

# Evontree: Ontology Rule-Guided Self-Evolution of Large Language Models

Mingchen Tu<sup>1</sup>, Zhiqiang Liu<sup>1</sup>, Juan Li<sup>1</sup>, Liangyurui Liu<sup>2</sup>, Junjie Wang<sup>3</sup>,  
Lei Liang<sup>3</sup>, and Wen Zhang<sup>1\*</sup>

<sup>1</sup> Zhejiang University, Hangzhou, China

<sup>2</sup> University of Electronic Science and Technology of China, Chengdu, China

<sup>3</sup> Ant Group, Hangzhou, China

{mingchentz,zhiqiangliu,zhang.wen}@zju.edu.cn

**Abstract.** Large language models (LLMs) have demonstrated exceptional capabilities across multiple domains by leveraging massive pre-training and curated fine-tuning data. However, in data-sensitive fields such as healthcare, the lack of high-quality, domain-specific training corpus hinders LLMs’ adaptation for specialized applications. Meanwhile, domain experts have distilled domain wisdom into ontology rules, which formalize relationships among concepts and ensure the integrity of knowledge management repositories. Viewing LLMs as implicit repositories of human knowledge, we propose Evontree—a novel framework that leverages a small set of high-quality ontology rules to systematically extract, validate, and enhance domain knowledge within LLMs, without requiring extensive external datasets. Specifically, Evontree extracts domain ontology from raw models, detects inconsistencies using two core ontology rules, and reinforces the refined knowledge via self-distilled fine-tuning. Extensive experiments on medical QA benchmarks with Llama3-8B-Instruct and Med42-v2 demonstrate consistent outperformance over both unmodified models and leading supervised baselines, achieving up to a 3.7% improvement in accuracy. These results confirm the effectiveness, efficiency, and robustness of our approach for low-resource domain adaptation of LLMs.

**Keywords:** Ontology · Large Language Models · Self-Evolution.

## 1 Introduction

In recent years, large language models (LLMs) have achieved remarkable performance across a broad spectrum of domains. This success is driven not only by pre-training on large and diverse datasets, but also by fine-tuning with carefully curated, synthesized, or annotated domain data. For example, domain-specific LLMs like BioBERT [9], SciBERT [1], and BloombergGPT [21] have demonstrated that scaling pre-training to massive, specialized corpora is crucial for superior downstream outcomes. More recent adaptation methods, such as Aloe [5],

---

\* Corresponding Author.

which leverages hundreds of thousands of medical QA pairs, and Med42-v2 [2], instruction-tuned with over a billion biomedical tokens, further underscore that large-scale domain data remains a key prerequisite for optimal model performance.

However, in data-sensitive fields such as healthcare and finance—where privacy requirements are stringent—acquiring large amounts of domain-specific data poses significant challenges. This raises the question: can a small but highly valuable set of supervision signals be leveraged to efficiently boost model performance?

We further observe that, alongside large-scale corpora, human experts have long distilled valuable knowledge rules and ontology rules that describe relationships between data and constrain their validity. Although few in number, these rules can substantially enhance the rationality and quality of knowledge management systems, making them critical for professional domain data governance and knowledge system optimization.

Currently, there is growing consensus in the community that LLMs can serve as implicit knowledge bases, with internalized human knowledge through pre-training stage[20]. Building on this insight, various technologies have been developed for knowledge editing and enhancement of LLMs[12,20]. This leads us to consider: can we leverage a limited set of domain-specific rules that have been accumulated over time to validate, revise, and improve the implicit knowledge within LLMs, thereby boosting their performance in low-resource domains?

Building on the established effectiveness of domain ontology knowledge in enhancing model performance[12,22], we propose a framework that exploits ontology rules to enhance LLM’s performance in specialized domains with extremely limited human supervision.

Specifically, our approach proceeds in three main steps. (1) We explicitly extract the implicit domain knowledge embedded within the LLM, particularly focusing on subclass and synonym relationships among domain concepts. (2) We introduce two ontology rules (as shown in Table 1) to externally detect inconsistencies in the extracted model’s knowledge. (3) Finally, we re-inject the revised knowledge back into the model via fine-tuning, thereby addressing missing ontology knowledge. In this manner, our method does not require access to large quantities of external domain data. Instead, we utilize only 2 high-quality ontology rules to validate and supplement the knowledge extracted from the model, and then reinject the improved knowledge to the model, resulting in enhanced domain abilities.

We conducted experiments within the medical domain to validate the effectiveness of our pipeline. Results on the raw model Llama3-8B-Instruct demonstrate consistent improvements across all medical QA benchmarks (MedMCQA [15], MedQA [7], PubMedQA [8]), with average improvements of 3.1% over raw model and 0.9% over best baseline, which is trained on supervised data. It is noteworthy that Med42-v2, which have already been fine-tuned on large domain-specific corpus, gained average 3.7% over its raw model across all datasets, and 1.1% over

best baseline which is based on large supervised corpus, showing the effectiveness and robustness of our framework. Our main contributions are as follows:

- We highlight the challenges involved in fine-tuning LLMs in data-scarce domains, and we are the first to utilize ontology rules to solve this data-scarce challenge.
- We propose a novel framework that includes explicit extraction of implicit ontology knowledge from language models, ontology rule-based examination and refinement of llm’s implicit knowledge, and re-injection of the refined knowledge back into the model by self-distilled fine-tuning.
- We validate the effectiveness and robustness of our framework on various medical benchmark datasets, as well as on two models Llama3-8B-Instruct and Med42-v2. With only two ontology rules, we achieve substantial improvements in model performance, outperforming post-training methods that rely on large-scale supervised data, such as TaxoLlama[14] and OntoTune[12]. Extensive evaluation, including generalization and ablation studies, further demonstrate the robustness of our framework.

## 2 Preliminary

### 2.1 Ontology

Ontology is a type of structured framework that captures concepts, their inter-connections, and rules within a specific domain which enables a shared understanding of a domain’s knowledge. It has been widely applied in the semantic web and knowledge management systems. Three core components in ontologies are: (1) Concepts: representing entities or categories within a domain. For example, in a medical ontology, concepts might include "Cell", "Symptom" and "Treatment". (2) Relationships: Relationships define how concepts are interconnected. The most common and important relationships in ontologies are Hyponymy (Is-subclass-of) and Synonymy (Is-synonym-of). Hyponymy represents a hierarchical, subclass relationship. For instance, "Muscle Cell" is a subclass of "Cell". Synonymy indicates that two concepts are semantically equivalent. For example, "Muscle Cell" and "Muscle Fiber" are synonyms. (3) Axioms (Rules): Ontologies are equipped with built-in rules which enable automated reasoning and consistency checking within knowledge graphs. For example: If (Concept A, Is-subclass-of, Concept B) and (Concept B Is-subclass-of Concept C), then it logically follows that (Concept A, Is-subclass-of, Concept C). However, if the ontology also includes (Concept A Is-Not-A-subclass-of, Concept C), this creates a conflict with the previously inferred relationship. Such rules allow ontologies to automatically detect and resolve inconsistencies, ensuring the integrity of the knowledge graph. This capability is particularly valuable in large-scale knowledge bases, where manual verification would be impractical.

**Table 1.** Ontology Rules Used.

ID	Premise	Conclusion
R1	$(x, \text{SynonymOf}, y) \wedge (y, \text{SubclassOf}, z)$	$\Rightarrow (x, \text{SubclassOf}, z)$
R2	$(x, \text{SubclassOf}, y) \wedge (y, \text{SubclassOf}, z)$	$\Rightarrow (x, \text{SubclassOf}, z)$

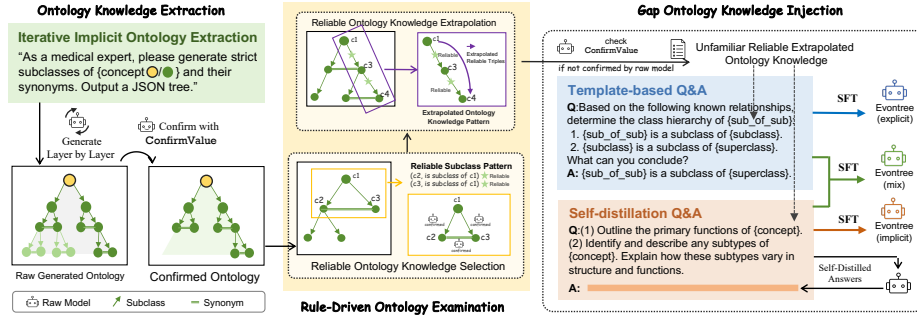
## 2.2 Perplexity

Model’s perplexity plays a critical role in our calculation of metric **ConfirmValue**. This metric quantifies the uncertainty of a probabilistic model in its predictions, where lower perplexity values correspond to higher prediction accuracy, while higher values indicate poorer performance. Formally, perplexity is defined as the exponential of the cross-entropy between the true distribution:

$$\text{Perplexity} = 2^{H(p,q)} \quad (1)$$

where  $H(p, q)$  is the cross-entropy between the true distribution between the true distribution  $p$  and the model’s predicted distribution  $q$ . In our framework, we leverage next-token prediction perplexity (e.g., [14], [10]) to design our metric **ConfirmValue**, which evaluates model’s confirmation towards certain ontology triple relationship, in order to prevent hallucination problems caused by model’s one-time generation.

## 3 Methodology



**Fig. 1.** The overview of Evontree.

### 3.1 Overview

As illustrated in Figure 1, our method comprises the following key steps. First, we employ predefined prompts to systematically elicit the original ontology knowledge from the raw model, layer by layer. To mitigate hallucination that may result from one-shot generation, we introduce a more robust metric **ConfirmValue**,

to quantify the model’s confidence in extracted ontology triples. Ultimately, for each root concept, we construct a tree-structured ontology—validated by the model itself—which is considered as its internal ontology knowledge. The second step is Rule-Driven Ontology Examination. The objective here is to leverage ontology rules to identify and correct inconsistencies within the model’s internal ontology knowledge. In our framework, the ontology knowledge targeted for correction comprises facts that the model is expected to know, but has not yet acquired. In practice, we first collect ontology knowledge confirmed by the model, then apply the ontology rules to extrapolate additional ontology facts. We subsequently use the metric **ConfirmValue** to assess whether the model has truly mastered these extrapolated facts. A critical aspect of this process is that, although the first step yields a considerable body of ontology knowledge from the model, some of this knowledge may be factually incorrect yet highly confirmed by the model. To avoid propagating such errors by using such erroneous knowledge serving as a basis for extrapolation, we design the Reliable Ontology Knowledge Selection module (see Section 3.3), which ensures that only factually reliable knowledge is used for extrapolating reliable new triples. Finally, newly extrapolated ontology triples deemed reliable are further evaluated based on their **ConfirmValue** scores. If the model’s **ConfirmValue** for these triples falls below a defined threshold, we designate them as gap triples. During the fine-tuning stage, we use these gap triples to synthesize training questions and generate self-distilled training data from the raw model. This process enables us to inject credible but unfamiliar ontology knowledge into the model, thereby enhancing its domain capability.

### 3.2 Ontology Knowledge Extraction

#### Generate Ontology Tree

As a medical expert, please generate strict subclasses of {concept} and their synonyms. Output a JSON tree like below:

```
{ "{c}": {
  "description": "",
  "subclasses": [{
    "name": "",
    "description": "",
    "synonyms": ["", ""]
  ]
}}
```

**Raw Triples Extraction:** Firstly, we use 15 root concepts to initiate the ontology extraction process, including Antibiotic, Bacterium, Cell, Enzyme, Fungus, Hormone, Tissue, Vertebrate, Virus, Vitamin, Chemical, Inorganic Chemical, Organic Chemical, Infectious Disease and Non-Infectious Disease. In order to ensure a sufficient volume of gap ontology knowledge for fine-tuning, we restrict these root concepts’ hierarchy depth  $\geq 3$ .

During the raw generation, for each concept  $c$  (initiated by root concept we manually designated, then iteratively replaced by generated last-layer concept), we feed the model the prompt as following.

After completing the iterative generation process, we obtain one ontology tree for each root concept, spanning at least three hierarchical layers. The initiated root concept is the apex of the ontology tree, and the underlying nodes are sub-concepts and synonym concepts generated layer-by-layer by the model itself. Edges of the tree denote the two most prevalent and fundamental ontology relations, subclass and synonym relationships.

**ConfirmValue Computation:** To reduce hallucination from one-time generation of models, we design the metrics **ConfirmValue**, a causal-perplexity(introduced in Section 2.2)-based confidence score, to evaluate model’s confirmation extent towards a single ontology triple. Taking raw triple  $(A, \text{SynonymOf}, B)$  as an example, we probe the model with the prompt as following.

Calculate ConfirmValue

Please determine if the statement is true or false, then answer with True or False: 'A' is an exact synonym of 'B'. Answer:

For each triple  $t$  under prompt  $p$ , we construct two completions:(1)  $S_{\text{True}} = p + \text{"True"}$ , (2)  $S_{\text{False}} = p + \text{"False"}$ , and compute the perplexity over the “Answer: True/False” token span. The **ConfirmValue** is defined as

$$\text{ConfirmValue}(t, p) = \frac{\text{sign}(\text{PPL}_{\text{False}} - \text{PPL}_{\text{True}})}{\min(\text{PPL}_{\text{True}}, \text{PPL}_{\text{False}})}. \quad (2)$$

To robustly gauge model’s confirmation towards certain ontology triple, we craft five paraphrased prompts for synonym relations and four for subclass relations. A triple is deemed model-confirmed only when every associated **ConfirmValue** surpasses a threshold  $\tau$ .

**Confirmation Threshold Setting:** The threshold  $\tau$  is obtained for each model and each prompt. For each model–prompt pair, we regard the one-time generated output (true / false) as a binary decision signal and **ConfirmValue** as its continuous confidence score. Using the  $\approx 10$  k raw subclass and synonym ontology triples from raw generation, we sweep  $\tau \in [0, 1]$  and compute the true-positive and false-positive rates. The optimal threshold  $\tau^*$  is the cut-off that maximizes the Youden index  $J = \text{TPR} - \text{FPR}$ , i.e. the point where the binary decision and continuous confidence are best aligned. If the model designates certain ontology triple with **ConfirmValue** higher than threshold  $\tau$ , then we deem that this ontology triple is "confirmed" by the raw model. If the model assigns a **ConfirmValue** exceeding the threshold  $\tau$  to a given ontology triple, we regard that triple as "confirmed" by the raw model.

### 3.3 Rule-Driven Ontology Examination

**Reliable Ontology Knowledge Selection:** While the preceding steps yield triples that are confirmed by the model, such confirmation alone does not en-

sure factual reliability. To guarantee that only robust knowledge is employed for downstream extrapolation, we must avoid arbitrarily selecting triples that may be erroneous. Therefore, we apply ontology rule R1 to identify closed triangles consisting of two subclass triples and one synonym triple, where the mutual structure allows each edge to corroborate the others. Any subclass triple that forms part of such a triangle is considered reliable. This rule leverages both ontological constraints and our prior belief that mislabelled instances within the model’s pre-training corpus are rare; therefore, the likelihood of three mutually supporting but erroneous statements co-occurring is extremely low. Consequently, only reliable subclass triples are retained for extrapolation in the subsequent step.

**Ontology Knowledge Extrapolation:** We take the reliable subclass triples  $\mathcal{T}_{\text{rel}}$  from last step and apply ontology rule R2 to generate new extrapolated subclass triples  $\mathcal{T}_{\text{extrapolated}}$ . For example, given reliable triples  $(D, \text{subClassOf}, C)$  and  $(C, \text{subClassOf}, A)$ , rule R1 yields  $(D, \text{subClassOf}, A)$ .

**Gap Ontology Triples Selection:** For every  $t_{\text{extrapolated}} \in \mathcal{T}_{\text{extrapolated}}$ , we recompute its **ConfirmValue** using the same prompt set as the raw ontology extraction step. Crucially, we retain only those triples whose **ConfirmValue** is below the threshold  $\tau^*$ , which means those triples the model are not familiar with before. These low-ConfirmValue triples are labelled  $\mathcal{T}_{\text{gap}}$ , which are supposed to be injected into the model during the subsequent model’s fine-tuning stage.

### 3.4 Gap Ontology Knowledge Injection

**Explicit Injection:** We introduce three injecting strategies for incorporating these credible but previously unfamiliar triples into the model. The simplest and most straightforward approach is to leverage our reliable and extrapolated triples to construct explicit reasoning chains, which are then used to synthesize question-answer pairs; the generation template is illustrated in the main figure 1.

**Implicit Injection:** To mitigate homogeneity in synthesized QA data from single-prompt reasoning chains, we follow prior work [12] to generate more natural and diverse training corpus. We append the ontology chain to pre-defined, concept-specific question templates in order to guide model to produce concept-aware higher-quality answers. We then finetune the model with the instruction and self-distilled output pairs back into the model, ensuring the ontology knowledge and concept-aware, self-distilled knowledge is reinforced to the model to continually enhance model’s domain capability.

The two pre-define concept-centric question templates are shown in below.

#### Generate Corpus

(1) Outline the primary functions of {concept}. (2) Identify and describe any subtypes of {concept}. Explain how these subtypes vary in structure and function.

For each gap triple  $(D, \text{SubclassOf}, A)$  and its derivation chain  $(D, \text{SubclassOf}, C), (C, \text{SubclassOf}, A)$ , we instantiate the placeholders with the three involved concepts  $A, C, D$  for template (1), and with  $A, C$  for template (2).

We hypothesise that appending the derivation chain as a *hint* after each template enables the model to produce higher-quality answers. The hint template are shown in following substantiated example. Given the weak triple (Skeletal Muscle Fiber, SubclassOf, Cell) supported by two reliable chains:(i) (Skeletal Muscle Fiber, SubclassOf, Muscle Cell), (Muscle Cell, SubclassOf, Cell). (ii) (Skeletal Muscle Fiber, SubclassOf, Myocyte), (Myocyte, SubclassOf, Cell). we generate the following instruction for chain (ii):

#### Hint with Ontology Chain

Outline the primary functions of Skeletal Muscle Fiber. You can consider these relationships as follows, but please ignore them if they are unnecessary: Skeletal Muscle Fiber is a subclass of Myocyte, and Myocyte is a subclass of Cell.

The same process is repeated for every supporting chain.

**Mixed Injection** As a further variant, we simultaneously inject knowledge via both implicit and explicit ways, which is termed mixed injection.

## 4 Experiments

We designed four research questions to verify the effectiveness of our method.

- **RQ1:** Can our method find high-quality ontology triples rightly?
- **RQ2:** Can our method effectively improve raw model’s domain capability without introducing any external corpora?
- **RQ3:** How does Evontree impact raw model’s generalization performance?
- **RQ4:** Does every key component of Evontree make a meaningful contribution?

### 4.1 Experiment Settings

**Datasets** Following prior work [12,5,2], we adopt 3 widely-used medical-domain benchmark datasets, including PubMedQA [8], MedQA [7] and MedMCQA [15] for evaluation. We exclude USMLE-Step1/2/3 because evaluation results can easily fluctuate due to their limited size.

**Baselines** We adopt the same baseline as prior ontology-based work[12]. Considering computational limits, we evaluate our method on this mainstream architecture Llama-3-8B-Instruct. To test whether the approach can enhance an already strong model, we additionally include Med42-v2—a Llama-3-8B variant that has been extensively fine-tuned for medical tasks and is widely used in the healthcare community.



**Implementation** We employ the Low Rank Adaptation(LoRA) technique to fine-tune raw model based on the LLaMA-Factory [25] framework. During the SFT stage, we use fp32 and a learning rate of 5e-5, training for 3 epochs with a cosine scheduler, a batch size per device initialized to 8 and gradient accumulation of 2.

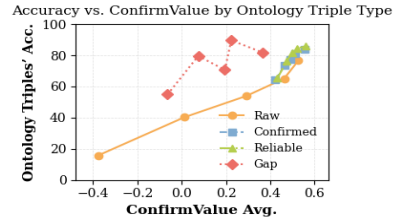
**Evaluation** We evaluate the effectiveness of our approach from two perspectives. On the one hand, we employ GPT-4o-mini and DeepSeek-V3 as strong supervisory models to assess the quality of the ontology triples generated at each stage of our pipeline. On the other hand, we measure our model’s performance on both medical-domain QA datasets and general capability benchmarks, following the evaluation protocols established in prior work [12].

## 4.2 Characteristics of Ontology Triples (RQ1)

As shown in Table 2 and Figure 2, the four ontology-triple types exhibit markedly different characteristics. (1) As for the raw ontology triples, which are extracted directly from the raw model, achieve the lowest accuracy across all categories. Nevertheless, the positive correlation between accuracy and **ConfirmValue** demonstrates that the original model already possesses a rather strong intrinsic alignment with factual correctness. (2) As for the confirmed triples, they are obtained by imposing a **ConfirmValue** threshold; consequently, their **ConfirmValues** begin at a high baseline. (3) Reliable ontology triples are further distilled from the confirmed set. Figure 2 shows that they span the entire **ConfirmValue** spectrum while achieving a consistently higher average accuracy than their confirmed counterparts. This confirms that our ontology-rule-based automatic labelling strategy can improve correctness with virtually no external supervision. (4) Gap triples, by design, are bounded above by the lower **ConfirmValue** limit of confirmed triples. Unlike the other three categories, gap triples exhibit no discernible correlation between **ConfirmValue** and accuracy; for a given accuracy level, their **ConfirmValue** is markedly lower than that of confirmed or reliable triples. This pattern aligns perfectly with our hypothesis: we have successfully isolated ontology knowledge that the model has not yet internalised, yet which is nevertheless highly reliable.

Triple Type	Relation	Num	ConfirmValue Avg.	Acc.
Raw	synonym	95,814	0.0537	54.90
Confirmed	synonym	30,532	0.4738	87.51
Raw	subclass	58,309	0.2094	47.79
Confirmed	subclass	24,276	0.4896	72.62
Reliable	subclass	14,169	0.4981	74.81
Extrapolated	subclass	8,661	0.4449	80.29
Gap	subclass	1,396	0.1595	74.57

**Table 2.** Ontology path statistics and quality evaluation.



**Fig. 2.** Relationship between accuracy and **ConfirmValue**.

### 4.3 Evaluation on Medical Datasets (RQ2)

As shown in Table 3, among the three Evontree variants, the “explicit” ontology injection approach does not effectively improve model performance, whereas the more natural “implicit” and “mix” variants yield significant enhancements. This trend is further validated when comparing TaxoLLaMA (which uses explicit ontology injection) and OntoTune (which adopts implicit injection). Notably, Evontree not only achieves substantial performance gains for general-purpose LLMs (e.g., LLaMA3-8B-Instruct), but also further boosts domain-specific models like Med42-v2, reaching improvements of 3.1% and 3.7%, respectively. Compared with existing LLaMA3-8B-based domain models, Evontree attains state-of-the-art results even without introducing additional external supervision data, outperforming baselines by an average of 1.1%. This indicates that large-scale, fragmented corpora alone cannot reliably achieve effective domain alignment, and may even introduce conflicting knowledge that increases model confusion. We conclude that the advantage of Evontree lies in its targeted identification of effective gap triples within the model’s knowledge system. Through fine-grained automated selection, Evontree ensures that only information highly compatible with the model’s existing knowledge structure is injected. This mechanism enables the model to self-reflect and precisely locate knowledge gaps, evolving itself based on reliable information while effectively avoiding unnecessary noise and conflicts. Compared to traditional data augmentation and domain adaptation approaches, Evontree prioritizes quality over quantity in knowledge supplementation, promoting greater consistency and completeness in the model’s knowledge base.

### 4.4 General Capabilities and Safety Evaluation (RQ3)

**General Capabilities Evaluation.** We evaluate the general capabilities of our model variants compared to the raw model on the MMLU, TriviaQA, and ARC datasets. Specifically, MMLU is assessed using the LLaMAFactory framework, while ARC and TriviaQA are evaluated using the OpenCompass tool in gen mode. As shown in Table 4, under the zero-shot setting, all three variants of our method exhibit no significant degradation in general capabilities. The best-performing variant shows an average performance drop of only 0.15% compared to the raw model, and LLaMA3-Evontree(mix) even achieves marginal improvements on certain test sets. Furthermore, under the supervised fine-tuning (SFT) evaluation setting, our models continue to demonstrate similar trends. Compared to med42-llama3-8b, our variants maintain stable performance and in some cases achieve notable gains, highlighting the robustness and potential benefits of our approach.

**Safety Evaluation.** Following prior work, we also assess whether our method introduces any safety risks to the model. We use harmful instructions from the AdvBench dataset to measure the proportion of safe responses, reported as the “Raw Safe” metric. Then, we append an inducing suffix to the harmful prompts to encourage unsafe behavior and measure the proportion of safe outputs under

**Table 3.** Results of the medical domain QA in the zero-shot and supervised fine-tuning (on evaluation) setting. The best results are highlighted in bold, while the second best are underlined. Results are taken from OntoTune [12], and we evaluate our method in the same way.

Setting	Model	MedQA	MedMCQA	PubMedQA	Average
zero-shot	LLaMA3 8B [4]	51.7	51.7	70.3	57.9
	TaxoLLaMA [14]	50.5	46.1	73.4	56.7
	OntoTune <sub>st</sub> [12]	51.5	56.7	72.0	<u>60.1</u>
	OntoTune <sub>dpo</sub> [12]	<b>53.3</b>	<u>57.2</u>	65.5	58.7
	OntoTune <sub>st+dpo</sub> [12]	51.9	56.7	66.3	58.3
	LLaMA3 8B-Evontree (explicit)	43.4	48.6	<b>76.4</b>	56.1
	LLaMA3 8B-Evontree (implicit)	<u>52.7</u>	54.4	72.9	60.0
	LLaMA3 8B-Evontree (mix)	51.0	<b>57.4</b>	<u>74.7</u>	<b>61.0</b>
	$\Delta$ Improvement over raw model	+1.0%	+5.7%	+6.1%	+3.1%
	$\Delta$ Improvement over best baseline	-0.6%	+0.2%	+3.0%	+0.9%
SFT (on evaluation)	LLaMA3 8B [4]	56.4	53.9	77.2	62.5
	Aloe [5]	51.1	56.8	75.4	61.1
	Med42-v2 [2]	57.8	58.1	74.6	63.5
	TaxoLLaMA [14]	55.9	57.5	77.6	63.7
	OntoTune <sub>st</sub> [12]	<u>58.4</u>	60.4	78.6	65.8
	OntoTune <sub>dpo</sub> [12]	58.3	<u>60.7</u>	<b>79.4</b>	<u>66.1</u>
	OntoTune <sub>st+dpo</sub> [12]	58.2	60.5	<u>78.9</u>	65.9
	Med42-v2-Evontree (explicit)	57.1	58.6	71.1	62.3
	Med42-v2-Evontree (implicit)	<b>60.3</b>	<b>62.4</b>	<u>78.9</u>	<b>67.2</b>
	Med42-v2-Evontree (mix)	57.2	57.7	74.9	63.3
	$\Delta$ Improvement over raw model	+2.5%	+4.3%	+4.3%	+3.7%
	$\Delta$ Improvement over best baseline	+1.9%	+1.7%	-0.5%	+1.1%

this adversarial setting, reported as the "Jailbreak Safe" metric. As shown in Table 4, the Evontree(implicit) and Evontree(mix) variants show no significant drop in safety compared to the raw model, and in fact, when built upon Med42, these two variants even lead to notable improvements. In contrast, the Evontree(explicit) variant performs worse in both the LLaMA3 and Med42 settings. We hypothesize that, similar to the observations in TaxoLLaMA, the introduction of rigid template-based Q&A formats may negatively impact the model's safety performance.

#### 4.5 Ablation Study (RQ4)

Variant	gpt deepseek Avg.		
w/o reliable triple sel.	58.4	53.0	55.7
Evontree	75.3	73.8	74.6

**Table 5.** Impact of Reliable-Triple Selection on Extrapolated Triple Accuracy.

Variant	MedQA	MedMCQA	PubMedQA	Avg.
w/o reliable triple sel.	<u>52.4</u>	<b>54.7</b>	72.3	<u>59.8</u>
w/o gap triple sel.	52.6	54.1	<u>72.5</u>	59.7
w/o ontology injecting	51.1	53.9	<b>73.1</b>	59.4
LLaMA3-Evontree(implicit)	<b>52.7</b>	<u>54.4</u>	<u>72.9</u>	<b>60.0</b>

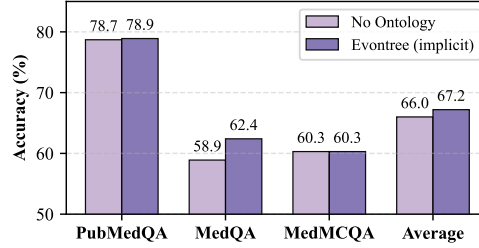
**Table 6.** Ablation Study of Different Modules.

To provide a comprehensive evaluation of our method, Evontree, we conducted an ablation study in Table 6 to validate the contributions of key components of Evontree. Specifically, we focused on assessing the impact of the reliable

**Table 4.** Results of general capabilities and safety evaluation. The results with the least decrease compared to the original model are highlighted in bold, while the second least are underlined.

Model	MMLU					ARC			TriviaQA	Advbench		
	STEM	Social Sciences	Humanities	Other	Average	ARC	C	ARC	E	-	Raw	Safe Jailbreak Safe
LLaMA3 8B [4]	56.83	76.61	60.81	74.10	66.49	78.64	92.77	64.81	97.50	96.35		
TaxoLLaMA [14]	55.96	73.74	56.92	69.43	63.29	72.88	89.24	63.12	94.04	73.27		
OntoTune <sub>ft</sub> [12]	56.47	75.73	<b>61.85</b>	<u>73.02</u>	<b>66.31</b>	78.31	<u>91.89</u>	<b>64.07</b>	94.04	<u>92.69</u>		
OntoTune <sub>po</sub> [12]	56.33	75.33	59.93	<b>73.64</b>	65.70	78.98	<b>92.06</b>	<u>63.96</u>	90.58	77.88		
OntoTune <sub>ft+ppo</sub> [12]	55.67	75.17	<u>61.79</u>	72.71	65.93	78.98	<b>92.06</b>	<u>63.96</u>	90.58	84.81		
LLaMA3-Evontree (implicit)	<u>56.89</u>	76.41	61.38	71.99	66.16	79.66	89.59	63.00	<u>95.38</u>	<b>96.15</b>		
LLaMA3-Evontree (mix)	<b>56.99</b>	<b>76.80</b>	61.62	71.71	<u>66.28</u>	<u>80.34</u>	90.30	63.22	95.19	<b>96.15</b>		
LLaMA3-Evontree (explicit)	<u>56.89</u>	<u>76.60</u>	60.85	72.24	66.08	<b>82.03</b>	91.71	63.42	<b>97.50</b>	83.65		
Aloe [5]	55.67	<u>76.24</u>	58.91	72.25	65.10	75.25	86.95	63.03	62.50	34.23		
Med42-v2 [2]	56.59	<u>76.24</u>	59.91	72.67	65.72	<b>82.37</b>	<b>92.59</b>	<b>65.19</b>	83.85	60.19		
Med42-v2-Evontree (implicit)	<b>56.86</b>	<b>76.37</b>	59.79	72.79	<u>65.80</u>	80.34	91.36	63.54	<u>86.54</u>	60.77		
Med42-v2-Evontree (mix)	56.46	<u>76.24</u>	<b>60.09</b>	<u>72.86</u>	<u>65.80</u>	<u>80.68</u>	<u>92.06</u>	63.71	<b>89.23</b>	<b>65.77</b>		
Med42-v2-Evontree (explicit)	<u>56.79</u>	76.11	<u>59.98</u>	<b>72.92</b>	<b>65.82</b>	<u>80.68</u>	91.18	<u>64.33</u>	75.38	45.19		

**Performance Comparison: Ontology Injection Effectiveness**



**Fig. 3.** Ablation Study of Med42-v2.

ontology knowledge selection, the gap triple selection and ontology injecting. By isolating these components respectively, we aim to understand their individual and collective contributions to the overall performance of Evontree. We designed three ablation variants of Evontree, each purposefully modified to test the significance of specific components:

**Ablation on Reliable Ontology Knowledge Selection.** The module of Reliable Ontology Knowledge Selection is designed to make sure that we extrapolate new triples based on ontology triples that are more reliable in factuality. To eliminate this module, we extrapolate ontology triples directly from model-confirmed ontology triples. We randomly sampled 10,000 triples from the full set of triples extrapolated from confirmed triples. To control for confounding variables, we employ a sliding-window procedure to draw a subset that matches the gap triples in both size and mean ConfirmValue. This design allows us to examine (i) whether removing the reliable-triple selection step reduces the accuracy of the extrapolated triples, and (ii) whether such drop in injected triples’ accuracy will subsequently degrade model’s downstream task performance. for question (i), we employ both gpt4o-mini and deepseek-v3 to check the accuracy of ontology triples, as shown in Table 5, although their ConfirmValue average have been controlled the same, the accuracy of triples are different largely, with the ablation variant lagging 18.9% from gap triples. This observation illustrates that the reliable triple selection module can effectively discriminate reliable triples

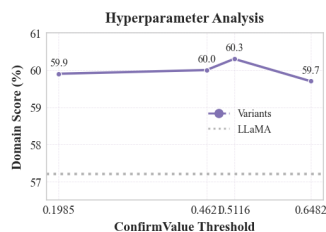
from hallucinated ones. for question (i), we evaluate ablated variant on medical benchmarks, the results are shown in Table 6. We can observe that the average score declined when eliminated the module of reliable triple selection, which illustrate the key role of reliable triples’ annotation.

**Ablation on Gap Triple Selection.** Gap-Triple Selection serves as a precision gate: only those extrapolated ontology triples whose `ConfirmValue` falls below a threshold are injected. To ablate this module, we skip the `ConfirmValue` filter and inject all extrapolated triples. As shown in Table 6, the ablated model declines on every benchmark, confirming that targeted injection is essential. We hypothesize that indiscriminately injecting all extrapolated triples perturbs the training distribution and triggers catastrophic forgetting, whereas restricting injection to the model’s blind spots (i.e., low-`ConfirmValue` gap triples) preserves performance while still expanding its factual coverage. We show further analysis of gapping threshold in 4.6.

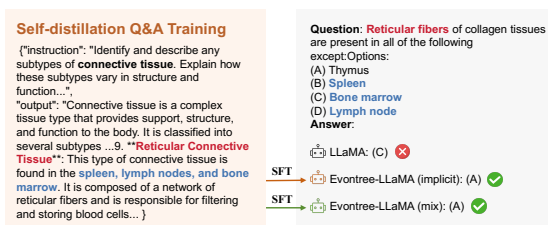
**Ablation on Ontology Injecting.** Instead of inject ontology knowledge into models by implicit or explicit way, we study that if model can improve themselves by simply based on the training questions. To eliminate this module, we delete ontology hint in the last of the question sentence, and only provide self-distilled model with the original templated questions. As a result ( Shown in Table 6), ablated model decline the most among all the ablated models, which indicate that ontology triples are really necessary for models to generate higher-quality answers. The same ablation implemented to Med42-v2 illustrate the same decline (Shown in Figure 3), which emphasizes the critical value of high-quality injected ontology knowledge.

## 4.6 Hyperparameter Sensitivity Analysis.

After extrapolating new ontology triples from the reliable triples, a natural question arises: should we use all of them, or only a subset? In our main model, we use the original confirmation threshold—used to cut off gap triples. However, to assess the sensitivity of this choice, we conduct a hyperparameter sensitivity experiments in this section, sweeping over different gap triple confirmation thresholds; the results are reported in Table 4. From the figure we can observed that with the threshold smaller than our chosen value (0.4621), domain performance degrades. While moving to the right, performance first improves and then declines. This yields two key observations: (1) Although the optimal threshold remains uncertain, every variant models trained from randomly sampled threshold consistently outperforms the raw model, demonstrating the overall robustness of our approach to hyper-parameter variation. (2) Based on the prior experiments, we can hypothesize that three factors may influence variant model’s performance: the factual accuracy of the selected ontology triples, the model’s prior familiarity (`ConfirmValue`) and the size of the whole injecting triples. However, owing to resource and time constraints, we will further investigate into it fully in future work.



**Fig. 4.** Hyperparameter Analysis.



**Fig. 5.** Case Study of Evontree. The example is selected from MedMCQA dataset.

## 4.7 Case Study

We also present an intuitive example illustrating how our pipeline enhances performance on domain-specific downstream tasks. As shown in Figure 5, the downstream multiple-choice question asks: “Reticular fibers of collagen tissues are present in all of the following except: (A) Thymus (B) Spleen (C) Bone marrow (D) Lymph node.” Since the model has been trained on a self-distilled corpus derived from carefully selected ontology triples, it has internalized that spleen, lymph nodes, and bone marrow are canonical sites of reticular connective tissue, while the thymus mainly consists of epithelial reticular cells rather than collagenous reticular fibers. The overlap between the training QA pairs and the required knowledge for answering the exam question (i.e., identifying organs with reticular fibers and distinguishing the thymus as the exception) highlights the effectiveness of our ontology rule-guided approach in supplementing crucial domain knowledge. This targeted alignment between ontology-based knowledge integration and downstream evaluation demonstrates the significant improvement of model capabilities in specialized domains.

## 5 Related Work

### 5.1 Domain-specific Large Language Models

Recent advances in large language models (LLMs) have demonstrated remarkable capabilities across general domains, yet their application to specialized fields often requires domain-specific adaptation [23]. Among the approaches [18,11] that have emerged to enhance LLMs for technical domains, early efforts to inject domain knowledge into language models relied on pre-training from scratch on specialized corpora. For instance, BioBERT [9] and SciBERT [1] are pre-trained on PubMed and scientific articles, respectively. While BloombergGPT [21] scaled the approach to tens of billions of parameters, confirming that domain corpus size and quality remain critical drivers of downstream performance. Rather than training from scratch, more recent work explores lightweight adaptation of general LLMs. Parameter-efficient fine-tuning (PEFT) methods [24], including LoRA [6] and adapters [16], successfully specialize open-source LLMs like LLaMA [4] for medical dialogue or financial reasoning with minimal compute.

These methods always rely on large-scale domain-specific datasets for optimization [5,2]. For example, Aloe [5] uses 348K medical QA pairs from 20+ sources, 430K synthetic medical QA pairs and 122K high-quality general-domain samples for fine-tuning. Med42-v2 [2] was instruction-tuned on 1 B tokens drawn from 58 open biomedical corpora, including MedQA, USMLE- and MMLU-style exams, PubMed abstracts, clinical dialogues, case reports, and medical textbooks.

In contrast to these methods, our approach does not rely on any external data resources. Instead, we extract the model’s inherent ontology knowledge and apply two fundamental ontology rules to systematically refine and reinforce it. This paradigm achieves strong domain adaptation in LLMs with maximal efficiency and zero reliance on external annotated corpora.

## 5.2 Ontology-Enhanced LLMs

Recent research has sought to improve the domain capabilities of large language models (LLMs) by integrating externally supervised ontology knowledge. These methods typically incorporate high-quality, curated ontological data to improve LLMs’ ability to represent and reason about domain concepts. For example, structured knowledge from ontologies [22], which are systematically developed within specific fields, has been integrated into LLM training to provide rich hierarchical organization and semantic constraints. Ronzano and Nanavati [17] apply a contrastive learning framework to embed medical ontology information into LLMs, thereby enhancing biomedical concept representation. Moskvoretskii et al. employ WordNet [13] to fine-tune LLaMA [19] for instruction alignment, resulting in TaxoLLaMA [14]. Similarly, Liu et al. [12] propose OntoTune, which aligns LLMs with the SNOMED CT [3] medical ontology using in-context learning and external supervision.

In contrast, our method refrains from directly incorporating external ontology knowledge, relying exclusively on two constraining ontology rules. We first extract the LLM’s inherent ontology knowledge, then automatically detect and rectify inconsistencies using these rules, and finally reintegrate the refined ontology knowledge into the model. This entire process is accomplished without reliance on external ontology databases.

## 6 Conclusion

In this work, we address the challenge of adapting large language models to data-scarce domains by leveraging the unique value of domain ontology rules. Our framework enables explicit extraction and validation of LLMs’ implicit knowledge, using only two carefully chosen ontology rules to systematically detect and correct inconsistencies. By reinjecting revised knowledge into the model, we substantially improve performance on medical QA tasks, outperforming baselines dependent on large supervised datasets. Extensive experiments across two representative models and several medical benchmarks validate the robustness and effectiveness of our paradigm. This study highlights the potential of ontology

rule-driven supervision as a practical and powerful solution for enhancing LLMs in highly specialized domains where conventional data-centric approaches are limited by privacy or scarcity. Future work can explore broader application of our framework to other professional domains and further enrichment of ontology rule-guided knowledge editing techniques.

## References

1. Beltagy, I., Lo, K., Cohan, A.: Scibert: A pretrained language model for scientific text. In: EMNLP/IJCNLP (1). pp. 3613–3618. Association for Computational Linguistics (2019)
2. Christophe, C., Kanithi, P.K., Raha, T., Khan, S., Pimentel, M.A.: Med42-v2: A suite of clinical llms. arXiv preprint arXiv:2408.06142 (2024)
3. Donnelly, K.: Snomed-ct: The advanced terminology and coding system for ehealth. *Studies in Health Technology & Informatics* **121**, 279 (2006)
4. Dubey, A., Jauhri, A., Pandey, A., et al.: The llama 3 herd of models. CoRR **abs/2407.21783** (2024)
5. Gururajan, A.K., Lopez-Cuena, E., Bayarri-Planas, J., et al.: Aloe: A family of fine-tuned open healthcare llms. arXiv preprint arXiv:2405.01886 (2024)
6. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. In: ICLR. OpenReview.net (2022)
7. Jin, D., Pan, E., Oufattole, N., Weng, W.H., Fang, H., Szolovits, P.: What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences* p. 6421 (Jul 2021). <https://doi.org/10.3390/app11146421>, <http://dx.doi.org/10.3390/app11146421>
8. Jin, Q., Dhingra, B., Liu, Z., Cohen, W., Lu, X.: Pubmedqa: A dataset for biomedical research question answering. Cornell University - arXiv, Cornell University - arXiv (Sep 2019)
9. Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H., Kang, J.: Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinform.* **36**(4), 1234–1240 (2020)
10. Li, M., Zhang, Y., Li, Z., Chen, J., Chen, L., Cheng, N., Wang, J., Zhou, T., Xiao, J.: From quantity to quality: Boosting LLM performance with self-guided data selection for instruction tuning. In: NAACL-HLT. pp. 7602–7635. Association for Computational Linguistics (2024)
11. Liu, K., Chen, Z., Fu, Z., Zhang, W., Jiang, R., Zhou, F., Chen, Y., Wu, Y., Ye, J.: Structure-aware domain knowledge injection for large language models. In: ACL (1). pp. 29443–29464. Association for Computational Linguistics (2025)
12. Liu, Z., Gan, C., Wang, J., et al.: Ontotune: Ontology-driven self-training for aligning large language models. In: Proceedings of the ACM on Web Conference 2025. pp. 119–133 (2025)
13. Miller, G.A.: Wordnet: a lexical database for english. *Communications of the ACM* **38**(11), 39–41 (1995)
14. Moskvoretskii, V., Neminova, E., Lobanova, A., Panchenko, A., Nikishina, I.: Taxollama: Wordnet-based model for solving multiple lexical semantic tasks. CoRR **abs/2403.09207** (2024)
15. Pal, A., Umapathi, L.K., Sankarasubbu, M.: Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In: Conference on health, inference, and learning. pp. 248–260. PMLR (2022)



16. Pfeiffer, J., Rücklé, A., Poth, C., Kamath, A., Vulic, I., Ruder, S., Cho, K., Gurevych, I.: Adapterhub: A framework for adapting transformers. In: EMNLP (Demos). pp. 46–54. Association for Computational Linguistics (2020)
17. Ronzano, F., Nanavati, J.: Towards ontology-enhanced representation learning for large language models (2024)
18. Song, Z., Yan, B., Liu, Y., Fang, M., Li, M., Yan, R., Chen, X.: Injecting domain-specific knowledge into large language models: A comprehensive survey. CoRR **abs/2502.10708** (2025)
19. Touvron, H., Lavril, T., Izacard, G., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)
20. Wang, M., Yao, Y., Xu, Z., Qiao, S., Deng, S., Wang, P., Chen, X., Gu, J., Jiang, Y., Xie, P., Huang, F., Chen, H., Zhang, N.: Knowledge mechanisms in large language models: A survey and perspective. In: EMNLP (Findings). pp. 7097–7135. Association for Computational Linguistics (2024)
21. Wu, S., Irsoy, O., Lu, S., Dabrovolski, V., Dredze, M., Gehrmann, S., Kambadur, P., Rosenberg, D.S., Mann, G.: Bloomberggpt: A large language model for finance. CoRR **abs/2303.17564** (2023)
22. Xiao, G., Calvanese, D., Kontchakov, R., Lembo, D., Poggi, A., Rosati, R., Zakharyashev, M.: Ontology-based data access: A survey. In: Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence (Jul 2018). <https://doi.org/10.24963/ijcai.2018/777>, <http://dx.doi.org/10.24963/ijcai.2018/777>
23. Xie, Y., Aggarwal, K., Ahmad, A.: Efficient continual pre-training for building domain specific large language models. In: ACL (Findings). pp. 10184–10201. Association for Computational Linguistics (2024)
24. Xu, L., Xie, H., Qin, S.J., Tao, X., Wang, F.L.: Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. CoRR **abs/2312.12148** (2023)
25. Zheng, Y., Zhang, R., Zhang, J., Ye, Y., Luo, Z., Feng, Z., Ma, Y.: Llamafactory: Unified efficient fine-tuning of 100+ language models. arXiv preprint arXiv:2403.13372 (2024)