# Structure-aware Domain Knowledge Injection for Large Language Models

**Anonymous ACL submission**

## Abstract

This paper introduces a pioneering methodology, termed *StructTuning*, to efficiently transform foundation Large Language Models (LLMs) into domain specialists. It significantly reduces the training corpus needs to a mere **5%** while achieving an impressive **100%** of traditional knowledge injection performance. Motivated by structured human education, we propose a novel two-stage strategy for knowledge injection and alignment: *Structure-aware Continual Pre-Training* (SCPT) and *Structure-aware Supervised Fine-Tuning* (SSFT). In the SCPT phase, we automatically extract the domain knowledge taxonomy and reorganize the training corpora, enabling LLMs to effectively link textual segments to targeted knowledge points within the taxonomy. In the SSFT phase, we explicitly prompt models to elucidate the underlying knowledge structure in their outputs, leveraging the structured domain insight to address practical problems. Our ultimate method was extensively evaluated across model architectures and scales on LongBench and MMed-Bench datasets, demonstrating superior performance against other knowledge injection methods. We also explored our method's scalability across different training corpus sizes, laying the foundation to enhance domain-specific LLMs with better data utilization. Code is available at this anonymous URL: https://anonymous.4open.science/r/StructTuning/.

## 1 Introduction

Large language models (LLMs) have recently seen extensive deployment across various applications (Vaswani et al., 2017; Achiam et al., 2023; Jiang et al., 2023; Bi et al., 2024). When adapting foundational models (*e.g.*, Llama series (Touvron et al., 2023a,b; Dubey et al., 2024)) to specialized AI assistants in distinct domains, developers usually employ two techniques to enhance LLMs' proficiency: retrieval-augmented generation (RAG) (Lewis et al., 2020) and domain knowledge injection (Gururangan et al., 2020). While RAG effectively utilizes an external knowledge base to augment information, the retrieval process's inherent noise poses challenges to generating reliable responses, especially in scenarios requiring logical reasoning where there is a semantic gap between the user's query and the knowledge base.(Zhang et al., 2023; Chen et al., 2023). Thus, another avenue tries to inject new knowledge to LLMs via training techniques (Gu et al., 2021; Hu et al., 2021; Mecklenburg et al., 2024).

Continual pre-training (Sun et al., 2020; Ibrahim et al., 2024) is preferred for new, domain-specific knowledge injection (Cui et al., 2023; Wang et al., 2023b; Qiu et al., 2024). However, it often entails resource-intensive training on billions of tokens from the internet to learn fragmented knowledge points, rather than absorbing structured knowledge from a few domain-specific textbooks (Jin et al., 2020). For example, MMedLM (Qiu et al., 2024) curates 25.5B tokens to derive a medical model, and DeepSeek-Coder (Guo et al., 2024) uses 2T tokens for coding adaptation. The common failure to learn effectively from limited textbook content has been attributed to insufficient data diversity (Zhu and Li, 2023a), which however violates the observation during the human education process in Fig. 1: students gain knowledge by sequentially studying from textbooks, reviewing knowledge points and structures, and applying this knowledge through proper exercises. Here, all the new data to learn are textbooks (structured content) and exercising examples (question-answering pairs), and students just adopt their world knowledge to memorize, understand, and apply the knowledge to become domain experts (Krathwohl, 2002; Yu et al., 2023).

As educating human students, we propose to inject structured domain knowledge into LLMs via two steps: *Structure-aware Continual Pre-Training* (SCPT) and *Structure-aware Supervised Fine-Tuning* (SSFT).
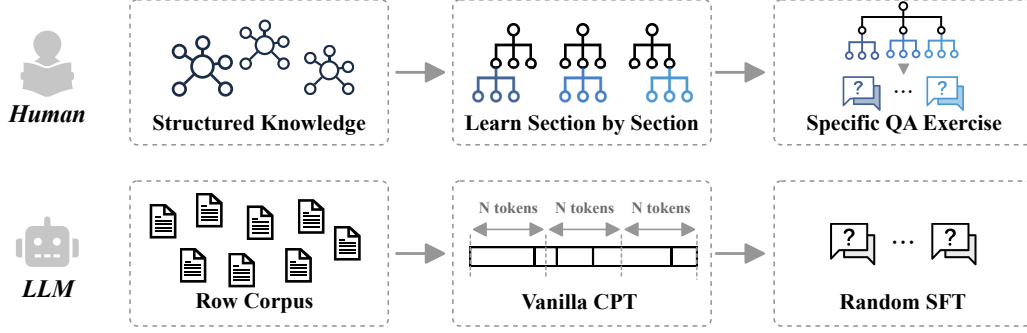
Figure 1: **Discrepancy between human education and vanilla LLM adaptation**. Human students learn structured knowledge through textbooks section by section, with particular exercises on related knowledge points. Traditional LLM adaptation continually pre-trains on data chunks from randomly concatenated text segments, with aimless supervised fine-tuning for conversation alignment. The inherent knowledge structure is ignored.

In the SCPT stage, we argue that high-quality textbook data (as well as regular corpora from the Internet) can adequately infuse domain knowledge (Gunasekar et al., 2023), where the organization of training corpora is crucial. In conventional paradigms (Fig. 1), text corpora are simply concatenated and divided into chunks of 2048 (Qiu et al., 2024) or 4096 (Guo et al., 2024) tokens, while the inherent semantic structure (*e.g.*, catalogs of textbooks) is discarded. Instead, we view each chunk as a knowledge point and automatically extract domain knowledge taxonomy from the whole corpus. Subsequently, LLMs are trained to predict textual content (corresponding to a knowledge point) *under the condition of* the knowledge path within the domain structure, linking individual training chunks with the entire knowledge architecture. Finally, models are asked to memorize the entire structure to review the whole domain knowledge system.

In the SSFT stage, the goal shifts from knowledge injection to enabling LLMs to recall and utilize their acquired knowledge to tackle real-world challenges. We explicitly elicit knowledge paths in LLMs' responses, as a beacon for models to targeted information retrieval or logical reasoning for reliable responses. To this end, we derive a scalable strategy to generate question-answer pairs as practice exercises by open-sourced LLMs or API, such as LLaMA3 (Dubey et al., 2024) and GPT4 (Achiam et al., 2023). In the scenarios with existing QA pairs like MMedBench (Qiu et al., 2024), we retrieve the related knowledge structure and content, instructing LLaMA3 to provide explanations from questions to answers based on the knowledge paths. For datasets lacking specific QA samples like LongBench (Bai et al., 2023b), we randomly select knowledge paths from the domain taxonomy and prompt LLaMA3 to craft question-answer-explanation triplets for training exercises.

Our ultimate approach *StructTuning* has been extensively evaluated across different model architectures and sizes. In particular, we first examine their capability to recall the knowledge injected through open-ended QA on the LongBench (Bai et al., 2023b) dataset, then assess the application of injected knowledge to address real-world issues through multiple-choice QA on MMedBench (Qiu et al., 2024). Both evaluations underscore the superiority of StructTuning, surpassing other SOTA domain knowledge injection methods (Cheng et al., 2023; Zhang et al., 2024). Remarkably, we achieve a **50%** improvement in knowledge injection compared to SOTA MMedLM2 in the medical domain, using only **0.3%** of the training data requirement. Furthermore, StructTuning exhibits good scalability, achieving comparable performance with only **5%** of the training data. These findings reveal the superiority of our method for enhancing domain-specific AI assistants with more efficient data utilization.

Our contribution is summarized as follows:

- We proposed a novel two-stage training strategy, SCPT and SSFT, to inject domain knowledge into LLMs by preserving and utilizing the inherent structure of the training corpus.

- We developed a scalable data construction framework to generate structure-aware training samples from original corpora to facilitate the SCPT and SSFT stages.

- We conducted extensive investigations on our StructTuning strategy on various data and model settings, and comprehensively illustrate our superiority in knowledge injection.
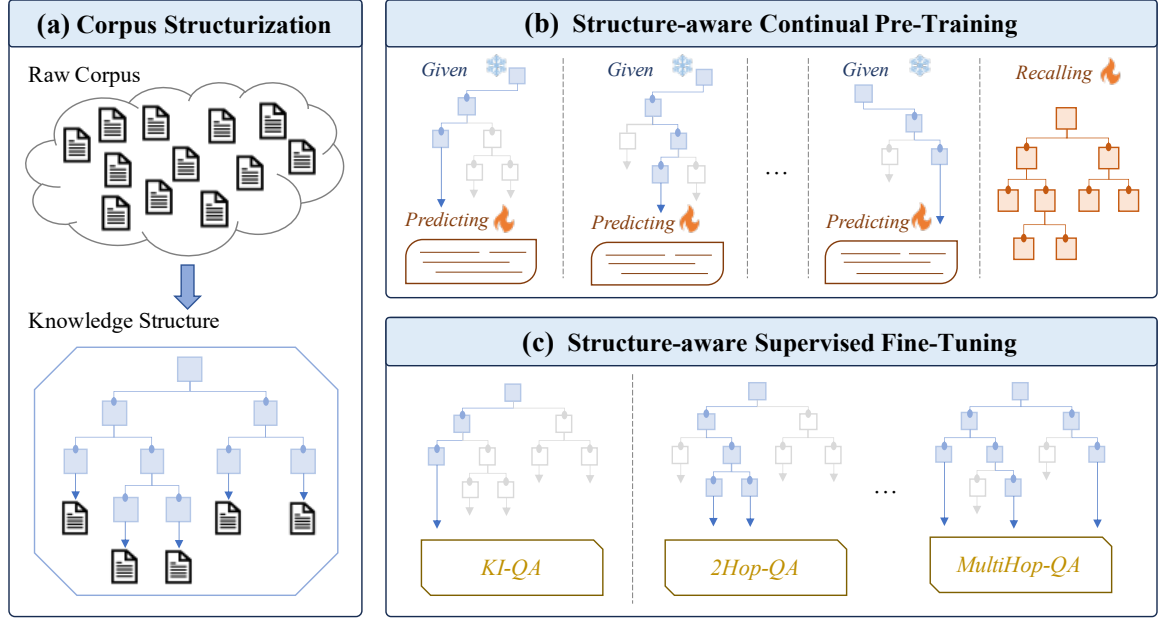
2

Figure 2: **Framework for structure-aware knowledge injection**. We extract the inherent knowledge structure in the training corpus, and associate training chunks to corresponding knowledge points. Models are continually pre-trained on data chunks in the condition of the knowledge structure, and fine-tuned with supervised QA samples to elicit their learned knowledge to solve knowledge-intensive (KI) and 2- or multi-hop questions in the real world.

## 2 Related Works

Here we briefly discuss the closely related works. A detailed discussion can be found in Appendix C.

**Domain Adaptation for LLMs**. To address the domain adaptation problem, a pre-trained model will be continually pre-trained (CPT) with domain-specific content (Sun et al., 2020; Xu et al., 2023b), and fine-tuned with supervised instruction-response pairs (SFT) to keep advancing interactive capabilities (Mecklenburg et al., 2024; Qiu et al., 2024). This paradigm is validated efective in dynamic fields like medicine (Wang et al., 2023b; Qiu et al., 2024) and coding (Roziere et al., 2023; Guo et al., 2024). Our study builds upon this CPT-SFT framework, innovating with SCPT-SSFT strategies to efficiently and effectively infuse domain knowledge with the inherent structure hierarchy.

**Structure-aware Knowledge Aggregation**. In conventional paradigms, researchers extract entity-relation-entity triplets from texts to construct knowledge graphs (Pan et al., 2024), to enhance LLMs's factual knowledge and logical reasoning(Zhang et al., 2022; Wen et al., 2023). Here, each node corresponds to either a specific entity or an abstract concept, lacking the capability to present an informative and self-contained *knowledge point*. This paper extends to structure-aware knowledge aggregation on existing training corpora, injecting the whole domain knowledge structure into LLMs' by linking training samples to corresponding knowledge points and reasoning paths.

**Data Augmentation and Synthesis**. Traditional methods aim to artificially expand the training dataset size (Xu et al., 2023a; Mukherjee et al., 2023) or generate entirely new samples to adapt LLMs to specific tasks (Tang et al., 2024). Yet, they often overlook the structured nature of domain knowledge, while the aimlessly generated samples may also lack diversity (Ovadia et al., 2023; Mecklenburg et al., 2024) and cannot cover the domain knowledge points (Mecklenburg et al., 2024; Tang et al., 2024). By contrast, our SSFT design is an innovative departure to address the challenge of retaining and utilizing the structured knowledge inherent in domain-specific content.

## 3 Methodology

Fig. 2 depicts our StructTuning methodology to inject domain knowledge into pre-trained LLMs using the inherent knowledge structure. With curated domain corpora (typically a few textbooks), we first extract the knowledge structure, and associate text chunks to corresponding knowledge paths and points (Sec. 3.1). Then, we design a two-stage training strategy to inject the highly structured domain knowledge into language models by mimicking the human education process, comprising the SCPT (Sec. 3.2) and SSFT (Sec. 3.3) techniques.
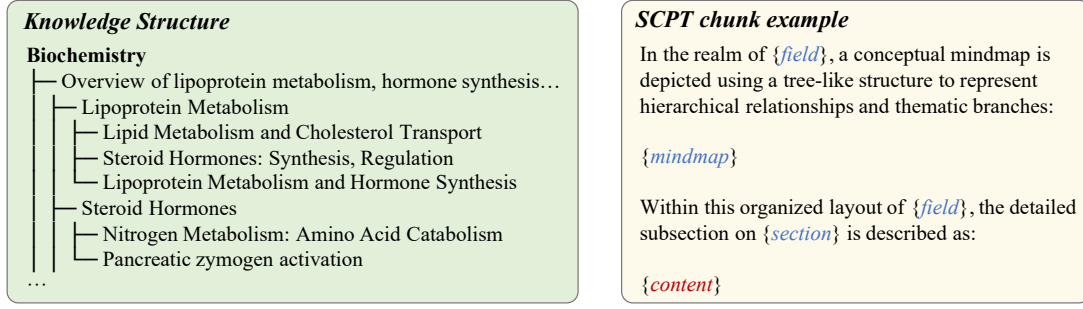
| Knowledge Structure | SCPT chunk example |
|---|---|
| **Biochemistry**<br>├── Overview of lipoprotein metabolism, hormone synthesis...<br>│  ├── Lipoprotein Metabolism<br>│  │  ├── Lipid Metabolism and Cholesterol Transport<br>│  │  ├── Steroid Hormones: Synthesis, Regulation<br>│  │  └── Lipoprotein Metabolism and Hormone Synthesis<br>│  ├── Steroid Hormones<br>│  │  ├── Nitrogen Metabolism: Amino Acid Catabolism<br>│  │  └── Pancreatic zymogen activation<br>... | In the realm of {*field*}, a conceptual mindmap is depicted using a tree-like structure to represent hierarchical relationships and thematic branches:<br><br>{*mindmap*}<br><br>Within this organized layout of {*field*}, the detailed subsection on {*section*} is described as:<br><br>{*content*} |

Figure 3: **Left**: extracted knowledge structure. **Right**: template to bridge mindmap structure and textual contents.

## 3.1 Knowledge Structure Extraction

For web-crawled corpus, previous data pre-processing focuses on quality assessment for individual documents (Bi et al., 2024), while the meta-info of knowledge structures (*e.g.*, the table content for a textbook) is usually neglected or filtered out, and all we have are those sequentially arranged text segments (*e.g.*, page-by-page-chunked content). As shown in Fig. 2 (a), we aim to extract (or, recover) the knowledge structure from the raw corpus for subsequent domain knowledge injection.

First, we use spaCy[1] to split the content from a textbook at the paragraph-level, and merge the sentences to form training chunks within a maximum size (*e.g.*, 2048 tokens (Qiu et al., 2024)). After that, we prompt the advanced Llama3-70B (Dubey et al., 2024) model to summarize the title for each chunk, where the textual content with the abstractive title jointly contributes to a "knowledge point".

Then, we aggregate knowledge points and extract the inherent structure hierarchy by leveraging advanced language models. Inspired by Liu et al. (2024a), we take the title list to instruct a specifically developed 7B model to identify the inherent knowledge structure (as exemplified in Fig. 3) within the text chunks, and Fig. A1 showcases how to deal with non-textbook data. The whole process does not involve human annotation, which reduces the cost and makes our method scalable for larger domain training corpora.

In particular, Appendix B.1 and Appendix B.4 verify that our specialized 7B model can identify sufficiently precise knowledge structure for effective and efficient domain adaptation, as more powerful LLMs like LLaMA3-70B (Dubey et al., 2024) and GPT-3.5 (Brown et al., 2020) cannot bring significant enhancement while largely increase the inference costs.

---

[1] https://github.com/explosion/spaCy

## 3.2 Structure-aware Continual Pre-Training

In conventional knowledge injection, training corpora are randomly concatenated and chunked into text segments without distinguishing the original content, making models only able to absorb domain knowledge emerging in data diversity (Ovadia et al., 2023; Mecklenburg et al., 2024; Qiu et al., 2024). In this section, we present another solution to inject knowledge from limited text corpora by leveraging the highly abstractive and exhaustive domain knowledge structures for continual pre-training.

We first transform the knowledge structure into natural languages using the same mindmap template (Wen et al., 2023) in Fig. 3 (left), and prepend it to each training chunk, forcing LLMs to memorize the textual content (knowledge points) in the condition of the associated knowledge path in the structure hierarchy. 20 diversified templates (see Fig. A5) are collected from GPT-4 (Achiam et al., 2023) to bridge mindmap structures and training chunks, one of which is displayed in Fig. 3 (right). The prepended mindmap, as well as the template, does not produce auto-regressive loss. Losses are only calculated in the *content* part. Formally, we turn the original language modeling in vanilla CPT to conditioned modeling (Keskar et al., 2019) in our SCPT stage:

$$p(\boldsymbol{x}^k) = \prod_{i=1}^{n} p(x_i^k|x_{<i}^k) \implies$$
$$p(\boldsymbol{x}^k|\boldsymbol{s}^k) = \prod_{i=1}^{n} p(x_i^k|x_{<i}^k, \boldsymbol{s}^k) \tag{1}$$

where $p(\boldsymbol{x}^k)$ models the probability distribution for the $k$-th chunk $\boldsymbol{x}^k = (x_1^k, \cdots, x_n^k)$ via the chain rule of probability (Bengio et al., 2000) on each token $x_i^k$, and $\boldsymbol{s}^k$ denotes the associated knowledge mindmap. Appendix B.5 extensively investigates the effectiveness of our SCPT strategy.

4

**(a) Knowledge-Intensive QA Generation**

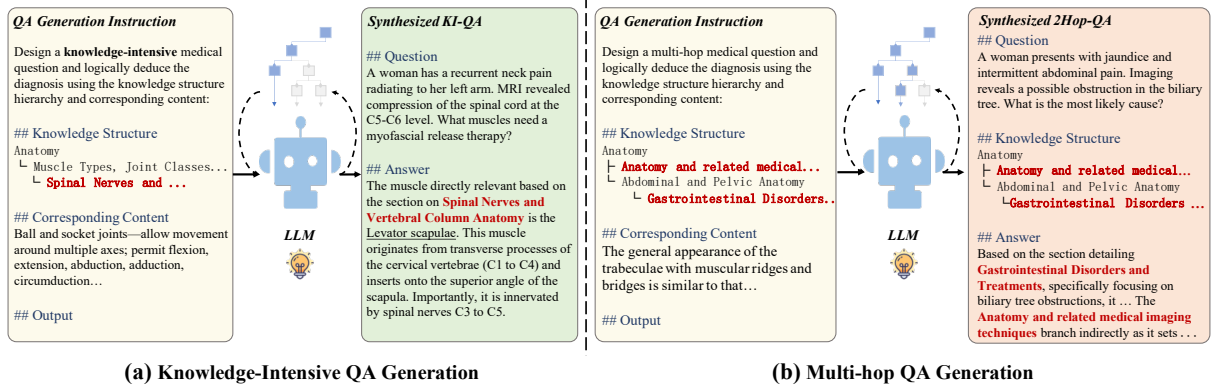**(b) Multi-hop QA Generation**

Figure 4: **QA samples synthesized for SSFT**. We instruct Llama3-70B to generate (a) knowledge-intensive and (b) multi-hop questions and derive the diagnosis answers with explicit reasoning.

As Fig. 2 (b) shows, after traversing the $m$ knowledge points in extracted structures, models are asked to recall the whole knowledge hierarchy, *i.e.*, to model the composed probability distribution:

$$p(\bar{s}) = \prod_{k=1}^{m} p(s^k) \qquad (2)$$

In SCPT, we mimic the human education process to inject knowledge into LLMs in a section-by-section manner, and replay the entire knowledge structure for the models to review and summarize the learned domain knowledge. These two steps iteratively alternate throughout training epochs.

### 3.3 Structure-aware Supervised Fine-Tuning

Traditional supervised fine-tuning aims to align the (continually) pre-trained models to interactive Chat-Bots through question-answering exercises (Cui et al., 2023; Qiu et al., 2024). Previous studies focus on enlarging the quantity and enhancing the diversity of training syntheses (Xu et al., 2023a; Mukherjee et al., 2023; Liu et al., 2024b) but neglect the highly structured domain knowledge. In contrast, our structure-aware supervised fine-tuning (SSFT) technique aims to elicit models' structured knowledge learned during SCPT, adapting LLMs to interactive and reliable domain experts.

Fig. 2 (c) illustrates SSFT samples synthesis guided by domain knowledge structures. First, we use the random walk algorithm to create knowledge paths with 1 to $l$ branches in the original mindmap (the illustration of knowledge paths and branches is displayed in Fig. A2). For paths linking to a single knowledge point, we use the corresponding text content to prompt Llama3-70B (Dubey et al., 2024) to generate knowledge-intensive QA pairs. For paths with two or more branches, we prompt

Llama3-70B with the knowledge path and textual contents to synthesize 2- or multi-hop QA samples, which require specific reasoning along the knowledge structure to derive from questions to answers. Fig. 4 presents several examples.

For every synthesized QA sample ($z$), we will prepend the relevant mindmap hierarchy to the answer, and add a CoT prompt in the question to construct another type of QA data ($z'$) for SFT alignment. This design explicitly elicits the learned knowledge in models' responses, teaching them how to apply the structured knowledge to address real-world problems. We use the two types of QA samples for training, as recommended by Qiu et al. (2024). During testing, we use the vanilla question as input to efficiently gather models' answers to calculate accuracy, and take the CoT prompt to probe to what extent LLMs can memorize and leverage the injected knowledge to answer the questions.

Integrating with SCPT and SSFT, our StructTuning approach translates into remarkable efficacy and efficiency in domain knowledge injection, as comprehensively evaluated in subsequent sections.

## 4 Experiments

We design extensive evaluations of our StructTuning through several experiments on two benchmarks. First, we investigate the free-form question-answering task on the LongBench (Bai et al., 2023b) dataset, so as to verify the **memorization and understanding** of injected knowledge (the answer can be directly found in training corpora). Then, we delve into the multi-choice question-answering task on MMedBench (Qiu et al., 2024), to explore how LLMs **apply** the injected knowledge in basic medicine to determine the real-world diagnosis for patients with logical reasoning.

Table 1: **Recall of Closed-Book QA (CBQA) on LongBench** (Bai et al., 2023b). The **best** and <u>secondary</u> results are marked in **bold** and <u>underlines</u>, respectively. The base model is Llama2-7B.

| Task | Adaptation | SingleDoc-QA | | | MultiDoc-QA | | | | Average |
|------|------------|--------|------|--------|------|------|-------|------|---------|
| | | Qasper | MFQA | MFQAzh | HpQA | 2Wiki | Musiq | Duzh | |
| CBQA | CPT+SFT | <u>20.7</u> | 35.3 | <u>20.6</u> | 29.9 | 32.1 | 18.9 | 12.0 | 24.2 |
| | SCPT+SFT | 18.8 | <u>42.5</u> | 17.7 | <u>35.7</u> | <u>36.4</u> | <u>20.5</u> | <u>15.3</u> | <u>26.7</u> |
| | SCPT+SSFT | **30.5** | **44.6** | **24.3** | **40.8** | **42.0** | **21.8** | **16.8** | **31.5** |

## 4.1 Preliminary Free-form QA Investigation

**Datasets and Tasks.** Seven subsets with 1,350 test examples from LongBench (Bai et al., 2023b) are utilized to evaluate closed-book question-answering evaluation, where the answers can be directly found in corresponding passages. The 14K reading-comprehension passages are used for knowledge injection via CPT/SCPT, and another 2,700 QA samples are generated for SFT/SSFT. Details are described in Appendix A.1.

**Evaluation Metrics.** Here we first report the *recall* (Zhu and Li, 2023b) for models' outputs against ground-truth answers to quantify the knowledge memorization degree. In Appendix B.9 we also evaluate F1-score for a thorough comparison.

**Investigated Models.** We mainly investigate Llama2-7B (Touvron et al., 2023b) to compare the knowledge injection performance.

**Implementation Details.** We train all models on LongBench passages for 3 epochs using a batch size of 128, and train for 1 epoch on synthetic SFT data to avoid overfitting. The learning rate is 2e-5.

**Main Results.** We first try to inject passage content into LLMs with a conventional CPT+SFT paradigm for the baseline, and use CoT instructions during testing to elicit models' memorized knowledge in their responses. However, Tab. 1 indicates such an injection approach is ineffective, as the knowledge recall is only 24.2%.

On the other hand, our SCPT strategy achieved a higher knowledge recall of 26.7%. It implies the model has successfully associated the relevant passages with their entire knowledge structure for the given question, especially on multi-doc QAs that require more complex information retrieval and reasoning on multiple passages. Furthermore, our SSFT technique continues to improve the knowledge recall to **31.5%**, which indicates the vanilla SFT strategy can only regularize LLMs' response styles, while our SSFT could teach LLMs to utilize their knowledge (injected in the SCPT stage) to answer corresponding questions. Fig. 5 showcases a representative example.

*Question*

Regarding Multilingual dialect detection and classification using deep learning, what models were used for dialect identification?

*Model Response*

Here is the knowledge structure regarding Multilingual dialect detection and classification using deep learning:
Multilingual dialect detection and classification using deep learning
├─ Exploring multilingual dialect identification using deep learning techniques
├─ Design and implementation of a neural network for language identification
├─ Proposed method for multilingual dialect identification using CNNs and LSTMs
├─ Evaluation and analysis of a dialect identification model using X-NMT architecture
└─ Analysis of a multilingual dialect identification system using NMT models

The question asks about the specific models employed in the study for dialect identification. By recalling the passage's description of the "proposed method," we can identify that **BERT, RoBERTa, and XLM-RoBERTa** were used as these three architectures were combined to form the X-NMT model

*Ground-Truth*

BERT, RoBERTa, ELECTRA, and XLM-RoBERTa.

Figure 5: An example for structure-aware responses.

In Appendix B.8, we also validate relatively good memorization of injected knowledge structures, further emphasizing our SCPT's efficacy.

## 4.2 In-depth Multi-choice QA Evaluation

**Datasets and Tasks.** We take several corpus sizes from MMedC (Qiu et al., 2024) for CPT/SCPT, with 45K QA data from MMedBench (Qiu et al., 2024)'s training set for SFT/SSFT. Models are evaluated on six multi-choice subsets from MMed-Bench, where LLMs should make real-world patient diagnoses with adequate reasoning on medical knowledge. Detailed setup is in Appendix A.2.

**Evaluation Metrics.** We follow the default setting to calculate the accuracy on six language subsets and the averaged scores. Metrics are computed by lexical exact-matching on models' responses, rather than maximum token probabilities.

**Investigated Models.** We extend the investigation across model architectures and scales, including Llama2-7B/13B (Touvron et al., 2023b), InternLM2-7B (Zheng et al., 2024), and Llama3-8B (Dubey et al., 2024). Other popular medical LLMs (Han et al., 2023; Wu et al., 2024; Qiu et al., 2024) are also included for a thorough comparison.

**Implementation Details.** Following Qiu et al. (2024), models are first trained for 3 epochs on medical corpora with a learning rate of 2e-5, and

6

Table 2: **Multiple-choice evaluation on MMedBench** (Qiu et al., 2024). We report separate accuracies across six languages, with "Average" denoting the mean score. "#Token" denotes required data for knowledge injection.

| Model | English | Chinese | Japanese | French | Russian | Spanish | Average | #Token |
|---|---|---|---|---|---|---|---|---|
| ChatDoctor (Yunxiang et al., 2023) | 43.52 | 43.26 | 25.63 | 18.81 | 62.50 | 43.44 | 39.53 | - |
| PMC-LLaMA (Wu et al., 2024) | 47.53 | 42.44 | 24.12 | 20.74 | 62.11 | 43.29 | 40.04 | - |
| MedAlpaca (Han et al., 2023) | 46.74 | 44.80 | 29.64 | 21.06 | 59.38 | 45.00 | 41.11 | - |
| Llama2-7B (Touvron et al., 2023b) | 43.36 | 50.29 | 25.13 | 20.90 | 66.80 | 47.10 | 42.26 | - |
| InternLM2-7B (Zheng et al., 2024) | 57.27 | 77.55 | 47.74 | 41.00 | 68.36 | 59.59 | 58.59 | - |
| Llama3-8B (AI, 2024) | 63.86 | 78.23 | 48.24 | 50.80 | 71.48 | 64.15 | 62.79 +0.00 | - |
| Llama3+MMed (Qiu et al., 2024) | <u>66.06</u> | **79.25** | **61.81** | **55.63** | <u>75.39</u> | 68.38 | **67.75** +4.96 | 25.5B |
| Llama3+StructTuning (**Ours**) | **66.77** | 77.44 | 53.27 | 51.61 | 74.61 | <u>68.49</u> | 65.36 +2.57 | **76M** |
| Llama3+StructTuning (**Ours**) | 65.36 | <u>79.04</u> | <u>56.28</u> | <u>55.47</u> | **80.47** | **69.80** | <u>67.74</u> +4.95 | <u>1.2B</u> |



Figure 6: Knowledge injection's scalability.

Table 3: **Generalization to various model architectures and sizes.** Here we use 76M training corpus.

| Model | Size | Adaptation | Accuracy |
|---|---|---|---|
| InternLM2 | 7B | CPT+SFT | 58.59 |
| | | **SCPT+SSFT** | **63.05** |
| Llama2 | 7B | CPT+SFT | 42.26 |
| | | **SCPT+SSFT** | **51.04** |
| Llama2 | 13B | CPT+SFT | 48.33 |
| | | **SCPT+SSFT** | **54.50** |

then fine-tuned for 1 epoch to avoid overfitting. The detailed setup is displayed in Appendix A.2.

**Main Results.** The results in Tab. 2 demonstrate the promising enhancement achieved by our Struct-Tuning technique largely outperforming the previous domain-specific LLMs like PMC-LLaMA (Wu et al., 2024) and MedAlpaca (Han et al., 2023). Notably, our structure-aware knowledge injection approach, using merely 76M tokens (**0.3%**) curated from medical textbooks, achieves over **50%** performance (2.57% *v.s.* 4.96%) against the state-of-the-art MMedLM (Qiu et al., 2024) method, which is trained on the entire MMedC (Qiu et al., 2024) corpora of 25.5B tokens. After scaling up the training tokens to 1.2B (around **5%**), our method makes nearly **100%** improvements to the average accuracy, significantly reducing the training cost of traditional knowledge injection approaches.

**Approach's Scalability.** We further curate a series of training corpora sizes to investigate our method's scalability in-depth: 30M, 76M, 132M, 250M, and 1.2B, which respectively take around 0.1%, 0.3%, 0.5%, 1%, and 5% of 25.5B tokens. The vanilla CPT-SFT paradigm and our SCPT-SSFT strategy are thoroughly compared across those data settings. According to Fig. 6, our method consistently surpasses the vanilla paradigm by a large margin, emphasizing the efficacy and effi-

ciency of domain knowledge injection. In particular, we fit two performance-ratio scaling curves from the data points in Fig. 6 as:

$$p_v \approx -0.04(\log r)^2 + 13.3 \log r + 100.0$$
$$p_s \approx -1.11(\log r)^2 + 7.63 \log r + 133.0 \quad (3)$$

where $p_v$ and $p_s$ denote the relative performance enhancement (%) for vanilla and structure-aware knowledge injection, and $r$ is the corpus ratio.

In fact, we use the scaling law derived from 0.1%, 0.3%, 0.5%, and 1% points to predict that achieving 100% performance would require 5% of the data corpus, which has been confirmed in Tab. 2. On the other hand, it also indicates our method may lead to 133% enhancement with a further 100% comprehensive data utilization, further validating the effectiveness and scalability of our method.

**Approach's Generalization.** In Tab. 3, we also validated the performance on Llama2 (Touvron et al., 2023b) and InternLM2 (Zheng et al., 2024) model series by using 76M tokens for knowledge injection. Our method leads to consistently significant improvements on InternLM2-7B (+4.46%), Llama2-7B (+8.78%) and Llama2-13B (+6.17%) backbone models, further demonstrating the generalizability and scalability of our Struct-Tuning across model architectures and sizes. Detailed results are presented in Tab. A6.

**Ablation Studies.** We perform a comprehensive ablation study on MMedBench's English subset in Tab. 4. As suggested by Qiu et al. (2024), we use

Table 4: **Ablation studies with Llama2-7B on the English subset of MMedBench.**

| Adaptation | English | Chinese | Japanese | French | Russian | Spanish | Average |
|---|---|---|---|---|---|---|---|
| SFT | 44.54 | 32.81 | **26.63** | 15.27 | 53.91 | 42.30 | 35.91 |
| CPT + SFT | 46.27 | 32.57 | 26.13 | 17.36 | 50.00 | 40.63 | 35.49 |
| **SCPT** + SFT | 46.50 | 32.14 | 20.10 | 18.17 | 53.91 | 39.97 | 35.13 |
| **SCPT** + **SSFT** | **49.96** | 32.63 | 22.11 | 17.52 | 51.17 | 41.28 | 35.78 |
| **SCPT** + **SSFT*** | 49.10 | **33.92** | 18.33 | **27.14** | **57.42** | **43.73** | **38.27** |
| RAG | 38.12 | 29.22 | 22.61 | 23.34 | 53.91 | 36.47 | 33.95 |
| AdaptLLM (Cheng et al., 2023) | 46.79 | 33.80 | 20.60 | 14.15 | 53.12 | 42.34 | 35.03 |
| RAFT (Zhang et al., 2024) | 43.60 | 32.34 | 21.11 | 14.95 | 50.39 | 42.16 | 34.09 |

the English textbooks (Jin et al., 2020) (26M tokens) to compare vanilla and structure-aware CPT, paired with corresponding SFT strategies. In particular, "SFT" uses vanilla SFT with 10K QA samples from MMedBench's training split, while "SSFT" applies structure-aware SFT on the same questions, enhancing answers with Llama3-70B-generated knowledge explanations (Sec. 3.3). "SSFT*" further includes 8K additional structure-aware QA pairs, totaling 18K training entries. Training hyperparameters align with the main experiment.

In Tab. 4, the CPT+SFT paradigm improves accuracy by 1.73%, while SCPT with vanilla SFT achieves a higher 46.50%. Combining SCPT with SSFT boosts performance significantly (49.96% v.s. 44.54%), highlighting the importance of structured knowledge elicitation. Adding 8K extra QA pairs ("SSFT*") further improves performance across five subsets, demonstrating a surprising cross-lingual knowledge transfer (Lai et al., 2023; Qin et al., 2024). After SSFT, LLMs effectively use knowledge injected in one language to solve problems in others, surpassing traditional SFT. Additional comparisons in Appendix B.6 confirm that structure-aware syntheses can enhance knowledge application better than random syntheses.

In addition, we observe that the commonly used RAG (Lewis et al., 2020) strategy does not bring significant advantages to the MMedBench evaluation. The main reason lies in the gap between the pre-training corpus (comprising official knowledge statements from textbooks) and evaluated QA samples (originating from practical diagnosis records). Knowledge injection by (S)CPT and (S)SFT shows more advantages in this situation. In-depth investigations can be found in Appendix B.7.

**Comparision with Other Methods.** We also compare two advanced knowledge injection methods in Tab. 4 to further demonstrate our StructTuning's efficacy: (1) AdaptLLM (Cheng et al., 2023): domain knowledge injection by appending reading comprehension QAs to each CPT chunk, and (2) RAFT (Zhang et al., 2024): improving LLM's robustness to domain-specific retrieval-augmented generation using noisy retrieval-augmented SFT samples. According to the experimental results, AdaptLLM (Cheng et al., 2023) brings negligible improvement in the final performance (e.g., 46.79% v.s. 46.27% on the English subset), indicating such a chunk-level reading comprehension augmentation during CPT cannot help LLMs capture the entire structured domain knowledge. Concurrently, RAFT (Zhang et al., 2024) causes even worse performance, since the retrieval process introduces too many unrelated chunks and hurts LLM's QA judgments, especially when there exists a significant gap between user query and knowledge chunks in the medical diagnosis scenario.

## 5 Conclusion

This work pioneers in incorporating structure-aware methodologies to enhance domain knowledge injection into large language models. Through a novel SCPT-SSFT paradigm, we have set a new precedent for adapting LLMs to specialized domains, and the promising and scalable results underscore the viability and potential of our method. We hope to inspire further research in efficient and effective domain adaptation in the LLM community, moving a step closer to models that can truly emulate human intelligence.

**Limitation.** Our two-stage strategy introduces added computational complexity, where taxonomy extraction and data reorganization are required in the SCPT phase, and extra QA syntheses are optionally applied in the SSFT stage. In Appendix B, we provide further discussion with extensive empirical experiments. Despite the additional computational overhead introduced, our method achieves greater overall benefits and can reduce the reliance on large-scale LLMs (e.g., 70B models). We will delve into the investigations in future work.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Meta AI. 2024. Introducing meta llama 3: The most capable openly available llm to date.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023a. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. 2023b. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.

Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. 2000. A neural probabilistic language model. *Advances in neural information processing systems*, 13.

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Hung-Ting Chen, Fangyuan Xu, Shane A Arora, and Eunsol Choi. 2023. Understanding retrieval augmentation for long-form question answering. *arXiv preprint arXiv:2310.12150*.

Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. *arXiv preprint arXiv:2402.03216*.

Daixuan Cheng, Shaohan Huang, and Furu Wei. 2023. Adapting large language models via reading comprehension. In *The Twelfth International Conference on Learning Representations*.

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 preprint arXiv:2306.16092*.

Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. 2021. A dataset of information-seeking questions and answers anchored in research papers. *arXiv preprint arXiv:2105.03011*.

Kuicai Dong, Derrick Goh Xin Deik, Yi Quan Lee, Hao Zhang, Xiangyang Li, Cong Zhang, and Yong Liu. 2024. Multi-view content-aware indexing for long document retrieval. *arXiv preprint arXiv:2404.15103*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23.

Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*.

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. 2024. Deepseek-coder: When the large language model meets programming–the rise of code intelligence. *arXiv preprint arXiv:2401.14196*.

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*.

Tianyu Han, Lisa C Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K Bressem. 2023. Medalpaca–an open-source collection of medical conversational ai models and training data. *arXiv preprint arXiv:2304.08247*.

Wei He, Kai Liu, Jing Liu, Yajuan Lyu, Shiqi Zhao, Xinyan Xiao, Yuan Liu, Yizhong Wang, Hua Wu, Qiaoqiao She, et al. 2017. Dureader: a chinese machine reading comprehension dataset from real-world applications. *arXiv preprint arXiv:1711.05073*.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*.

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*.

Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L Richter, Quentin Anthony, Timothée Lesort, Eugene Belilovsky, and Irina Rish. 2024. Simple and scalable strategies to continually pre-train large language models. *arXiv preprint arXiv:2403.08763*.

9

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *arXiv preprint arXiv:2009.13081*.

Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.

David R Krathwohl. 2002. A revision of bloom's taxonomy: An overview. *Theory into practice*, 41(4):212–218.

Viet Dac Lai, Nghia Trung Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. 2023. Chatgpt beyond english: Towards a comprehensive evaluation of large language models in multilingual learning. *arXiv preprint arXiv:2304.05613*.

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*, 33:9459–9474.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Kai Liu, Zhihang Fu, Chao Chen, Wei Zhang, Rongxin Jiang, Fan Zhou, Yaowu Chen, Yue Wu, and Jieping Ye. 2024a. Enhancing llm's cognition via structurization. *arXiv preprint arXiv:2407.16434*.

Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinmeng Rao, Steven Zheng, Daiyi Peng, Diyi Yang, Denny Zhou, et al. 2024b. Best practices and lessons learned on synthetic data for language models. *arXiv preprint arXiv:2404.07503*.

Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024c. Dora: Weight-decomposed low-rank adaptation. *arXiv preprint arXiv:2402.09353*.

Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. 2020. K-bert: Enabling language representation with knowledge graph. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 2901–2908.

Robert L Logan IV, Nelson F Liu, Matthew E Peters, Matt Gardner, and Sameer Singh. 2019. Barack's wife hillary: Using knowledge-graphs for fact-aware language modeling. *arXiv preprint arXiv:1906.07241*.

Nick Mecklenburg, Yiyou Lin, Xiaoxiao Li, Daniel Holstein, Leonardo Nunes, Sara Malvar, Bruno Silva, Ranveer Chandra, Vijay Aski, Pavan Kumar Reddy Yannam, et al. 2024. Injecting new knowledge into large language models via supervised fine-tuning. *arXiv preprint arXiv:2404.00213*.

Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*.

Oded Ovadia, Menachem Brief, Moshik Mishaeli, and Oren Elisha. 2023. Fine-tuning or retrieval? comparing knowledge injection in llms. *arXiv preprint arXiv:2312.05934*.

Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*.

Libo Qin, Qiguang Chen, Yuhang Zhou, Zhi Chen, Yinghui Li, Lizi Liao, Min Li, Wanxiang Che, and Philip S Yu. 2024. Multilingual large language model: A survey of resources, taxonomy and frontiers. *arXiv preprint arXiv:2404.04925*.

Pengcheng Qiu, Chaoyi Wu, Xiaoman Zhang, Weixiong Lin, Haicheng Wang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2024. Towards building multilingual language model for medicine. *arXiv preprint arXiv:2402.13963*.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.

Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D Manning. 2024. Raptor: Recursive abstractive processing for tree-organized retrieval. In *The Twelfth International Conference on Learning Representations*.

Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. 2020. Ernie 2.0: A continual pre-training framework for language understanding. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8968–8975.

Zhengyang Tang, Xingxing Zhang, Benyou Wang, and Furu Wei. 2024. Mathscale: Scaling instruction tuning for mathematical reasoning. In *Forty-first International Conference on Machine Learning*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

10

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Musique: Multi-hop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, et al. 2023a. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity. *arXiv preprint arXiv:2310.07521*.

Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. 2023b. Huatuo: Tuning llama model with chinese medical knowledge. *arXiv preprint arXiv:2304.06975*.

Jianing Wang, Qiushi Sun, Xiang Li, and Ming Gao. 2023c. Boosting language models reasoning with chain-of-knowledge prompting. *arXiv preprint arXiv:2306.06427*.

Zekun Moore Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Man Zhang, et al. 2023d. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. *arXiv preprint arXiv:2310.00746*.

Yilin Wen, Zifeng Wang, and Jimeng Sun. 2023. Mindmap: Knowledge graph prompting sparks graph of thoughts in large language models. *arXiv preprint arXiv:2308.09729*.

Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Weidi Xie, and Yanfeng Wang. 2024. Pmc-llama: toward building open-source language models for medicine. *Journal of the American Medical Informatics Association*, page ocae045.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023a. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*.

Yan Xu, Mahdi Namazifar, Devamanyu Hazarika, Aishwarya Padmakumar, Yang Liu, and Dilek Hakkani-Tür. 2023b. Kilm: Knowledge injection into encoder-decoder language models. *arXiv preprint arXiv:2302.09170*.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.

Jifan Yu, Xiaozhi Wang, Shangqing Tu, Shulin Cao, Daniel Zhang-Li, Xin Lv, Hao Peng, Zijun Yao, Xiaohan Zhang, Hanming Li, et al. 2023. Kola: Carefully benchmarking world knowledge of large language models. *arXiv preprint arXiv:2306.09296*.

Li Yunxiang, Li Zihan, Zhang Kai, Dan Ruilong, and Zhang You. 2023. Chatdoctor: A medical chat model fine-tuned on llama model using medical domain knowledge. *arXiv preprint arXiv:2303.14070*.

Taolin Zhang, Chengyu Wang, Nan Hu, Minghui Qiu, Chengguang Tang, Xiaofeng He, and Jun Huang. 2022. Dkplm: decomposable knowledge-enhanced pre-trained language model for natural language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11703–11711.

Tianjun Zhang, Shishir G Patil, Naman Jain, Sheng Shen, Matei Zaharia, Ion Stoica, and Joseph E Gonzalez. 2024. Raft: Adapting language model to domain specific rag. *arXiv preprint arXiv:2403.10131*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. 2023. Siren's song in the ai ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219*.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.

Cai Zheng, Cao Maosong, and et al Haojiong. 2024. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*.

Zeyuan Allen Zhu and Yuanzhi Li. 2023a. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*.

Zeyuan Allen Zhu and Yuanzhi Li. 2023b. Physics of language models: Part 3.2, knowledge manipulation. *arXiv preprint arXiv:2309.14402*.

11

## A  Implementation Details

### A.1  Detailed Setup on LongBench

**Dataset Composition.** To focus on the investigation of knowledge injection, we choose 7 subsets from LongBench (Bai et al., 2023b) across single- and multi-document QA tasks in English and Chinese, and the remaining synthetic or code-orientated tasks are eliminated:

- **Single-Doc QA.** For single-document QA, we take three subsets from LongBench: (1) *Qasper* (Dasigi et al., 2021), featured by question-answering over NLP technical papers and annotated by NLP practitioners; (2) *MultiFieldQA* (Bai et al., 2023b), manually curated from multiple data sources and annotated by Ph.D. students; and (3) *MultiFieldQA-zh*, the Chinese split also provided by Bai et al. (2023b), covering multiple Chinese scenarios. *MultiFieldQA* contains 150 Context-Question-Answer triplets to test, and the others subsets include 200 pieces of test samples respectively.

- **Multi-Doc QA.** Multi-document QA requires LLMs to extract and combine information from multiple documents to derive the answer, which is generally more challenging than single-doc QA. We take four multi-hop QA datasets: (1) *HotpotQA* (Yang et al., 2018), containing 2-hop questions written by native speakers given two related paragraphs; (2) *2WikiMultihopQA* (Ho et al., 2020), involving up to 5-hop questions synthesized through manually designed templates on Wikipedia passages; (3) *MuSiQue* (Trivedi et al., 2022), carefully composed with up to 4-hop reasoning on an increased number of supporting and distracting context evidence; and (4) *Dureader* (He et al., 2017), developed based on Baidu Search and Baidu Zhidao and filtered by Bai et al. (2023b) to reduce the data noise. Each subset has 200 test samples.

In Single-Doc QA, we extract knowledge structures for each single passage; in Multi-Doc QA, we identify the knowledge structure across multiple passages for each test sample. There are ultimate 1350 *question-answer-passage(s)-(knowledge)structure* quadruples to evaluate knowledge injection approaches on LongBench.

**SFT Data-Synthesis.** We query Llama3-70B to generate 2,700 QA examples and remove those with over 0.5 F1-Score similarity to test samples to prevent data leakage. During inference, when the model can generate correct answers (corresponding to specific knowledge points) that haven't been seen during the SFT stage, we can ensure the knowledge is injected at the CPT stage and SFT only enhances the instruction-following capability. In practice, merely 13 out of 2700 (around 0.5%) synthetic data have over 0.5 F1-Score and are thus filtered out from the SFT data.

Tab. A1 statistics the semantic similarity (measured by BERTScore (Zhang et al., 2020)) between generated and GT questions and answers, and the results emphasize there is no knowledge leakage in the generated SFT data (they share poor semantic similarity across questions, answers, and QAs).

Table A1: Similarity statistics on synthetic SFT data and LongBench's test samples.

| Target | Question | Answer | Question-Answer |
|---|---|---|---|
| BERTScore | 0.277 | 0.106 | 0.093 |

### A.2  Detailed Setup on MMedBench

**Data for Evaluation.** The Multilingual Medical Benchmark (MMedBench) (Qiu et al., 2024) represents a comprehensive and diverse multilingual medical Question and Answering (QA) benchmark designed to evaluate models' capabilities of understanding and processing medical content.

MMedBench's robust dataset extends across 6 languages (*i.e.*, English, Chinese, Japanese, French, Russian, and Spanish) and 21 medical fields, which include, but are not limited to, Internal Medicine, Biochemistry, Pharmacology, Psychiatry, and many others. It provides 45,048 training pairs and 8,518 testing pairs for diverse learning and testing scenarios. The training split is specifically designed for domain-specific finetuning of large language models (LLMs), while the entire testing set allows for a precise assessment of multi-choice question-answering performance. Statistics on six languages are displayed in Tab. A2. Notably, the benchmark includes scenarios where questions may have multiple correct answers (*i.e.*, in Japanese and French subsets), introducing additional complexity for the model evaluation process.

**Data for Continual Pre-Training.** To investigate high-quality domain knowledge injection for LLMs and the scalability of injection methods, we curate a series of training corpus sizes from the 25.5B MMedC (Qiu et al., 2024) dataset, including 0.1%, 0.3%, 0.5%, 1%, and 5%, which respectively takes 30M, 76M, 132M, 250M, and 1.2B tokens.

Table A2: Sample statistics on MMedBench's QA data.

| Split | English | Chinese | Japanese | French | Russian | Spanish | Total |
|---|---|---|---|---|---|---|---|
| Train | 10,178 | 27,400 | 1,590 | 2,171 | 1,052 | 2,657 | 45,048 |
| Test | 1,273 | 3,426 | 199 | 622 | 256 | 2,742 | 8,518 |

Table A3: Sample statistics on different training data settings.

| Stage | Ratio | English | Chinese | Japanese | French | Russian | Spanish | Total |
|---|---|---|---|---|---|---|---|---|
| CPT | 0.1% | 6.9M | 8.8M | - | 4.1M | 4.9M | 5.4M | 30.1M |
| CPT | 0.3% | 26.1M | 21.5M | - | 8.1M | 10.3M | 10.1M | 76.1M |
| CPT | 0.5% | 35.9M | 27.2M | 4.5M | 14.0M | 24.9M | 26.1M | 132.6M |
| CPT | 1.0% | 44.1M | 39.8M | 34.2M | 45.1M | 47.9M | 38.4M | 249.5M |
| CPT | 5.0% | 169.4M | 227.8M | 119.8M | 232.7M | 235.5M | 226.1M | 1.2B |
| SFT | - | 18.8K | 39.1K | 1.6K | 5.3K | 5.9K | 7.5K | 78.2K |

We sorted the document-level text content (including textbooks and other corpus from websites and wikipedia) by length in descending order and progressively included more samples to expand the training dataset at larger scales, where Tab. A3 provides detailed statistics. In particular, as MMedC does not provide English textbooks in its released data (due to copyright issues), we collect 18 English textbooks (Jin et al., 2020) as an alternative for English medical knowledge injection, since they have a common OCR source with MMedC (Qiu et al., 2024). These English textbooks take around 21.5M tokens.

**Knowledge Structure Extraction.** For textbooks, we split the data to chunks (knowledge points) within 3072 tokens, and use our specifically developed 7B-size LLM to extract the structured medical knowledge system. For non-textbook data (for instance, MMedC does not provide Japanese textbooks), a clustering-based technique (Sarthi et al., 2024) is adopted to recursively build knowledge structures from fragmented text segments. Fig. A1 presents an example to illustrate the two kinds of knowledge structure extraction processes, Tab. A4 displays the overall statistics on extracted knowledge structures for the final 1.2B tokens.

**Data for Supervised Fine-Tuning.** As introduced in Sec. 3.3, we prompt Llama3-70B (Dubey et al., 2024) to build the structure-aware answer explanations on top of the raw SFT samples in MMedBench's training split, and generate extra QA pairs by traversing the extracted knowledge structure from textbooks. The final quantity statistics are presented in Tab. A3. Note that the 70B-size model is not necessary to synthesize QAs and explanations, since we have provided relevant text sources (retrieved from the training corpus) to sup-

Table A4: Knowledge structure on 1.2B training corpus.

| Lang. | Book | Chapter | Section | KnowledgePoints |
|---|---|---|---|---|
| 6 | 2933 | 14,411 | 23,239 | 180,793 |

plement medical knowledge. In this way, LLMs only need to perform in-context comprehension, rather than generate new QAs based on their own knowledge, which can also be achieved by smaller models like Qwen2.5-7B (Yang et al., 2024).

### A.3 Terminology Explanation

**Knowledge Structures.** We extract the domain knowledge structure for each textbook, where Fig. 3 presents an example, and combine the medical knowledge for six languages in MMedBench. As the (S)CPT corpus for Japanese is collected from Wikipedia rather than textbooks, we derive a single knowledge structure for Japanese medicine.

**Knowledge Paths and Branches.** Fig. A2 shows an example of how we define the knowledge paths and branches of the extracted knowledge structure for SSFT data synthesis.

- A **path** means a knowledge path from the domain summary (*e.g.*, *Biochemistry*) to specific knowledge points (*e.g.*, *Lipid Metabolism and Cholesterol Transport*): "Biochemistry – Overview of lipoprotein metabolism, hormone synthesis – Lipoprotein Metabolism – Lipid Metabolism and Cholesterol Transport".

- A **branch** means the knowledge branch of the tree structure. If a question is related to two knowledge points (*e.g.*, *Lipid Metabolism and Cholesterol Transport* and *Pancreatic zymogen activation*) at different branches of the knowledge tree, the knowledge path contains two branches, which becomes the right-bottom part of Fig. A2.

13

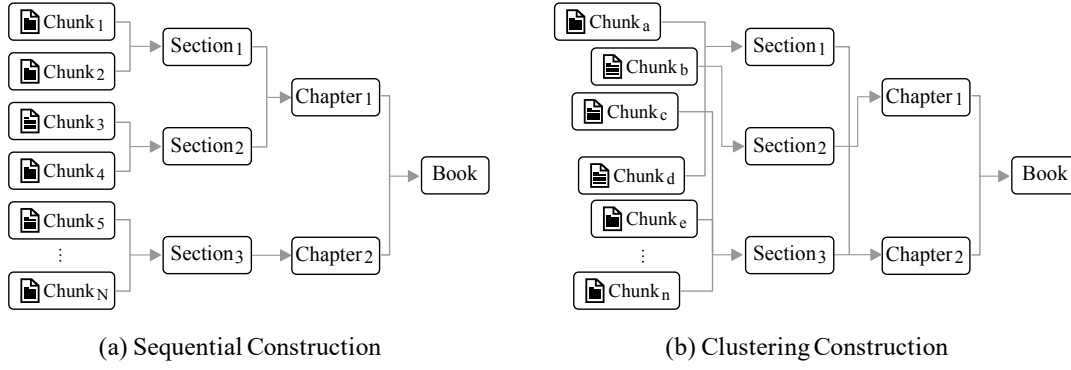(a) Sequential Construction  (b) Clustering Construction

Figure A1: Knowledge structure extraction from (a) sequential chunks (*e.g.*, from textbooks) by our specialized 7B model and (b) separated trunks (*e.g.*, from websites) by clustering-based methods (Sarthi et al., 2024). Here the terms of "Section", "Chapter", and "book" are just examples to help illustrate the knowledge structure.
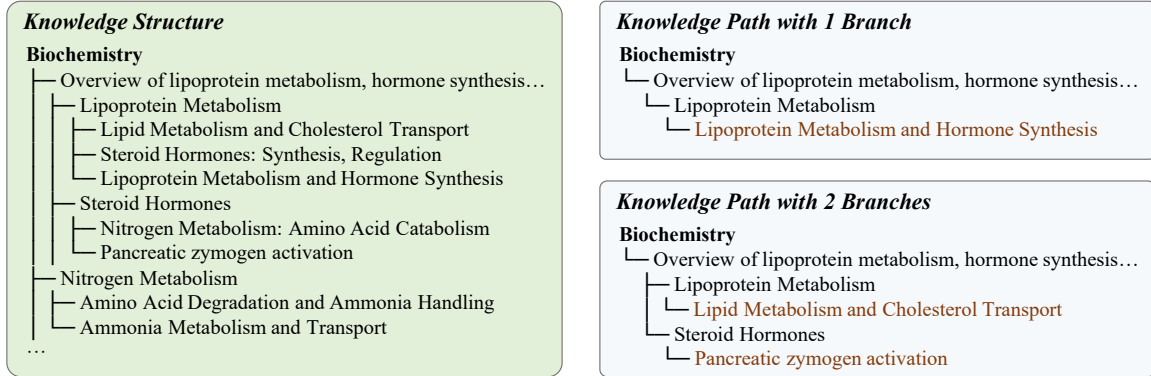


Figure A2: Definition and example of knowledge paths and branches.

## A.4 Resource Requirement

We use 8 NVIDIA A100-80G GPUs to train all the models, and leverage 1-2 NVIDIA A100-80G GPUs for inference.

## B Additional Experiments

### B.1 Knowledge Structure Extraction

Extracting domain knowledge structure is a prerequisite for subsequent knowledge injection (including both SCPT and SSFT) for language models. In Sec. 3.1, we propose a bottom-up strategy to re-chunk the texts from domain textbooks, summarize a title for each chunk, and send the title list to a specialized 7B model to derive the knowledge structure. The prompt template is displayed in Appendix D.

In fact, we argue that due to the language diversity, a perfectly recovered table of contents of textbooks is unnecessary for domain knowledge injection. A reasonable knowledge structure is sufficient enough. In Tab. A5, we individually adopt few-shot GPT-3.5-Turbo (Brown et al., 2020) and LLaMA3-70B (Dubey et al., 2024) models to extract medical knowledge structure from 18 English textbooks (Jin et al., 2020) (with 26.1M tokens)

Table A5: Comparison of models to extract knowledge structures on 26M English corpus.

| Model | Improvement | Time Cost | Extra Fee |
|---|---|---|---|
| GPT-3.5-Turbo | +3.58 | - | 15$ |
| LLaMA3-70B | +3.72 | 1.5h | - |
| Ours-7B | +3.69 | **0.2h** | - |

for subsequent knowledge injection (the backbone LLM is LLaMA2-7B (Touvron et al., 2023b)). Although they present 3.6%-3.7% enhancement on MMedBench (Qiu et al., 2024)'s English test set (denoted as "improvement"), leveraging GPT-3.5-Turbo and LLaMA3-70B is either expensive or time-consuming. GPT-3.5-Turbo costs around 15 dollars to process 26M tokens, while LLaMA3-70B takes around 1.5 hours on 2 A100-80G GPUs, of which both limit scaling data pre-processing for structure-aware knowledge injection.

Inspired by Liu et al. (2024a), we distilled the knowledge structure extraction capability from giant LLMs to a LLaMA2-7B model via SFT. In particular, we instruct LLaMA3-70B to generate 22K training examples (pairs of raw knowledge points and extracted knowledge structures) from Wikipedia, and train a LLaMA2-7B model at a

14

Table A6: Structure-aware knowledge injection to Llama2 (Touvron et al., 2023b) model series.

| Model | English | Chinese | Japanese | French | Russian | Spanish | **Average** |
|-------|---------|---------|----------|--------|---------|---------|-------------|
| InternLM2-7B | 57.27 | 77.55 | 47.74 | 41.00 | 68.36 | 59.59 | 58.59 |
| **+Ours** | **60.80** | **79.19** | **50.75** | **45.34** | **75.39** | **66.85** | **63.05** |
| Llama2-7B | 43.36 | 50.29 | 25.13 | 20.90 | 66.80 | 47.10 | 42.26 |
| **+Ours** | **49.41** | **65.15** | **36.68** | **35.21** | **69.14** | **50.62** | **51.04** |
| Llama2-13B | 51.37 | 57.97 | 32.66 | 25.08 | 69.92 | 52.99 | 48.33 |
| **+Ours** | **53.02** | **68.30** | **37.78** | **41.71** | **70.70** | **55.51** | **54.50** |

Table A7: Model evaluation on comment benchmarks.

| Injection | MMLU | C-Eval | TruthfulQA | WinoGrande | ARC_c |
|-----------|------|--------|------------|------------|-------|
| Before | 0.46 | 0.35 | 0.49 | 0.52 | 0.51 |
| After | 0.43 | 0.35 | 0.43 | 0.51 | 0.56 |

Table A8: Training costs for knowledge injection.

| Paradigm | Corpus | Improve. | Preprocess | Total |
|----------|--------|----------|------------|-------|
| CPT+SFT | 100% | 100% | - | >30d |
| SCPT+SSFT | 0.3% | 50% | 0.6h | 4.5h |
| SCPT+SSFT | 5% | 100% | 9.7h | **3d** |

batch size of 128 and a learning rate of 2e-5 for 1 epoch. After utilizing the specialized 7B model to identify the knowledge structure in medical textbooks, as shown in Tab. A5, the results translate to comparable performance on structure-aware knowledge injection. Meanwhile, the inference cost significantly decreases to 0.3 hours, which is more scalable to handle a larger domain corpus.

## B.2 Evaluation of Approach's Generalization

In addition to the Llama3 models in previous experiments, we also investigate the generalization ability of our SCPT+SSFT paradigm on Llama2 (Touvron et al., 2023b) and InternLM2 (Zheng et al., 2024) model series. As shown in Tab. A6, our method leads to consistently significant improvements on InternLM2-7B (+4.46%), Llama2-7B (+8.78%) and Llama2-13B (+6.17%) backbone models. The results further demonstrate the generalizability and scalability of our StructTuning strategy across model architectures and sizes.

## B.3 Investigation on Model Overfitting

Here we provide a supplementary evaluation on common benchmarks in Tab. A7, which indicates the current methodology of structure-aware knowledge injection does not significantly hurt LLM's general capabilities, and even brings slight enhancement to the ARC_c benchmark (maybe the knowledge structure enhances the reasoning ability). As previous research states, the issue of overfitting can be further mitigated by incorporating common QA examples, and we leave this in our future work.

## B.4 Comparison on Training Costs

In Tab. A8, we quantify the total training cost on 8 A100-80G GPUs. According to Qiu et al. (2024), the conventional CPT+SFT paradigm

on 25.5B medicine corpus takes more than 30 days to derive the SOTA MMedLM model. In our SCPT+SSFT framework, although the preprocessing (*i.e.*, knowledge structure extraction) introduces an extra 0.6 hours to process 0.3% data (around 76M tokens), the total training process only costs 4.5 hours. As suggested in Fig. 6, when 5% training data is leveraged for knowledge injection to achieve 100% improvement, the overall cost is limited to 3 days, much less than the CPT+SFT approach with more than a month. Those analyses further demonstrate the efficacy and efficiency of our structure-aware knowledge injection framework.

## B.5 Ablation on Structured Knowledge Injection

During the Structure-aware Continual Pre-Training (SCPT) stage, we proposed to learn specific text chunks (knowledge points) in the condition of the mindmap inputs (knowledge structures), in order to relate the knowledge points to corresponding structure nodes. In this section, we conduct a series of ablation studies to investigate the design efficacy. The vanilla CPT+SFT paradigm is adopted as the comparison baseline, where the Llama2-7B model is trained with CPT and SFT data on the English subset of MMedBench, while tested on all subsets across six languages. The hyper-parameter settings follow the main experiment in our manuscript. The empirical results are presented in Tab. A9.

First, we investigate the choice of formatting template to convert the knowledge structure to a mindmap condition. In particular, we try to fix the template to convert all knowledge mindmaps for SCPT, and randomly select two templates to repeat

Table A9: Ablation studies of SCPT on MMedBench subsets. The base model is Llama2-7B.

| Adaptation | English | Chinese | Japanese | French | Russian | Spanish | Average |
|---|---|---|---|---|---|---|---|
| CPT+SFT | 46.27 | 32.57 | 26.13 | 17.36 | 50.00 | 40.63 | 35.49 |
| Ours-FixTmpl1 | 48.27 | 32.86 | 20.61 | 23.70 | 56.17 | 42.56 | 37.36 |
| Ours-FixTmpl2 | 48.10 | 32.99 | 21.23 | 23.97 | 55.97 | 43.10 | 37.56 |
| Ours-RemoveL1 | 47.90 | 32.90 | 20.11 | 24.63 | 57.10 | 43.43 | 37.68 |
| Ours-RemoveL2 | 48.05 | 33.63 | 21.62 | 24.11 | 57.49 | 43.13 | <u>38.00</u> |
| Ours-RemoveL3 | 48.47 | 33.14 | 20.15 | 23.33 | 56.86 | 43.65 | 37.60 |
| Ours-NTPLoss | <u>48.99</u> | 33.15 | 20.57 | 25.31 | 56.78 | 42.94 | 37.96 |
| **Ours-Full** | **49.10** | 33.92 | 18.33 | 27.14 | 57.42 | 43.73 | **38.27** |

the experiment. According to Tab. A9, fixed SCPT templates lead to inferior performance against randomly choosing the template from the diversified 20 template pool. This is consistent with Zhu and Li (2023a)'s observation, that text rewriting can provide better knowledge augmentation for large language models.

Then, we explore the impact of the extracted knowledge structure itself. In MMedBench, a 3-layer knowledge structure (follow the *chapter-section-subsection* hierarchy) is constructed for each textbook, and we respectively remove the 1st (chapter), 2nd (section), and 3rd (subsection) layer of the hierarchy during knowledge injection. As Tab. A9 shows, removing the top layer (chapter) leads to the worst performance of 47.90%, because the remaining knowledge points cannot effectively relate to each other without the organization of the top layer. On the other hand, removing the bottom layer (subsection) performs slightly better on the English subset (because of the controlled structure-information lost), but hinders the cross-language knowledge utilization on the remaining subsets (*e.g.*, 37.60% on average across six languages).

Finally, we revisit the modeling choice of the mindmap-conditioning learning. Specifically, we try to turn the conditional modeling $p(\boldsymbol{x}^k|\boldsymbol{s}^k)$ back to complete next-token prediction $p(\boldsymbol{x}^k, \boldsymbol{s}^k)$ (the next-token prediction loss is computed on mindmap condition as well). According to Tab. A9, the performance is slightly inferior to our full version of SCPT strategy (*e.g.*, 37.96% *v.s.* 38.27% on Average). Therefore, we reserve conditional modeling for our SCPT stage.

## B.6 Ablation on SFT Data Synthesis

In Sec. 4.2, we compared our structure-aware knowledge injection with conventional CPT+SFT

paradigm on MMedBench. On its English subset, we ablated the training components of our method, and found that the newly synthesized 8K SSFT data (by traversing the extracted knowledge structure) can inspire LLMs' cross-language capability to apply the learned structured knowledge to solve practical diagnosis problems. Here, we follow Liu et al. (2024b) to randomly generate another 8K QA pairs for SFT alignment for further comparison, denoted as "SFT*". We randomly sample medical texts and instruct Llama3-70B (Dubey et al., 2024) for data synthesis, without the knowledge structure provided. Tab. A10 indicates that "SFT*" brings slight enhancement to the English test subset, but the average accuracy drops to 34.51% instead. The results further demonstrate our method's efficacy in the application of the injected, structured domain knowledge.

## B.7 In-depth Comparison on Retrieval-Augmented Generation

In Sec. 4.2, we briefly compare RAG adaptation and injection-based approaches in the MMedBench dataset, and this section provides more implementation details and further investigations on the popular retrieval-augmented generation approach.

**Experimental Settings.** On the implementation of the RAG baseline, we utilize the BAAI/bge-m3 (Chen et al., 2024) embedding model for dense retrieval, due to its state-of-the-art and multilingual semantic retrieval ability. For the experiments in Tab. 4, we take the same 26M English CPT data as the knowledge base, re-chunk the data corpus for every 512 tokens, and retrieve top-3 related chunks as context inputs for LLM's generation process. The retrieval process is implemented using the LlamaIndex[2] framework.

**Additional Experiments.** We also conduct

---

[2]https://www.llamaindex.ai/

Table A10: Comparison of SFT data synthesis strategies on MMedBench. The backbone LLM is the same Llama2-7B model after SCPT on English textbooks.

| SFT synthesis | English | Chinese | Japanese | French | Russian | Spanish | Average |
|---|---|---|---|---|---|---|---|
| - | 46.50 | 32.14 | **20.10** | 18.17 | 53.91 | 39.97 | 35.13 |
| SFT* | 47.13 | 32.49 | 16.58 | 16.72 | 51.95 | 42.16 | 34.51 |
| **SSFT*** | **49.10** | **33.92** | 18.33 | **27.14** | **57.42** | **43.73** | **38.27** |

Table A11: Ablation on the hyper-parameter settings for the RAG baseline.

| ChunkSize | 256 | | | 512 | | | 1024 | | |
|---|---|---|---|---|---|---|---|---|---|
| RetrieveNum | 10 | 5 | 3 | 5 | 3 | 2 | 3 | 2 | 1 |
| Accuracy | 35.08 | 37.67 | 38.04 | 36.42 | **38.12** | 37.99 | 35.00 | 36.89 | 38.07 |

Table A12: Attempts to integrate hybrid-search and reranker models.

| Hybrid | Reranker | Accuracy |
|---|---|---|
| × | × | **38.12** |
| √ | × | 37.97 |
| × | √ | 37.52 |
| √ | √ | 37.75 |

a variety of experiments to evaluate the hyper-parameters for the RAG baseline. As shown in Tab. A11, changing the chunk size and retrieved chunk number cannot bring any significant benefits. The core reason lies in the gap between user query and retrieved chunks. In particular, user queries contain many descriptive and quantitative sentences and numbers (such as the example in Fig. A3, "They enrolled 800 patients in the study, half of which have breast cancer".), and may even talk about an entirely new thing that has not been recorded in the knowledge base.

Furthermore, we also try to use the hybrid (dense+sparse) search strategy and larger rerank model (BAAI/bge-reranker-v2-m3 (Chen et al., 2024)) to enhance the retrieval quality. However according to the results in Tab. A12, the semantic gap between user queries and retrieved chunks still exists. Introducing the hybrid search and rerank model even gets worse performance (*e.g.*, the keyword *age* may be considered a key factor for hybrid search, but it cannot help to derive the answer of test sensitivity).

**Conclusion.** RAG may assist in some knowledge-intensive tasks for information-seeking, but will encounter problems when there exists a significant semantic gap between user query and retrieved documents. MMedBench is a typical scenario, where LLMs are asked to derive medical diagnoses with proper reasoning according to the descriptions of patients or medical examinations.

In this case, the retrieval process introduces too many unrelated chunks and hurts LLM's QA judgments. Fig. A3 provides an example where the retrieved chunks are actually unrelated to the complicated user query (the user asks about the analysis of a given research study, but the retrieved documents contain several keywords, *e.g.*, *age*, while having nothing to do with the *blood test study*.)

### B.8 Knowledge Memorization on LongBench

Here, we also use lexical ROUGE-L (Lin, 2004) and semantic BERTScore (Zhang et al., 2020) to quantify the memorization of injected knowledge structures, by comparing the mindmap in models' responses (as Fig. 5 displays) with ground-truth answers. The results in Tab. A13 indicate a relatively good memorization of the injected knowledge mindmap, emphasizing the efficacy of our SCPT strategy.

Table A13: MindMap Recall

| F1-Score | BERTScore |
|---|---|
| 0.61 | 0.87 |

### B.9 F1-Score Evaluation on LongBench

In Sec. 4.1, we mainly follow Zhu and Li (2023b) to investigate the memorization and understanding of injected knowledge by calculating the knowledge recall in models' responses. Here we report the F1-score measure over the Closed-Book QA (CBQA) settings for a thorough comparison. Note that here we use the vanilla question prompt to obtain concise answers, instead of the CoT prompt used in Sec. 4.1 to elicit models' memorized knowledge. The evaluated models are the same as Sec. 4.1.

In Tab. A14, we report the Closed-Book QA (CBQA) baseline by traditional CPT+SFT to inject passage contents into model parameters, and

Figure A3: An example of retrieved document/chunk based on a given query.

Table A14: F1 Score evaluation of Closed-Book QA (CBQA) tasks on the LongBench (Bai et al., 2023b) dataset. The best results are marked in **bold**, and the secondary results are marked with <u>underlines</u>. The backbone model is Llama2-7B (Touvron et al., 2023b).

| Task | Adaptation | SingleDoc-QA | | | MultiDoc-QA | | | | Average |
|------|-----------|--------|------|--------|------|-------|-------|------|---------|
| | | Qasper | MFQA | MFQAzh | HpQA | 2Wiki | Musiq | Duzh | |
| CBQA | CPT+SFT | <u>16.8</u> | <u>23.1</u> | 13.2 | <u>21.3</u> | 19.1 | <u>10.4</u> | 13.4 | <u>16.8</u> |
| | SCPT+SFT | 15.2 | 21.5 | <u>15.2</u> | 14.9 | <u>19.8</u> | 5.6 | <u>14.1</u> | 15.2 |
| | SCPT+SSFT | **19.7** | **23.5** | **19.5** | **26.4** | **24.1** | **12.2** | **15.4** | **20.1** |

supplement the experiment of our SCPT+SSFT technique for comparison. According to the results shown in Tab. A14, our SCPT+SSFT approach successfully boosts the closed-book QA performance to 20.1% on average. The results are consistent with Sec. 4.1, which jointly demonstrate the effectiveness of structure-aware knowledge injection for large language models.

## C  Detailed Related Work

**Domain Adaptation for LLMs**. While pre-trained LLMs possess promising capabilities, their performance is often hampered by the scope and recency of their training data, which particularly affects smaller models in downstream applications (Zhao et al., 2023; Wang et al., 2023a). Continual Pre-Training (CPT) addresses this by perpetually updating a pre-trained model with domain-specific content (Sun et al., 2020; Xu et al., 2023b). , with parameter-efficient tuning methods devised to curtail training costs (Hu et al., 2021; Liu et al., 2024c). To keep pace with the latest information, models can be fine-tuned with supervised instruction-response pairs (SFT), thus staying current with the advancing knowledge landscape (Mecklenburg et al., 2024; Qiu et al., 2024). Existing literature confirms that combining CPT and SFT is effective for LLMs to remain precise and up-to-date in dynamic fields like medicine (Wang et al., 2023b; Qiu et al., 2024) and coding (Roziere et al., 2023; Guo et al., 2024). Our study builds upon this CPT-SFT framework, innovating with SCPT-SSFT strategies to efficiently and effectively infuse domain knowledge with the inherent structure hierarchy.

**Conditional Language Modeling**. The idea of continual pre-training language models on domain corpus in the condition of the knowledge structure is mainly inspired by CTRL (Keskar et al., 2019). Keskar et al. (2019) demonstrates the effectiveness of steering text generation through control codes (one or two words) that signify the desired genre, style, or task. In the era of LLM, system prompt plays a similar role in controlling models' responses to adapt to different needs and functionalities, such as role-playing, language style transfer, task setting, and behavior setting (Brown et al., 2020; Wang et al., 2023d; Bai et al., 2023a). Our SCPT approach extends the control codes or system prompts to domain-specific knowledge structures, so as to guide the learning process and tailor the model's output more closely to specialized fields.

**Structure-aware Knowledge Aggregation**. Knowledge structure has been widely explored in the recent LLM community. In conventional paradigms, researchers extract entity-relation-entity triplets from texts to construct knowledge graphs (Pan et al., 2024), to enhance LLMs's factual knowledge and logical reasoning by feature aggregation (Liu et al., 2020; Zhang et al., 2022), prompt engineering (Wen et al., 2023; Wang

et al., 2023c), information searching (Logan IV et al., 2019), training data synthesis (Tang et al., 2024), etc. In these cases, each node corresponds to either a specific entity or an abstract concept, lacking the capability to present an informative and self-contained *knowledge point*. Some works have recently related a piece of descriptive text to a knowledge point, and constructed the knowledge structure for LLMs' retrieval-augmented generation (Sarthi et al., 2024; Dong et al., 2024), where the top-to-down retrieval provides precise information-seeking paths along the knowledge structure. In this paper, we extend the structure-aware knowledge aggregation to LLMs' training phase, injecting the whole domain knowledge structure into LLMs' by linking training samples to corresponding knowledge points and reasoning paths.

**Data Augmentation and Synthesis**. Due to the lack of high-quality datasets, data augmentation has emerged as a promising solution to mimic real-world patterns (Liu et al., 2024b). Traditional methods aim to artificially expand the training dataset size (Xu et al., 2023a; Mukherjee et al., 2023) or generate entirely new samples that could help models learn better or adapt to specific tasks (Tang et al., 2024). Yet, they often overlook the structured nature of domain knowledge, and the aimlessly generated samples may also lack diversity (Ovadia et al., 2023; Mecklenburg et al., 2024), leading to potentially suboptimal training outcomes when applied for domain adaptations (Mecklenburg et al., 2024; Tang et al., 2024). By contrast, our SSFT design is an innovative departure to address the challenge of retaining and utilizing the structured knowledge inherent in domain-specific content.

## D  Prompt Template for Knowledge Structure Extraction

Fig. A4 displays the prompt template to query our specialized 7B model to extract knowledge structure on given knowledge points, which introduces the task definition, detailed instruction, and output formats to illustrate the process.

You are a sophisticated AI expert in Natural Language Processing (NLP), with the specialized capability to deconstruct complex sentences and map their semantic structure. Your task is to analyze the given sentences to extract and represent the intrinsic semantic hierarchy systematically.

Follow this approach to ensure clarity and utility in your analysis:
**1. **Comprehension**:** Begin with a thorough reading to understand the overarching theme of the input sentences.
**2. **Defining Scope**:** Summarize the central theme to establish the scope of the semantic analysis.
**3. **Aspect Breakdown**:** Identify the core aspects of the discussion. For any aspect with additional layers, delineate "SubAspects" and repeat as necessary for complex structures. Each aspect or subaspect should be highly summarized and self-contained.
**4. **Mapping**:** Assign sentence numbers to their respective aspects or subaspects, indicating where in the text they are addressed.

Structure your analysis in a YAML format according to this template, and ensure the format is clean, well-organized, and devoid of extraneous commentary:
```yaml
Scope: <central theme summary>
Aspects:
 - AspectName: <main aspect>
   SentenceRange:
     start: <start sentence number>
     end: <end sentence number>
   SubAspects:
 - AspectName: <subaspect>
   SentenceRange:
     start: <start sentence number>
     end: <end sentence number>
   # Recursively repeat "SubAspects" structure as needed
 # Adjust "SubAspect" entries as needed
# Adjust "Aspect" entries as needed
```

Now, analyze the provided sentences with the structured analytical process, and output your analysis in the structured YAML format. NOTE: each aspect or subaspect should be highly summarized and self-contained, which covers at least two sentences, except for introduction or conclusion aspects.

## Content
```
{title_list}
```

## Analysis

Figure A4: Prompt template for knowledge structure identification.

```
(
"In the realm of `{field}`, a conceptual mindmap is depicted using a tree-like structure "
"to represent hierarchical relationships and thematic branches:\n\n"
"```\n{mindmap}\n```\n\n"
"Within this organized layout of `{field}`, the detailed subsection on `{section}` is described as:\n\n"
),
(
"The area of `{field}` unfolds into a rich and detailed structure, encapsulating a diverse array of topics and their interconnections. "
"These topics are organized in a manner that reflects their relationships and thematic relevance to one another, depicted through a
structured diagram:\n\n"
"```\n{mindmap}\n```\n\n"
"Within this elaborate organization, the concept of `{section}` serves as a detailed exploration into a specific element of `{field}`:\n\n"
),
(
"The `{field}` sector is structured through a complex network of concepts and categories, "
"as reflected in the following outlined representation:\n\n"
"```\n{mindmap}\n```\n\n"
"Zooming in on a discrete element of this intellectual landscape, the topic tagged as `{section}` "
"covers specific subject matter related to `{field}`:\n\n"
),
(
"Exploring the `{field}`, structured insights reveal a network of thematic areas. "
"The essence is captured in a concise diagram:\n\n"
"```\n{mindmap}\n```\n\n"
"A closer look at the portion labeled `{section}` unveils a segment rich in detail, contributing "
"to the broader understanding of `{field}`:\n\n"
),
(
"`{field}` encompasses a diverse array of themes, organized for clarity. "
"The visual schema below illustrates this organization:\n\n"
"```\n{mindmap}\n```\n\n"
"Investigating `{section}` furnishes insight into a selected theme within `{field}`, enriching the overall comprehension:\n\n"
),
(
"Contextualizing within the broader spectrum of `{field}`, the organizational structure is delineated as follows:\n\n"
"```\n{mindmap}\n```\n\n"
"Delving into `{section}`, an integral component of the `{field}` fabric, enriches the grasp of the thematic variety and depth.\n\n"
),
(
"Within the expansive knowledge area of `{field}`, an organizational schema is represented as:\n\n"
"```\n{mindmap}\n```\n\n"
"Exploring `{section}` reveals a critical facet of `{field}`, offering insights into its thematic diversity and detail.\n\n"
),
(
"The discipline of `{field}` is encapsulated by a series of interlinked concepts, mapped out as:\n\n"
"```\n{mindmap}\n```\n\n"
"The segment labeled `{section}` delves into a particular topic within `{field}`, "
"illuminating a slice of the broader intellectual landscape:\n\n"
),
(
"Navigating through `{field}`, one encounters a structured depiction of knowledge as illustrated below:\n\n"
"```\n{mindmap}\n```\n\n"
"Within this schema, `{section}` serves as a gateway to a distinct area of interest, "
"shedding light on specific sections of `{field}`:\n\n"
),
(
"Diving into the `{field}` landscape, a coherent outline presents itself, showcasing the interconnectedness of its themes:\n\n"
"```\n{mindmap}\n```\n\n"
"Focusing on the section of `{section}`, it serves as a focal point into nuanced exploration within the vast `{field}` territory:\n\n"
),
(
"The sphere of `{field}` unfolds as a network of insights and principles, outlined for comprehensive understanding:\n\n"
"```\n{mindmap}\n```\n\n"
"The exploration of `{section}` unveils a segment pivotal to the fabric of `{field}`, providing a perceiving lens:\n\n"
),
(
"As we chart the terrain of `{field}`, a constellation of concepts emerges, graphically represented as follows:\n\n"
"```\n{mindmap}\n```\n\n"
"Focusing on the component marked as `{section}`, we uncover layers within `{field}` that resonate with both breadth and depth, offering a
panoramic view into the diverse thought processes and methodologies encapsulated within.\n\n"
),
(
"`{field}` is organized into various key areas, as shown in the diagram below:\n\n"
"```\n{mindmap}\n```\n\n"
"`{section}` highlights a core area, integral for understanding the overall structure of `{field}`:\n\n"
),
(
"The structure of `{field}` is detailed below:\n\n"
"```\n{mindmap}\n```\n\n"
"A deeper understanding of `{field}` can be achieved by examining `{section}`, a vital element of its framework:\n\n"
),
(
"Overview of `{field}`'s foundational structure is as follows:\n\n"
"```\n{mindmap}\n```\n\n"
"Exploring `{section}` reveals its crucial role in comprehending the comprehensive schema of `{field}`:\n\n"
),
(
"`{field}` encompasses a range of interconnected topics, illustrated in the diagram below:\n\n"
"```\n{mindmap}\n```\n\n"
"The examination of `{section}` provides insight into how key concepts within `{field}` are interrelated:\n\n"
),
(
"Key elements within `{field}` can be organized as follows:\n\n"
"```\n{mindmap}\n```\n\n"
"Investigating the component of `{section}` is essential for grasping the complex dynamics in the `{field}` realm:\n\n"
),
(
"The `{field}` includes various components as detailed in the following structure:\n\n"
"```\n{mindmap}\n```\n\n"
"Focusing on `{section}` offers an opportunity to explore one of the numerous elements that comprise the `{field}`:\n\n"
),
(
"Within the scope of `{field}`, multiple dimensions unfold as depicted below:\n\n"
"```\n{mindmap}\n```\n\n"
"Delving into `{section}` contributes to a broader understanding of the diverse elements that construct the landscape of `{field}`:\n\n"
),
(
"Comprehensive knowledge of `{field}` can be achieved by examining its individual components, as depicted below:\n\n"
"```\n{mindmap}\n```\n\n"
"An exploration of `{section}` sheds light on its unique contribution to the `{field}`:\n\n"
)
```

Figure A5: Full template pool for mindmap conversion with 20 diversified templates.