with informative entities," in *Proc. 57th Ann. Meet. Assoc. Comput. Linguistics.*, 2019, pp. 1441–1451.

[32] P. He, X. Liu, J. Gao, and W. Chen, "Deberta: Decoding-enhanced bert with disentangled attention," in *International Conference on Learning Representations*, 2021.

[33] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, pp. 1–9, 2019.

[34] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, "Xlnet: Generalized autoregressive pretraining for language understanding," in *Advances in neural information processing systems*, 2019.

[35] T. Shen, Y. Mao, P. He, G. Long, A. Trischler, and W. Chen, "Exploiting structured knowledge in text via graph-guided representation learning," in *Proc. 2020 Conf. Empirical Methods Nat. Lang. Process.*, 2020, p. 8980–8994.

[36] W. Fedus, B. Zoph, and N. Shazeer, "Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity," *The Journal of Machine Learning Research*, vol. 23, no. 1, pp. 5232–5270, 2022.

[37] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. Dehghani, S. Brahma *et al.*, "Scaling instruction-finetuned language models," *arXiv:2210.11416*, 2022.

[38] A. Zeng, X. Liu, Z. Du, Z. Wang, H. Lai, M. Ding, Z. Yang, Y. Xu, W. Zheng, X. Xia *et al.*, "Glm-130b: An open bilingual pre-trained model," *arXiv:2210.02414*, 2022.

[39] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, and A. Askell, "Language models are few-shot learners," in *Adv. Neural Inform. Process. Syst.*, 2020, pp. 1877–1901.

[40] N. Du, Y. Huang, A. M. Dai, S. Tong, D. Lepikhin, Y. Xu, M. Krikun, Y. Zhou, A. W. Yu, O. Firat, B. Zoph, L. Fedus, M. P. Bosma, Z. Zhou, T. Wang, E. Wang, K. Webster, M. Pellat, K. Robinson, K. Meier-Hellstern, T. Duke, L. Dixon, K. Zhang, Q. Le, Y. Wu, Z. Chen, and C. Cui, "GLaM: Efficient scaling of language models with mixture-of-experts," in *Proc. 39th Int. Conf. Machine Learning*, 2022, pp. 5547–5569.

[41] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, and A. Ray, "Training language models to follow instructions with human feedback," in *Adv. Neural Inform. Process. Syst.*, 2022, pp. 27 730–27 744.

[42] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann *et al.*, "Palm: Scaling language modeling with pathways," *arXiv:2204.02311*, 2022.

[43] R. Thoppilan, D. De Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H.-T. Cheng, A. Jin, T. Bos, L. Baker, and Y. Du, "Lamda: Language models for dialog applications," *arXiv:2201.08239*, 2022.

[44] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin *et al.*, "Opt: Open pre-trained transformer language models," *arXiv:2205.01068*, 2022.

[45] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv:2302.13971*, 2023.

[46] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto, "Stanford alpaca: An instruction-following llama model," 2023.

[47] OpenAI, "Gpt-4 technical report," *arXiv:2303.08774*, 2023.

[48] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv:2307.09288*, 2023.

[49] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, "A survey of large language models," *arXiv:2303.18223*, 2023.

[50] J.-W. Lu, C. Guo, X.-Y. Dai, Q.-H. Miao, X.-X. Wang, J. Yang, and F.-Y. Wang, "The chatgpt after: Opportunities and challenges of very large scale pre-trained models," *Acta Autom. Sin.*, vol. 49, no. 4, pp. 705–717, 2023.

[51] P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei, "Deep reinforcement learning from human preferences," in *Adv. Neural Inf. Process. Syst.*, 2017.

[52] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing," *ACM Comput. Surv.*, vol. 55, no. 9, pp. 1–35, 2023.

[53] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv:1707.06347*, 2017.

[54] J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. Chi, Q. Le, and D. Zhou, "Chain of thought prompting elicits reasoning in large language models," *arXiv:2201.11903*, 2022.

[55] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg *et al.*, "Sparks of artificial general intelligence: Early experiments with gpt-4," *arXiv:2303.12712*, 2023.

[56] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui, "A survey for in-context learning," *arXiv:2301.00234*, 2022.

[57] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, "Finetuned language models are zero-shot learners," *arXiv:2109.01652*, 2021.

[58] Y. Qin, S. Hu, Y. Lin, W. Chen, N. Ding, G. Cui, Z. Zeng, Y. Huang, C. Xiao, and C. Han, "Tool learning with foundation models," *arXiv:2304.08354*, 2023.

[59] K. Yang, Y. Tian, N. Peng, and K. Dan, "Re3: Generating longer stories with recursive reprompting and revision," in *Proc. 2022 Conf. Empir. Methods Nat. Lang. Process.*, 2022, p. 4393–4479.

[60] W. Yu, C. Zhu, Z. Li, Z. Hu, Q. Wang, H. Ji, and M. Jiang, "A survey of knowledge-enhanced text generation," *ACM Comput. Surv.*, vol. 54, no. 11, pp. 1–18, 2022.

[61] X. Li, R. Zhao, Y. K. Chia, B. Ding, S. Joty, S. Poria, "Chain of knowledge: A framework for grounding large language models with structured knowledge bases," *arXiv:2305.13269*, 2023.

[62] E. Yanai, K. Nora, R. Shauli, F. Amir, R. Abhilasha, M. Marius, B. Yonatan, S. Hinrich, and G. Yoav, "Measuring causal effects of data statistics on language model's 'factual' predictions," *arXiv:2207.14251*, 2022.

[63] J. Sun, C. Xu, L. Tang, S. Wang, C. Lin, Y. Gong, H.-Y. Shum, and J. Guo, "Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph," *arXiv:2307.07697*, 2023.

[64] S. Abulhair and H. He, "Language models are greedy reasoners: A systematic formal analysis of chain-of-thought," *arXiv:2210.01240*, 2022.

[65] T. A. Chang and B. K. Bergen, "Language model behavior: A comprehensive survey," *arXiv:2303.11504*, 2023.

[66] S. Wang, Z. Wei, J. Xu, and Z. Fan, "Unifying structure reasoning and language model pre-training for complex reasoning," *arXiv:2301.08913*, 2023.

[67] X. Zhang, A. Bosselut, M. Yasunaga, H. Ren, P. Liang, C. D. Manning, and J. Leskovec, "Greaselm: Graph reasoning enhanced language models for question answering," *arXiv:2201.08860*, 2022.

[68] R. Ding, X. Han, and L. Wang, "A unified knowledge graph augmentation service for boosting domain-specific nlp tasks," in *Find. Assoc. Comput. Linguist.: ACL 2023*, 2023, pp. 353–369.

[69] J. Baek, A. F. Aji, and A. Saffari, "Knowledge-augmented language model prompting for zero-shot knowledge graph question answering," in *Proc. 1st Workshop Nat. Lang. Reasoning Struct. Expl.*, 2023, pp. 78–106.

[70] S. Gehman, S. Gururangan, M. Sap, Y. Choi, and N. A. Smith, "RealToxicityPrompts: Evaluating neural toxic degeneration in language models," in *Find. Assoc. Comput. Linguist.: EMNLP 2020*, 2020, pp. 3356–3369.

[71] R. Zellers, A. Holtzman, E. Clark, L. Qin, A. Farhadi, and Y. Choi, "TuringAdvice: A generative and dynamic evaluation of language use," in *Proc. 2021 Conf. North Am. Chapter Assoc. Comput. Linguist.: Hum. Lang. Technol.*, 2021, pp. 4856–4880.

[72] V. Swamy, A. Romanou, and M. Jaggi, "Interpreting language models through knowledge graph extraction," *arXiv:2111.08546*, 2021.

[73] W. Liu, P. Zhou, Z. Zhao, Z. Wang, Q. Ju, H. Deng, and P. Wang, "K-bert: Enabling language representation with knowledge graph," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 2901–2908.

[74] T. Sun, Y. Shao, X. Qiu, Q. Guo, Y. Hu, X. Huang, and Z. Zhang, "CoLAKE: Contextualized language and knowledge embedding," in *Proc. 28th Int. Conf. Comput. Linguistics*, 2020, pp. 3660–3670.

[75] Y. Zhang, J. Lin, Y. Fan, P. Jin, Y. Liu, and B. Liu, "Cn-hit-it. nlp at semeval-2020 task 4: Enhanced language representation with multiple knowledge triples," in *Proc. 14th Workshop Semant. Eval.*, 2020, pp. 494–500.

[76] I. Yamada, A. Asai, H. Shindo, H. Takeda, and Y. Matsumoto, "LUKE: Deep contextualized entity representations with entity-aware self-attention," in *Proc. 2020 Conf. Empir. Methods Nat. Lang. Process.*, 2020, pp. 6442–6454.

[77] N. Poerner, U. Waltinger, and H. Schütze, "E-BERT: Efficient-yet-effective entity embeddings for BERT," in *Find. Assoc. Comput. Linguist.: EMNLP 2020*, 2020, pp. 803–818.

[78] C. Rosset, C. Xiong, M. Phan, X. Song, P. Bennett, and S. Tiwary, "Knowledge-aware language model pretraining," *arXiv:2007.00655*, 2020.

[79] X. Liu, D. Yin, J. Zheng, X. Zhang, P. Zhang, H. Yang, Y. Dong, and J. Tang, "Oag-bert: Towards a unified backbone language model for academic knowledge services," in *Proc. 28th ACM SIGKDD Conf. Knowl. Discov. Data Min.*, 2022, pp. 3418–3428.

[80] T. Zhang, C. Wang, N. Hu, M. Qiu, C. Tang, X. He, and J. Huang, "Dkplm: Decomposable knowledge-enhanced pre-trained language model for natural language understanding," in *Proc. AAAI Conf. Artif. Intell.*, 2022, pp. 11 703–11 711.

[81] Z.-X. Ye, Q. Chen, W. Wang, and Z.-H. Ling, "Align, mask and select: A simple method for incorporating commonsense knowledge into language representation models," *arXiv:1908.06725*, 2019.

[82] W. Chen, Y. Su, X. Yan, and W. Y. Wang, "KGPT: Knowledge-grounded pre-training for data-to-text generation," in *Proc. 2020 Conf. Empir. Methods Nat. Lang. Process.*, 2020, pp. 8635–8648.

[83] R. Logan, N. F. Liu, M. E. Peters, M. Gardner, and S. Singh, "Barack's wife hillary: Using knowledge graphs for fact-aware language modeling," in *Proc. 57th Annu. Meet. Assoc. Comput. Linguist.*, 2019, pp. 5962–5971.

[84] M. Sap, R. Le Bras, E. Allaway, C. Bhagavatula, N. Lourie, H. Rashkin, B. Roof, N. A. Smith, and Y. Choi, "Atomic: An atlas of machine commonsense for if-then reasoning," in *Proc. 33rd AAAI Conf. Artif. Intell. & 31st Innov. Appl. Artif. Intell. Conf. & 9th AAAI Symp. Educ. Adv. Artif. Intell.*, 2019, p. 3027–3035.

[85] X. Wang, T. Gao, Z. Zhu, Z. Zhang, Z. Liu, J. Li, and J. Tang, "Kepler: A unified model for knowledge embedding and pre-trained language representation," *Trans. Assoc. Comput. Linguist.*, vol. 9, pp. 176–194, 2021.

[86] W. Xiong, J. Du, W. Y. Wang, and V. Stoyanov, "Pretrained encyclopedia: Weakly supervised knowledge-pretrained language model," *arXiv:1912.09637*, 2019.

[87] Y. Sun, S. Wang, S. Feng, S. Ding, C. Pang, J. Shang, J. Liu, X. Chen, Y. Zhao, Y. Lu *et al.*, "Ernie 3.0: Large-scale knowledge enhanced pre-training for language understanding and generation," *arXiv:2107.02137*, 2021.

[88] B. He, D. Zhou, J. Xiao, X. Jiang, Q. Liu, N. J. Yuan, and T. Xu, "BERT-MK: Integrating graph contextualized knowledge into pre-trained language models," in *Find. Assoc. Comput. Linguist.: EMNLP 2020*, 2020, pp. 2281–2290.

[89] Y. Su, X. Han, Z. Zhang, Y. Lin, P. Li, Z. Liu, J. Zhou, and M. Sun, "Cokebert: Contextual knowledge selection and embedding towards enhanced pre-trained language models," *AI Open*, vol. 2, pp. 127–134, 2021.

[90] Y. Sun, Q. Shi, L. Qi, and Y. Zhang, "JointLK: Joint reasoning with language models and knowledge graphs for commonsense question answering," in *Proc. 2022 Conf. North Am. Chapter Assoc. Comput. Linguist.: Hum. Lang. Technol.*, 2022, pp. 5049–5060.

[91] P. Zhong, D. Wang, and C. Miao, "Knowledge-enriched transformer for emotion detection in textual conversations," in *Proc. 2019 Conf. Empir. Methods Nat. Lang. Process. & 9th Int. Joint Conf. Nat. Lang. Process.*, 2019, pp. 165–176.

[92] Q. Liu, D. Yogatama, and P. Blunsom, "Relational memory-augmented language models," *Trans. Assoc. Comput. Linguist.*, vol. 10, pp. 555–572, 2022.

[93] M. Yasunaga, H. Ren, A. Bosselut, P. Liang, and J. Leskovec, "QA-GNN: Reasoning with language models and knowledge graphs for question answering," in *Proc. 2021 Conf. North Am. Chapter Assoc. Comput. Linguist.: Hum. Lang. Technol.*, Online, 2021, pp. 535–546.

[94] L. He, S. Zheng, T. Yang, and F. Zhang, "KLMo: Knowledge graph enhanced pretrained language model with fine-grained relationships," in *Find. Assoc. Comput. Linguist.: EMNLP 2021*, 2021, pp. 4536–4542.

[95] M. E. Peters, M. Neumann, R. Logan, R. Schwartz, V. Joshi, S. Singh, and N. A. Smith, "Knowledge enhanced contextual word representations," in *Proc. 2019 Conf. Empir. Methods Nat. Lang. Process. & 9th Int. Joint Conf. Nat. Lang. Process.*, 2019, pp. 43–54.

[96] D. Yu, C. Zhu, Y. Yang, and M. Zeng, "Jaket: Joint pre-training of knowledge graph and language understanding," in *Proc. AAAI Conf. Artif. Intell.*, 2022, pp. 11630–11638.

[97] Y. Liu, Y. Wan, L. He, H. Peng, and S. Y. Philip, "Kg-bart: Knowledge graph-augmented bart for generative commonsense reasoning," in *Proc. AAAI Conf. Artif. Intell.*, 2021, pp. 6418–6425.

[98] R. Wang, D. Tang, N. Duan, Z. Wei, X. Huang, G. Cao, D. Jiang, and M. Zhou, "K-adapter: Infusing knowledge into pre-trained models with adapters," in *Proc. Joint Conf. 59th Annu. Meet. Assoc. Comput. Linguist. and 11th Int. Joint Conf. Nat. Lang. Process.*, 2021, p. 1405–1418.

[99] A. Lauscher, O. Majewska, L. F. R. Ribeiro, I. Gurevych, N. Rozanov, and G. Glavaš, "Common sense or world knowledge? investigating adapter-based knowledge injection into pretrained transformers," in *Proc. DeeLIO: 1st Workshop Knowl. Extract. Integr. Deep Learn. Archit.*, 2020, pp. 43–49.

[100] Q. Lu, D. Dou, and T. H. Nguyen, "Parameter-efficient domain knowledge integration from multiple sources for biomedical pre-trained language models," in *Find. Assoc. Comput. Linguist.: EMNLP 2021*, 2021, pp. 3855–3865.

[101] G. Lu, H. Yu, Z. Yan, and Y. Xue, "Commonsense knowledge graph-based adapter for aspect-level sentiment classification," *Neurocomput.*, vol. 534, pp. 67–76, 2023.

[102] Y. Levine, B. Lenz, O. Dagan, O. Ram, D. Padnos, O. Sharir, S. Shalev-Shwartz, A. Shashua, and Y. Shoham, "SenseBERT: Driving some sense into BERT," in *Proc. 58th Annu. Meet. Assoc. Comput. Linguist.*, 2020, pp. 4656–4667.

[103] Y. Qin, Y. Lin, R. Takanobu, Z. Liu, P. Li, H. Ji, M. Huang, M. Sun, and J. Zhou, "ERICA: Improving entity and relation understanding for pre-trained language models via contrastive learning," in *Proc. 59th Ann. Meet. Assoc. Comput. Ling. Int. Jt. Conf. Nat. Lang. Process.*, 2021, pp. 3350–3363.

[104] P. Ke, H. Ji, S. Liu, X. Zhu, and M. Huang, "SentiLARE: Sentiment-aware language representation learning with linguistic knowledge," in *Proc. 2020 Conf. Empir. Methods Nat. Lang. Process.*, Online, 2020, pp. 6975–6988.

[105] Y. Sun, S. Wang, Y. Li, S. Feng, H. Tian, H. Wu, and H. Wang, "Ernie 2.0: A continual pre-training framework for language understanding," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 8968–8975.

[106] M. Yasunaga, A. Bosselut, H. Ren, X. Zhang, C. D. Manning, P. S. Liang, and J. Leskovec, "Deep bidirectional language-knowledge graph pretraining," in *Adv. Neural Inform. Process. Syst.*, 2022, pp. 37309–37323.

[107] H. Hayashi, Z. Hu, C. Xiong, and G. Neubig, "Latent relation language models," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 7911–7918.

[108] M. Kang, J. Baek, and S. J. Hwang, "KALA: knowledge-augmented language model adaptation," in *Proc. 2022 Conf. North Am. Chapter Assoc. Comput. Linguist.: Hum. Lang. Technol.*, 2022, pp. 5144–5167.

[109] Q. Xie, J. A. Bishop, P. Tiwari, and S. Ananiadou, "Pre-trained language models with domain knowledge for biomedical extractive summarization," *Knowl. Based Syst.*, vol. 252, p. 109460, 2022.

[110] B. Y. Lin, X. Chen, J. Chen, and X. Ren, "KagNet: Knowledge-aware graph networks for commonsense reasoning," in *Proc. 2019 Conf. Empir. Methods Nat. Lang. Process. & 9th Int. Joint Conf. Nat. Lang. Process.*, 2019, pp. 2829–2839.

[111] H. Fei, Y. Ren, Y. Zhang, D. Ji, and X. Liang, "Enriching contextualized language model from knowledge graph for biomedical information extraction," *Brief. Bioinform.*, vol. 22, no. 3, p. bbaa110, 2021.

[112] T.-Y. Chang, Y. Liu, K. Gopalakrishnan, B. Hedayatnia, P. Zhou, and D. Hakkani-Tur, "Incorporating commonsense knowledge graph in pretrained models for social commonsense tasks," in *Proc. DeeLIO: 1st Workshop Knowl. Extract. Integr. Deep Learn. Archit.*, Nov. 2020, pp. 74–79.

[113] B. R. Andrus, Y. Nasiri, S. Cui, B. Cullen, and N. Fulda, "Enhanced story comprehension for large language models through dynamic document-based knowledge graphs," in *Proc. AAAI Conf. Artif. Intell.*, 2022, pp. 10436–10444.

[114] J. Wang, W. Huang, Q. Shi, H. Wang, M. Qiu, X. Li, and M. Gao, "Knowledge prompting in pre-trained language model for natural language understanding," in *Proc. 2022 Conf. Empir. Methods Nat. Lang. Process.*, 2022, pp. 3164–3177.

[115] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," in *Adv. Neural Inf. Process. Syst.*, 2013.

[116] A. Wettig, T. Gao, Z. Zhong, and D. Chen, "Should you mask 15% in masked language modeling?" *arXiv:2202.08005*, 2022.

[117] Z. Bi, N. Zhang, Y. Jiang, S. Deng, G. Zheng, and H. Chen, "When do program-of-thoughts work for reasoning?" *arXiv:2308.15452*, 2023.

[118] T. Schick and H. Schütze, "Exploiting cloze-questions for few-shot text classification and natural language inference," in *Proc. 16th Conf. Eur. Chapter Assoc. Comput. Linguist.*, 2021, pp. 255–269.

[119] N. Bian, X. Han, B. Chen, and L. Sun, "Benchmarking knowledge-enhanced commonsense question answering via knowledge-to-text transformation," in *Proc. AAAI Conf. Artif. Intell.*, 2021, pp. 12574–12582.

[120] J. Giorgi, X. Wang, N. Sahar, W. Y. Shin, G. D. Bader, and B. Wang, "End-to-end named entity recognition and relation extraction using pre-trained language models," *arXiv:1912.13415*, 2019.

[121] Z. Yuan, Y. Liu, C. Tan, S. Huang, and F. Huang, "Improving biomedical pretrained language models with knowledge," in *Proc. BioNLP 2021 workshop*, 2021, pp. 180–190.

[122] D. Seyler, T. Dembelova, L. Del Corro, J. Hoffart, and G. Weikum, "A study of the importance of external knowledge in the named entity recognition task," in *Proc. 56th Ann. Meet. Assoc. Comput. Linguistics.*, 2018, pp. 241–246.

[123] Q. He, L. Wu, Y. Yin, and H. Cai, "Knowledge-graph augmented word representations for named entity recognition," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 7919–7926.

[124] Y. Song, W. Zhang, Y. Ye, C. Zhang, and K. Zhang, "Knowledge-enhanced relation extraction in chinese emrs," in *Proc. 2022 5th Int. Conf. Mach. Learn. Nat. Lang. Process.*, 2023, p. 196–201.

[125] J. Li, Y. Katsis, T. Baldwin, H.-C. Kim, A. Bartko, J. McAuley, and C.-N. Hsu, "Spot: Knowledge-enhanced language representations for information extraction," in *Proc. 31st ACM Int. Conf. Inf. Knowl. Manage.*, 2022, p. 1124–1134.

[126] A. Roy and S. Pan, "Incorporating medical knowledge in BERT for clinical relation extraction," in *Proc. 2021 Conf. Empir. Methods Nat. Lang. Process.*, 2021, pp. 5357–5366.

[127] J. Zhou, J. Tian, R. Wang, Y. Wu, W. Xiao, and L. He, "SentiX: A sentiment-aware pre-trained model for cross-domain sentiment analysis," in *Proc. 28th Int. Conf. Comput. Linguist.*, 2020, pp. 568–579.

[128] Q. Wang, X. Cao, J. Wang, and W. Zhang, "Knowledge-aware collaborative filtering with pre-trained language model for personalized review-based rating prediction," *IEEE Trans. Knowl. Data Eng.*, pp. 1–13, 2023.

[129] S. Liang, J. Shao, D. Zhang, J. Zhang, and B. Cui, "Drgi: Deep relational graph infomax for knowledge graph completion," *IEEE Trans. Knowl. Data Eng.*, vol. 35, no. 3, pp. 2486–2499, 2023.

[130] Q. Lin, R. Mao, J. Liu, F. Xu, and E. Cambria, "Fusing topology contexts and logical rules in language models for knowledge graph completion," *Inf. Fusion*, vol. 90, pp. 253–264, 2023.

[131] W. Li, R. Peng, and Z. Li, "Knowledge graph completion by jointly learning structural features and soft logical rules," *IEEE Trans. Knowl. Data Eng.*, vol. 35, no. 3, pp. 2724–2735, 2023.

[132] A. Ghanbarpour and H. Naderi, "An attribute-specific ranking method based on language models for keyword search over graphs," *IEEE Trans. Knowl. Data Eng.*, vol. 32, no. 1, pp. 12–25, 2020.

[133] Z. Hu, Y. Xu, W. Yu, S. Wang, Z. Yang, C. Zhu, K.-W. Chang, and Y. Sun, "Empowering language models with knowledge graph reasoning for open-domain question answering," in *Proc. 2022 Conf. Empir. Methods Nat. Lang. Process.*, 2022, pp. 9562–9581.

[134] J. Guan, F. Huang, Z. Zhao, X. Zhu, and M. Huang, "A knowledge-enhanced pretraining model for commonsense story generation," *Trans. Assoc. Comput. Linguist.*, vol. 8, pp. 93–108, 2020.

[135] H. Ji, P. Ke, S. Huang, F. Wei, X. Zhu, and M. Huang, "Language generation with multi-hop reasoning on commonsense knowledge graph," in *Proc. 2020 Conf. Empir. Methods Nat. Lang. Process.*, 2020, pp. 725–736.

[136] X. Yang and I. Tiddi, "Creative storytelling with language models and knowledge graphs," in *Proc. CIKM 2020 Workshops*, 2020.

[137] L. Du, X. Ding, T. Liu, and B. Qin, "Learning event graph knowledge for abductive reasoning," in *Proc. 59th Annu. Meet. Assoc. Comput. Linguist. and 11th Int. Joint Conf. Nat. Lang. Process.*, 2021, pp. 5181–5190.

[138] B. AlKhamissi, M. Li, A. Celikyilmaz, M. Diab, and M. Ghazvininejad, "A review on language models as knowledge bases," *arXiv:2204.06031*, 2022.

[139] A. Talmor, Y. Elazar, Y. Goldberg, and J. Berant, "olmpics-on what language model pre-training captures," *Trans. Assoc. Comput. Linguist.*, vol. 8, pp. 743–758, 2020.

[140] C. Wang, S. Liang, Y. Zhang, X. Li, and T. Gao, "Does it make sense? and why? a pilot study for sense making and explanation," in *Proc. 57th Ann. Meet. Assoc. Comput. Linguistics.*, 2019, pp. 4020–4026.

[141] M. Sung, J. Lee, S. Yi, M. Jeon, S. Kim, and J. Kang, "Can language models be biomedical knowledge bases?" in *Proc. 2021 Conf. Empir. Methods Nat. Lang. Process.*, 2021, pp. 4723–4734.

[142] R. Zhao, F. Zhao, G. Xu, S. Zhang, and H. Jin, "Can language models serve as temporal knowledge bases?" in *Find. Assoc. Comput. Linguist.: EMNLP 2022*, 2022, pp. 2024–2037.

[143] N. Kassner, P. Dufter, and H. Schütze, "Multilingual lama: Investigating knowledge in multilingual pretrained language models," *arXiv:2102.00894*, 2021.

[144] T. Pires, E. Schlinger, and D. Garrette, "How multilingual is multilingual BERT?" in *Proc. 57th Ann. Meet. Assoc. Comput. Linguistics.*, 2019, pp. 4996–5001.

[145] B. Cao, H. Lin, X. Han, F. Liu, and L. Sun, "Can prompt probe pretrained language models? understanding the invisible risks from a causal view," *arXiv:2203.12258*, 2022.

[146] N. Kandpal, H. Deng, A. Roberts, E. Wallace, and C. Raffel, "Large language models struggle to learn long-tail knowledge," *arXiv:2211.08411*, 2022.

[147] S. Hao, B. Tan, K. Tang, H. Zhang, E. P. Xing, and Z. Hu, "Bert-net: Harvesting knowledge graphs from pretrained language models," *arXiv:2206.14268*, 2022.

[148] A. Bosselut, H. Rashkin, M. Sap, C. Malaviya, A. Celikyilmaz, and Y. Choi, "COMET: Commonsense transformers for automatic knowledge graph construction," in *Proc. 57th Ann. Meet. Assoc. Comput. Linguistics.*, 2019, pp. 4762–4779.

[149] O. Ram, Y. Levine, I. Dalmedigos, D. Muhlgay, A. Shashua, K. Leyton-Brown, and Y. Shoham, "In-context retrieval-augmented language models," *arXiv:2302.00083*, 2023.

[150] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu *et al.*, "Palm-e: An embodied multimodal language model," *arXiv:2303.03378*, 2023.

[151] Y. Hou, G. Fu, and M. Sachan, "Understanding the integration of knowledge in language models with graph convolutions," *arXiv:2202.00964*, 2022.

[152] H. Schuff, H.-Y. Yang, H. Adel, and N. T. Vu, "Does external knowledge help explainable natural language inference? automatic evaluation vs. human ratings," in *Proc. 4th BlackboxNLP Workshop on Analyz. Interpr. Neural Networks NLP*, 2021, pp. 26–41.

[153] Z. Chen, A. K. Singh, and M. Sra, "Lmexplainer: a knowledge-enhanced explainer for language models," *arXiv:2303.16537*, 2023.