

RAG+: Enhancing Retrieval-Augmented Generation with Application-Aware Reasoning

Yu Wang^{1,2}, Shiwan Zhao³, Ming Fan², Zhihu Wang¹, Yubo Zhang¹,
Xicheng Zhang², Zhengfan Wang², Heyuan Huang¹, Ting Liu²

¹Huawei Technologies Ltd. ²Xi'an Jiaotong University ³Nankai University

Abstract

The integration of external knowledge through Retrieval-Augmented Generation (RAG) has become foundational in enhancing large language models (LLMs) for knowledge-intensive tasks. However, existing RAG paradigms often overlook the cognitive step of applying knowledge, leaving a gap between retrieved facts and task-specific reasoning. In this work, we introduce **RAG+**, a principled and modular extension that explicitly incorporates application-aware reasoning into the RAG pipeline. RAG+ constructs a dual corpus consisting of knowledge and aligned application examples, created either manually or automatically, and retrieves both jointly during inference. This design enables LLMs not only to access relevant information but also to apply it within structured, goal-oriented reasoning processes. Experiments across mathematical, legal, and medical domains, conducted on multiple models, demonstrate that RAG+ consistently outperforms standard RAG variants, achieving average improvements of 3–5%, and peak gains up to 7.5% in complex scenarios. By bridging retrieval with actionable application, RAG+ advances a more cognitively grounded framework for knowledge integration, representing a step toward more interpretable and capable LLMs.

1 Introduction

Large language models (LLMs) have demonstrated strong performance across a broad range of natural language processing tasks (Rong et al., 2025; Tang et al., 2025). To further enhance their capabilities, Retrieval-Augmented Generation (RAG) has become a widely adopted framework. By equipping LLMs with access to external knowledge sources, RAG enables the dynamic retrieval of up-to-date information at inference time, significantly improving performance in knowledge-intensive scenarios (Li et al., 2025; Mostafa et al., 2025).

However, existing RAG methods often focus on lexical or semantic similarity when retrieving

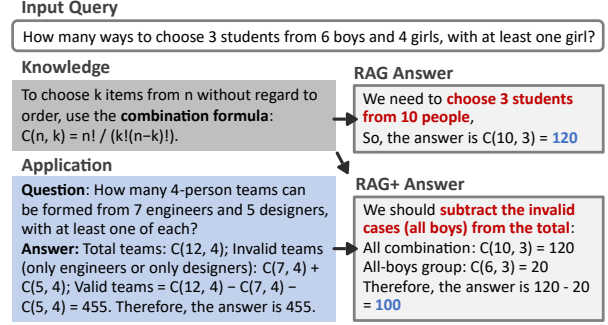


Figure 1: An Illustrative Case of RAG and RAG+: Knowledge Only vs. Knowledge with Application.

knowledge, paying little attention to how the retrieved content should be applied in downstream tasks. While effective in factual recall and open-domain question answering, RAG frequently underperforms on domain-specific reasoning tasks (Lin et al., 2025; Hayashi et al., 2025; Ammann et al., 2025), where solving complex problems requires not only relevant information but also reasoning about how to use it to arrive at a solution, as shown in Figure 1.

RA-DIT (Lin et al., 2023) fine-tunes the retriever and generator in a dual-instruction manner, aligning retrieval more closely with what the model needs to generate accurate responses. Other recent RAG extensions tackle this issue by decomposing reasoning tasks into smaller steps and retrieving relevant knowledge for each substep (Singh et al., 2025; Xiong et al., 2025; Zihao et al., 2024). However, they provide limited guidance on applying retrieved knowledge, which hampers performance in procedural reasoning tasks that require understanding both the process and the underlying concepts.

This limitation reflects insights from educational psychology. Bloom’s Taxonomy identifies “applying” knowledge as a distinct cognitive skill that goes beyond simple recall (Bloom, 2010). Similarly, cognitive architectures like ACT-R distinguish between declarative memory (facts) and pro-

cedural memory (skills), suggesting that coupling factual knowledge with procedural examples enhances performance on complex tasks (Anderson et al., 2004). Consistently, Re-TASK (Wang et al., 2024b) introduces the concept of capability items, emphasizing that accomplishing domain-specific tasks requires jointly leveraging domain knowledge and task-specific skills.

Motivated by these insights, we propose RAG+, a simple yet effective extension to the RAG framework that enhances reasoning by bridging retrieval and generation with an application-aware stage. Instead of retrieving only relevant knowledge, RAG+ additionally retrieves examples that demonstrate how the knowledge is applied in practice—such as structured reasoning chains and stepwise solutions—to ground the model’s output in task-relevant usage and to improve reasoning accuracy.

RAG+ builds a dual corpus: one stores domain knowledge, and the other containing automatically generated application instances aligned to each fact. At inference time, the system retrieves relevant knowledge based on the input query and then fetches aligned application examples to provide practical context. This design encourages the model not only to recall facts but also to produce outputs that follow grounded reasoning patterns based on prior usage. RAG+ is modular and retrieval-agnostic; it can be integrated into any existing RAG pipeline without changing the model architecture or requiring additional fine-tuning.

We evaluate RAG+ across three reasoning-intensive domains: mathematics, medicine, and legal, using four representative RAG variants: vanilla RAG, Answer-First RAG, GraphRAG, and Rerank RAG. Experiments show that RAG+ consistently improves performance. Notably, Qwen2.5-72B improves from 76.5% to 87.5% on legal prediction, and LLaMA3.3-70B rises from 78.2% to 86.5% on medical QA. Even smaller models like DS-Qwen-7B benefit, demonstrating the broad effectiveness of application-aware augmentation. These results demonstrate that bridging retrieval and reasoning through application-aware augmentation can yield substantial gains, especially in domains where accurate reasoning is essential.

Our contributions are as follows:

- We identify a key limitation of existing RAG systems: the lack of an application-aware step that impeding reasoning in complex tasks.
- We introduce RAG+, a modular and plug-and-

play extension that jointly retrieves factual knowledge and its aligned usage examples to better support downstream reasoning.

- We validate the effectiveness of RAG+ through comprehensive experiments, demonstrating consistent gains across multiple domains and retrieval strategies.

2 Related Work

Retrieval-Augmented Generation (RAG) has become a widely adopted framework for enhancing large language models (LLMs), especially in knowledge-intensive tasks. A standard RAG pipeline includes query formulation, corpus construction, retrieval, and answer generation. Prior work has sought improvements at each stage.

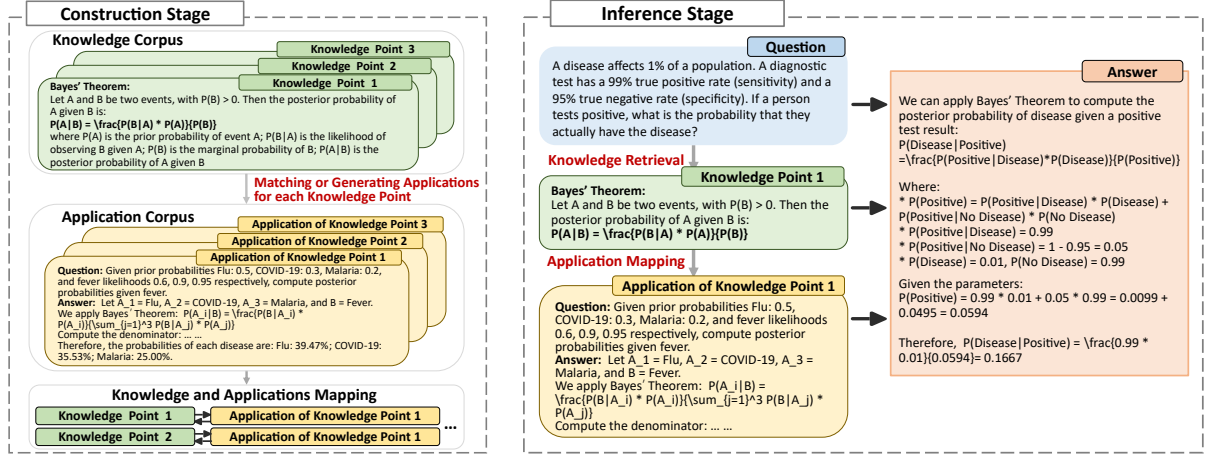
For retrieval quality, Rewrite-Retrieve-Read (Ma et al., 2023) rewrites queries to better match document style. R^2AG (Ye et al., 2024) uses a retrieval-aware encoder to highlight key signals, while Query2Doc (Wang et al., 2023a) expands queries into pseudo-documents to clarify intent. On the corpus side, GraphRAG (Edge et al., 2024) integrates knowledge graphs to support entity-centric reasoning. To enhance retrieval precision, Reranking (Wang et al., 2024a) and filtering techniques (Wang et al., 2023b; Pickett et al., 2025) further refine retrieved content.

Beyond these enhancements, recent work has tackled more complex reasoning tasks through modular and agent-based RAG variants. Agentic RAG (Singh et al., 2025) and RAG-Gym (Xiong et al., 2025) decompose tasks into subtasks handled by specialized agents or workflows. OPEN-RAG (Islam et al., 2024) introduces agent-based decomposition and context selection, while RAT (Zihao et al., 2024) uses Chain-of-Thought prompting to support step-wise retrieval.

Despite these efforts, most RAG methods remain optimized for fact-centric tasks such as open-domain QA. In reasoning-centric domains such as mathematics, retrieving the right facts is only the first step. The model must also understand how to apply them to reach a specific goal. Our work addresses this gap by introducing an application-aware step that explicitly guides how retrieved knowledge is used.

3 Methods

We propose RAG+, a principled extension of Retrieval-Augmented Generation (RAG) that incor-



(a) Construction Stage: Building an Application Corpus Aligned with the Knowledge Corpus. (b) Inference Stage: Retrieving Aligned Knowledge and Application Examples to Generate Output.

Figure 2: The overview framework of RAG+: (a) construction of the application corpus and (b) inference with retrieved knowledge and applications.

porates the explicit application of retrieved knowledge to improve reasoning. While prior RAG frameworks focus on retrieving relevant knowledge for downstream tasks, they often overlook explicitly guiding models on how to utilize this knowledge in reasoning. RAG+ addresses this limitation by introducing an explicit step that applies the retrieved knowledge through aligned application examples, illustrating its practical use.

The RAG+ pipeline consists of two stages: (a) a construction stage, where an application corpus is built and aligned with the knowledge corpus; and (b) an inference stage, in which both knowledge and corresponding applications are retrieved to form a comprehensive prompt for response generation, as shown in Figure 2.

3.1 Construction Stage

The construction stage aims to build an application corpus A aligned with an existing knowledge corpus K , as shown in Figure 2(a). For each knowledge item $k \in K$, an application example $a \in A$ is either retrieved or generated to demonstrate the practical use of k . These examples bridge the gap between passive knowledge access and task-oriented reasoning.

Depending on domain characteristics and data availability, we consider two complementary strategies for constructing application examples: *application generation* and *application matching*.

Application Generation: In many domains, while structured knowledge corpora exist, corresponding application examples remain scarce or

incomplete. To address this gap, we leverage powerful LLMs to automatically generate application examples. This process produces a structured application corpus aligned with the knowledge base, facilitating application-aware reasoning. We categorize knowledge items into two types based on their inherent nature to ensure generation of relevant and task-appropriate applications:

Conceptual knowledge comprises static, descriptive information, such as definitions, theoretical explanations, or descriptions of entities and principles. The corresponding applications generally involve comprehension tasks, contextual interpretations, or analogies that elucidate meaning and deepen understanding.

Procedural knowledge refers to dynamic, actionable information including problem-solving strategies, inference rules, and step-by-step methods. Its associated applications are demonstrated through worked examples, reasoning chains, or practical problem-solving instances where the knowledge is actively applied.

Figure 3 illustrates examples of these knowledge types alongside their generated applications. By aligning each knowledge item with a representative application, the constructed corpus enables downstream systems not only to retrieve relevant information but also to engage in more effective, application-aware reasoning.

Guided by the prior classification of knowledge items into conceptual and procedural types, we design tailored prompting strategies to elicit task-appropriate applications: comprehension or con-

textualization tasks for conceptual knowledge and worked examples or reasoning chains for procedural knowledge.

Application Matching: In domains where real-world cases naturally exemplify the use of specific knowledge, each knowledge item is paired with one or more application instances drawn from authentic scenarios, serving as grounded demonstrations that contextualize and concretize the corresponding knowledge. To establish these pairings, we first perform category alignment by assigning both problems and knowledge items to broad categories using powerful LLMs with temperature sampling and self-consistency voting, followed by manual refinement to ensure accuracy. Next, within each category, relevance selection is conducted by prompting the model to identify the most pertinent knowledge entries for each problem.

This two-stage process yields a many-to-many mapping between knowledge points and application examples, ensuring comprehensive and grounded coverage of relevant applications. For knowledge points lacking matched real-world applications, we supplement them with automatically generated examples as described above, thereby maintaining the completeness and robustness of the application corpus.

These two strategies enable the construction of an application corpus that supports application-aware reasoning across diverse domains.

3.2 Inference Stage

During inference, given a test query, RAG+ first retrieves relevant knowledge items from the knowledge corpus using any retrieval method (e.g., dense retrieval, reranking), as shown in Figure 2(b). For each retrieved knowledge item k , its corresponding application example a , pre-aligned during the construction stage, is retrieved from the application corpus. The pair (k, a) is then incorporated into a predefined prompt template that guides the model with both factual information and procedural cues. This prompt is subsequently fed into the language model for final answer generation.

RAG+ is retrieval-agnostic and can be seamlessly integrated into any existing RAG pipeline. Since the knowledge-application alignment is established offline, no modifications to retrieval or generation models are needed at inference. This modular design allows RAG+ to function as a plug-and-play enhancement for diverse reasoning tasks.

Conceptual Knowledge

Anatomy includes those structures that can be seen grossly (without the aid of magnification) and microscopically (with the aid of magnification). Typically, when used by itself, the term anatomy tends to mean gross or macroscopic anatomy—that is, the study of structures that can be seen without using a microscopic.....

An Application of Conceptual Knowledge

Question: What is the relationship between gross anatomy and microscopic anatomy in the study of human body structures?

Answer Choices:

- A) Gross anatomy and microscopic anatomy are unrelated fields of study.
- B) Gross anatomy deals with structures visible to the naked eye, while microscopic anatomy involves the study of cells and tissues using a microscope.
- C) Microscopic anatomy is a subset of gross anatomy, focusing on larger structures.
- D) Gross anatomy is used only in medical practice, while microscopic anatomy is used in research.

Correct Answer: B) Gross anatomy deals with structures visible to the naked eye, while microscopic anatomy involves the study of cells and tissues using a microscope.

(a) An Example of Conceptual Knowledge and Corresponding Applications.

Procedural Knowledge

Euler's Theorem for Homogeneous Functions: If $P(x_1, \dots, x_m)$ is a homogeneous polynomial of degree d , then: $\nabla P(x_1, \dots, x_m) \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} = d \cdot P(x_1, \dots, x_m)$

An Application of Procedural Knowledge

Question: Let $P(x, y) = x^3 + 3x^2y + 3xy^2 + y^3$. Use Euler's theorem to compute $x \frac{\partial P}{\partial x} + y \frac{\partial P}{\partial y}$

Correct Answer:

$$\frac{\partial P}{\partial x} = 3x^2 + 6xy + 3y^2, \quad \frac{\partial P}{\partial y} = 3x^2 + 6xy + 3y^2$$

Then apply Euler's theorem:

$$\begin{aligned} x \frac{\partial P}{\partial x} + y \frac{\partial P}{\partial y} &= x(3x^2 + 6xy + 3y^2) + y(3x^2 + 6xy + 3y^2) \\ &= 3x^3 + 6x^2y + 3xy^2 + 3x^2y + 6xy^2 + 3y^3 \\ &= 3(x^3 + 3x^2y + 3xy^2 + y^3) \\ &= 3P(x, y) \end{aligned}$$

(b) An Example of Procedural Knowledge and Corresponding Applications.

Figure 3: Examples of aligning different types of knowledge with application instances.

4 Experiment Setup

4.1 Baseline

To assess the effectiveness of the proposed RAG+ framework, we compare it against several representative RAG-based baselines, each embodying a distinct approach to utilizing retrieved information. RAG (Lewis et al., 2020) is the standard framework that retrieves relevant documents based on the input query and generates responses conditioned on both the query and the retrieved content. Answer-First RAG (AFRAG) first generates a candidate answer from the query, which is then used to retrieve supporting evidence. The final output is produced using both the original query and the retrieved context. GraphRAG (Edge et al., 2024) incorporates structured knowledge via knowledge graphs to facilitate multi-hop reasoning and improve contextual relevance. Rerank RAG re-ranks the top- k retrieved documents by a large language model and selects the top three for answer generation to enhance query-context alignment.

Table 1: Accuracy (%) of different models on the MathQA dataset with and without application-level augmentation. “+” indicates the use of application-level augmentation.

| Methods | GLM4-9B | Qwen2.5-7B | DS-Qwen-7B | Qwen2.5-14B | DS-Qwen-32B | Qwen2.5-72B | LLaMA3.3-70B |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Baseline | 46.51 | 58.84 | 24.19 | 66.98 | 80.00 | 69.07 | 69.07 |
| RAG | 47.21 | 64.65 | 24.42 | 73.49 | 82.79 | 70.47 | 71.16 |
| RAG+ | 52.09 | 65.58 | 26.28 | 74.67 | 84.65 | 73.72 | 71.86 |
| AFRAG | 48.14 | 63.51 | 22.56 | 66.98 | 82.09 | 71.86 | 70.23 |
| AFRAG+ | 51.16 | 64.42 | 23.95 | 70.00 | 83.95 | 76.05 | 71.86 |
| GraphRAG | 33.95 | 56.98 | 27.21 | 69.07 | 82.79 | 73.02 | 68.37 |
| GraphRAG+ | 36.51 | 59.77 | 33.72 | 69.77 | 83.49 | 72.56 | 69.00 |
| Rerank RAG | 48.21 | 56.05 | 26.05 | 71.40 | 80.46 | 73.26 | 74.65 |
| Rerank RAG+ | – | 56.28 | 32.09 | 78.90 | 83.26 | 77.21 | 76.74 |

In the Rerank RAG setup, the same model performs both reranking and generation. Smaller models (e.g., GLM4-9B, DS-Qwen-7B) often fail to comply with reranking instructions and instead generate answers directly. This issue remains despite prompt tuning and leads to missing results on the MathQA and MedQA datasets. Larger models (14B and above), however, reliably perform reranking without this problem.

4.2 Datasets

RAG+ is evaluated on three reasoning-intensive domains: mathematics, legal, and medicine. The MathQA dataset is constructed from publicly available educational resources and is paired with a custom mathematical knowledge corpus. For legal, the sentencing prediction dataset from CAIL 2018 (Xiao et al., 2018; Zhong et al., 2018) is used, with a knowledge corpus composed of statutes from the Criminal Law of China. For medicine, the MedQA dataset (Jin et al., 2020) is employed, together with a curated medical corpus from (Xiong et al., 2024) relevant to clinical reasoning.

Because the legal and medical corpora lack sufficient real-world applications, we use automatic generation methods for their application corpora. Specifically, generated applications in legal reflect case rulings from the Chinese Criminal Law corpus, while in medicine, knowledge items are categorized before generating aligned applications based on a clinical knowledge base and the MedQA dataset.

In contrast, the math corpus includes authentic application instances, enabling us to employ an application matching approach. This approach employs a two-stage filtering process: first, category alignment assigns both problems and knowledge items to broad categories using Qwen2.5-72B with temperature sampling and self-consistency voting, followed by manual refinement to improve

accuracy. Second, relevance selection prompts the model to identify the most pertinent knowledge entries within each category. To maintain corpus completeness, knowledge points without matched real-world applications (under 10%) are supplemented with automatically generated examples.

The resulting corpus pairs each knowledge item with one or more applications, enabling RAG+ to retrieve both during inference. All prompts used are detailed in Appendix C.

4.3 Models

Nine conversational models from the Qwen (Qwen et al., 2025; Yang et al., 2025), LLaMA (Dubey et al., 2024), DeepSeek (DeepSeek-AI, 2025), and ChatGLM (Team GLM et al., 2024) series were evaluated. Detailed configurations of the evaluated models are provided in Appendix A.1.

5 Results

In this section, the effectiveness of RAG+ is demonstrated across multiple models on three reasoning-intensive datasets spanning distinct domains. Performance trends are also analyzed across different model scales within the Qwen2.5 series. Ablation studies examine the effects of incorporating only application examples or using larger models for reranking. Finally, case studies illustrate how RAG+ enhances complex reasoning.

5.1 Mathematics Domain

Table 1 reports the model performance on the MathQA dataset with and without application-level augmentation. Most retrieval methods demonstrate effectiveness, though some negatively impact performance. Nevertheless, almost all augmented variants outperform their non-augmented counterparts.

Notably, Qwen2.5-14B achieves a substantial improvement of over 7.5% with Rerank RAG+,

Table 2: Accuracy (%) of different models on sentencing prediction tasks with and without application-level augmentation. LLaMA3.1-8B* denotes the Chinese version, and “+” indicates the use of application-level augmentation.

| Methods | LLaMA3.1-8B* | DS-Qwen-7B | DS-Qwen-32B | QwQ-32B | Qwen3-32B | Qwen2.5-72B | LLaMA3.3-70B |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Baseline | 29.00 | 53.00 | 80.50 | 80.00 | 73.00 | 73.00 | 51.50 |
| RAG | 36.00 | 65.50 | 85.50 | 81.50 | 83.00 | 76.50 | 70.50 |
| RAG+ | 41.00 | 67.50 | 85.50 | 86.00 | 82.50 | 83.00 | 76.00 |
| AFRAG | 27.50 | 65.50 | 85.00 | 82.50 | 76.00 | 85.00 | 41.50 |
| AFRAG+ | 33.00 | 68.00 | 86.50 | 83.00 | 77.50 | 86.50 | 52.50 |
| GraphRAG | 36.50 | 42.00 | 81.50 | 76.00 | 68.50 | 64.00 | 38.50 |
| GraphRAG+ | 46.00 | 47.50 | 81.50 | 77.50 | 75.00 | 64.00 | 52.00 |
| Rerank RAG | 33.00 | 60.00 | 82.00 | 83.50 | 80.50 | 77.50 | 77.50 |
| Rerank RAG+ | 34.00 | 61.00 | 82.50 | 83.50 | 82.00 | 87.50 | 77.50 |

Table 3: Accuracy (%) of different models on the MedQA dataset with and without application-level augmentation. “+” indicates the use of application-level augmentation.

| Methods | LLaMA3.1-8B | Qwen2.5-7B | DS-Qwen-7B | DS-Qwen-32B | QwQ-32B | Qwen2.5-72B | LLaMA3.3-70B |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Baseline | 57.80 | 41.80 | 32.60 | 80.00 | 80.40 | 73.80 | 78.20 |
| RAG | 63.00 | 59.20 | 34.60 | 79.00 | 80.20 | 75.00 | 80.20 |
| RAG+ | 63.60 | 57.60 | 40.20 | 80.20 | 80.80 | 75.40 | 81.40 |
| AFRAG | 56.40 | 53.40 | 32.20 | 78.20 | 81.20 | 76.40 | 82.40 |
| AFRAG+ | 57.00 | 57.20 | 34.60 | 78.60 | 82.20 | 77.40 | 83.00 |
| Rerank RAG | 60.00 | 58.60 | 35.20 | 79.80 | 80.60 | 76.40 | 81.00 |
| Rerank RAG+ | 63.40 | 61.40 | — | 80.20 | 81.40 | 78.20 | 85.60 |

while DS-Qwen-7B showing gains of 6.5% and 6.0% with GraphRAG+ and Rerank RAG+, respectively. GLM4-9B and Qwen2.5-72B show consistent gains between 2.8% and 4.8% across multiple RAG+ variants. In contrast, AFRAG and GraphRAG tend to be less effective on smaller models, such as Qwen2.5-14B with GraphRAG, likely because these methods depend heavily on a model’s ability to interpret complex relational structures and integrate them into reasoning, which smaller models often lack.

Additionally, GraphRAG without application augmentation can lead to performance drops. For instance, Qwen2.5-72B shows a slight decline with plain GraphRAG, possibly due to its emphasis on entity definitions and relations may not align well with mathematical tasks that require solution-oriented knowledge such as formulas.

Overall, most models achieve accuracy improvements between 2.5% and 6.5%, while the overall range spans from approximately 0.7% to 7.5%. These findings highlight the value of application-aware augmentation in reasoning tasks that require more than factual knowledge.

5.2 Legal and Medicine Domain

Table 2 and Table 3 show the performance of various RAG-based methods on the sentencing predic-

tion task and the MedQA dataset across different models. Application-level augmentation consistently improves accuracy over both base models and standard RAG variants.

In the legal domain, Qwen2.5-72B achieves 87.5% accuracy with Rerank RAG+, a 10% gain over its non-augmented version. DS-Qwen-32B and QwQ-32B also show notable improvements with RAG+ and AFRAG+, demonstrating the effectiveness of application-aware augmentation. In contrast, GraphRAG alone underperforms, especially for smaller models like DS-Qwen-7B and LLaMA3.1-8B, likely due to its focus on entity-level information. However, combining GraphRAG with application augmentation significantly improves results, highlighting the need for task-aligned retrieval. LLaMA3.3-70B performs well across all methods but shows marginal gains with application augmentation, indicating diminishing returns from retrieval augmentation as model size increases.

On the MedQA dataset, Rerank RAG+ yields the best performance for most models, especially the larger ones. For example, LLaMA3.3-70B reaches 85.6%, surpassing its baseline (81%) and Rerank RAG (81.0%) methods. Smaller models like Qwen2.5-7B and LLaMA3.1-8B also benefit, with gains of 2.2% and 3.4%, respectively. AFRAG

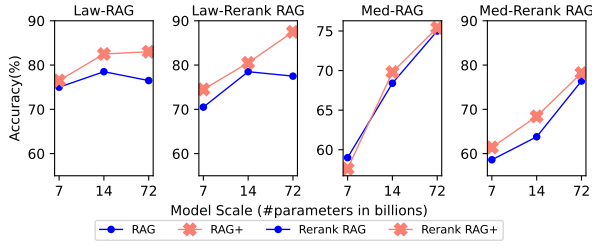


Figure 4: Performance comparison on sentencng prediction tasks and the MedQA dataset using Qwen2.5 series models across different scales with and without application examples.

and its augmented version provide steady improvements, showing that grounding abstract medical knowledge through applications enhances reasoning. While standard RAG offers a solid boost, application-level augmentation is essential for maximizing performance across model sizes.

These results demonstrate the consistent effectiveness of application-aware augmentation in enhancing RAG-based methods, benefiting both large and small models when paired with effective retrieval strategies.

5.3 Effect of Model Scale

As shown in Figure 4, all methods show consistent performance gains as model size increases from 7B to 14B and then to 72B, reflecting the enhanced reasoning capabilities of larger models. Notably, approaches augmented with application-level examples achieve larger improvements than their non-augmented counterparts.

In the legal domain, accuracy rises steadily with model scale when application augmentation is used. In the medicine domain, while all models benefit, the performance gains from scaling are less pronounced. These results indicate that application-guided retrieval scales especially well with model size, further boosting knowledge integration and reasoning in complex, domain-specific tasks.

5.4 Impact of Reranking Model

To investigate the suboptimal performance of Rerank RAG on smaller models, we conducted an ablation study by replacing the initial reranker with a stronger model, Qwen2.5-72B, to assess its impact on downstream accuracy.

As shown in Table 4 and Table 5, we isolate the effect of reranking quality by varying the reranker while keeping other components fixed. Stronger reranking consistently improves accuracy across

Table 4: Accuracy (%) on sentencng prediction tasks across different RAG methods. “Rerank (72B)” denotes reranked by Qwen2.5-72B, generated by base model.

| Methods | Qwen2.5-7B | Qwen2.5-14B |
|-------------------|--------------|--------------|
| Baseline | 46.00 | 74.00 |
| RAG | 75.00 | 78.50 |
| RAG+ | 76.50 | 82.50 |
| Rerank RAG | 70.50 | 78.50 |
| Rerank RAG+ | 74.50 | 81.00 |
| Rerank (72B) RAG | 72.50 | 80.00 |
| Rerank (72B) RAG+ | 83.50 | 86.00 |

Table 5: Accuracy (%) on the MedQA dataset across different RAG methods.

| Methods | Qwen2.5-7B | Qwen2.5-14B |
|-------------------|--------------|--------------|
| Baseline | 41.80 | 68.80 |
| RAG | 59.20 | 68.40 |
| RAG+ | 57.60 | 69.80 |
| Rerank RAG | 58.60 | 63.80 |
| Rerank RAG+ | 61.40 | 68.40 |
| Rerank (72B) RAG | 60.80 | 69.80 |
| Rerank (72B) RAG+ | 61.80 | 70.00 |

both tasks, with the largest gains in the Rerank (72B) RAG+ setting, where application-augmented reranking achieves the best performance. Notably, its improvement over non-augmented reranking exceeds that between Rerank RAG+ and Rerank RAG, highlighting the joint benefits of reranking strength and application-level guidance. These results underscore the value of cross-model collaboration, where large models enhance retrieval to support more efficient inference in smaller ones.

5.5 Effect of Application Examples Only

We conducted experiments to assess the contribution of application examples without their associated knowledge on the legal domain. Results are shown in Table 6. Incorporating only application examples leads to performance improvements over the base models in most cases. This indicates that application-level guidance alone provides a meaningful benefit. However, this setting still underperforms compared to the full configuration where both knowledge and application examples are included, referred to as RAG+ and AFRAG+.

These findings suggest that while application examples alone can partially support reasoning, likely by providing structural cues or task-specific formulations, the explicit inclusion of knowledge remains essential for achieving optimal performance. There-

| | |
|-----------------|--|
| Math Question A | Find a polynomial $P(x)$ of degree 3 or less whose graph passes through the four data points $(-2,8)$, $(0,4)$, $(1,2)$, $(3,-2)$. A) $P(x) = 2x^2 - 3x + 1$, B) $P(x) = x^3 - 2x^2 + x - 1$, C) $P(x) = 3x^2 + 2x - 4$, D) $P(x) = 4 - 2x$ |
| RAG Answer | To find a polynomial $P(x)$ of degree 3 or less that passes through the four points $(-2,8)$, $(0,4)$, $(1,2)$, and $(3,-2)$, we use polynomial interpolation. We'll construct it using the Lagrange interpolation formula. For each point (x_i, y_i) , we define a corresponding Lagrange basis polynomial $\ell_i(x)$ and we construct interpolating polynomial then: $P(x) = \sum_{i=1}^4 y_i \ell_i(x)$ For the given points, we calculate each of the Lagrange basis polynomials: $\ell_1(x) = \frac{(x-1)(x-3)(x-0)}{(x+2)(x-1)(x-3)(x-0)}$; $\ell_2(x) = \frac{(x+2)(x-1)(x-3)(x-0)}{(x+2)(x-1)(x-3)(x-0)}$ Next, we construct the polynomial $P(x)$: $P(x) = 8 \ell_1(x) + 4 \ell_2(x) + 2 \ell_3(x) - 2 \ell_4(x)$ Substituting the expressions for the basis polynomials and simplifying, we get: $P(x) = \frac{(x+2)(x-1)(x-3)(x-0)}{(x+2)(x-1)(x-3)(x-0)} + \frac{(x+2)(x-1)(x-3)(x-0)}{(x+2)(x-1)(x-3)(x-0)} - \frac{(x+2)(x-1)(x-3)(x-0)}{(x+2)(x-1)(x-3)(x-0)} - \frac{(x+2)(x-1)(x-3)(x-0)}{(x+2)(x-1)(x-3)(x-0)} = x^3 - 2x^2 + x - 1$. Therefore, the answer is B) $P(x) = x^3 - 2x^2 + x - 1$ |
| RAG+ Answer | To find the polynomial $P(x)$ of degree 3 or less that passes through the points $(-2, 8)$, $(0, 4)$, $(1, 2)$, and $(3, -2)$, we use Newton's divided differences method. We begin by calculating the divided differences step by step: First-order differences: $f[-2, 0] = \frac{4 - 8}{0 - (-2)} = -2$, $f[0, 1] = \frac{2 - 4}{1 - 0} = -2$, $f[1, 3] = \frac{-2 - 2}{3 - 1} = -2$ Second-order differences: $f[-2, 0, 1] = \frac{-2 - (-2)}{1 - (-2)} = 0$, $f[0, 1, 3] = \frac{-2 - (-2)}{3 - 0} = 0$ Third-order difference: $f[-2, 0, 1, 3] = \frac{0 - 0}{3 - (-2)} = 0$ From this, we get the coefficients for the Newton form of the polynomial: Constant term: $f[-2] = 8$; First-order term: $f[-2, 0] = -2$; Higher-order terms are zero. So the polynomial is: $P(x) = 8 + (-2)(x + 2) = 4$. Therefore, the answer is D) $P(x) = 4 - 2x$ |

Figure 5: Comparison of output reasoning paths between RAG and RAG+ in the math domain using the Qwen2.5-72B model.

Table 6: Accuracy (%) on legal tasks across different RAG methods. “—” indicates application-only retrieval; “+” indicates the inclusion of retrieved knowledge and application.

| Methods | Qwen2.5-7B | Qwen2.5-14B | DS-Qwen-32B |
|----------|--------------|--------------|--------------|
| Baseline | 46.00 | 74.00 | 80.50 |
| RAG- | 65.00 | 76.50 | 84.50 |
| RAG | 75.00 | 78.50 | 85.50 |
| RAG+ | 76.50 | 82.50 | 85.50 |
| AFRAG- | 34.50 | 74.50 | 83.00 |
| AFRAG | 40.50 | 81.00 | 85.00 |
| AFRAG+ | 47.00 | 85.50 | 86.50 |

fore, combining knowledge with its applications produces the most effective results.

5.6 Case Study

To evaluate the practical effectiveness of RAG+, we present a case study on the mathematics domain, as illustrated in Figure 5. The answer generated by RAG correctly identifies the interpolation method as Lagrange interpolation but fails to execute it accurately due to the complexity of intermediate symbolic expressions. While the approach is mathematically valid, errors in deriving the basis polynomials lead to an incorrect final result. In comparison, the Newton divided differences method, although less commonly emphasized in retrieval-based settings, provides a more transparent and step-by-step procedure. Its recursive computation of coefficients reduces algebraic errors and produces the correct polynomial. This suggests that even when the correct method is retrieved, symbolic reasoning may fail due to execution errors, highlighting the need for verification mechanisms alongside retrieval.

6 Discussion

To enable application-aware retrieval, we propose a lightweight method for constructing application-specific corpora from a knowledge base, typically generating one to two examples per item. The corpus grows linearly with the number of items and introduces no retrieval overhead, as each item is directly paired with an application example. In experiments, the initial corpus sizes were 223 KB (math), 528 KB (legal), and 99,382 KB (medicine). Generating the legal corpus with Qwen2.5-72B on eight 64 GB NPUs took around six hours, which is acceptable given this scale. Final sizes reached 612 KB, 868 KB, and 105,558 KB, respectively, closely matching the sizes of the underlying knowledge corpora. The method is scalable and supports efficient incremental updates.

7 Conclusion

In this work, we introduce RAG+, a framework that integrates application-level augmentation into retrieval-augmented generation. Through comprehensive experiments across diverse domains and model scales, we demonstrate that incorporating application examples consistently leads to performance improvements. RAG+ outperforms baselines across tasks and RAG variants, highlighting the value of structured application in leveraging retrieved knowledge for reasoning. Our results indicate that retrieval alone is insufficient—effective alignment and application of retrieved knowledge are crucial. Future work may explore more advanced application strategies and tighter integration to further enhance reasoning in LLMs.

Limitations

While RAG+ consistently improves performance, it also has several limitations. First, constructing a high-quality application corpus can be resource-intensive, especially in domains with limited annotated data. Automated generation depends heavily on large language models, which may introduce errors or oversimplify complex reasoning.

Second, RAG+ assumes a strong alignment between knowledge and application pairs, but mismatches can occur—particularly when retrieved knowledge is noisy or incomplete—leading to incorrect or misleading reasoning.

Finally, our current approach focuses on enhancing reasoning via application-level augmentation, but does not directly address retrieval quality or efficiency, which remain critical to overall performance. Future work should explore joint optimization of retrieval and application generation, as well as better handling of uncertainty and ambiguity in retrieved content.

Ethics Statement

All datasets and models used in this study are open-source, and the licenses of the models have been clearly specified. The prompts used in the experiments, along with the full experimental environment and all relevant parameters, are provided to ensure reproducibility. The entire study can be replicated using widely available large model API frameworks. Furthermore, all manual operations involved in the process were carried out solely by the authors, without the involvement of any external collaborators or paid contributors. This guarantees both the transparency and independence of the research.

References

- L. Ammann, S. Ott, C. R. Landolt, and M. P. Lehmann. 2025. [Securing rag: A risk assessment and mitigation framework](#). *arXiv preprint arXiv:2505.08728*.
- John R Anderson, Daniel Bothell, Michael D Byrne, Scott Douglass, Christian Lebiere, and Yulin Qin. 2004. An integrated theory of the mind. *Psychological review*, 111(4):1036.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*.
- Benjamin Samuel Bloom. 2010. *A taxonomy for learning, teaching, and assessing: A revision of Bloom's taxonomy of educational objectives*. Longman.
- 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](#).
- DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#).
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, and et al. 2024. [The llama 3 herd of models](#).
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitan, Robert Osazuwa Ness, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*.
- K. Hayashi, H. Kamigaito, and S. Kouda. 2025. [Iterkey: Iterative keyword generation with llms for enhanced retrieval augmented generation](#). *arXiv preprint arXiv:2505.08450*.
- Shayekh Bin Islam, Md Asib Rahman, K S M Tozammel Hossain, Enamul Hoque, Shafiq Joty, and Md Rizwan Parvez. 2024. [Open-RAG: Enhanced retrieval augmented reasoning with open-source large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 14231–14244, Miami, Florida, USA. Association for Computational Linguistics.
- 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*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- 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.
- Z. Li, Z. Wang, W. Wang, K. Hung, and H. Xie. 2025. [Retrieval-augmented generation for educational application: A systematic survey](#). *Computers and Education: Artificial Intelligence*.
- K.H. Lin, T.H. Kao, L.C. Wang, C.T. Kuo, and P.C.H. Chen. 2025. [Benchmarking large language models gpt-4o, llama 3.1, and qwen 2.5 for cancer genetic variant classification](#). *NPJ Precision Oncology*.

- Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Richard James, Pedro Rodriguez, Jacob Kahn, Gergely Szilvasy, Mike Lewis, et al. 2023. Ra-dit: Retrieval-augmented dual instruction tuning. In *The Twelfth International Conference on Learning Representations*.
- Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. [Query rewriting in retrieval-augmented large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5303–5315, Singapore. Association for Computational Linguistics.
- S. Mostafa, M. K. Abdel-Aziz, and M. S. Elbamby. 2025. [Rag-enabled intent reasoning for application-network interaction](#). *arXiv preprint arXiv:2505.09339*.
- Marc Pickett, Jeremy Hartman, Ayan Kumar Bhowmick, Raquib ul Alam, and Aditya Vempaty. 2025. [Better rag using relevant information gain](#).
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. [Qwen2.5 technical report](#).
- L. Rong, Y. Zhang, P. Tiwari, and M. Yu. 2025. [Begoni-agpt: Cultivating the large language model to be an exceptional k-12 english teacher](#). *Neural Networks*.
- Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talaei Khoei. 2025. [Agentic retrieval-augmented generation: A survey on agentic rag](#).
- H. Tang, H. Yan, and R. Song. 2025. [Synthetic sentiment cue enhanced graph relation-attention network for aspect-level sentiment analysis](#). *IEEE Access*.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, and et al. 2024. [Chatglm: A family of large language models from glm-130b to glm-4 all tools](#).
- Liang Wang, Nan Yang, and Furu Wei. 2023a. Query2doc: Query expansion with large language models. *arXiv preprint arXiv:2303.07678*.
- Xiaohua Wang, Zhenghua Wang, Xuan Gao, Feiran Zhang, Yixin Wu, Zhibo Xu, Tianyuan Shi, Zhengyuan Wang, Shizheng Li, Qi Qian, Ruicheng Yin, Changze Lv, Xiaoqing Zheng, and Xuanjing Huang. 2024a. [Searching for best practices in retrieval-augmented generation](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17716–17736, Miami, Florida, USA. Association for Computational Linguistics.
- Zhihu Wang, Shiwan Zhao, Yu Wang, Heyuan Huang, Sitao Xie, Yubo Zhang, Jiaxin Shi, Zhixing Wang, Hongyan Li, and Junchi Yan. 2024b. Re-task: Revisiting llm tasks from capability, skill, and knowledge perspectives. *arXiv preprint arXiv:2408.06904*.
- Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. 2023b. Learning to filter context for retrieval-augmented generation. *arXiv preprint arXiv:2311.08377*.
- Chaojun Xiao, Haoxi Zhong, Zhipeng Guo, Cunchao Tu, Zhiyuan Liu, Maosong Sun, Yansong Feng, Xianpei Han, Zhen Hu, Heng Wang, and Jianfeng Xu. 2018. [Cail2018: A large-scale legal dataset for judgment prediction](#).
- Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. Benchmarking retrieval-augmented generation for medicine. *arXiv preprint arXiv:2402.13178*.
- Guangzhi Xiong, Qiao Jin, Xiao Wang, Yin Fang, Haolin Liu, Yifan Yang, Fangyuan Chen, Zhixing Song, Dengyu Wang, Minjia Zhang, Zhiyong Lu, and Aidong Zhang. 2025. Rag-gym: Optimizing reasoning and search agents with process supervision. *arXiv preprint arXiv:2502.13957*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. 2025. [Qwen3 technical report](#).
- Fuda Ye, Shuangyin Li, Yongqi Zhang, and Lei Chen. 2024. [R²AG: Incorporating retrieval information into retrieval augmented generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 11584–11596, Miami, Florida, USA. Association for Computational Linguistics.
- Haoxi Zhong, Chaojun Xiao, Zhipeng Guo, Cunchao Tu, Zhiyuan Liu, Maosong Sun, Yansong Feng, Xianpei Han, Zhen Hu, Heng Wang, and Jianfeng Xu. 2018. [Overview of cail2018: Legal judgment prediction competition](#).
- Wang Zihao, Liu Anji, Lin Haowei, Li Jiaqi, Ma Xiaojian, and Liang Yitao. 2024. Rat: Retrieval augmented thoughts elicit context-aware reasoning in long-horizon generation. *arXiv preprint arXiv:2403.05313*.

A Experimental Details

In this study, we employ a range of hardware devices, including units of 32 GB Ascend 910B4, 64 GB Ascend 910B3, and 32 GB Tesla V100 PCIe units. For large-scale models exceeding 70 billion parameters (e.g., LLaMA3.3-70B and Qwen2.5-72B), we deploy eight Ascend 910B4 devices (32 GB each). For 32B-scale models (e.g., Qwen2.5-32B, Qwen3-32B, QwQ-32B, and DS-Qwen-32B), two Ascend 910B4 devices are deployed. Models with approximately 14 billion parameters (e.g., Qwen2.5-14B) are run on a single 64 GB Ascend 910B3, while models around the 8B scale (e.g., Llama3.1-8B, Qwen2.5-7B, DS-Qwen-7B, and ChatGLM4-9B) are executed on a single 32 GB Ascend 910B4. For models used in retrieval-augmented generation (RAG) components (e.g., bg3-m3 and bge-reranker-v2-m3), we utilize a 32 GB Tesla V100 PCIe. All deployments are carried out using the vLLM inference framework (Kwon et al., 2023).

For models generating long-form outputs (i.e., QwQ-32B, DS-Qwen-7B, and DS-Qwen-32B), we set the decoding temperature to 1 and top_p set to 1 to encourage output diversity. For all other models, we apply a deterministic decoding strategy with a temperature of 0 and top_p of 1. Each experiment is conducted three times, and we report the average performance across these runs.

To implement RAG functionality, we utilize the Dify framework. Specifically, Dify is used for corpus uploading, with a segmentation length of 600 tokens. The text is pre-segmented to ensure that the divisions are structurally coherent and suitable for retrieval.

A.1 Model sources and licenses

This study employs nine publicly available conversational models: LLaMA-Chinese-8B-Instruct (sourced from ModelScope) and LLaMA3.1-8B, LLaMA3.3-70B, Qwen2.5-7B, Qwen2.5-72B, QwQ-32B, DeepSeek-R1-Distill-Qwen-7B, DeepSeek-R1-Distill-Qwen-32B, and ChatGLM4-9B (all obtained from Hugging Face). Additionally, the BGE-M3 embedding model (Chen et al., 2024) is utilized for generating text embeddings. Details regarding model sources and licensing are provided in Table 7.

The LLaMA-Chinese-8B-Instruct model is exclusively used for the sentencing prediction task in the legal domain, as this task is conducted in

Chinese. For the remaining two tasks, LLaMA3.1-8B is employed. The QwQ-32B model is excluded from the numerical analysis task due to the high complexity of the questions, which frequently lead to excessively long outputs and potential infinite loops. These issues significantly increase inference time and GPU memory usage, thereby exceeding available computational resources.

The BGE-M3 model is employed to generate embeddings for corpus texts, while the BGE-reranker-v2-m3 model is employed to reorder candidate documents based on their relevance to the user query. This reranking step enhances the overall quality of the semantic retrieval process.

A.2 Baseline Details

We evaluate nine configurations: Baseline, RAG (Lewis et al., 2020), RAG+, Answer-First RAG, Answer-First RAG+, GraphRAG (Edge et al., 2024), GraphRAG+, Rerank RAG, and Rerank RAG+.

The RAG configuration follows the original implementation by Lewis et al. (Lewis et al., 2020), and is implemented using the Dify framework. We employ the BGE-M3 model to embed both corpus chunks and key terms for similarity matching, and the BGE-Reranker-v2-M3 model to rerank the retrieved candidates. High-quality indices are constructed for retrieval, with each corpus chunk limited to 800 tokens. For all RAG-based configurations, the top three most relevant text chunks are retrieved for each query.

Baseline: The Baseline configuration directly uses the large language model (LLM) to answer questions without incorporating any external knowledge retrieval. It serves as the baseline for evaluating the effectiveness of various RAG-based enhancements.

RAG: Vanilla RAG directly applies the retrieval-augmented generation method. The matching fields are the question and its answer options. Retrieved content is incorporated into a prompt template and passed to the model to generate an answer.

Answer-First RAG: Following the approach by Asai et al. (Asai et al., 2023), we adopt a more direct implementation. Given a question Q , the model M is first prompted to generate a preliminary answer A^* , which is then used as a query to retrieve relevant content C . The final answer A is generated by prompting the model with both Q and C .

Table 7: Models, sources and licenses used in this work.

| Models | URL | Licenses |
|--------------------|---|--------------------|
| LLaMA3-Chinese-8B | https://www.modelscope.cn/models/FlagAlpha/LLaMA3-Chinese-8B-Instruct/summary | Apache License 2.0 |
| LLaMA3.1-8B | https://huggingface.co/meta-llama/Meta-Llama-3.1-8B-Instruct | llama3.1 license |
| LLaMA3.3-70B | https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct | llama3.3 license |
| Qwen2.5-7B | https://huggingface.co/Qwen/Qwen2.5-7B-Instruct | Apache License 2.0 |
| Qwen2.5-72B | https://huggingface.co/Qwen/Qwen2.5-72B-Instruct | Qwen license |
| QWQ-32B | https://huggingface.co/Qwen/QwQ-32B | Apache License 2.0 |
| DS-Qwen-7B | https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-7B | MIT License |
| DS-Qwen-32B | https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-32B | MIT License |
| ChatGLM4-9b | https://huggingface.co/THUDM/glm-4-9b-chat | glm-4-9b License |
| bg3-m3 | https://huggingface.co/BAAI/bge-m3 | MIT License |
| bge-reranker-v2-m3 | https://huggingface.co/BAAI/bge-reranker-v2-m3 | Apache License 2.0 |

GraphRAG: Based on [Edge et al. \(Edge et al., 2024\)](#), GraphRAG integrates structured knowledge into a prompt for answer generation. We use Qwen2.5-72B for entity extraction, relationship extraction, claim identification, summary generation, and community report creation. The embedding model is set to “cl100k_base”, while “nomic/text-embedding-nomic-embed-text-v1.5@q4_k_m” is used for graph construction. We adopt a chunk size of 1200 tokens and a maximum cluster size of 10, and use “local_search” for graph retrieval. Only the top relevant chunk is returned during inference. Due to the large size of the MedQA corpus in the medical domain, we omit GraphRAG for this task, as graph construction and retrieval are prohibitively expensive. However, for the NA and legal domains (each corpus approximately 1000 KB), GraphRAG is feasible and conducted accordingly.

Rerank RAG: Rerank RAG is a multi-step retrieval process. First, vanilla RAG retrieves the top ten candidates. These are then passed, along with the question, to the target model for re-ranking. The top three items from this reranked list are used to construct the final prompt for answer generation. This process allows the model to participate in the retrieval pipeline, enabling more informed selection of relevant knowledge.

RAG+ and Answer-First RAG+: In these variants, each retrieved knowledge item is mapped to a corresponding application. Both the knowledge and its application are inserted into the prompt template to guide the model’s answer generation.

GraphRAG+: GraphRAG+ extends GraphRAG

by enriching retrieved knowledge with application content. Since constructing dedicated applications for GraphRAG-extracted entities and relations would be labor-intensive, we instead reuse applications from the original corpus via fuzzy matching. Once a match is found, both the application and the matched knowledge—along with the original GraphRAG output—are included in the prompt.

Rerank RAG+: Rerank RAG+ extends Rerank RAG by requiring the model to output results in a specific format, enabling application mapping. However, models frequently fail to follow this format, often providing direct answers or returning too few knowledge items. This inconsistency complicates the parsing process. To address this, we increase the number of runs and re-prompt the model to ensure successful output parsing. This approach proves effective across most models in our experiments, with the exception of GLM4-9B on the Numerical Analysis dataset.

A.3 Dataset Details

A.3.1 Math QA in Mathematic Domain

We conduct experiments in the domain of Numerical Analysis. Most existing mathematics datasets are relatively simple—often achieving over 90% accuracy with current models—and primarily emphasize reasoning rather than Retrieval-Augmented Generation (RAG). To better suit our research objectives, we constructed a new Numerical Analysis dataset from scratch.

We first collected a range of publicly available online resources, supplemented with original ques-

tions and corresponding answers. The dataset includes both standard QA items and multiple-choice questions with varying numbers of answer options. Only the correct answers were retained. For compound questions that inquire about multiple values (e.g., "What is the value of y when $x = 0$ or $x = 1$?"), we decomposed them into separate single-answer questions, such as "What is the value of y when $x = 0$ " and "What is the value of y when $x = 1$?".

Next, we transformed all QA pairs into a multiple-choice format by prompting GPT to generate three additional plausible but incorrect answer options. This step expands the test set while simultaneously constraining the model’s generative space.

Finally, we manually reviewed and refined the knowledge points, solution demonstrations, and generated answer choices to ensure overall quality and accuracy. All annotations and validations were conducted solely by the authors without external assistance.

The final dataset consists of 430 test questions in the domain of Numerical Analysis.

A.3.2 Sentencing Prediction in Legal Domain

We use data from the CAIL 2018 dataset (Xiao et al., 2018; Zhong et al., 2018), selecting 200 samples from our training set. We focus exclusively on questions related to Article 234 of the Criminal Law, which concerns sentencing for intentional injury. All questions were converted to a multiple-choice format, with answer options corresponding to sentencing durations: less than three years, three to ten years, and more than ten years.

A.3.3 MedQA in Medicine Domain

We use the MedQA dataset curated by Jin et al.. We randomly sampled 500 examples from the dataset to serve as our training set. All selected items are in a multiple-choice format.

A.4 Corpus Details

A.4.1 Mathematics

The knowledge points in the Numerical Analysis (NA) corpus are collected from various online sources. These include definitions, theorems, lemmas, factual statements, and methods, and are typically concise. Due to inconsistent formatting, we used Qwen2.5-72B to extract key knowledge points and normalize them into a unified format, as illustrated in Figure 19.

After processing, we obtained a total of 816 knowledge points. We examined their lengths with a predefined chunk size of 800 tokens. No knowledge point exceeded this threshold, so no additional splitting was required.

We followed the Dify framework to construct the knowledge base: processed knowledge points were uploaded using the pre-segmented chunks, embedded using the BGE-M3 model, and indexed for high-quality retrieval. During inference, the top three most relevant chunks were retrieved.

A.4.2 Legal

The legal corpus is extracted from the Criminal Law of the People’s Republic of China. Chunking was based on the natural structure of the legal text, treating each article as a single knowledge point, resulting in 451 items. The corpus is divided into two major sections: the General Provisions (101 items), which provide conceptual knowledge such as definitions and factual descriptions, and the Specific Provisions, which provide solution knowledge in the form of sentencing guidelines.

No further processing was applied. We used a chunk size of 800 tokens and followed the same Dify-based uploading, embedding, and retrieval pipeline as in the mathematics domain.

A.4.3 Medicine

The medical corpus is sourced from Xiong et al. (Xiong et al., 2024). We use only the text-book portion of the corpus, where knowledge is already structured into discrete knowledge points, each representing a self-contained fragment.

The corpus spans 18 subjects and includes 64,117 knowledge points, with a total size of 99,382 KB. Each knowledge point, averaging less than 600 words, was treated as a single chunk without further segmentation. The knowledge encompasses both conceptual and procedural content.

We used a chunk size of 800 tokens, consistent with other domains, and uploaded the data to the Dify system using the same embedding and retrieval procedure.

B Experimental Results

B.1 Complete Results

Accuracy is used as the evaluation metric. We conducted three independent inference runs for all experiments and calculated the average results. The complete results are shown in Table 8 - 10.

Table 8: Accuracy (%) of different models on the sentencing prediction task with and without application-level augmentation. LLaMA3.1-8B* denotes the Chinese version.

| Methods | LLaMA3.1-8B* | DS-Qwen-7B | DS-Qwen-32B | QwQ-32B | Qwen3-32B | Qwen2.5-72B | LLaMA3.3-70B |
|-------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| Baseline | 29.00 \pm 0.0 | 53.00 \pm 2.9 | 80.50 \pm 1.8 | 80.00 \pm 2.6 | 73.00 \pm 0.0 | 73.00 \pm 0.0 | 51.50 \pm 0.0 |
| RAG | 36.00 \pm 0.0 | 65.50 \pm 1.7 | 85.50 \pm 1.8 | 81.50 \pm 2.5 | 83.00 \pm 0.0 | 76.50 \pm 0.0 | 70.50 \pm 0.0 |
| RAG+ | 41.00 \pm 0.0 | 67.50 \pm 1.8 | 85.50 \pm 2.1 | 86.00 \pm 2.4 | 82.50 \pm 0.0 | 83.00 \pm 0.0 | 76.00 \pm 0.0 |
| AFRAG | 27.50 \pm 0.0 | 65.50 \pm 2.2 | 85.00 \pm 2.6 | 82.50 \pm 2.9 | 76.00 \pm 0.0 | 85.00 \pm 0.0 | 41.50 \pm 0.0 |
| AFRAG+ | 33.00 \pm 0.0 | 68.00 \pm 2.3 | 86.50 \pm 2.5 | 83.00 \pm 2.9 | 77.50 \pm 0.0 | 86.50 \pm 0.0 | 52.50 \pm 0.0 |
| GraphRAG | 36.50 \pm 0.0 | 42.00 \pm 2.8 | 81.50 \pm 2.2 | 76.00 \pm 2.3 | 68.50 \pm 0.0 | 64.00 \pm 0.0 | 38.50 \pm 0.0 |
| GraphRAG+ | 46.00 \pm 0.0 | 47.50 \pm 1.6 | 81.50 \pm 2.5 | 77.50 \pm 2.7 | 75.00 \pm 0.0 | 64.00 \pm 0.0 | 52.00 \pm 0.0 |
| Rerank RAG | 33.00 \pm 0.0 | 60.00 \pm 2.1 | 82.00 \pm 1.8 | 83.50 \pm 2.0 | 80.50 \pm 0.0 | 77.50 \pm 0.0 | 77.50 \pm 0.0 |
| Rerank RAG+ | 34.00 \pm 0.0 | 61.00 \pm 1.1 | 82.50 \pm 1.9 | 83.50 \pm 2.1 | 82.00 \pm 0.0 | 87.50 \pm 0.0 | 77.50 \pm 0.0 |

Table 9: Accuracy (%) of different models on the MedQA dataset with and without application-level augmentation.

| Methods | LLaMA3.1-8B | Qwen2.5-7B | DS-Qwen-7B | DS-Qwen-32B | QwQ-32B | Qwen2.5-72B | LLaMA3.3-70B |
|-------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| Baseline | 57.80 \pm 0.0 | 41.80 \pm 0.0 | 32.60 \pm 1.2 | 80.00 \pm 0.8 | 80.40 \pm 1.2 | 73.80 \pm 0.0 | 78.20 \pm 0.0 |
| RAG | 63.00 \pm 0.0 | 59.20 \pm 0.0 | 34.60 \pm 1.6 | 79.00 \pm 1.0 | 80.20 \pm 0.9 | 75.00 \pm 0.0 | 80.20 \pm 0.0 |
| RAG+ | 63.60 \pm 0.0 | 57.60 \pm 0.0 | 40.20 \pm 0.5 | 80.20 \pm 0.5 | 80.80 \pm 0.4 | 75.40 \pm 0.0 | 81.40 \pm 0.0 |
| AFRAG | 56.40 \pm 0.0 | 53.40 \pm 0.0 | 32.20 \pm 1.3 | 78.20 \pm 0.9 | 81.20 \pm 1.6 | 76.40 \pm 0.0 | 82.40 \pm 0.0 |
| AFRAG+ | 57.00 \pm 0.0 | 57.20 \pm 0.0 | 34.60 \pm 1.2 | 78.60 \pm 0.9 | 82.20 \pm 1.2 | 77.40 \pm 0.0 | 83.00 \pm 0.0 |
| Rerank RAG | 60.00 \pm 0.0 | 58.60 \pm 0.0 | 35.20 \pm 0.8 | 79.80 \pm 0.5 | 80.60 \pm 0.9 | 76.40 \pm 0.0 | 81.00 \pm 0.0 |
| Rerank RAG+ | 63.40 \pm 0.0 | 61.40 \pm 0.0 | — | 80.20 \pm 1.0 | 81.40 \pm 0.7 | 78.20 \pm 0.0 | 85.60 \pm 0.0 |

Based on the results in Table 11, incorporating application examples alone improves performance over the base models in several cases, particularly for Qwen2.5-14B and Qwen2.5-72B. However, combining both knowledge and application examples generally yields better results, as shown by the superior accuracy of RAG+ and AFRAG+ across most models.

B.2 Model Scale

Table 12 presents the performance of Qwen2.5 models of different scales on the sentencing prediction task in the legal domain, and Table 13 shows their performance on the MedQA dataset.

B.3 Case Study

We present two case studies in Figures 6 and 7, comparing the outputs of Qwen3-32B and DeepSeek-R1-Distill-Qwen-7B under RAG and RAG+ configurations on the sentencing prediction and MedQA dataset, respectively. The results demonstrate that merely retrieving external knowledge, as done in standard RAG, is often insufficient: models may still make reasoning mistakes or misuse the retrieved content. In contrast, with the integration of the application module in RAG+, the models are able to apply the retrieved knowledge more appropriately, leading to correct predictions. These qualitative results further support the

effectiveness of RAG+ in enhancing knowledge utilization during inference.

C Prompts

The prompt configurations used across three domains—legal, medical, and mathematical—under the Base, RAG, and RAG+ settings are shown in Figures 8 to 16. Specifically, Figures 8, 9, and 10 illustrate the prompt designs for the legal domain, while similar configurations for the medical, and mathematical domains are provided in Figures 11 to 13, and 14 to 16, respectively. These templates clearly reflect how retrieved knowledge is introduced and applied in each configuration, enabling consistent comparison across domains and setups.

C.1 Dataset

The prompt templates used to generate application examples for different domains are shown in Figures 17, 18, and 19, corresponding to the legal, medical, and mathematical domains, respectively. These templates are designed to guide the model in producing domain-specific knowledge applications that align with the downstream tasks.

C.2 Examples of Knowledge and Applications

Examples of the retrieved and applied knowledge in different domains are shown in Figures 20 and 21 for the mathematical domain, Figures 22 and 23

Table 10: Accuracy (%) of different models on the Math dataset with and without application-level augmentation.

| Methods | GLM4-9B | Qwen2.5-7B | DS-Qwen-7B | Qwen2.5-14B | DS-Qwen-32B | Qwen2.5-72B | LLaMA3.3-70B |
|-------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| Baseline | 46.51 \pm 0.0 | 58.84 \pm 0.0 | 24.19 \pm 0.7 | 66.98 \pm 0.0 | 80.00 \pm 0.6 | 69.07 \pm 0.0 | 69.07 \pm 0.0 |
| RAG | 47.21 \pm 0.0 | 64.65 \pm 0.0 | 24.42 \pm 1.0 | 73.49 \pm 0.0 | 82.79 \pm 1.2 | 70.47 \pm 0.0 | 71.16 \pm 0.0 |
| RAG+ | 52.09 \pm 0.0 | 65.58 \pm 0.0 | 26.28 \pm 0.7 | 74.67 \pm 0.0 | 84.65 \pm 1.3 | 73.72 \pm 0.0 | 71.86 \pm 0.0 |
| AFRAG | 48.14 \pm 0.0 | 63.51 \pm 0.0 | 22.56 \pm 0.6 | 66.98 \pm 0.0 | 82.09 \pm 0.5 | 71.86 \pm 0.0 | 70.23 \pm 0.0 |
| AFRAG+ | 51.16 \pm 0.0 | 64.42 \pm 0.0 | 23.95 \pm 0.8 | 70.00 \pm 0.0 | 83.95 \pm 0.6 | 76.05 \pm 0.0 | 71.86 \pm 0.0 |
| GraphRAG | 33.95 \pm 0.0 | 56.98 \pm 0.0 | 27.21 \pm 0.8 | 69.07 \pm 0.0 | 82.79 \pm 0.6 | 73.02 \pm 0.0 | 68.37 \pm 0.0 |
| GraphRAG+ | 36.51 \pm 0.0 | 59.77 \pm 0.0 | 33.72 \pm 0.8 | 69.77 \pm 0.0 | 83.49 \pm 0.5 | 72.56 \pm 0.0 | 69.00 \pm 0.0 |
| Rerank RAG | 48.21 \pm 0.0 | 56.05 \pm 0.0 | 26.05 \pm 0.9 | 71.40 \pm 0.0 | 80.46 \pm 0.7 | 73.26 \pm 0.0 | 74.65 \pm 0.0 |
| Rerank RAG+ | — | 56.28 \pm 0.0 | 32.09 \pm 0.6 | 78.90 \pm 0.0 | 83.26 \pm 0.5 | 77.21 \pm 0.0 | 76.74 \pm 0.0 |

Table 11: Performance comparison of application-only retrieval vs knowledge-enhanced retrieval on the MedQA dataset.

| Methods | Qwen2.5-14B | Qwen3-32B | DS-Qwen-32B | Qwen2.5-72B |
|-----------|--------------|--------------|--------------|--------------|
| Baseline | 68.80 | 84.40 | 79.00 | 73.80 |
| RAG | 68.40 | 82.60 | 81.00 | 75.00 |
| RAG-app | 70.60 | 83.60 | 78.40 | 76.60 |
| RAG-plus | 69.80 | 84.20 | 82.20 | 75.20 |
| AFRAG | 69.00 | 83.00 | 80.20 | 76.40 |
| AFRAG-app | 70.60 | 82.42 | 77.20 | 77.80 |
| AFRAG+ | 69.40 | 83.60 | 80.60 | 77.40 |

for the legal domain, and Figures 24 and 25 for the medical domain. Each pair of figures illustrates the difference between knowledge produced under the RAG and RAG+ configurations, highlighting how RAG+ promotes more targeted and applicable knowledge generation in support of downstream reasoning.

Table 12: Accuracy (%) Comparison of Qwen2.5 Models of Varying Scales on Sentencing Prediction.

| Methods | Qwen2.5-7B | Qwen2.5-14B | Qwen2.5-72B |
|-------------|--------------|--------------|--------------|
| Baseline | 46.00 | 74.00 | 58.00 |
| RAG | 75.00 | 78.50 | 76.50 |
| RAG+ | 76.50 | 82.50 | 83.00 |
| AF-RAG | 40.50 | 81.00 | 85.00 |
| AF-RAG+ | 47.00 | 85.50 | 85.50 |
| GraphRAG | 47.00 | 68.00 | 64.00 |
| GraphRAG+ | 59.00 | 79.00 | 64.00 |
| Rerank RAG | 70.50 | 78.50 | 77.50 |
| Rerank RAG+ | 74.50 | 81.00 | 87.50 |

Table 13: Accuracy (%) Comparison of Qwen2.5 Models of Varying Scales on the MedQA Dataset.

| Methods | Qwen2.5-7B | Qwen2.5-14B | Qwen2.5-72B |
|-------------|--------------|--------------|--------------|
| Baseline | 41.80 | 68.80 | 73.80 |
| RAG | 59.20 | 68.40 | 75.00 |
| RAG+ | 57.60 | 69.80 | 75.20 |
| AF-RAG | 53.40 | 69.00 | 76.40 |
| AF-RAG+ | 57.20 | 69.40 | 77.40 |
| Rerank RAG | 58.60 | 63.80 | 76.40 |
| Rerank RAG+ | 61.40 | 68.40 | 78.20 |

| | |
|------------------|--|
| Legal Question C | <p>Select the most appropriate sentence interval from the three options A, B, and C (A: less than or equal to 36 months; B: greater than 36 months and less than or equal to 120 months; C: greater than 120 months) for the criminal Li:</p> <p>It was found that at about 1:00 on April 29, 2016, the defendant Chen and Chen No. 2 (a minor, handled in another case) and others went to an Internet Cafe on Binjiang to surf the Internet. They learned that Wen was beaten by Liu due to a debt dispute. The defendant Chen went up to Liu for trouble, and had an argument with Chen No. 3 who was traveling with Liu. The defendant Chen and Chen No. 2 beat Chen No. 3 together. After the two sides had another argument, Chen No. 2 returned to his residence and brought a small Japanese machete, an iron baseball bat and three stainless steel hollow steel pipes. The defendant Chen beat Liu with a baseball bat, Chen No. 2 with a steel pipe, and Wen with a machete. After being beaten, Liu fled in the direction of the Bund Square, and the defendant Chen and three others continued to chase and beat Liu. According to the appraisal of the Physical Evidence Appraisal Office of the Public Security Bureau of Yunyang County, Chongqing, the degree of injury of Liu's facial skin lacerations was in the category of first-grade minor injuries, and the degree of injury of the remaining scars and bleeding of the left ring nail bed was in the category of minor injuries. The appraisal opinion was that Liu's injury was a first-degree minor injury.</p> |
| RAG Answer | <p>According to the facts of the case provided, the defendant Chen and others (including the minor Chen 2) participated in the attack on Liu with weapons (baseball bats) in a group conflict, resulting in Liu's facial skin laceration reaching the level of first -degree minor injury, and the rest of the injuries were minor injuries. According to Article 234 of the Criminal Law of the People's Republic of China, for the intentional injury to the body of another person, if the person is slightly injured, he may be sentenced to fixed -term imprisonment of not more than three years, detention or control. However, in this case, Chen participated in the fight with weapons. Although Liu's injuries have not yet reached the standard of serious injury, considering the circumstances of the armed fight, he should be punished more severely.</p> <p>Taking the above circumstances into consideration, although Liu's injury is a first -degree minor injury, considering the fact of the armed fight, it is necessary to impose a heavier punishment; at the same time, Chen has the circumstances of surrendering him self, so he can be given a lighter punishment. Therefore, considering the specific circumstances, the sentence for Chen should be more than three years, but not more than ten years.</p> <p>Therefore, based on the above analysis, the answer is B.</p> |
| RAG+ Answer | <p>According to the facts provided, the defendant Chen participated in the beating of Liu, causing Liu to suffer a minor injury of the first degree. According to Article 234 of the Criminal Law of the People's Republic of China, if someone intentionally injures another person's body and causes minor injuries to others, he shall be sentenced to fixed -term imprisonment of not more than three years, detention or control in accordance with the law.</p> <p>In the trial of a specific case, it is also necessary to consider the defendant's mitigating or aggravating circumstances. According to the information provided, Chen has the following circumstances that can be given a lighter punishment: 1. Chen was arrested by his father with the assistance of the public security organs, indicating that he did not have a strong desire to evade legal sanctions. 2. After being arrested, Chen truthfully confessed the facts of the crime and had a confession. 3. Chen was not the ringleader and did not personally carry a weapon when committing the crime.</p> <p>Combining the above circumstances, according to the provisions of the Criminal Law of the People's Republic of China and relevant judicial interpretations, criminals with confessions, assistance in arrests, etc. can be given a lighter punishment in accordance with the law.</p> <p>Therefore, Chen's sentence of less than three years' imprisonment is in line with the nature and circumstances of his behavior, and is a lenient sentence within the scope of the law. In summary, the answer is A.</p> |

Figure 6: Comparison of Qwen3-32B Outputs under RAG and RAG+ Configurations on Sentencing Prediction.

| | |
|---------------------|---|
| Medicine Question B | <p>A 67-year-old man presents to his primary care physician for erectile dysfunction. He states that for the past month he has been unable to engage in sexual intercourse with his wife despite having appropriate sexual desire. He also endorses deep and burning buttock and hip pain when walking, which is relieved by rest. The patient states that he does not have erections at night or in the morning. His past medical history is notable for diabetes, coronary artery disease, and hypertension, and he has a 40 pack-year smoking history. Physical exam is notable for weak lower extremity and femoral pulses. Which of the following is the most specific etiology of this patient's symptoms?</p> <p>A: Anxiety, B: Aortoiliac atherosclerosis, C: Spinal stenosis, D: Vascular claudication</p> |
| RAG Answer | <p>Okay, so I'm trying to figure out the etiology of this 67-year-old man's symptoms.</p> <p>First, the patient presents with erectile dysfunction, specifically inability to engage in sexual intercourse for the past month, despite having the sexual desire. He also mentions deep and burning buttock and hip pain that's relieved by rest. There's no erections at night or in the morning.</p> <p>His medical history includes diabetes, coronary artery disease, and hypertension, plus a 40 pack-year smoking history. The physical exam shows weak lower extremity and femoral pulses, which suggests possible peripheral vascular issues.</p> <p>Looking at the options: A is anxiety, B is aortoiliac atherosclerosis, C is spinal stenosis, and D is vascular claudication.</p> <p>Anxiety (A) could cause ED, but the symptoms here are more physical. The presence of diabetes, hypertension, and coronary artery disease points towards cardiovascular issues. The physical exam findings, like weak pulses, support a vascular problem.</p> <p>Aortoiliac atherosclerosis (B) is a type of cardiovascular disease where the aorta and ilies are narrowed, sometimes causing claudication. Vascular claudication (D) is when blood flow becomes intermittent, often due to narrowing or stenosis in the arteries, leading to pain on exertion, which matches the patient's buttock pain.</p> <p>Spinal stenosis (C) typically affects nerve root pressure, causing numbness or tingling, not necessarily pain on the buttocks or hip. Plus, the physical exam doesn't mention numbness or tingling, just pain that's relieved by rest.</p> <p>The patient's age is 67, which is above 40, making vascular conditions more likely. His smoking history and presence of diabetes and hypertension also support a cardiovascular etiology.</p> <p>Putting it all together, the symptoms and medical history point most strongly to vascular claudication, which is caused by narrowing in the aortoiliac arteries.</p> <p>The patient's symptoms, including erectile dysfunction, pain on exertion, and peripheral pulses, are consistent with vascular issues. His medical history and physical exam support aortoiliac atherosclerosis causing claudication. Therefore, the answer is D.</p> |
| RAG+ Answer | <p>Alright, let's tackle this question step by step. So, we have a 67-year-old man presenting with erectile dysfunction. He's been experiencing this for a month, says he can't have sex with his wife. He mentions deep and burning buttock and hip pain that gets better when he rests. He doesn't have erections at night or in the morning. His medical history includes diabetes, coronary artery disease, hypertension, and a 40 pack-year smoking history. The physical exam shows weak lower extremities and femoral pulses.</p> <p>First, I need to figure out the possible causes of his symptoms. Erectile dysfunction can be due to either psychological or physical factors. But the options given are all physical conditions: anxiety, aortoiliac atherosclerosis, spinal stenosis, and vascular claudication.</p> <p>Given his symptoms, the buttock and hip pain that improves with rest suggests something affecting the lower extremities. His physical exam shows weak lower extremity and femoral pulses, which is a classic sign of peripheral artery disease, or PAD. PAD is often caused by atherosclerosis in the carotid arteries, leading to reduced blood flow. Vascular claudication is a type of PAD that causes intermittent, but his symptoms are more of pain and weakness, which points more towards atherosclerosis.</p> <p>His diabetes and hypertension are known to contribute to atherosclerosis. The presence of aortoiliac atherosclerosis would indeed cause narrowing of the aortae, leading to reduced blood flow to the lower extremities. Spinal stenosis would cause pain and possibly numbness, but the presence of weak pulses and pain with rest suggests something in the lower extremities rather than the spine.</p> <p>Anxiety is a possible cause, but the question is asking for the most specific etiology. The physical exam and symptoms point more to a physical issue. His smoking history is significant for atherosclerosis, which aligns with aortoiliac atherosclerosis.</p> <p>Spinal stenosis is less likely because the pain is localized to the buttocks and hips, not involving the spinal area. Vascular claudication is possible but typically presents with pain that is intermittent, not constant but worsening with exercise. The fact that his pulses are weak and he has reduced blood flow supports atherosclerosis rather than claudication.</p> <p>So, considering all the factors, the most specific etiology here seems to be aortoiliac atherosclerosis, which is option B. Therefore, the correct answer should be B: Aortoiliac atherosclerosis.</p> |

Figure 7: Comparison of DeepSeek-R1-Distill-Qwen-7B model Outputs under RAG and RAG+ Configurations on MedQA.

| |
|---|
| <p>## Role: You are a judge, please complete the user-specified task.</p> <p>## Task: According to the following input, the criminal XXX will be sentenced to sentence.</p> <p>##Format: Choose from A, B, and C. (A: less than or equal to 36 months; B: greater than 36 months and less than or equal to 120 months; C: greater than 120 months) Select the most appropriate sentence interval from. Please fill in the answer between [Sentence] and [eoa], for example, [Sentence] C [eoa].)</p> <p>## Input: [Put The Input Fact Here]</p> |
|---|

Figure 8: An Example of the Prompt Template Used for the Sentencing Prediction Task under the Base Configuration.

Role:
You are a judge, please complete the user-specified task based on the understanding of the <Reference Law>.

Reference Law
[Put the Reference Law Here]

Task:
According to the following input, the criminal XXX will be sentenced to sentence.

##Format:
Choose from A, B, and C. (A: less than or equal to 36 months; B: greater than 36 months and less than or equal to 120 months; C: greater than 120 months) Select the most appropriate sentence interval from. Please fill in the answer between [Sentence] and [eoal], for example, [Sentence] C [eoal].)

Input:
[Put The Input Fact Here]

Figure 9: An Example of the Prompt Template Used for the Sentencing Prediction Task under the RAG Configuration.

Role:
You are a judge, please complete the user-specified task based on the understanding of the <Law> and the apply of the <Law Application> in the <Reference Law>.

Reference Law:
Law:
[Put the Law Here]

Law Application
[Put the Application of Law Here]

Task:
According to the following input, the criminal XXX will be sentenced to sentence.

##Format:
Choose from A, B, and C. (A: less than or equal to 36 months; B: greater than 36 months and less than or equal to 120 months; C: greater than 120 months) Select the most appropriate sentence interval from. Please fill in the answer between [Sentence] and [eoal], for example, [Sentence] C [eoal].)

Input:
[Put The Input Fact Here]

Figure 10: An Example of the Prompt Template Used for the Sentencing Prediction Task under the RAG+ Configuration.

You are a decision-evaluation assistant. Please help me solve the medicine question. Please answer the question step by step in a XML format with the reasoning step enclosed with the tag <Think></Think> and the answer option enclosed with the tag <Answer></Answer>. You must choose one correct option among A-D.

Here is the question: [Put The Questions Here]
Options: [Put the Options Here]
Answer:

Figure 11: An Example of the Prompt Template Used for the MedQA Dataset under the Base Configuration.

You are a decision-evaluation assistant. Your task is to solve the given question based on the Reference Knowledge. There are some knowledge to help you solve the problem:

```
<Reference>
  <Knowledge>
    [Put the Knowledge Here]
  </Knowledge>

  <Knowledge>
    [Put the Knowledge Here]
  </Knowledge>

  <Knowledge>
    [Put the Knowledge Here]
  </Knowledge>
</Reference>
```

Now use the reference knowledge provided for guidance (but do not be strictly bound by them), please answer the question step by step in a XML format with the reasoning step enclosed with the tag <Think></Think> and the answer option enclosed with the tag <Answer></Answer>. You must choose one correct option among A -D.

Here is the question: [Put the Questions Here]
Options: [Put the Options Here]
Answer:

Figure 12: An Example of the Prompt Template Used for the MedQA Dataset under the RAG Configuration.

You are a decision-evaluation assistant. Your task is to solve the given question based on the Reference Knowledge. There are some knowledge to help you solve the problem:

```
<Reference>
  <Knowledge Point>
    <Knowledge>
      [Put One Knowledge Here]
    </Knowledge>
    <Application>
      [Put The Application of the Knowledge Here]
    </Application>
  </Knowledge Point>

  <Knowledge Point>
    <Knowledge>
      [Put One Knowledge Here]
    </Knowledge>
    <Application>
      [Put The Application of the Knowledge Here]
    </Application>
  </Knowledge Point>
</Reference>
```

Now use the reference knowledge and applications provided for guidance (but do not be strictly bound by them), please answer the question step by step in a XML format with the reasoning step enclosed with the tag <Think></Think> and the answer option enclosed with the tag <Answer></Answer>. You must choose one correct option among A -D.

Here is the question:[Put the Questions Here]
Options: [Put the Options Here]
Answer:

Figure 13: An Example of the Prompt Template Used for the MedQA Dataset under the RAG+ Configuration.

Please help me solve the numerical analysis question. Please answer the question step by step in a XML format with the reasoning step enclosed with the tag <Think></Think> and the answer option enclosed with the tag <Answer></Answer>. You must choose one correct option among A -D.

Here is the question: [Put the Questions Here]
Options: [Put the Options Here]
Answer:

Figure 14: An Example of the Prompt Template Used for the Math Task under the Base Configuration.

Please help me solve the numerical analysis question. Here are some reference knowledge that might help you solve the question:

```

<Reference>
  <Knowledge>
    [Put the Knowledge Here]
  </Knowledge>

  <Knowledge>
    [Put the Knowledge Here]
  </Knowledge>

  <Knowledge>
    [Put the Knowledge Here]
  </Knowledge>
</Reference>

```

Now use the reference knowledge provided for guidance (but do not be strictly bound by them), please answer the question step by step in a XML format with the reasoning step enclosed with the tag <Think></Think> and the answer option enclosed with the tag <Answer></Answer>. You must choose one correct option among A -D.

Here is the question: [Put the Questions Here]
Options:[Put the Options Here]
Answer:

Figure 15: An Example of the Prompt Template Used for the Math Task under the RAG Configuration.

Please help me solve the numerical analysis question. Here are some reference knowledge and applications of the knowledges that might help you solve the question; each Knowledge includes one Knowledge Point and some Application Examples:

```

<Reference>
  <Knowledge Point>
    <Knowledge>
      [Put One Knowledge Here]
    </Knowledge>
    <Application>
      [Put The Application of the Knowledge Here]
    </Application>
  </Knowledge Point>

  <Knowledge Point>
    <Knowledge>
      [Put One Knowledge Here]
    </Knowledge>
    <Application>
      [Put The Application of the Knowledge Here]
    </Application>
  </Knowledge Point>
</Reference>

```

Now use the reference knowledge and applications provided for guidance (but do not be strictly bound by them), please answer the question step by step in a XML format with the reasoning step enclosed with the tag <Think></Think> and the answer option enclosed with the tag <Answer></Answer>. You must choose one correct option among A -D.

Here is the question: [Put the Questions Here]
Options: [Put the Options Here]
Answer:

Figure 16: An Example of the Prompt Template Used for the Math Task under the RAG+ Configuration.

I will provide you with a text and ask you to generate corresponding examples based on the required knowledge types. If it is factual knowledge or conceptual knowledge, please generate examples or explanations corresponding to the knowledge. If it is procedural knowledge, please generate an example of using that procedural knowledge to complete a task. All examples should be presented in the form of questions or multiple-choice questions.

Here is the text: [Put the Knowledge here]

Figure 17: The prompt template of generating the applications for the knowledge in legal domain.

I'll give you three different types of knowledge. Please help me generate examples of knowledge as required.
 For factual knowledge and conceptual knowledge, please help me give an example of that conceptual knowledge, or generate an example of understanding of that knowledge.
 For procedural knowledge, please help me generate an example of using that procedural knowledge.

All examples must be question-and-answer or multiple-choice questions. Generate an example for each knowledge that meets the requirements. Use ----- to separate the examples.

Here is the text: [Put the Knowledge here]

Figure 18: The prompt template of generating the applications for the knowledge in medicine domain.

I will provide you with a section. I need your help to extract key knowledge points and examples. Here are the specific tasks :

1. Extract Knowledge Points:
 - Include definitions, theorems, lemmas, and problem-solving methods extracted from plain text.
 - Use tags like `` and ``, `` and `` to enclose each knowledge point. Each point should be ended with tags like ``.
 - For definitions, theorems, and lemmas, extract them directly. List the name of the knowledge point first, followed by its detailed content.
 - For problem-solving methods or problem-solving demonstrations, provide the name of the method, steps to extract more general methods for accomplishing such tasks, and the core points of each step. The steps should be general and not tied to a specific instance.
2. Extract Examples:
 - Use tags like `` and ``, `` and `` to enclose each example. Each demonstration should first provide the corresponding knowledge point, which is represented by tag like ``.
 - Some examples may correspond to multiple knowledge points. If this is the case, list the same demonstration for each relevant knowledge point. For example, if a demonstration corresponds to more than one knowledge point, list all the points.
 - The demonstration may be given in methods, examples, or plain text. If it is provided in a method or plain text, the example corresponds to the knowledge points extracted from the method. If it is given in an example, it corresponds to the knowledge points near the example or mentioned earlier.
 - Each demonstration should contain explicit questions and outputs, where the outputs include problem solving steps and answers.
 - Do not miss any examples mentioned in the text.

Here is the text: [Put the Knowledge here]

Figure 19: The prompt template of generating the applications for the knowledge in mathematic domain.

```
<Knowledge>
Name: Romberg Integration Algorithm
Content: The Romberg Integration algorithm involves the following steps:
1. Initialize  $R_{11} = \frac{b-a}{2}(f(a) + f(b))$ .
2. For  $j = 2, 3, \dots$ : - Calculate  $h_j = \frac{b-a}{2^{j-1}}$ . - Compute  $R_{j1} = \frac{1}{2}R_{j-1,1} + h_j \sum_{i=1}^{2^{j-2}} f(a + (2i-1)h_j)$ .
3. For  $k = 2, \dots, j$ : - Compute  $R_{jk} = \frac{4^{k-1}R_{j,k-1} - R_{j-1,k-1}}{4^{k-1} - 1}$ .
4. Repeat until the desired accuracy is achieved.
</Knowledge>
```

```
<Knowledge>
Name: Romberg Integration
Content: Romberg Integration is an efficient method for calculating definite integrals that can be extended by adding data until the required accuracy is attained. It applies extrapolation to the composite Trapezoid Rule to achieve higher-order accuracy.
</Knowledge>
```

```
<Knowledge>
Name: Romberg Integration Tableau
Content: The Romberg Integration tableau is a lower triangular matrix used to organize the extrapolated values. The general  $(j, k)$ th entry is given by:  $R_{jk} = \frac{4^{k-1}R_{j,k-1} - R_{j-1,k-1}}{4^{k-1} - 1}$ .
</Knowledge>
```

Figure 20: Examples of the knowledge in Mathematic Domain.

<KnowledgePoint>
 Let $(P(x))$ be the interpolating polynomial of $((x_1, f(x_1)), \dots, (x_n, f(x_n)))$ where the (x_i) are distinct. Then: (a) $(P(x) = f(x_1) + f(x_1, x_2)(x - x_1) + f(x_1, x_2, x_3)(x - x_1)(x - x_2) + \dots + f(x_1, x_2, \dots, x_n)(x - x_1)(x - x_2)\dots(x - x_{n-1}))(b)$
 For $(k > 1)$, $(f(x_1, \dots, x_k) = \frac{f(x_2, \dots, x_k) - f(x_1, \dots, x_{k-1})}{x_k - x_1})$.
 </KnowledgePoint>
 <ApplicationExample>
 What should you do if you are interpolating a set of data points with a high-degree polynomial and notice extreme oscillations?
 A) Increase the degree of the polynomial to get a better fit.
 B) Use a different interpolation method, such as spline interpolation, to avoid the Runge phenomenon.
 C) Add more data points to the set to reduce the oscillations
 Answer: B) Use a different interpolation method, such as spline interpolation, to avoid the Runge phenomenon
 </ApplicationExample>

<KnowledgePoint>
 Let $(P(x_1, \dots, x_m))$ be a homogeneous polynomial of degree (d) . Then $(\nabla P \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} = dP)$. Proof: Differentiating with respect to (c) yields $(x_1 P_{x_1}(cx_1, \dots, cx_m) + \dots + x_m P_{x_m}(cx_1, \dots, cx_m) = dc^{d-1}P(x_1, \dots, x_m))$ using the multivariable chain rule. Evaluating at $(c = 1)$ results in the desired conclusion.
 </KnowledgePoint>
 <ApplicationExample>
 Question: Use the Backward Difference Equation with Newton iteration to solve Fisher's equation with homogeneous Neumann boundary conditions $(\begin{cases} u_t = D * u_{xx} + u(1 - u) \\ u(x, 0) = \sin(\pi x) \end{cases} \text{ for } 0 \leq x \leq 1, u_x(0, t) = 0 \text{ for all } t \geq 0)$
 Answer: The discretization retraces the derivation that was carried out for Burgers' equation: $(\frac{w_{i+1}^j - 2w_i^j + w_{i-1}^j}{h^2} + w_i^j(1 - w_i^j))$ or $((1 + 2\sigma - k(1 - w_i^j))w_i^j - \sigma(w_{i+1}^j + w_{i-1}^j) - w_i^j(j-1) = 0)$. This results in the nonlinear equations $(F_i(z_1, \dots, z_m) = (1 + 2\sigma - k(1 - z_i))z_i - \sigma(z_{i+1} + z_{i-1}) - w_i^j(j-1) = 0)$ to solve for the $(z_i = w_i^j)$ at the (j) th time step. The first and last equations will establish the Neumann boundary conditions: $(F_1(z_1, \dots, z_m) = \frac{-3z_0 + 4z_1 - z_2}{2h} = 0)$ $(F_m(z_1, \dots, z_m) = \frac{-z_{m-2} + 4z_{m-1} - 3z_m}{2h} = 0)$ The Jacobian (DF) has the form $(\begin{bmatrix} -3 & 4 & -1 \\ -\sigma & 1 + 2\sigma - k + 2kw_2 & -\sigma \\ -\sigma & 1 + 2\sigma - k + 2kw_3 & -\sigma & \dots & \dots & \dots & \dots & \dots & \dots \\ -\sigma & 1 + 2\sigma - k + 2kw_{m-1} & -\sigma & -1 & 4 & -3 \end{bmatrix})$ After altering the function (F) and Jacobian (DF) , the Newton iteration implemented in Program 8.7 can be used to solve Fisher's equation. Lemma 8.11 can be used to separate the degree 1 and 2 parts of (DF) . Neumann boundary conditions are also applied, as shown in the code fragment below: $(DF1 = \text{diag}(1 - k + 2\sigma, \text{ones}(m, 1)) + \text{diag}(-\sigma, \text{ones}(m - 1, 1), 1); DF1 = DF1 + \text{diag}(-\sigma, \text{ones}(m - 1, 1), -1); DF2 = \text{diag}(2 * k * w1); DF = DF1 + DF2; F = -w(:, j) + (DF1 + DF2 / 2) * w1; DF(1, :) = [-3 \ 4 \ -1 \ \text{zeros}(1, m - 3)]; F(1) = DF(1, :) * w1; DF(m, :) = [\text{zeros}(1, m - 3) \ -1 \ 4 \ -3]; F(m) = DF(m, :) * w1;$ Figure 8.20 shows approximate solutions of Fisher's equation with $(D = 1)$ that demonstrate the tendency to relax to the attracting equilibrium $(u(x, t) \equiv 1)$. Of course, $(u(x, t) \equiv 0)$ is also a solution of (8.69) with $(f(u) = u(1 - u))$, and will be found by the initial data $(u(x, 0) = 0)$. Almost any other initial data, however, will eventually approach $(u = 1)$ as (t) increases.
 </ApplicationExample>

Figure 21: Examples of the knowledge and the applications in Mathematic Domain.

<Knowledge>
 Article 234: Whoever intentionally harms the body of another person shall be sentenced to fixed-term imprisonment of not more than three years, criminal detention, or public surveillance. Whoever commits the crime mentioned in the preceding paragraph and causes serious injury to another person shall be sentenced to fixed-term imprisonment of not less than three years but not more than ten years; Whoever causes death or serious disability by using particularly cruel means shall be sentenced to fixed-term imprisonment of not less than ten years, life imprisonment, or death penalty. If there are other provisions in this law, they shall be followed.
 </Knowledge>

Figure 22: Examples of the knowledge in legal Domain.

<KnowledgePoint>

Article 234: Whoever intentionally harms the body of another person shall be sentenced to fixed-term imprisonment of not more than three years, criminal detention, or public surveillance. Whoever commits the crime mentioned in the preceding paragraph and causes serious injury to another person shall be sentenced to fixed-term imprisonment of not less than three years but not more than ten years; Whoever causes death or serious disability by using particularly cruel means shall be sentenced to fixed-term imprisonment of not less than ten years, life imprisonment, or death penalty. If there are other provisions in this law, they shall be followed.

</KnowledgePoint>

<ApplicationExample>

Question: What criteria should be used to determine the applicable penalty measures when dealing with cases of intentional injury?

Options: A) Based solely on the severity of the injury; B) Only in accordance with legal provisions; C) Based on specific circumstances and legal regulations; D) According to the victim's request.

Answer: C) According to specific circumstances and legal regulations

</ApplicationExample>

Figure 23: Examples of the knowledge and the applications in legal Domain.

<Knowledge>

What is anatomy? Anatomy includes those structures that can be seen grossly (without the aid of magnification) and microscopically (with the aid of magnification). Typically, when used by itself, the term anatomy tends to mean gross or macroscopic anatomy—that is, the study of structures that can be seen without using a microscope. Microscopic anatomy, also called histology, is the study of cells and tissues using a microscope. Anatomy forms the basis for the practice of medicine. Anatomy leads the physician toward an understanding of a patient's disease, whether he or she is carrying out a physical examination or using the most advanced imaging techniques. Anatomy is also important for dentists, chiropractors, physical therapists, and all others involved in any aspect of patient treatment that begins with an analysis of clinical signs. The ability to interpret a clinical observation correctly is therefore the endpoint of a sound anatomical understanding.

</Knowledge>

<Knowledge>

Each of these approaches has benefits and deficiencies. The regional approach works very well if the anatomy course involves cadaver dissection but falls short when it comes to understanding the continuity of an entire system throughout the body. Similarly, the systemic approach fosters an understanding of an entire system throughout the body, but it is very difficult to coordinate this directly with a cadaver dissection or to acquire sufficient detail. The anatomical position

</Knowledge>

Figure 24: Examples of the knowledge in Medicine Domain.

<KnowledgePoint>

What is anatomy? Anatomy includes those structures that can be seen grossly (without the aid of magnification) and microscopically (with the aid of magnification). Typically, when used by itself, the term anatomy tends to mean gross or macroscopic anatomy—that is, the study of structures that can be seen without using a microscope. Microscopic anatomy, also called histology, is the study of cells and tissues using a microscope. Anatomy forms the basis for the practice of medicine. Anatomy leads the physician toward an understanding of a patient's disease, whether he or she is carrying out a physical examination or using the most advanced imaging techniques. Anatomy is also important for dentists, chiropractors, physical therapists, and all others involved in any aspect of patient treatment that begins with an analysis of clinical signs. The ability to interpret a clinical observation correctly is therefore the endpoint of a sound anatomical understanding.

</KnowledgePoint>

<ApplicationExample>

Question: What is the relationship between gross anatomy and microscopic anatomy in the study of human body structures?

Answer Choices: A) Gross anatomy and microscopic anatomy are unrelated fields of study. B) Gross anatomy deals with structures visible to the naked eye, while microscopic anatomy involves the study of cells and tissues using a microscope. C) Microscopic anatomy is a subset of gross anatomy, focusing on larger structures. D) Gross anatomy is used only in medical practice, while microscopic anatomy is used in research.

Correct Answer: B) Gross anatomy deals with structures visible to the naked eye, while microscopic anatomy involves the study of cells and tissues using a microscope.

</ApplicationExample>

<KnowledgePoint>

Each of these approaches has benefits and deficiencies. The regional approach works very well if the anatomy course involves cadaver dissection but falls short when it comes to understanding the continuity of an entire system throughout the body. Similarly, the systemic approach fosters an understanding of an entire system throughout the body, but it is very difficult to coordinate this directly with a cadaver dissection or to acquire sufficient detail.

</KnowledgePoint>

<ApplicationExample>

Question: If you were to use a systemic approach to study the human body, which of the following sequences would be correct?

Answer Choices: A) Study the thorax, then the abdomen, then the pelvis, and so on. B) Study the cardiovascular system, then the nervous system, then the skeletal system, and so on. C) Study the thorax, then the cardiovascular system, then the abdomen, and so on. D) Study the thorax, then the upper limb, then the lower limb, and so on.

Correct Answer: B) Study the cardiovascular system, then the nervous system, then the skeletal system, and so on.

</ApplicationExample>

Figure 25: Examples of the knowledge and the applications in Medicine Domain.