BLADE: Enhancing Black-box Large Language Models with Small Domain-Specific Models

Haitao Li DCST, Tsinghua University Quan Cheng Laboratory liht22@mails.tsinghua.edu.cn

Qian Dong DCST, Tsinghua University Quan Cheng Laboratory dq22@mails.tsinghua.edu.cn Qingyao Ai* DCST, Tsinghua University Quan Cheng Laboratory aiqy@tsinghua.edu.cn

Zhijing Wu Beijing Institute of Technology zhijingwu.bit.edu.cn Jia Chen DCST, Tsinghua University Quan Cheng Laboratory chenjia0831@gmail.com

Yiqun Liu DCST, Tsinghua University Zhongguancun Laboratory yiqunliu@tsinghua.edu.cn

Chong Chen Huawei Cloud BU chenchong55@huawei.com

ABSTRACT

Large Language Models (LLMs) like ChatGPT and GPT-4 are versatile and capable of addressing a diverse range of tasks. However, general LLMs, which are developed on open-domain data, may lack the domain-specific knowledge essential for tasks in vertical domains, such as legal, medical, etc. To address this issue, previous approaches either conduct continuous pre-training with domainspecific data or employ retrieval augmentation to support general LLMs. Unfortunately, these strategies are either cost-intensive or unreliable in practical applications. To this end, we present a novel framework named BLADE, which enhances Black-box LArge language models with small Domain-spEcific models. BLADE consists of a black-box LLM and a small domain-specific LM. The small LM preserves domain-specific knowledge and offers specialized insights, while the general LLM contributes robust language comprehension and reasoning capabilities. Specifically, our method involves three steps: 1) pre-training the small LM with domainspecific data, 2) fine-tuning this model using knowledge instruction data, and 3) joint Bayesian optimization of the general LLM and the small LM. Extensive experiments conducted on public legal and medical benchmarks reveal that BLADE significantly outperforms existing approaches. This shows the potential of BLADE as an effective and cost-efficient solution in adapting general LLMs for vertical domains.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Qi Tian Huawei Cloud BU tian.qi1@huawei.com

KEYWORDS

Large Language Models, Domain Adaptation, Bayesian Optimization

ACM Reference Format:

Haitao Li, Qingyao Ai, Jia Chen, Qian Dong, Zhijing Wu, Yiqun Liu, Chong Chen, and Qi Tian. 2024. BLADE: Enhancing Black-box Large Language Models with Small Domain-Specific Models. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 11 pages. https://doi.org/10.1145/nnnnnnn.nnnnnnn

1 INTRODUCTION

Recently, large language models (LLMs) have attracted considerable attention in both academia and industry [15, 40, 55]. These models, driven by expansive neural networks and trained on extensive data sets, exhibit remarkable ability in comprehending and generating natural language. The wide application of LLMs trained with open-domain data, denoted in this paper as *General LLMs*, has profoundly impacted various aspects of daily life and professional environments.

Despite their superior capabilities, large language models often face challenges in vertical domains (e.g., medicine, legal) where in-depth, domain-specific knowledge is crucial [7, 9]. For instance, as shown in this paper, ChatGPT exhibits suboptimal performance in Chinese legal question-answering tasks due to its limited knowledge of the Chinese legal system. Therefore, how to adapt general LLMs for domain-specific applications has become an important problem for the research community [3, 9].

Existing methods for adapting general LLMs to specific domains can be broadly divided into two main categories: domain data continuous pre-training and retrieval augmentation. Continuous pre-training involves infusing domain knowledge into general LLMs by training them on a domain-specific corpus [2, 44]. While straightforward, this paradigm requires direct access to large-scale domain data and LLM parameters, which are not available in many conditions. Also, even with access to general LLM parameters and sufficient domain-specific data, directly tuning a general LLM (such

^{*}Corresponding author

as GPT-4) can be prohibitively expensive and poses a risk of over-fitting. Aware of these challenges, researchers propose retrieval augmentation as a new paradigm, aiming to enhance general LLMs by leveraging their in-context learning ability [45]. It involves first using a text retriever to find relevant content from the domain corpus, which is then incorporated into the LLM's input to aid in understanding domain-specific knowledge. However, there may exist two problems in this paradigm. First, retrievers primarily rely on exact matches or semantic similarity, lacking inferential capabilities. This limitation means they may not always retrieve documents that fully address specific queries. Second, the retrieved documents may include irrelevant or misleading information, which could adversely affect the performance of LLMs.

When humans face questions in new domains, besides taking classes (i.e., continuous pre-training) or conducting online searches via platforms like Google (i.e., retrieval augmentation), a more direct and practical approach is to seek advice from experts possessing domain-specific knowledge. With this idea in mind, we present BLADE, a novel paradigm where the general LLM is viewed as a black box and the small domain-specific LM (#parameters < 3B) is added as a tuneable module. BLADE leverages the superior language comprehension and logical reasoning capabilities of the general LLM, while also incorporating the domain-specific expertise and precision provided by the smaller, domain-focused LM. This approach includes Domain-specific Pretraining (DP) of the smaller LM and introduces two strategies: Knowledge Instruction Tuning (KIT) and Bayesian Prompted Optimization (BPO). Knowledge Instruction Tuning leverages general LLMs to generate pseudo data, which refines the smaller LM, enabling it to generate answers tailored to specific queries. Then, the Bayesian Prompted Optimization aligns the output of small LMs with general LLMs using derivative-free optimization on soft embeddings.

To validate the effectiveness of BLADE, we conduct experiments on question-answering tasks in legal and medical fields, which require in-depth knowledge and strong reasoning abilities. The experiments show that BLADE can improve the performance of diverse general LLMs across legal and medical benchmarks. Compared with existing retrieval augmentation methods, the domain-specific small LM can generate more in-depth, comprehensive, and contextually appropriate external knowledge. This capability significantly improves the application of general LLMs in specialized domains. Overall, the principal contributions of this paper can be summarized as follows:

- (1) We introduce BLADE, a new framework for adapting general LLMs to specific domains. This method involves training a smaller, domain-specific Language Model (LM) that aids the general LLM in excelling at tasks within its respective domain.
- (2) We propose Knowledge Instruction Tuning (KIT) and Bayesian Prompted Optimization (BPO). These strategies enhance the small LM's adaptation to general LLMs, resulting in improved performance.
- (3) We conduct extensive experiments on public legal and medical benchmarks. Our method significantly outperforms existing approaches that rely on continuous pre-training or retrieval augmentation in both professional fields.

2 RELATED WORK

2.1 Large Language Models

Recently, the emergence of LLMs has attracted substantial attention and research efforts across various fields, primarily owing to their superior linguistic capability. The foundation for modern LLMs is the Transformer architecture [10, 48], which is fundamental in major language models like BERT [13] and GPT [42]. BERT [13] achieves promising results by pre-training on large corpus and fine-tuning on downstream tasks [31, 32, 52]. Similarly, GPT [42] introduces auto-regressive language modeling, generating text sequentially based on previously produced tokens. Both BERT and GPT series models require extensive pre-training on a massive corpus to acquire general linguistic knowledge, leveraging tasks such as masked language modeling [13], next token prediction [42], etc. However, infusing domain-specific knowledge during the pretraining phase of LLMs poses challenges due to the substantial data requirements. In specific domains, the availability of sufficient data for effective pre-training is often limited, thus constraining LLMs' applicability. This limitation is particularly acute in fields with stringent data privacy regulations, such as legal and medical domains. Another crucial training stage for LLMs is Supervised Fine-Tuning (SFT), which involves training LLMs using task-specific datasets with labeled examples. This stage adapts the general linguistic knowledge acquired during pre-training to specific tasks, such as sentiment analysis [17], text classification [16, 18], and dialogue systems [14, 41]. Despite the effectiveness of SFT, high-quality taskspecific annotation is usually costly. This cost is further amplified by the scale of parameters in LLMs, making the adaptation of LLMs to specific domains during this stage financially prohibitive. In this paper, we investigate how to leverage existing data to adapt LLMs to specific domains effectively, without the need of additional LLM pre-training or extra annotations.

2.2 Domain adaptation of LLMs

The domain adaptation of LLMs is an extensively researched field. Researchers have explored methods such as continuous pre-training, and retrieval augmentation to improve the performance of LLMs in a specific domain [2, 9, 22, 44, 45]. The most intuitive approach is continuous pre-training a language model on domain-specific corpora. Previous research has focused on data selection [2, 22] as well as adjusting or extending tokenizers [44] to achieve superior performance in the target domain. Despite being effective, fully training these models is often impractical due to extensive computational costs. Additionally, high-quality LLMs can only accessed through the inference API as black boxes. One possible alternative is to answer domain-specific questions by retrieving relevant information from a specific knowledge base [6, 29, 45]. Retrieval augmentation has been shown to be effective in improving performance on various tasks. For instance, RETRO [6] modifies the model architecture to incorporate retrieved text. Furthermore, RE-PLUG [45] treats the language model as a black box and enhances it using a retrieval model. Additionally, recent research has explored substituting traditional document retrievers with large language model generators [34, 47, 54]. In this paper, we focus on employing smaller language models to tackle domain adaptation challenges

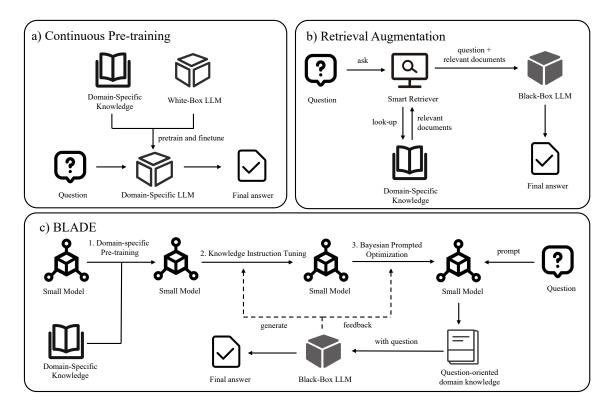


Figure 1: Comparison of the workflow of BLADE with existing domain adaptation methods. There are three steps in BLADE: (1) Domain-specific Pre-training imparts domain knowledge to the small LM. (2) Knowledge Instruction Tuning, which enhances the small LM's ability to follow instructions, thereby sharpening its capacity to produce precise, question-specific knowledge. (3) Bayesian Prompted Optimization contributes to aligning the output of small LM with the comprehension of black-box LLM.

and further propose methods to adapt the small models to the large ones.

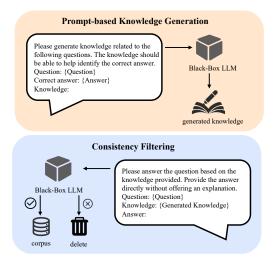


Figure 2: The process of generating data for Knowledge Instruction Tuning. Only knowledge that can help the blackbox LLM correctly answer a question is reserved.

3 METHOD

In this section, we describe our approach in detail. First, we introduce the overview of BLADE. After that, we elaborate on Domain-specific Pre-training (DP), Knowledge Instruction Tuning (KIT), and Bayesian Prompted Optimization (BPO).

3.1 Overview

Figure 1 illustrates the comparison between BLADE and existing domain adaptation methods. Compared to the paradigm of continuous pre-training and retrieval augmentation, BLADE solves domain-specific problems through a collaborative approach between general black-box LLMs and small white-box LMs. General black-box LLMs, such as GPT-4 and GLM-130B, excel in reasoning and inference but often present challenges in terms of cost and feasibility for fine-tuning in downstream applications. Conversely, small white-box LMs may lack sufficient reasoning ability, but can easily update and memorize domain-specific knowledge. In the BLADE framework, when faced with a domain-specific query, the small LM initially generates knowledge tailored to the question. Following this, the general LLM synthesizes this information to generate a comprehensive answer.

The training objectives for the small white-box LMs should satisfy two properties. Firstly, these models must effectively memorize domain-specific knowledge. Secondly, they need to effectively communicate with the general black-box LLMs. To achieve this, we

introduce two methodologies: Knowledge Instruction Tuning and Bayesian Prompted Optimization. These techniques markedly improve the interaction and collaboration capabilities of the smaller LMs with the general LLMs.

BLADE has several advantages. First, small LMs have reduced size and can be easily trained to memorize new knowledge. This separation of knowledge memorization from reasoning capabilities aids in better safeguarding private data. Moreover, unlike the shallow interactions (e.g., inner product) between questions and documents in modern dense retrieval models, the small LM in BLADE generates deeper, question-specific knowledge through intricate token-level cross-attention. We believe that BLADE presents a promising approach for adapting general LLMs to specialized domains, offering a solution that is both effective in performance and cost-efficient.

3.2 Domain-specific Pre-training (DP)

A variety of studies have shown that pre-trained language models implicitly contain substantial knowledge [26, 28, 36, 49]. This knowledge can be elicited from the language model through instructional prompts [36]. In light of this, we incorporate domain-specific knowledge via Domain-specific Pre-training (DP). Specifically, given domain-specific unsupervised text $T = \{t_1, t_2,, t_n\}$, we optimize the model by maximizing the following training objective:

$$G(T) = \sum_{i} \log P\left(t_{i} \mid t_{i-k}, \dots, t_{i-1}; \Theta\right),\,$$

where Θ is the parameter of the model. P is the conditional probability of generating the current token based on the previous tokens. After Domain-specific Pre-training, the domain-specific knowledge is effectively encoded into the parameters of the smaller language model.

3.3 Knowledge Instruction Tuning (KIT)

While DP enhances small LMs with domain-specific knowledge, its principal purpose is to extend the existing text, rather than actively generating valuable knowledge tailored to the specific queries. To tackle this limitation, we introduce Knowledge Instruction Tuning (KIT) which focuses on enhancing the instruction-following abilities of the small LM. This process enables the small LM to leverage its inherent knowledge for particular tasks, aligning the model with more practical applications.

To be specific, KIT consists of three components: Prompt-based Knowledge Generation, Consistency Filtering, and Instruction Tuning. Figure 2 show the process of generating data of KIT. During Prompt-based Knowledge Generation, we combine question-answer pairs from the training dataset with task-specific prompts. Subsequently, a general black-box LLM is employed to generate knowledge that assists in accurately answering the questions. The process entails instructing the model using precise answers and promoting reverse reasoning to deduce the underlying reasons. To enhance the trustworthiness of the generated knowledge, retrieval models can be employed to supply relevant information. Prompt-based Knowledge Generation is based on the intuition that the general black-box LLM always generates knowledge in a format and style that aligns with its own preferences. Similar to people,

we tend to convey our knowledge in a manner that is readily comprehensible to us.

Afterwards, we improve the quality of the generated knowledge by ensuring round-trip consistency. This involves verifying that the general black-box LLM can accurately respond to queries based on the knowledge it produces. Consistency Filtering has been shown to be effective for query generation in QA tasks and various information retrieval tasks [12, 30]. The instruction prompts are structured as depicted in Figure 2.

We filter the generated knowledge and only retain the information that can help the LLM to answer questions correctly. The refined data is then used to fine-tune the small LM, employing the same prompt template as in the Prompt-based Knowledge Generation stage. After KIT, the small LM acquires the ability to generate specific knowledge to the questions.

3.4 Bayesian Prompted Optimization (BPO)

In this section, we try to align the output of the small LM with the understanding of the general LLM. This approach draws inspiration from human society, where pre-collaborative training often leads to better effectiveness. Different from previous work [8, 46], our objective is not to train better black-box general LLMs or to manually craft task-specific prompts. Instead, the primary aim here is to fine-tune the small LM so that it aligns with the general LLM's preferences.

Figure 3 illustrates the detailed process of Bayesian Prompted Optimization (BPO) where only soft embeddings are trainable. To be specific, the primary optimization objective is to enhance the performance of the general $\mathrm{LLM}f(\cdot)$ on domain-specific tasks. Consider an example (X,Y) from the dataset \mathcal{T}_t . k represents the domain knowledge that is specific to the query X. In our framework, knowledge k is generated by the small domain-specific $\mathrm{LM}\ g(\cdot)$. $h(\cdot,\cdot)$ is defined as the evaluation metric for output f(k,X) and ground truth Y. For example, in multiple-choice tasks, $h(\cdot,\cdot)$ can be accuracy. The optimization objective is to maximize the performance with appropriate knowledge, i.e.,

$$\max_{L} \mathbb{E}_{(X,Y) \sim \mathcal{T}_t} h(f([k;X]), Y), \text{ s.t. } k = g(X),$$

As discussed in [8], the above problems are combinatorial optimization with complicated structural constraints. Since $f(\cdot)$ is a black-box model, traditional optimization via backpropagation is not feasible. Therefore, we apply derivative-free optimization to refine the soft prompt p_h on small model $g(\cdot)$. Specifically, we concatenate n soft tokens $p_{h_1:h_n} \in \mathbb{R}^D$ with input queries X as inputs to the small model, facilitating the generation of domain knowledge $k = g(p_{h_1:h_n}, X)$. Therefore, our objective is to identify the optimal soft prompt:

$$p_h^* = \arg \max_{p_h \in \mathbb{R}^D} \mathbb{E}_{(X,Y) \sim \mathcal{T}_t} h(f([g(p_h, X); X]), Y).$$

Although the original optimization problem is transformed into a feasible continuous optimization problem, derivative-free black-box optimization remains challenging due to the high dimensionality of the optimized soft prompt. To address this problem, we attempt to optimize a lower dimensional vector $\boldsymbol{p} \in \mathbb{R}^d$ where $d \ll D$ and apply a random projection $A \in \mathbb{R}^{d \times D}$ to project \boldsymbol{p} into the

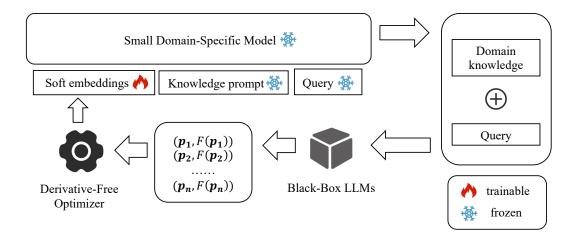


Figure 3: Illustration of the Bayesian Prompted Optimization where only soft embeddings are trainable. F(p) is the objective score corresponding to soft embedding p. In each iteration, the derivative-free optimizer explores new soft embedding based on previous evaluation scores. The knowledge prompt is consistent with the instruction used in the Prompt-based Knowledge Generation stage.

original space. This is feasible based on (1) pre-trained language models possess low intrinsic dimensionality, indicating that the minimal reparameterization required for effective optimization is significantly lower than the full parameter space, as supported by findings in [24, 46]. (2) According to Johnson-Lindenstrauss Lemma [27], the random projection is distance-preserving. This means that the kernel similarity is consistent before and after the random projection. Thus, the optimization objective is transformed into the following formula:

$$\boldsymbol{p^*} = \arg\max_{A\boldsymbol{p} \in \mathbb{R}^D} \mathbb{E}_{(X,Y) \sim \mathcal{T}_t} h(f([g(A\boldsymbol{p},X);X]),Y).$$

We employ Bayesian optimization (BO) [21] to handle the above optimization. BO is an effective technique for updating the posterior distribution of the objective function by iteratively incorporating new sample points. To be specific, we first define the objective function as $F(p) = \mathbb{E}_{(X,Y) \sim \mathcal{T}_t} h(f([g(Ap,X);X]),Y)$. Then, we employ the Gaussian Process (GP) as the prior to estimate the distribution of $F(\cdot)$, i.e.,

$$F(\mathbf{p}) \sim GP(\mu, \sigma^2),$$

where μ is the mean function and σ^2 is the variance function. This GP can be updated iteratively as the optimization process progresses, incorporating new observations to better approximate the true function and reduce uncertainty w.r.t. its behavior. For each p_i , we can obtain a score $F(p_i)$. Let D denote all collected data in previous BO steps, i.e., $D = \{(p_1, F(p_1)), ..., (p_n, F(p_n))\}$. Then the μ and σ^2 of the GP can be updated as follows:

$$\begin{split} \mu(\boldsymbol{p}) &= c(\boldsymbol{p}, \boldsymbol{P}) \left(\boldsymbol{C} + \sigma_n^2 \boldsymbol{I} \right)^{-1} F(\boldsymbol{P}), \\ \sigma^2(\boldsymbol{p}) &= c(\boldsymbol{p}, \boldsymbol{p}) - c(\boldsymbol{p}, \boldsymbol{P}) \left(\boldsymbol{C} + \sigma_n^2 \boldsymbol{I} \right)^{-1} c(\boldsymbol{P}, \boldsymbol{p}), \end{split}$$

where $P = [p_1, ..., p_n]$, $c(\cdot, \cdot)$ is the covariance function and C is the covariance matrix of P. σ_n is the noise variance. I represents the identity matrix. p_1 is randomly initialized. After finishing an

iteration, we employ Expected Improvement (EI) to find the next p_{n+1} . The Expected Improvement (EI) is a popular acquisition function that balances exploration and exploitation. It quantifies the potential improvement over the current best observed value. Formally, the next soft prompt p_{n+1} is defined as follows:

$$\boldsymbol{p}_{\boldsymbol{n}+1} \in \arg\max_{\boldsymbol{p} \in \mathbb{R}^d} \mathbb{E}_{F(\boldsymbol{p}) \sim GP(\mu, \sigma^2)} \left[\max \left\{ 0, F(\boldsymbol{p}) - \max_{i \in [n]} F\left(\boldsymbol{p}_i\right) \right\} \right],$$

In practice, we employ an evolutionary search algorithm known as CMA-ES [23] as the optimization method to identify the most effective soft prompts. When obtaining p_{n+1} , we evaluate performance $F(p_{n+1})$ of BLADE on the training batch. Subsequently, the pair $(p_{n+1}, F(p_{n+1}))$ is incorporated into the D to update μ and σ^2 . This process is performed iteratively until convergence (effectiveness gains less than a threshold for a given number of steps) or reaches the maximum iteration number.

4 EXPERIMENT

In this section, we first introduce our experimental setup, including datasets and metrics, baselines, and implementation details. Then, we report experimental results to demonstrate the effectiveness of BLADE.

4.1 Datasets and Metrics

We conduct extensive experiments with question-answering datasets in the legal and medical domains. Experimental datasets are as follows:

- JEC-QA [58] is the largest Chinese multiple-choice dataset in the legal domain. The legal questions in JEC-QA require high-level reasoning ability and are divided into two types: Knowledge-Driven Questions (KD-questions) and Case-Analysis Questions (CA-questions). There are 26,365 questions in JEC-QA, of which 5,289 of them comprising the test set. It's worth noting that the number of correct options for each question is uncertain.

- CaseHOLD [57] is an English multiple-choice dataset with the purpose of identifying the relevant holding of a cited case. It contains 53,000+ multiple choice questions with 3,600 questions in the test set
- MLEC-QA [33] is the the largest-scale Chinese multi-choice biomedical QA dataset. This dataset contains five subsets: Clinic(Cli), Stomatology (Sto), Public Health (PH), Traditional Chinese Medicine (TCM), and Traditional Chinese Medicine Combined with Western Medicine (CWM), all of them collected from the National Medical Licensing Examination in China. There are 136,236 questions in MLEC-QA, each presenting five options with one correct answer.

We report the zero-shot performance of various LLMs on these datasets. Answers are extracted from model predictions by regular matching. Accuracy serves as the primary evaluation metric. When the correct answer contains more than one option, the model prediction is considered correct only if it exactly matches the golden answer.

4.2 Baselines

We adopt four groups of baselines for comparison: General LLMs, Legal-specific LLMs, Medical-specific LLMs, and Retrieval-augmented LLMs.

- 4.2.1 General LLMs. We consider a range of multilingual general LLMs: ChatGLM-6B [19], ChatGLM2-6B [19], Baichuan-7B/13B-Chat [5], Baichuan2-7B-Chat/13B-Chat [5], Qwen-7B-Chat [4], Chat-GPT [20].
- 4.2.2 Legal-specific LLMs. Legal-specific LLMs are further fine-tuned in the legal corpus to improve the understanding of the law. LaywerLLaMA [25], LexiLaw 1 , ChatLaw-13B/33B [11], are considered for the evaluation. The LawyerLLaMA model is developed based on the Chinese-LLaMA-13B 2 with a combination of general and legal instructions. LexiLaw is fine-tuned based on ChatGLM-6B with legal datasets. Additionally, ChatLaw-13B is refined based on Ziya-LLaMA-13B-v1 [56], and ChatLaw-33B is fine-tuned with Anima-33B 3 .
- 4.2.3 Medical-specific LLMs. For Medical-specific LLMs, Taiyi [38] and Zhongjing [53] are selected as baseline models. Zhongjing is developed through a series of pre-training and fine-tuning processes based on the Ziya-LLaMA-13B-v1 [56]. The model's initial stage, referred to as Zhongjing_base, involves pre-training on an extensive medical corpus. This is followed by Zhongjing_sft, a version that undergoes multiple rounds of supervised fine-tuning based on Zhongjing_base. Taiyi is a bilingual fine-tuned large language model, specifically designed for diverse biomedical tasks, and is continuously trained on the Qwen-7B.
- 4.2.4 Retrieval-augmented LLMs. To fully evaluate the effectiveness of our proposed method, we also compare it with Retrieval-augmented LLMs. We utilize the following retrieval models as baseline: BGE-base [51], M3E-base [50], GTE-base [35], piccolo-base 4 .

These models are advanced text embedding systems capable of converting natural language into dense embeddings suitable for retrieval purposes.

4.3 Pretraining and Implementation Details

We implement BLADE with BLOOMZ_1b7 [1] as the small LM since BLOOMZ is a multi-language model with diverse sizes. To construct the Chinese legal pre-training corpus, we collect legal articles, legal books, legal cases, and other resources from official websites ⁵. The English legal pre-training corpus is derived from the English division of MultiLegalPile [39]. In the medical domain, our corpus comprises medical Wikipedia entries and various medical texts. We pre-train the small LMs for 6 epochs using AdamW [37] optimizer, with a learning rate of 5e-5, batch size of 32, and linear schedule with a warmup ratio of 0.1.

For KIT, we utilized ChatGPT to generate knowledge data, incorporating three manual demonstrations for in-context learning during the Prompt-based Knowledge Generation stage. We finetune the small LM up to 10 epochs using the AdamW [37] optimizer, with a learning rate of 5e-6, batch size of 32, and linear schedule with warmup ratio 0.1. In the process of BPO, we utilize accuracy as the evaluation metric $h(\cdot, \cdot)$. The number of tokens in soft prompts is set to 5. The entries of the random projection matrix A are drawn from a uniform distribution between [-1, 1]. In the BO process, A is fixed. The dimensionality of p_l is set to 10. The maximum number of iterations is set to 50. All models except ChatGPT use greedy decoding with default settings. For each question, the small LM generates only one piece of relevant knowledge. All the experiments in this work are conducted on 8 NVIDIA Tesla A100 GPUs. To facilitate the reproductivity of our results, we will release the source code for our experiments after the reviewing phase.

4.4 Experiment Result

- 4.4.1 Main result in legal domain. Table 1 presents the results from the baselines and BLADE on the JEC-QA dataset. We derive the following observations from the experiment results.
- Legal-specific LLMs show relatively poor results. There is even some performance degradation after domain-specific fine-tuning. For example, LexiLaw underperforms compared to ChatGLM-6B. We hypothesize that although continuous tuning can enhance domain knowledge, it also significantly impacts the model's ability in prompt processing. This observation is also in line with Cheng et al [9]. Furthermore, it's worth noting that the legal-specific LLMs demonstrate substantial improvement when integrated with BLADE, implying that these models may not be fully leveraging the domain knowledge encoded in their parameters.
- BLADE consistently enhances performance across various models. For example, Baichuan-7B achieves 28.4% performance improvement, while ChatGPT achieves 31.3% performance improvement. This indicates that BLADE is applicable to diverse language models with different sizes.
- Overall, BLADE effectively utilizes domain knowledge without affecting the reasoning ability of the original model. It has achieved state-of-the-art results on the Chinese legal question-answering

¹https://github.com/CSHaitao/LexiLaw

²https://github.com/ymcui/Chinese-LLaMA-Alpaca

³https://github.com/lyogavin/Anima

 $^{^4} https://hugging face.co/sensenova/piccolo-base-zh\\$

⁵https://wenshu.court.gov.cn/

Table 1: Zero-shot test accuracy on JEC-QA dataset. BLADE achieves consistent improvements on two types of questions. The gain % shows the relative improvement of methods compared to the original language model. */** denotes that BLADE performs significantly better than the original language model at p < 0.05/0.01 level using the fisher randomization test [43]. Best results are marked bold.

Model	# Parameters	KD	-questions	CA-questions		All	
		Original	+BLADE	Original	+BLADE	Original	+BLADE
Legal Specific LLMs							
LaywerLLaMA	13B	9.76	12.94**(32.6%)	6.05	8.66**(43.1%)	7.45	10.26**(37.7%)
LexiLaw	6B	15.50	19.63**(26.6%)	14.35	18.07**(25.9%)	14.78	18.66**(26.5%)
ChatLaw-13B	13B	10.32	17.32**(67.8%)	5.03	8.08**(60.6%)	7.01	11.55**(64.8%)
ChatLaw-33B	33B	15.66	21.80**(39.2%)	17.01	20.46**(20.3%)	16.50	20.96**(27.0%)
General LLMs	General LLMs						
ChatGLM-6B	6B	17.08	21.19**(24.1%)	16.64	18.62**(11.9%)	16.81	19.58**(16.5%)
ChatGLM2-6B	6B	27.39	30.81**(12.5%)	24.09	26.34**(9.3%)	25.32	28.01**(10.6%)
Qwen-7B-Chat	7B	25.78	31.26**(21.2%)	24.52	25.07*(2.2%)	24.99	27.39**(9.6%)
Baichuan-7B	7B	15.31	21.80**(41.4%)	17.80	21.58**(21.2%)	16.86	21.66**(28.4%)
Baichuan-13B-Chat	13B	17.87	23.06**(14.1%)	19.19	21.71**(13.1%)	18.69	21.21**(13.4%)
Baichuan2-7B-Chat	7B	19.23	24.27**(26.2%)	19.53	21.73**(11.3%)	19.41	22.68**(16.8%)
Baichuan2-13B-Chat	13B	25.78	28.29**(9.73%)	21.80	24.22**(11.1%)	23.29	25.75**(10.5%)
ChatGPT	-	20.53	28.45**(38.6%)	18.70	23.67**(26.6%)	19.38	25.46**(31.3%)

Table 2: Overall Zero-shot performance on the English dataset CaseHOLD. */** denotes that BLADE performs significantly better than baselines at p < 0.05/0.01 level using the fisher randomization test. Best results are marked bold.

Model	CaseHOLD			
Model	Original	+BLADE		
ChatGLM2-6B	48.03	58.19**(21.1%)		
Qwen-7B-Chat	54.28	57.33**(5.6%)		
Baichuan2-7B-Chat	47.53	58.69**(23.5%)		
Baichuan2-13B-Chat	48.69	62.55**(28.5%)		
ChatGPT	62.58	64.78(3.5%)		

dataset, demonstrating its effective use of domain-specific information.

We also evaluate the performance of the several best general LLMs on the English dataset CaseHOLD. The performance comparisons are presented in Table 2. From the experimental results, we have the following findings:

- BLADE consistently shows improvements on the English dataset. Notably, Baichuan2-13B-Chat exhibits a more significant enhancement compared to ChatGLM2-6B and Baichuan2-7B-Chat. This suggests that larger models might derive greater benefits from the knowledge generated by the smaller LM. However, in the context of Chinese datasets, an inverse trend is noticed, implying that performance improvements may also depend on the intrinsic knowledge and in-context learning abilities of the models
- Another interesting observation is that the performance enhancement of ChatGPT on this dataset is not substantially high. This

may be because the original training data of ChatGPT has already incorporated a portion of legal pre-training data, thus reducing the additional advantage gained from external knowledge sources.

4.4.2 Main result in medical domain. Table 3 shows the performance of BLADE on the medical domain dataset MLEC-QA. From the experimental results, we have the following findings:

- Similar to the legal domain, the Medical-specific LLMs exhibit unsatisfactory performance. The challenge of integrating domain knowledge through continuous training without compromising the original capabilities of LLMs deserves further investigation.
- Both Zhongjing_base and Zhongjing_sft originally performed suboptimal in MLEC-QA. Surprisingly, under the guidance of generated knowledge, Zhongjing_sft showed superior performance than Zhongjing_base. This may indicate that supervised finetuning can improve the LLMs' ability to comprehend external knowledge.
- In the medical domain, BLADE demonstrates consistent performance improvements across all five subsets, with Bacihuan2-13B-Chat achieving the best performance. Overall, BLADE is proven to be able to maintain excellent performance in different evaluation tasks under multiple domains, which underscores its robust applicability in real-world settings.

4.4.3 Comparison with Retrieval Augmented LLMs. To further explore the effectiveness of BLADE, we compare it to the retrieval augmentation paradigm. More specifically, given a question, the retrieval model first retrieves the most relevant documents from a corpus. Then the general LLMs answer the question in the context of the relevant documents. We employ three different legal corpora, including legal_article, legal_book, and legal_all. The legal_article

Table 3: Overall Zero-shot performance on the medical dataset MLEC-QA. */** denotes that BLADE performs significantly better than baselines at p < 0.05/0.01 level using the fisher randomization test. Best results are marked bold.

Model		Cli	CWM		PH		Sto		TCM	
Model	Original	+BLADE	Original	+BLADE	Original	+BLADE	Original	+BLADE	Original	+BLADE
Medical Specific LLN	Medical Specific LLMs									
Zhongjing_base	15.58	35.74**	19.03	37.52**	16.55	36.98**	14.48	34.86**	17.41	36.65**
Zhongjing_sft	16.00	47.92**	18.50	49.64**	15.85	50.24**	15.76	46.12**	18.88	47.82**
Taiyi	43.42	49.72**	32.71	42.99**	35.11	45.63**	31.53	41.77**	32.83	43.65**
General LLMs	General LLMs									
ChatGLM-6B	30.04	53.42**	30.84	55.06**	30.47	55.66**	27.56	52.24**	32.96	53.64**
ChatGLM2-6B	48.86	60.20**	44.82	57.23**	44.39	59.75**	41.77	57.61**	46.12	55.72**
Qwen-7B-Chat	56.57	59.78*	52.59	58.20**	52.64	62.26**	49.33	57.39**	51.53	56.62**
Baichuan-7B	27.80	54.86**	25.19	56.03**	26.75	58.54**	22.34	50.34**	24.66	52.59**
Baichuan-13B-Chat	42.17	58.98**	45.27	56.59**	42.01	61.54**	38.52	56.42**	41.97	55.66**
Baichuan2-7B-Chat	51.10	59.99**	51.14	58.69**	50.00	62.45**	45.29	57.61**	51.79	56.82**
Baichuan2-13B-Chat	58.98	61.62*	54.39	58.79**	57.92	63.80**	50.39	57.84**	54.87	57.34*
ChatGPT	47.56	58.92**	38.69	57.91**	47.73	63.37**	43.32	57.58**	36.49	56.40**

Table 4: The performance comparison of BLADE and Retrieval-augmented LLMs on JEC-QA. The gain % shows the relative improvement of methods compared to the original language model. */** denotes that BLADE performs significantly better than the original language model at p < 0.05/0.01 level using the fisher randomization test. The best method in each column is marked in bold. The legal_pretrain corpus contains the entire contents of legal_article and legal_book.

Retrieval model	C	ChatGLM2-6B			ChatGPT		
Retrievai_model	Corpus	KD-questions(%)	CA-questions(%)	All(%)	KD-questions(%)	CA-questions(%)	All(%)
=	=	27.39	24.09	25.33	20.53	18.70	19.38
BGE-base	legal_article	28.51*(5.3%)	24.95(3.6%)	26.41*(4.3%)	26.73**(30.2%)	19.73*(5.5%)	22.36(15.4%)
BGE-base	legal_book	27.54(0.5%)	23.49(-2.4%)	25.01(-1.2%)	27.19**(32.4%)	19.55(4.5%)	22.42**(15.7%)
BGE-base	legal_all	30.11**(9.9%)	24.13(0.2%)	26.38*(4.1%)	27.75**(35.2%)	20.54**(9.8%)	23.25**(19.9%)
GTE-base	legal_article	27.09(-1.1%)	23.55(-2.2%)	24.88(-1.8%)	22.15**(7.8%)	19.04(1.8%)	20.21(4.3%)
GTE-base	legal_book	25.58(-6.6%)	22.84(-5.2%)	23.86(-5.7%)	21.90**(6.6%)	19.55(4.5%)	20.43(5.4%)
GTE-base	legal_all	25.43(-7.1%)	23.28(-3.4%)	24.09(-4.9%)	22.25**(8.3%)	19.10(2.1%)	20.28(4.6%)
M3E-base	legal_article	28.55*(4.2%)	24.77(2.8%)	26.19(3.4%)	26.03**(26.7%)	20.58**(10.1%)	22.63**(16.7%)
M3E-base	legal_book	27.74(1.3%))	24.77(2.8%)	25.88(2.2%)	26.28**(28.0%)	20.98**(12.2%)	22.97**(18.5%)
M3E-base	legal_all	30.56**(11.6%)	24.88(3.3%)	27.02*(6.6%)	28.20**(37.3%)	21.19**(13.3%)	23.82**(22.9%)
piccolo-base	legal_article	28.85*(5.3%)	23.46(-2.6%)	25.49(0.6%)	26.53**(29.2%)	20.46**(9.4%)	22.74**(17.3%)
piccolo-base	legal_book	29.20**(6.6%)	23.46(-2.6%)	25.61(1.1%)	26.18**(27.5%)	19.07(1.9%)	21.74**(12.17%)
piccolo-base	legal_all	30.72**(12.1%)	24.37(1.2%)	26.75*(5.6%)	28.29**(37.7%)	20.61**(10.2%)	23.50**(21.3%)
BLADE		30.81**(12.5%)	26.34**(9.3%)	28.01**(10.6%)	28.45**(38.6%)	23.67**(26.6%)	25.46**(31.3%)

Table 5: Impact of the number of retrieved documents on JEC-QA. The retrieved corpus is legal_all. Best results are marked bold.

Model	doc_num	KD-questions	CA-questions	All
-	0	27.39	24.09	25.33
M3E_base	1	30.56	24.88	27.02
M3E_base	3	30.71	24.67	26.93
M3E_base	5	30.36	25.29	27.19
M3E_base	7	29.75	24.28	26.33
$M3E_base$	9	29.63	24.40	26.36
BLADE	-	30.81	26.34	28.01

corpus contains all the Chinese legal provisions. Legal_book denotes the National Unified Legal Professional Qualification Examination Counseling Book, which consists of 15 topics and 215

chapters organized in a hierarchical manner. Legal_all corpus is consistent with the corpus in the DP phase, which contains all documents from legal_article and legal_book.

We select ChatGLM2-6B and ChatGPT, which have shown the best results in open-source and closed-source models respectively on the JEC-QA dataset, to conduct the experiment. To ensure a fair comparison, we exclusively use the top-1 document from the retrieval results and employ an identical prompt to BLADE, which also generates only a single document for analysis. The retrieval augmentation methods are implemented using LangChain ⁶. Table 4 demonstrates the comparison results. We have the following observations:

 Retrieval augmentation is proven to be effective in enhancing the performance of general LLMs in specific domains. However, its effectiveness is significantly influenced by the retrieval model and

 $^{^6} https://github.com/langchain-ai/langchain \\$

Query	Zhang has a watch. One day, his friend Li offers to buy it for 30,000 yuan, and Zhang immediately agrees. They sign a contract, with Li agreeing to pay in three installments over two months. After two months, Li has only paid 10,000 yuan. Zhang sues Li for breach of contract. During the trial, Li's parents claim that he has intermittent mental illness and was experiencing an episode when signing the contract, thus unable to be responsible for his actions. They apply to the court to declare Li as a person with limited capacity for civil conduct. How should the court handle this?
Retrieved documents	Article 367 of the Criminal Procedure Rules of the People's Procuratorate states that if a criminal suspect is suspected of having mental illness during the review for prosecution, the People's Procuratorate should conduct an assessment according to relevant rules. If the suspect's defender or relatives apply for an assessment due to suspected mental illness, the People's Procuratorate can also conduct an assessment according to the rules, with the costs borne by the applicant
BLADE	According to the Civil Procedure Law, the determination of whether a party has full civil conduct capacity requires methods like forensic assessment. If a party is declared as having limited capacity for civil conduct during litigation, the contract they signed may be deemed invalid. Before ruling on the contract's validity, it is necessary to determine Li's civil conduct capacity status through a special procedure. The court should suspend the original lawsuit and allow Li's parents to initiate a separate procedure to confirm Li's status as a person with limited civil conduct capacity. Based on the forensic assessment, the court will determine Li's status and make a corresponding judgment

Figure 4: Comparison of retrieved knowledge with that generated by BLADE.

the corpus. Consequently, not all retrieved knowledge contributes positively to the task at hand.

- Knowledge-Driven questions, focusing on the definition and explanation of legal concepts, tend to benefit more from the retrieved knowledge. However, Case-Analysis Questions, involving the analysis of real-life scenarios, may not see significant improvement from retrieved knowledge. This reflects the limitations of the retrieval augmentation paradigm, which lacks causal inference ability to identify question-specific knowledge.
- Regardless of Knowledge-Driven or Case Analysis questions, BLADE consistently provides stable enhancements and achieves optimal performance. When comparing overall performance with ChatGPT, the most effective retrieval model shows a 22.9% improvement, while BLADE achieves a notable 31.3% enhancement This success is attributed to our Knowledge Instruction Tuning and Bayesian Prompted Optimization strategy, which effectively supply the crucial knowledge that general LLMs require.

We further explore the impact of the number of retrieved documents on the performance of ChagtGLM2-6B. Specifically, we use M3E-base as the retrieval model and legal_all as the corpus because they achieve the best results in the retrieval augmentation paradigm. As shown in Table 5, when an appropriate number of documents are retrieved, there is a slight performance improvement due to more relevant documents being recalled. However, the performance of ChatGLM2-6B degrades when too many documents are retrieved, probably due to excessive noise introduced by the additional documents. In contrast, BLADE achieves the best results by generating only one piece of knowledge, suggesting its proficiency in producing more targeted and refined knowledge.

4.5 Ablation Studies

To better illustrate the effectiveness of our approach, we further conduct ablation studies on JEC-QA in zero-shot setting. Table 6 shows

Table 6: Ablation study on JEC-QA under zero-shot setting. The general LLM is ChatGLM2-6B. Best results are marked bold.

Small Model	KD-questions(%)	CA-questions(%)	All(%)
-	27.39	24.09	25.33
BLOOMZ_1b7	26.38	22.40	23.89
+ DP	26.87	23.63	24.85
+ DP & KIT	28.45	24.89	26.23
+ DP & KIT & BPO	30.81	26.34	28.01

Table 7: Impact of sizes on JEC-QA. The general LLM is ChatGLM2-6B. Best results are marked bold.

Small Model	KD-questions(%)	CA-questions(%)	All(%)
-	27.39	24.09	25.33
BLOOMZ_560m	29.05	24.92	26.47
BLOOMZ_1b1	29.80	25.52	27.13
BLOOMZ_1b7	30.81	26.34	28.01

the impact of different strategies. It's noticeable that while Domain-specific Pretraining successfully imparts domain knowledge to the small LM, it falls short in enabling instruction-following capabilities and in generating suitable knowledge, leading to a decrease in performance. With the integration of Knowledge Instruction Tuning, the small LM begins to offer beneficial knowledge. Bayesian Prompted Optimization further enhances the performance. The above experiments verify the effectiveness of each process within our approach.

4.6 Impact of Sizes

In this section, we aim to investigate the impact of the small LM's size. We conducted experiments on the JEC-QA dataset, utilizing

ChatGLM2-6B as the general model. Three versions of the small LM, namely BLOOMZ_560m, BLOOMZ_1b1, and BLOOMZ_1b7, were tested, each trained with the same training parameters and datasets. The results are shown in Table 7. We can observe that the small model with 560m parameters can also lead to performance gains. As the parameters of the small LM increase, the performance improvement brought by BLADE also increases. This phenomenon could be attributed to larger models' enhanced capability to generate more accurate and reliable knowledge.

4.7 Case Study

In this section, we conduct a case study to facilitate a clear understanding of the effectiveness of BLADE. Figure 4 illustrates the comparison of retrieved knowledge retrieved by M3E-base from the legal all corpus with the knowledge generated by BLADE. This question involves the assessment of civil conduct capacity in the context of a contract dispute. The appropriate legal procedure involves suspending the ongoing proceedings and initiating a specialized process by Li's parents to affirm Li's status as a person with limited civil capacity. The retrieval model returns the article about proceedings for people with mental illnesses, which fails to directly address the civil litigation process and the implications of limited civil capacity in contract disputes. BLADE's response is more accurate and directly relevant to the question. It correctly identifies the key issue - the civil litigation process concerning the assessment of civil conduct capacity in the context of a contract dispute. This case shows BLADE's strength in providing domainspecific, contextually appropriate responses. The domain-specific LM, trained on nuanced legal knowledge, is adept at interpreting the underlying legal implications of the described events. Therefore, BLADE can effectively bridge the gap between the specific details of an event and the relevant legal principles or precedents.

5 CONCLUSION

This paper proposes BLADE, a new framework for applying general large language models to new domains. At its core, BLADE employs small language models to assimilate and continually update domain-specific knowledge. The framework solves problems by realizing collaboration between general large language models and a small domain-specific model. It comprises three main stages: Domain-specific Pre-training, Knowledge Instruction Tuning, and Bayesian Prompted Optimization. Domain-specific Pre-training injects domain-specific knowledge into the small model. Knowledge Instruction Tuning activates the instruction-following capacity of the small model. Bayesian Prompted Optimization facilitates better alignment of the small model with the large model. Through extensive experiments on legal datasets, we find BLADE consistently demonstrates performance improvement across various language models with different sizes. In the future, we will investigate approaches to minimize hallucinations in small models and explore additional methods for joint optimization. A limitation is that our experiments are conducted only in multiple-choice datasets, the feasibility of our approach in generative tasks still deserves further investigation.

REFERENCES

[1] 2023. BLOOM: A 176B-Parameter Open-Access Multilingual Language Model. arXiv:2211.05100 [cs.CL]

- [2] Roee Aharoni and Yoav Goldberg. 2020. Unsupervised domain clusters in pretrained language models. arXiv preprint arXiv:2004.02105 (2020).
- [3] Md Adnan Arefeen, Biplob Debnath, and Srimat Chakradhar. 2023. LeanContext: Cost-Efficient Domain-Specific Question Answering Using LLMs. arXiv preprint arXiv:2309.00841 (2023).
- [4] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen Technical Report. arXiv:2309.16609 [cs.CL]
- [5] Baichuan. 2023. Baichuan 2: Open Large-scale Language Models. arXiv preprint arXiv:2309.10305 (2023). https://arxiv.org/abs/2309.10305
- [6] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*. PMLR, 2206–2240.
- [7] Ilias Chalkidis. 2023. ChatGPT may Pass the Bar Exam soon, but has a Long Way to Go for the LexGLUE benchmark. arXiv:2304.12202 [cs.CL]
- [8] Lichang Chen, Jiuhai Chen, Tom Goldstein, Heng Huang, and Tianyi Zhou. 2023. InstructZero: Efficient Instruction Optimization for Black-Box Large Language Models. arXiv preprint arXiv:2306.03082 (2023).
- [9] Daixuan Cheng, Shaohan Huang, and Furu Wei. 2023. Adapting Large Language Models via Reading Comprehension. arXiv:2309.09530 [cs.CL]
- [10] Zhumin Chu, Qingyao Ai, Yiteng Tu, Haitao Li, and Yiqun Liu. 2024. PRE: A Peer Review Based Large Language Model Evaluator. arXiv preprint arXiv:2401.15641 (2024)
- [11] Jiaxi Cui, Zongjian Li, Yang Yan, Bohua Chen, and Li Yuan. 2023. ChatLaw: Open-Source Legal Large Language Model with Integrated External Knowledge Bases. arXiv:2306.16092 [cs.CL]
- [12] Zhuyun Dai, Vincent Y Zhao, Ji Ma, Yi Luan, Jianmo Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith B Hall, and Ming-Wei Chang. 2022. Promptagator: Few-shot dense retrieval from 8 examples. arXiv preprint arXiv:2209.11755 (2022).
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018).
- [14] Qian Dong, Yiding Liu, Qingyao Ai, Haitao Li, Shuaiqiang Wang, Yiqun Liu, Dawei Yin, and Shaoping Ma. 2023. 13 Retriever: Incorporating Implicit Interaction in Pre-trained Language Models for Passage Retrieval. In Proceedings of the 32nd ACM International Conference on Information and Knowledge Management. 441–
- [15] Qian Dong, Yiding Liu, Qingyao Ai, Zhijing Wu, Haitao Li, Yiqun Liu, Shuaiqiang Wang, Dawei Yin, and Shaoping Ma. 2023. Aligning the Capabilities of Large Language Models with the Context of Information Retrieval via Contrastive Feedback. arXiv preprint arXiv:2309.17078 (2023).
- [16] Qian Dong, Yiding Liu, Suqi Cheng, Shuaiqiang Wang, Zhicong Cheng, Shuzi Niu, and Dawei Yin. 2022. Incorporating Explicit Knowledge in Pre-trained Language Models for Passage Re-ranking. arXiv preprint arXiv:2204.11673 (2022).
- [17] Qian Dong and Shuzi Niu. 2021. Latent Graph Recurrent Network for Document Ranking. In International Conference on Database Systems for Advanced Applications. Springer, 88–103.
- [18] Qian Dong and Shuzi Niu. 2021. Legal judgment prediction via relational learning. In Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval. 983–992.
- [19] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. GLM: General Language Model Pretraining with Autoregressive Blank Infilling. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 320–335.
- [20] Luciano Floridi and Massimo Chiriatti. 2020. GPT-3: Its nature, scope, limits, and consequences. Minds and Machines 30 (2020), 681–694.
- [21] Peter I Frazier. 2018. A tutorial on Bayesian optimization. arXiv preprint arXiv:1807.02811 (2018).
- [22] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. arXiv preprint arXiv:2004.10964 (2020).
- [23] Nikolaus Hansen. 2016. The CMA evolution strategy: A tutorial. arXiv preprint arXiv:1604.00772 (2016).
- [24] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021).
- [25] Quzhe Huang, Mingxu Tao, Zhenwei An, Chen Zhang, Cong Jiang, Zhibin Chen, Zirui Wu, and Yansong Feng. 2023. Lawyer LLaMA Technical Report. ArXiv abs/2305.15062 (2023).

- [26] Mandar Joshi, Kenton Lee, Yi Luan, and Kristina Toutanova. 2020. Contextualized representations using textual encyclopedic knowledge. arXiv preprint arXiv:2004.12006 (2020).
- [27] Jon M Kleinberg. 1997. Two algorithms for nearest-neighbor search in high dimensions. In Proceedings of the twenty-ninth annual ACM symposium on Theory of computing. 599–608.
- [28] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. Advances in Neural Information Processing Systems 33 (2020), 9459–9474.
- [29] Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel. 2020. Question and Answer Test-Train Overlap in Open-Domain Question Answering Datasets. arXiv:2008.02637 [cs.CL]
- [30] Patrick Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra Piktus, Pontus Stenetorp, and Sebastian Riedel. 2021. Paq: 65 million probably-asked questions and what you can do with them. Transactions of the Association for Computational Linguistics 9 (2021), 1098–1115.
- [31] Haitao Li, Qingyao Ai, Jia Chen, Qian Dong, Yueyue Wu, Yiqun Liu, Chong Chen, and Qi Tian. 2023. SAILER: Structure-Aware Pre-Trained Language Model for Legal Case Retrieval (SIGIR '23). Association for Computing Machinery, New York, NY, USA, 1035–1044. https://doi.org/10.1145/3539618.3591761
- [32] Haitao Li, Qingyao Ai, Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Zheng Liu, and Zhao Cao. 2023. Constructing Tree-based Index for Efficient and Effective Dense Retrieval. arXiv:2304.11943 [cs.IR]
- [33] Jing Li, Shangping Zhong, and Kaizhi Chen. 2021. MLEC-QA: A Chinese multichoice biomedical question answering dataset. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. 8862–8874.
- [34] Yuang Li, Yu Wu, Jinyu Li, and Shujie Liu. 2023. Prompting large language models for zero-shot domain adaptation in speech recognition. arXiv preprint arXiv:2306.16007 (2023).
- [35] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards General Text Embeddings with Multi-stage Contrastive Learning. arXiv:2308.03281 [cs.CL]
- [36] Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hannaneh Hajishirzi. 2021. Generated knowledge prompting for commonsense reasoning. arXiv preprint arXiv:2110.08387 (2021).
- [37] Ilya Loshchilov and Frank Hutter. 2018. Fixing weight decay regularization in adam. (2018).
- [38] Ling Luo, Jinzhong Ning, Yingwen Zhao, Zhijun Wang, Zeyuan Ding, Peng Chen, Weiru Fu, Qinyu Han, Guangtao Xu, Yunzhi Qiu, Dinghao Pan, Jiru Li, Hao Li, Wenduo Feng, Senbo Tu, Yuqi Liu, Zhihao Yang, Jian Wang, Yuanyuan Sun, and Hongfei Lin. 2023. Taiyi: A Bilingual Fine-Tuned Large Language Model for Diverse Biomedical Tasks. arXiv preprint arXiv:2311.11608 (2023).
- [39] Joel Niklaus, Veton Matoshi, Matthias Stürmer, Ilias Chalkidis, and Daniel E Ho. 2023. MultiLegalPile: A 689GB Multilingual Legal Corpus. arXiv preprint arXiv:2306.02069 (2023).
- [40] OpenAI. 2023. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL]
- [41] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. Advances in Neural Information Processing Systems 35 (2022), 27730–27744.
- [42] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).

- [43] Donald B Rubin. 1980. Randomization analysis of experimental data: The Fisher randomization test comment. Journal of the American statistical association 75, 371 (1980). 591–593.
- [44] Vin Sachidananda, Jason S Kessler, and Yi-An Lai. 2021. Efficient domain adaptation of language models via adaptive tokenization. arXiv preprint arXiv:2109.07460 (2021).
- [45] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. Replug: Retrieval-augmented black-box language models. arXiv preprint arXiv:2301.12652 (2023).
- [46] Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. 2022. Black-box tuning for language-model-as-a-service. In *International Conference on Machine Learning*. PMLR, 20841–20855.
- [47] Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. 2022. Recitation-augmented language models. arXiv preprint arXiv:2210.01296 (2022).
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in neural information processing systems. 5998–6008.
- [49] Wenya Wang, Vivek Srikumar, Hanna Hajishirzi, and Noah A Smith. 2022. Elaboration-generating commonsense question answering at scale. arXiv preprint arXiv:2209.01232 (2022).
- [50] He sicheng Wang Yuxin, Sun Qingxuan. 2023. M3E: Moka Massive Mixed Embedding Model.
- [51] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighof. 2023. C-Pack: Packaged Resources To Advance General Chinese Embedding. arXiv preprint arXiv:2309.07597 (2023).
- [52] Xiaohui Xie, Qian Dong, Bingning Wang, Feiyang Lv, Ting Yao, Weinan Gan, Zhijing Wu, Xiangsheng Li, Haitao Li, Yiqun Liu, and Jin Ma. 2023. T2Ranking: A Large-scale Chinese Benchmark for Passage Ranking. In Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval (, Taipei, Taiwan.) (SIGIR '23). Association for Computing Machinery, New York, NY, USA, 2681–2690. https://doi.org/10.1145/3539618.3591874
- [53] Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan. 2023. Zhongjing: Enhancing the Chinese Medical Capabilities of Large Language Model through Expert Feedback and Real-world Multi-turn Dialogue. arXiv:2308.03549 [cs.CL]
- [54] Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2022. Generate rather than retrieve: Large language models are strong context generators. arXiv preprint arXiv:2209.10063 (2022).
- [55] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. arXiv preprint arXiv:2210.02414 (2022).
- [56] Jiaxing Zhang, Ruyi Gan, Junjie Wang, Yuxiang Zhang, Lin Zhang, Ping Yang, Xinyu Gao, Ziwei Wu, Xiaoqun Dong, Junqing He, Jianheng Zhuo, Qi Yang, Yongfeng Huang, Xiayu Li, Yanghan Wu, Junyu Lu, Xinyu Zhu, Weifeng Chen, Ting Han, Kunhao Pan, Rui Wang, Hao Wang, Xiaojun Wu, Zhongshen Zeng, and Chongpei Chen. 2022. Fengshenbang 1.0: Being the Foundation of Chinese Cognitive Intelligence. CoRR abs/2209.02970 (2022).
- [57] Lucia Zheng, Neel Guha, Brandon R Anderson, Peter Henderson, and Daniel E Ho. 2021. When does pretraining help? assessing self-supervised learning for law and the casehold dataset of 53,000+ legal holdings. In Proceedings of the eighteenth international conference on artificial intelligence and law. 159–168.
- [58] Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020. JEC-QA: a legal-domain question answering dataset. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34. 9701–9708.