

# Reinforcement Learning Improves Traversal of Hierarchical Knowledge in LLMs

Anonymous ACL submission

## Abstract

Reinforcement learning (RL) is often credited with improving language model reasoning and generalization at the expense of degrading memorized knowledge. We challenge this narrative by observing that **RL-enhanced models consistently outperform their base and supervised fine-tuned (SFT) counterparts on pure knowledge recall tasks**, particularly those requiring traversal of hierarchical, structured knowledge (e.g., medical codes). We hypothesize these gains stem not from newly acquired data, but from improved procedural skills in navigating and searching existing knowledge hierarchies within the model parameters. To support this hypothesis, we show that structured prompting—which explicitly guides SFTed models through hierarchical traversal—recovers most of the performance gap (reducing 24pp to 7pp on MedConceptsQA for DeepSeek-V3/R1). We further find that while prompting improves final-answer accuracy, RL-enhanced models retain superior ability to recall correct procedural paths on deep-retrieval tasks. Finally our layer-wise internal activation analysis reveals that while factual representations (e.g., activations for the statement “code 57.95 refers to urinary infection”) maintain high cosine similarity between SFT and RL models, query representations (e.g., “what is code 57.95”) diverge noticeably, indicating that **RL primarily transforms how models traverse knowledge rather than the knowledge representation itself**.

## 1 Introduction

Large Language Models (LLMs) acquire vast parametric knowledge during pre-training, encoding facts, concepts, and their relationships across billions of parameters. Post-training techniques—including

supervised fine-tuning (SFT), Reinforcement Learning from Human Feedback (RLHF), and specialized reasoning-focused RL—are then applied to transform these base models into instruction-following agents capable of complex reasoning (Yu et al., 2025; Wang et al., 2025c; Bai et al., 2022). While these methods improve performance on reasoning benchmarks and user preference metrics, a growing body of evidence reveals a concerning trade-off known as the “alignment tax” (Lin et al., 2024; Askell et al., 2021; Sorensen et al., 2025): models sacrifice factual memorization capabilities to optimize for other objectives, leading to reduced performance on knowledge-intensive benchmarks (Yuan et al., 2024; Gekhman et al., 2024). However, existing work has primarily focused on direct factual recall tasks over unstructured knowledge, leaving a critical gap: *do these degradation patterns hold for all forms of parametric knowledge retrieval tasks?*

To address this question, we investigate tasks where retrieval demands navigating hierarchical structures encoded within the model’s parameters. Consider medical code lookup (Figure 3): to identify that ICD-9-CM code 57.95 refers to “Replacement of indwelling urinary catheter,” a model can attempt direct recall—often failing due to the vast code space—or systematically traverse the taxonomy (Chapter 11 → codes 57.0-57.99 → specific procedure). Surprisingly, reasoning-enhanced models outperform their base counterparts by 24 percentage points on MedConceptsQA, directly **challenging the conventional wisdom that RL sacrifices memorization for reasoning** (Ghosh et al., 2024; Chu et al., 2025). We hypothesize these models succeed through systematic hierarchical navigation rather than direct recall, proposing that **reinforcement learning enhances navigation of ex-**

091 **isting parametric knowledge rather**  
092 **than adding new factual content.**

093 To disentangle knowledge acquisition from  
094 navigation, we design three complemen-  
095 tary experiments. First, inspired by work  
096 showing prompt optimization can match  
097 RL gains (Agrawal et al., 2025; Khattab  
098 et al., 2023; Ziems et al., 2025), we de-  
099 velop structured prompting that explic-  
100 itly guides base models through hierar-  
101 chical traversal. If knowledge exists in  
102 base models, prompting should surface  
103 it. **Structured prompting reduces the**  
104 **24pp gap between DeepSeek-V3 and**  
105 **DeepSeek-R1 to 7pp, suggesting in-**  
106 **formation is present but inaccessi-**  
107 **ble without proper navigation** (Fig-  
108 ure 3, right-hand side). Second, to validate  
109 that improved traversal drives these gains,  
110 we introduce a complexity-stratified patent  
111 classification dataset and Path Matching  
112 Score metric measuring traversal accuracy.  
113 We show that as recall depth increases  
114 (from fewer than 3 hops to more than 5),  
115 **reasoning models demonstrate supe-**  
116 **rior path recall accuracy**, with the per-  
117 formance gap widening from 5pp to 9pp,  
118 demonstrating that reasoning models excel  
119 at complex hierarchical navigation (Ta-  
120 ble 5).

121 Third, to provide internal validation, we  
122 conduct layer-wise representational analy-  
123 sis inspired by work examining how post-  
124 training modifies internal model structure  
125 (Mukherjee et al., 2025; Skean et al., 2025;  
126 Huan et al., 2025; Agarwal et al., 2024).  
127 We extract layer-wise representations for  
128 matched query-answer pairs, comparing in-  
129 terrogative queries (e.g., “What is the med-  
130 ical code 57.95?”) versus declarative state-  
131 ments (e.g., “Code 57.95 refers to uri-  
132 nary catheter replacement”). We find a  
133 striking pattern (Figure 2): declarative  
134 statements maintain high cosine similarity  
135 (0.85-0.92) between base and RL models  
136 throughout most layers, while interroga-  
137 tive queries diverge substantially (similar-  
138 ity dropping to 0.65-0.73 in middle lay-  
139 ers). This asymmetry reveals that **RL and**  
140 **instruction tuning primarily trans-**  
141 **forms how models process questions**  
142 **while leaving factual knowledge rep-**  
143 **resentations intact**, consistent with our  
144 hypothesis that RL enhances navigation  
145 mechanisms rather than knowledge con-  
146 tent.

147 We further conduct ablation studies com-

#### Research Questions

1. **RQ1: Does explicit prompt-  
ing close the performance gap?**  
If instruction-tuned models contain  
the required knowledge, can struc-  
tured prompts that explicitly in-  
struct hierarchical traversal match  
the performance of RL-enhanced  
models?
2. **RQ2: Do reasoning mod-  
els navigate deeper hierarchies  
better?** On tasks requiring multi-  
step hierarchical traversal, do rea-  
soning models demonstrate supe-  
rior path accuracy beyond what  
prompting achieves?
3. **RQ3: How do internal rep-  
resentations differ?** Do reason-  
ing models transform how they en-  
code queries, factual knowledge, or  
both?

148 comparing distilled R1 models to R1 and base  
149 models (Kim et al., 2025; Chen et al.,  
150 2025a), finding that distilled models cap-  
151 ture only surface-level improvements with-  
152 out acquiring robust navigation capabili-  
153 ties—achieving intermediate performance  
154 on complex retrieval tasks. Structured  
155 prompting provides minimal gains for dis-  
156 tilled models, and layer-wise analysis re-  
157 veals greater representational changes than  
158 instruction-tuned variants, yet without im-  
159 proved deep-retrieval navigation.

160 Our findings carry important implica-  
161 tions: RL-enhanced models succeed not  
162 through expanded knowledge but through  
163 improved cognitive scaffolding—the abil-  
164 ity to systematically traverse structures al-  
165 ready encoded during pretraining, which is  
166 inline with recent work showing that RL  
167 surfaces intelligence (Huan et al., 2025; Wu  
168 et al., 2025). While our experiments fo-  
169 cus on two datasets (MedConceptsQA and  
170 IPC) and specific model families (Qwen2.5,  
171 DeepSeek, Mistral), the patterns suggest  
172 more efficient training paradigms sepa-  
173 rating knowledge acquisition (pretraining)  
174 from organization (post-training). We en-  
175 courage future work to investigate these  
176 phenomena across broader domains and  
177 develop RL mechanisms that explicitly op-  
178 timize for hierarchical navigation.

## 2 Experimental Methodology

179 Our investigation into how reinforcement  
180 learning enhances hierarchical knowledge  
181 traversal is guided by three research ques-  
182 tions:  
183

184 We address these questions through three  
185 complementary experiments. Section 2.1

demonstrates that structured prompting can induce hierarchical reasoning in instruction-tuned models, reducing the performance gap by up to 68%. Section 2.2 introduces retrieval tasks of varying complexity with a path matching metric, revealing that reasoning models excel particularly on deep-retrieval tasks requiring extensive hierarchical navigation. Section 2.3 presents layer-wise activation analysis showing that while factual representations remain largely unchanged, query processing diverges substantially between SFT and RL models, supporting our hypothesis that RL primarily enhances navigation mechanisms rather than knowledge content.

## 2.1 Hierarchical Navigation Through Structured Prompting

We investigate tasks requiring pure information recall without multi-step computation or logical deduction, to determine whether the performance gap between base and reasoning models can be mitigated through prompting strategies alone. Remarkably, structured prompting reduces the performance gap for 671B base models such as DeepSeek-V3 from 23.7 pp to 7.5 pp (a 68% gap reduction), demonstrating the effectiveness of our method.

### Datasets

- **MedConceptsQA:** A multiple-choice question answering dataset focused on biomedical and clinical concepts. The questions are designed to test factual recall of medical terminology, concept definitions, and their relationships, without reasoning over patient cases or performing calculations.
- **International Patent Classification (IPC):** A dataset consists of queries mapped to patent classification codes. The task requires identifying the correct category for a given technical description, relying on recalling standardized knowledge of patent domains rather than multi-step reasoning.

### Prompting

- **Direct QA Prompting:** This baseline requires the model to provide only a single-letter answer to each multiple-choice question without any explanation.

- **Standard CoT Prompting:** This template requests both a final answer and a supporting explanation, aiming to capture the model’s intrinsic reasoning without imposing any procedural constraints.
- **Structured Prompting:** We introduce hierarchical instructions that enforce systematic reasoning. This strategy involves a two-stage process: (1) recall the hierarchical structural breakdown of the relevant medical code or concept, and (2) systematically evaluate each option with justification before elimination. This approach tests our hypothesis that enforcing structured knowledge recall and stepwise elimination can reduce performance gaps (see Appendix B.1 for complete prompt templates).

**Models** We evaluate a diverse set of large language models, focusing on comparisons between base, instruction-tuned, reasoning, and distilled models. The first group includes instruction-tuned models, such as the Qwen2.5 family (7B, 14B, 32B, and 72B parameters) (Team et al., 2024) and Mistral-Small-3.1-24B-Instruct (Karamcheti et al., 2021), each paired with their respective base models. The second group consists of reasoning models, including QwQ-32B (reasoning-enhanced Qwen2.5-32B), DeepSeek-R1 (from DeepSeek-V3), Magistral (from Mistral-Small-3.1-24B), and the reasoning model of Qwen3-235B-A22B (Liu et al., 2024; Guo et al., 2025a; Yang et al., 2025a). The third group includes models distilled from DeepSeek-R1: Qwen2.5-Math-7B, Qwen2.5-32B, and Llama3.3-70B, each compared against their pre-distillation ones (Grattafiori et al., 2024). We sample from all models using a temperature of 0.8 and top-p of 0.7 across three independent runs. Performance is reported as both mean accuracy ( $\pm$  standard deviation) and majority-voted accuracy, where majority voting selects the most frequent answer among the three runs for each question.

## 2.2 Hierarchical Navigation Across Retrieval Complexity

While we previously conclude that reasoning models use hierarchical navigation that can be externalized through structured prompting, a fundamental question

Table 1: Stratification of the “Nearest Common Ancestor” task by retrieval complexity, defined by the number of unique ancestor nodes (traversals) recalled to find the common node.

Task Complexity	Traversals	Figure Example	Example
Memory-Light	< 3		<p><b>Question:</b> Nearest common ancestor of H04B 1/7075 and H04B 1/7083 is: A) H04B 1/707 B) H04B 1/7073 C) H04B 1/69 D) H04B</p> <p><b>Hierarchical Paths:</b></p> <ul style="list-style-type: none"> <li>• H04B 1/7075 → H04B 1/7073</li> <li>• H04B 1/7083 → H04B 1/7073</li> </ul> <p><b>Answer:</b> B</p>
Memory-Heavy	5+		<p><b>Question:</b> Nearest common ancestor of A01B 3/421 and A01B 15/06 is: A) A01B 3/00 B) A01B 15/00 C) A01B D) A01</p> <p><b>Hierarchical Paths:</b></p> <ul style="list-style-type: none"> <li>• A01B 3/421 → A01B 3/42 → A01B 3/40 → A01B 3/36 → A01B 3/00 → A01B</li> <li>• A01B 15/06 → A01B 15/04 → A01B 15/02 → A01B 15/00 → A01B</li> </ul> <p><b>Answer:</b> C</p>

remains: *do reasoning models merely execute these strategies more consistently, or are there tasks that they execute fundamentally better?* To address this, we need to analyze not just whether models retrieve correct answers but how they traverse knowledge hierarchies to reach those answers. Therefore, we extend the original IPC dataset to stratify it by retrieval complexity and introduce a new metric to measure path traversal quality. Subsequent results reveal that reasoning models show superior hierarchical traversal—an ability that emerges on complex tasks requiring deeper knowledge navigation.

**IPC Multi-Level Retrieval Dataset** As shown in Table 1, this expanded dataset tests basic structural knowledge, including identifying common ancestors of a given pair of nodes. The questions are categorized by retrieval complexity, defined as the total number of ancestor nodes that must be recalled along both hierarchical paths (excluding the initial query nodes) to reach the nearest common ancestor. This stratification allows us to isolate the effect of retrieval depth on model performance.

- **Memory-Light (ML)** tasks require recalling < 3 ancestor nodes total across both paths to reach the common ancestor.
- **Memory-Heavy (MH)** tasks demand recalling  $\geq 5$  ancestor nodes across both paths.

**Path Matching Score** To evaluate the quality of predicted hierarchical paths for IPC codes, we propose the path matching score, which combines two metrics:

- **F1-Score:** Measures precision and recall of hierarchical ancestor identification, defined as  $F_1 = \frac{2 \times P \times R}{P + R}$ , where  $P$  and  $R$  denote precision and recall over the set of hierarchical ancestors (Buckland and Gey, 1994).
- **Common Subsequence Score (CSS):** Evaluates structural integrity of sequential paths via the ratio of the Longest Common Subsequence (LCS) (Paterson and Dančák, 1994) between the predicted and true paths to the length of the true path: 
$$CSS = \frac{|\text{LCS}(\text{predicted}, \text{ground truth})|}{|\text{ground truth ancestors}|}$$

The path matching score combines both components via harmonic mean:  $\text{Path Matching} = \frac{2 \times F_1 \times \text{CSS}}{F_1 + \text{CSS}}$ . This metric captures both structural accuracy and hierarchical coherence in patent classification navigation.

**Models** To analyze the impact of retrieval complexity, we conduct a case study using the DeepSeek-V3 and R1 pair on our expanded IPC dataset. While a broader evaluation would be ideal, we select the DeepSeek pair due to their instruction-following capabilities suitable for a reliable analysis.

### 2.3 Hierarchical Navigation in Internal Representations

To investigate whether base and specialized models<sup>1</sup> possess equivalent knowledge

<sup>1</sup>Here “base models” refer to the foundation model from which “specialized models” (instruction-tuned/reasoning/distilled) variants are derived. We adopt this terminology throughout the section to clearly distinguish the two categories.

Table 2: Performance comparison of Instruct vs. Reasoning models on MedConceptsQA and IPC datasets. The first column indicates the dataset. Models are evaluated across three prompt templates (QA, CoT, Structured). Metrics shown are majority voting accuracy (Maj. Vote Acc.) and mean accuracy (Mean Acc.). Mean accuracy is reported as Mean Acc. (Std.), with the standard deviation in subscripted parentheses. For each model pair, a  $\Delta$  row shows the gap from the reasoning model for both Maj. Acc. (red) and Mean Acc. (green). Bold values indicate the best performance within each model pair.  $\Delta$  values are highlighted, with darker shades indicating larger gaps.

Dataset	Model	Model Type	Maj. Vote Acc.			Mean Acc.(Std.)		
			QA	CoT	Structured	QA	CoT	Structured
MedConceptsQA	Qwen2.5-32B	Instruct	0.379	0.475	0.469	0.371 <sub>(.012)</sub>	0.449 <sub>(.010)</sub>	0.454 <sub>(.007)</sub>
		Reasoning	<b>0.482</b>	<b>0.513</b>	<b>0.505</b>	<b>0.470</b> <sub>(.012)</sub>	<b>0.487</b> <sub>(.009)</sub>	<b>0.481</b> <sub>(.005)</sub>
		$\Delta$	<b>+0.103</b>	<b>+0.038</b>	<b>+0.036</b>	<b>+0.099</b>	<b>+0.038</b>	<b>+0.027</b>
	Qwen3-235B-A22B	Instruct	0.542	0.548	<b>0.631</b>	0.503 <sub>(.004)</sub>	0.528 <sub>(.005)</sub>	<b>0.589</b> <sub>(.007)</sub>
		Reasoning	<b>0.641</b>	<b>0.656</b>	0.580	<b>0.599</b> <sub>(.003)</sub>	<b>0.617</b> <sub>(.003)</sub>	0.554 <sub>(.008)</sub>
		$\Delta$	<b>+0.099</b>	<b>+0.108</b>	<b>-0.051</b>	<b>+0.096</b>	<b>+0.089</b>	<b>-0.035</b>
DeepSeek-V3	Instruct	0.541	0.632	0.717	0.551 <sub>(.014)</sub>	0.636 <sub>(.049)</sub>	0.701 <sub>(.026)</sub>	
	Reasoning	<b>0.778</b>	<b>0.790</b>	<b>0.792</b>	<b>0.830</b> <sub>(.006)</sub>	<b>0.774</b> <sub>(.013)</sub>	<b>0.775</b> <sub>(.026)</sub>	
	$\Delta$	<b>+0.237</b>	<b>+0.158</b>	<b>+0.075</b>	<b>+0.279</b>	<b>+0.138</b>	<b>+0.074</b>	
IPC Codes	Qwen2.5-32B	Instruct	0.759	0.754	0.774	0.759 <sub>(.007)</sub>	0.754 <sub>(.000)</sub>	0.774 <sub>(.007)</sub>
		Reasoning	<b>0.777</b>	<b>0.875</b>	<b>0.790</b>	<b>0.713</b> <sub>(.015)</sub>	<b>0.754</b> <sub>(.070)</sub>	<b>0.769</b> <sub>(.033)</sub>
		$\Delta$	<b>+0.018</b>	<b>+0.121</b>	<b>+0.016</b>	<b>-0.046</b>	<b>+0.000</b>	<b>-0.005</b>
	Qwen3-235B-A22B	Instruct	0.800	<b>0.846</b>	0.846	0.800 <sub>(.013)</sub>	<b>0.846</b> <sub>(.013)</sub>	0.846 <sub>(.013)</sub>
		Reasoning	<b>0.908</b>	0.877	<b>0.893</b>	<b>0.846</b> <sub>(.013)</sub>	0.836 <sub>(.026)</sub>	<b>0.851</b> <sub>(.015)</sub>
		$\Delta$	<b>+0.108</b>	<b>+0.031</b>	<b>+0.047</b>	<b>+0.046</b>	<b>-0.010</b>	<b>+0.005</b>
DeepSeek-V3	Instruct	0.831	<b>0.923</b>	0.877	0.846 <sub>(.000)</sub>	<b>0.882</b> <sub>(.007)</sub>	0.872 <sub>(.007)</sub>	
	Reasoning	<b>0.923</b>	0.892	<b>0.923</b>	<b>0.913</b> <sub>(.019)</sub>	0.867 <sub>(.026)</sub>	<b>0.903</b> <sub>(.007)</sub>	
	$\Delta$	<b>+0.092</b>	<b>-0.031</b>	<b>+0.046</b>	<b>+0.067</b>	<b>-0.015</b>	<b>+0.031</b>	

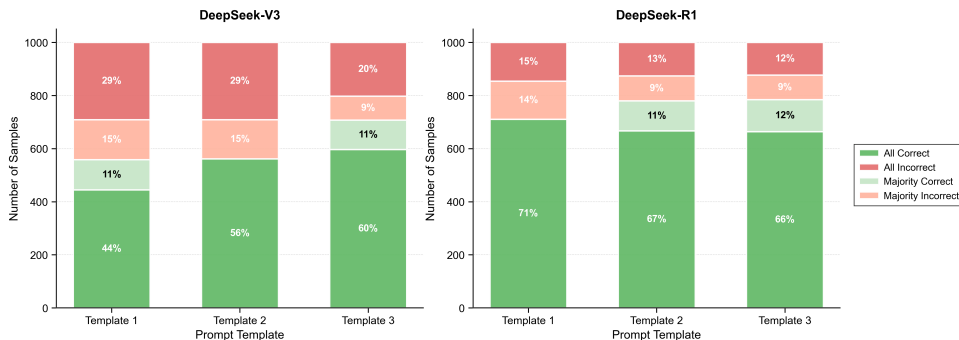


Figure 1: Comparative performance analysis of DeepSeek-V3 and DeepSeek-R1 across prompt strategies: direct question-answering (Template 1), chain-of-thought (Template 2), and structured prompting (Template 3) on MedConceptsQA dataset. Four categories are defined based on the number of correct votes across three independent runs: “All Incorrect” (0/3 correct), “Majority Incorrect” (1/3 correct), “Majority Correct” (2/3 correct), and “All Correct” (3/3 correct).

for hierarchical reasoning, we analyze their internal activations on MedConceptsQA using contrastive question-answer pairs. We conduct two complementary analyses: an *inter-model* comparison to show how enhancement modifies representations relative to the base model, and an *intra-model* comparison to trace how individual models transform questions into answers across layers. Our findings show that enhancement refines query processing while preserving factual knowledge.

**Probe Construction.** We construct probes from the MedConceptsQA dataset, which spans five medical vocabularies: ATC, ICD9CM, ICD10CM, ICD9PROC, and ICD10PROC. To ensure balanced

representation, we randomly sample 100 question-answer pairs from each vocabulary. Each probe consists of a factual question and its corresponding ground-truth answer, formatted as declarative statements. For example, a probe for medical code OQD2OZ from ICD10PROC takes the following form:

*Question:* What is the description of the medical code OQD2OZ in ICD10PROC?

*Answer:* The description of the medical code OQD2OZ in ICD10PROC is extraction of right pelvic bone, open approach.

We process questions and answers inde-

Table 3: Performance comparison of Base vs. Instruct models on MedConceptsQA and IPC datasets. The first column indicates the dataset. Models are evaluated across three prompt templates (QA, CoT, Structured). Metrics shown are majority voting accuracy (Maj. Vote Acc.) and mean accuracy (Mean Acc.). Mean accuracy is reported as Mean Acc.<sub>(Std.)</sub>, with the standard deviation in subscripted parentheses. For each model pair, an  $\Delta$  row shows the gap from the instruct model for both Maj. Acc. (red) and Mean Acc. (green). Bold values indicate the best performance within each model pair.  $\Delta$  values are highlighted, with darker shades indicating larger gaps. **This gap shrinks as we optimize the prompt, showing that the knowledge exists in the instruct model, it just needs to surface.**

Dataset	Model	Model Type	Maj. Vote Acc.			Mean Acc. <sub>(Std.)</sub>		
			QA	CoT	Structured	QA	CoT	Structured
MedConceptsQA	Qwen2.5-7B	Base	0.148	0.277	0.286	0.159 <sub>(.007)</sub>	0.239 <sub>(.036)</sub>	0.270 <sub>(.012)</sub>
		Instruct	<b>0.295</b>	<b>0.329</b>	<b>0.313</b>	<b>0.289<sub>(.006)</sub></b>	<b>0.316<sub>(.008)</sub></b>	<b>0.307<sub>(.015)</sub></b>
		$\Delta$	<b>+0.147</b>	<b>+0.052</b>	<b>+0.027</b>	<b>+0.130</b>	<b>+0.077</b>	<b>+0.037</b>
	Qwen2.5-14B	Base	0.335	0.332	0.386	0.316 <sub>(.015)</sub>	0.293 <sub>(.025)</sub>	0.372 <sub>(.007)</sub>
		Instruct	<b>0.395</b>	<b>0.420</b>	<b>0.420</b>	<b>0.385<sub>(.006)</sub></b>	<b>0.415<sub>(.007)</sub></b>	<b>0.409<sub>(.012)</sub></b>
		$\Delta$	<b>+0.060</b>	<b>+0.088</b>	<b>+0.034</b>	<b>+0.069</b>	<b>+0.122</b>	<b>+0.037</b>
	Qwen2.5-32B	Base	0.221	0.332	0.404	0.219 <sub>(.012)</sub>	0.260 <sub>(.071)</sub>	0.372 <sub>(.007)</sub>
		Instruct	<b>0.379</b>	<b>0.475</b>	<b>0.469</b>	<b>0.371<sub>(.012)</sub></b>	<b>0.449<sub>(.010)</sub></b>	<b>0.454<sub>(.007)</sub></b>
		$\Delta$	<b>+0.158</b>	<b>+0.143</b>	<b>+0.065</b>	<b>+0.152</b>	<b>+0.189</b>	<b>+0.082</b>
	Qwen2.5-72B	Base	0.443	0.351	0.468	0.389 <sub>(.005)</sub>	0.305 <sub>(.028)</sub>	0.418 <sub>(.008)</sub>
		Instruct	<b>0.546</b>	<b>0.520</b>	<b>0.546</b>	<b>0.519<sub>(.007)</sub></b>	<b>0.512<sub>(.005)</sub></b>	<b>0.537<sub>(.008)</sub></b>
		$\Delta$	<b>+0.103</b>	<b>+0.169</b>	<b>+0.078</b>	<b>+0.130</b>	<b>+0.207</b>	<b>+0.119</b>
Qwen2.5-7B	Base	0.463	0.436	<b>0.588</b>	0.349 <sub>(.040)</sub>	0.364 <sub>(.038)</sub>	<b>0.585<sub>(.038)</sub></b>	
	Instruct	<b>0.615</b>	<b>0.554</b>	0.574	<b>0.615<sub>(.025)</sub></b>	<b>0.554<sub>(.013)</sub></b>	0.574 <sub>(.015)</sub>	
	$\Delta$	<b>+0.152</b>	<b>+0.118</b>	<b>-0.014</b>	<b>+0.266</b>	<b>+0.190</b>	<b>-0.011</b>	
Qwen2.5-14B	Base	0.526	0.608	0.609	0.421 <sub>(.038)</sub>	0.492 <sub>(.033)</sub>	0.600 <sub>(.013)</sub>	
	Instruct	<b>0.708</b>	<b>0.691</b>	<b>0.718</b>	<b>0.708<sub>(.025)</sub></b>	<b>0.687<sub>(.029)</sub></b>	<b>0.718<sub>(.007)</sub></b>	
	$\Delta$	<b>+0.182</b>	<b>+0.083</b>	<b>+0.109</b>	<b>+0.287</b>	<b>+0.195</b>	<b>+0.118</b>	
IPC Codes	Qwen2.5-32B	Base	0.644	0.641	<b>0.777</b>	0.482 <sub>(.059)</sub>	0.503 <sub>(.038)</sub>	0.769 <sub>(.013)</sub>
	Instruct	<b>0.759</b>	<b>0.754</b>	0.774	<b>0.759<sub>(.007)</sub></b>	<b>0.754<sub>(.000)</sub></b>	<b>0.774<sub>(.007)</sub></b>	
	$\Delta$	<b>+0.115</b>	<b>+0.113</b>	<b>-0.003</b>	<b>+0.277</b>	<b>+0.251</b>	<b>+0.005</b>	

Table 4: Performance of distilled models compared to the **DeepSeek-R1 (reasoning model)**. Each cell for a distilled model shows its absolute score, followed in parentheses by the  $\Delta$  gap (reasoning - distilled).  $\Delta$  values for Maj. Vote Acc. are shaded red, and  $\Delta$  values for Mean Acc. are shaded green. Darker shades indicate a larger performance gap. All  $\Delta$  values are positive, showing the gap to the stronger R1 model.

Dataset	Model	Maj. Vote Acc. ( $\Delta$ vs. R1)			Mean Acc. <sub>(Std.)</sub> ( $\Delta$ vs. R1)		
		QA	CoT	Structured	QA	CoT	Structured
	<b>DeepSeek-R1 (Reasoning)</b>	<b>0.778</b>	<b>0.790</b>	<b>0.792</b>	<b>0.830<sub>(.006)</sub></b>	<b>0.774<sub>(.013)</sub></b>	<b>0.775<sub>(.026)</sub></b>
MedConceptsQA	Qwen2.5-Math-7B (Dist.)	0.296 (+0.482)	0.256 (+0.534)	0.282 (+0.510)	0.292 <sub>(.010)</sub> (+0.538)	0.250 <sub>(.017)</sub> (+0.524)	0.289 <sub>(.017)</sub> (+0.486)
	Qwen2.5-32B (Dist.)	0.375 (+0.403)	0.380 (+0.410)	0.447 (+0.345)	0.351 <sub>(.009)</sub> (+0.479)	0.369 <sub>(.005)</sub> (+0.405)	0.420 <sub>(.002)</sub> (+0.355)
	Llama3.3-70B (Dist.)	0.537 (+0.241)	0.633 (+0.157)	0.610 (+0.182)	0.495 <sub>(.002)</sub> (+0.335)	0.609 <sub>(.011)</sub> (+0.165)	0.596 <sub>(.012)</sub> (+0.179)
	<b>DeepSeek-R1 (Reasoning)</b>	<b>0.923</b>	<b>0.892</b>	<b>0.923</b>	<b>0.913<sub>(.019)</sub></b>	<b>0.867<sub>(.026)</sub></b>	<b>0.903<sub>(.007)</sub></b>
IPC Codes	Qwen2.5-32B (Dist.)	0.778 (+0.145)	0.730 (+0.162)	0.788 (+0.135)	0.754 <sub>(.038)</sub> (+0.159)	0.667 <sub>(.019)</sub> (+0.200)	0.780 <sub>(.019)</sub> (+0.123)
	Llama3.3-70B (Dist.)	0.785 (+0.138)	0.831 (+0.061)	0.815 (+0.108)	0.785 <sub>(.015)</sub> (+0.128)	0.785 <sub>(.041)</sub> (+0.082)	0.790 <sub>(.018)</sub> (+0.113)

pendently through each model to extract their respective layer-wise representations, enabling both inter-model and intra-model comparative analyses.

**Representation Extraction.** For a model with  $L$  layers and hidden dimension  $d$ , we extract the hidden state at the final token position for each layer  $\ell \in \{1, \dots, L\}$  as the layer’s representation vector  $\mathbf{h}_\ell \in \mathbb{R}^d$ . This representation attends to all preceding tokens, thereby capturing the full input context at that layer.

**Representation Analysis.** We quantify representational differences across and within models using *inter-model* and *intra-model* analyses:

- **Inter-Model Analysis.** By com-

paring the question-question (Q-Q) and answer-answer (A-A) representations between the base and specialized models, we assess how they differ at understanding query and retrieving factual knowledge.

- **Intra-Model Comparison.** This analysis investigates the internal transformation of information within a single model. By comparing a model’s question and answer representations layer by layer, we trace how internal activations evolve from encoding a problem to producing a solution.

**Comparison Metric** For each layer  $\ell \in \{1, \dots, L\}$ , we use cosine similarity, a mea-

Table 5: Comparison of structured prompting performance by task complexity for DeepSeek-R1 and DeepSeek-V3 models. Memory-Light tasks (1-2 hierarchical recalls); Memory-Heavy tasks (5+ hierarchical recalls). Bold values indicate the best performance for each metric within each complexity category. **As we move to retrieval heavier tasks with structure, the gap between path matching score of R1 and V3 increases.**

Task Complexity	Model	Accuracy (%)	Path Matching Score
Memory-Light	DeepSeek-R1	<b>44.8</b>	<b>0.681</b>
	DeepSeek-V3	37.9	0.627
Memory-Heavy	DeepSeek-R1	<b>67.7</b>	<b>0.597</b>
	DeepSeek-V3	<b>67.7</b>	0.503

sure of directional alignment, to define representation similarity:

$$d_{\cos}^{(a,b)}(\ell) = 1 - \frac{1}{N} \sum_{i=1}^N \frac{\mathbf{h}_{\ell}^{(a)}(i) \top \mathbf{h}_{\ell}^{(b)}(i)}{\|\mathbf{h}_{\ell}^{(a)}(i)\|_2 \|\mathbf{h}_{\ell}^{(b)}(i)\|_2}, \quad (1)$$

Here,  $\mathbf{h}_{\ell}^{(s)}(i)$  denotes the layer- $\ell$  hidden representation for probe  $i$  from a source  $s$ . The set of sources  $\mathcal{S} = \{Q^{\text{base}}, A^{\text{base}}, Q^{\text{specialized}}, A^{\text{specialized}}\}$  includes representations for both the question (Q) and answer (A) components from the base and specialized models. Pair  $(a, b)$  represents either inter-model (e.g.,  $Q^{\text{base}}$  vs  $Q^{\text{specialized}}$ ,  $A^{\text{base}}$  vs  $A^{\text{specialized}}$ ) or intra-model comparisons (e.g.,  $Q^{\text{base}}$  vs  $A^{\text{base}}$ ,  $Q^{\text{specialized}}$  vs  $A^{\text{specialized}}$ ). Results are reported per vocabulary using  $N = 100$  probes.

**Models** We compare Qwen2.5-32B (base) against three specialized variants: Qwen2.5-32B-Instruct (instruction-tuned), DeepSeek-R1-Distill-Qwen-32B (distilled), and QwQ-32B (reasoning). We select this 32B parameter family because it spans multiple enhancement methods while remaining computationally tractable for single-GPU inference. A supplementary analysis comparing variants of the Mistral-Small-24B family (base, instruct, and reasoning) is included in Appendix C.

## 3 Experimental Results

### 3.1 Hierarchical Navigation Through Structured Prompting

Hierarchical navigation and stepwise elimination strategies systematically narrow the accuracy gap between base models and

their reasoning-enhanced, or instruction-tuned versions across both MedConceptsQA and IPC code datasets. For example, on MedConceptsQA, structured prompting allows the Qwen3-235B Instruct model (Table 2) to outperform its reasoning counterpart, reversing a +0.108 majority vote accuracy gap (CoT) to a -0.051 advantage. Similarly, on the IPC dataset (Table 3), this prompting reduces the gap between the Qwen2.5-32B base and instruct models from +0.115 (QA) to -0.003. However, this effect is less pronounced for distilled models (Table 4), where the performance gap relative to the reasoning model remains substantial, even with structured prompts (e.g., Llama3.3-70B on MedConceptsQA, +0.182 gap).

To understand the mechanisms underlying structured prompting’s effectiveness, we examine response consistency patterns. Figure 1 presents results for DeepSeek-V3 and R1 across three independent runs under majority voting on MedConceptsQA. When transitioning from the baseline to the structured prompt, DeepSeek-V3 shows significant sample migration: questions initially categorized as “All Incorrect” and “Majority Incorrect” shift toward “Majority Correct” and “All Correct”. In contrast, R1 exhibits static distribution across these categories, suggesting it already operates near its ceiling. This redistribution in V3 indicates that explicit structural guidance improves the consistency of the model’s internal reasoning and that its underlying knowledge is sufficient. Therefore, the primary role of specialized post-training is not to introduce entirely novel knowledge, but rather to enhance the procedural consistency and strategic reasoning of existing knowledge structures.

### 3.2 Hierarchical Navigation Across Retrieval Complexity

Stratifying performance by retrieval complexity highlights a distinction between the base and reasoning models. Despite similar overall accuracy, R1 consistently achieves a higher path matching score, particularly on complex tasks such as common ancestor identification, suggesting it can correctly navigate the hierarchy step-by-step (Table 5). This is a deeper form of understanding that goes beyond simple memorization. Ultimately, R1 understands the process of navigating a knowledge hierarchy better

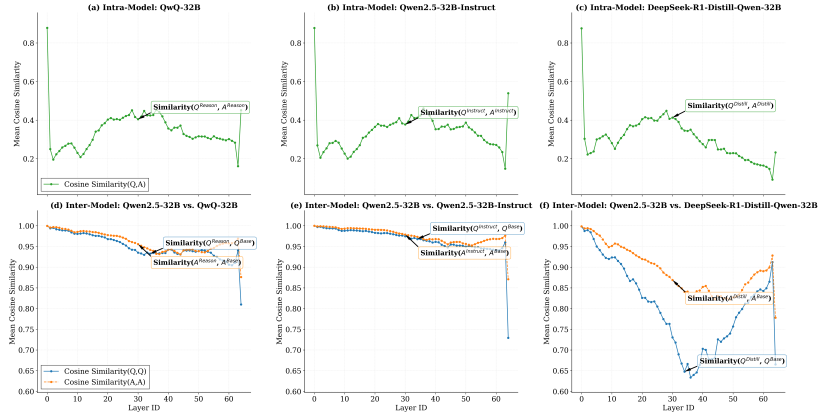


Figure 2: Layerwise Representation Similarity for ICD9PROC Vocabulary from MedConceptsQA. Plots compare last-token hidden state representations across layers (x-axis) using cosine similarity. Top Row (Intra-Model): Question vs. Answer representation similarity within QwQ-32B, Qwen2.5-32B-Instruct, and DeepSeek-R1-Distill-Qwen-32B. Bottom Row (Inter-Model): Similarity between the base model (Qwen2.5-32B) and each respective advanced model, comparing Question representations ( $Q_{Reason}$  vs.  $Q_{Base}$ ) and Answer representations ( $A_{Reason}$  vs.  $A_{Base}$ ) separately. **The representations of questions diverge more, specially in the last layer, compared to the answers. This hints at the knowledge being encoded similarly in base and reasoning models, but navigated differently.**

524 than the base model (V3), even when their  
525 final-answer accuracy is similar.

### 526 3.3 Hierarchical Navigation in 527 Internal Representations

528 **Intra-Model Representational Simi-**  
529 **larity.** Within each model, representa-  
530 tions for questions and answers are initially  
531 highly similar, but this similarity decreases  
532 in later layers, suggesting that the repre-  
533 sentations accumulate increasingly distinct  
534 features.

535 **Inter-Model Representational Simi-**  
536 **larity.** Instruction-tuned and reasoning  
537 models show strong directional alignment  
538 with the base model for both question  
539 and answer representations, whereas the  
540 distilled model shows much greater di-  
541 vergence (Figure 2(d-f)). Notably, ques-  
542 tion representations diverge more than an-  
543 swer representations across all specialized  
544 models, suggesting that performance gains  
545 arise primarily from refining question un-  
546 derstanding rather than reorganizing fac-  
547 tual knowledge.

## 548 4 Conclusion

549 This work challenges the view that rein-  
550 forcement learning enhances reasoning at  
551 the expense of memorization. We demon-  
552 strate that RL-enhanced models outper-  
553 form base counterparts by 24pp on hier-  
554 archical knowledge tasks, not through ac-

quiring new knowledge, but by improving  
555 navigation of existing structures. Struc-  
556 tured prompting reduces this gap to 7pp on  
557 simple tasks, yet reasoning models main-  
558 tain superior path traversal on complex  
559 deep-retrieval tasks (5pp to 9pp gap widen-  
560 ing). Layer-wise analysis reveals that  
561 RL transforms query processing (simi-  
562 larity drops to 0.65-0.73) while preserv-  
563 ing factual representations (0.85-0.92), con-  
564 firming that improvements stem from en-  
565 hanced navigation mechanisms rather than  
566 knowledge content changes. 567

Several open questions warrant investiga-  
568 tion. First, do similar navigation mecha-  
569 nisms underlie RL improvements on other  
570 structured reasoning tasks such as mathe-  
571 matical proof generation, code debugging,  
572 or multi-hop question answering? 573  
574 Second, can we develop RL objectives that  
575 explicitly optimize for hierarchical naviga-  
576 tion rather than relying on implicit emer-  
577 gence? 578  
579 Third, how do these findings ex-  
580 tend to knowledge domains with different  
581 structural properties—flat versus deeply  
582 nested hierarchies, dense versus sparse  
583 connectivity? 584  
585 Finally, can we design hy-  
586 brid approaches that combine the efficien-  
587 cy of structured prompting with the robust-  
588 ness of RL-trained navigation for practical  
589 deployment? Addressing these questions  
will deepen our understanding of how lan-  
guage models organize and access paramet-  
ric knowledge, ultimately enabling more  
capable and efficient reasoning systems. 589

590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

## Limitations

Limitations A key limitation of this work is that the experiments are restricted to a small set of structured knowledge domains, which may limit the generalizability of the findings to other tasks or less hierarchical settings.

There is a risk that structured prompting and path-based metrics overestimate true understanding, as models may learn superficial traversal patterns rather than robust, transferable navigation strategies.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, and 257 others. 2023. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.

Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos Garea, Matthieu Geist, and Olivier Bachem. 2024. On-policy distillation of language models: Learning from self-generated mistakes. In *The twelfth international conference on learning representations*.

Lakshya A Agrawal, Shangyin Tan, Dilara Soyly, Noah Ziem, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. 2024. [GEPA: Reflective prompt evolution can outperform reinforcement learning](#). *Preprint*, arXiv:2507.19457. Note: Citation key was xu2024gepa but first author is Agrawal.

Lakshya A Agrawal, Shangyin Tan, Dilara Soyly, Noah Ziem, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, and 1 others. 2025. [Gepa: Reflective prompt evolution can outperform reinforcement learning](#). *arXiv preprint arXiv:2507.19457*.

Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom

Brown, Jack Clark, and 3 others. 2021. [A general language assistant as a laboratory for alignment](#). *Preprint*, arXiv:2112.00861.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, and 32 others. 2022. [Constitutional AI: Harmlessness from AI feedback](#). *Preprint*, arXiv:2212.08073.

Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. 2024. The reversal curse: LLMs trained on "A is B" fail to learn "B is A". *arXiv preprint arXiv:2309.12288*.

Michael Buckland and Fredric Gey. 1994. The relationship between recall and precision. *Journal of the American society for information science*, 45(1):12–19.

Hardy Chen, Haoqin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. 2025a. Sft or rl? an early investigation into training rl-like reasoning large vision-language models. *arXiv preprint arXiv:2504.11468*.

Mingyang Chen, Linzhuang Sun, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Haofen Wang, Jeff Z Pan, Wen Zhang, Huajun Chen, and 1 others. 2025b. Learning to reason with search for llms via reinforcement learning. *arXiv preprint arXiv:2503.19470*.

Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. 2025. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*.

Matteo Gabburo, Nicolaas Paul Jedema, Siddhant Garg, Leonardo FR Ribeiro, and Alessandro Moschitti. 2024. Measuring retrieval complexity in question answering systems. *arXiv preprint arXiv:2406.03592*.

Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. Does fine-tuning llms on new knowledge encourage hallucinations? *arXiv preprint arXiv:2405.05904*.

Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Deepali Aneja, Zeyu Jin, Ramani Duraiswami, Dinesh Manocha, and 1 others. 2024. A closer look at the limitations of instruction tuning. *arXiv preprint arXiv:2402.05119*.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek

710	Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	
711		
712		
713		
714		
715	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	
716		
717		
718		
719		
720		
721		
722	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025b. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	
723		
724		
725		
726		
727		
728		
729	Maggie Huan, Yuetai Li, Tuney Zheng, Xiaoyu Xu, Seungone Kim, Minxin Du, Radha Poovendran, Graham Neubig, and Xiang Yue. 2025. Does math reasoning improve general llm capabilities? understanding transferability of llm reasoning. <i>arXiv preprint arXiv:2507.00432</i> .	
730		
731		
732		
733		
734		
735		
736	Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. <i>arXiv preprint arXiv:2503.09516</i> .	
737		
738		
739		
740		
741		
742		
743	Siddharth Karamcheti, Laurel Orr, Jason Bolton, Tianyi Zhang, Karan Goel, Avanika Narayan, Rishi Bommasani, Deepak Narayanan, Tatsunori Hashimoto, Dan Jurafsky, and 1 others. 2021. Mistral—a journey towards reproducible language model training.	
744		
745		
746		
747		
748		
749		
750	Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, and 1 others. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. <i>arXiv preprint arXiv:2310.03714</i> .	
751		
752		
753		
754		
755		
756		
757		
758	Minwu Kim, Anubhav Shrestha, Safal Shrestha, Aadim Nepal, and Keith Ross. 2025. Reinforcement learning vs. distillation: Understanding accuracy and capability in llm reasoning. <i>arXiv preprint arXiv:2505.14216</i> .	
759		
760		
761		
762		
763		
764	Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. 2023. Understanding the effects of rlhf on llm generalisation and diversity. <i>arXiv preprint arXiv:2310.06452</i> .	
765		
766		
767		
768		
769		
770		
771	Junliang Li, Yucheng Wang, Yan Chen, Yu Ran, Ruiqing Zhang, Jing Liu, Hua	
772		
	Wu, and Haifeng Wang. 2025. Knowledge-level consistency reinforcement learning: Dual-fact alignment for long-form factuality. <i>arXiv preprint arXiv:2509.23765</i> .	773
		774
		775
		776
	Yusheng Liao, Chaoyi Wu, Junwei Liu, Shuyang Jiang, Pengcheng Qiu, Haowen Wang, Yun Yue, Shuai Zhen, Jian Wang, Qianrui Fan, and 1 others. 2025. Ehr-r1: A reasoning-enhanced foundational language model for electronic health record analysis. <i>arXiv preprint arXiv:2510.25628</i> .	777
		778
		779
		780
		781
		782
		783
	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. <i>Preprint, arXiv:2305.20050</i> .	784
		785
		786
		787
		788
		789
	Yong Lin, Hangyu Lin, Wei Xiong, Shizhe Diao, Jianmeng Liu, Jipeng Zhang, Rui Pan, Haoxiang Wang, Wenbin Hu, Hanqing Zhang, Hanze Dong, Renjie Pi, Han Zhao, Nan Jiang, Heng Ji, Yuan Yao, and Tong Zhang. 2024. Mitigating the alignment tax of RLHF. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	790
		791
		792
		793
		794
		795
		796
		797
		798
		799
	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024. Deepseek-v3 technical report. <i>arXiv preprint arXiv:2412.19437</i> .	800
		801
		802
		803
		804
		805
	Leibo Liu, Oscar Perez-Concha, Anthony Nguyen, Vicki Bennett, and Louisa Jorm. 2022. Hierarchical label-wise attention transformer model for explainable icd coding. <i>Journal of biomedical informatics</i> , 133:104161.	806
		807
		808
		809
		810
		811
	Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. <i>arXiv preprint arXiv:2305.14251</i> .	812
		813
		814
		815
		816
		817
		818
	Sagnik Mukherjee, Lifan Yuan, Dilek Hakkani-Tür, and Hao Peng. 2025. Reinforcement learning finetunes small subnetworks in large language models. In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	819
		820
		821
		822
		823
		824
	Oded Ovadia, Menachem Brief, Moshik Mishaeli, and Oren Elisha. 2024. Fine-tuning or retrieval? comparing knowledge injection in LLMs. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 237–250.	825
		826
		827
		828
		829
		830
		831
	Mike Paterson and Vlado Dančik. 1994. Longest common subsequences. In <i>International symposium on mathematical foun-</i>	832
		833
		834

835	<i>dations of computer science</i> , pages 127–				
836	142. Springer.				
837	Hoang Phan, Xianjun Yang, Kevin Yao,				
838	Jingyu Zhang, Shengjie Bi, Xiaocheng				
839	Tang, Madian Khabisa, Lijuan Liu, and				
840	Deren Lei. 2025. Beyond reasoning gains:				
841	Mitigating general capabilities forgetting				
842	in large reasoning models. <i>arXiv preprint</i>				
843	<i>arXiv:2510.21978</i> .				
844	Laura Ruis, Maximilian Mozes, Juhan				
845	Bae, Siddhartha Rao Kamalakara, Dwarak				
846	Talupuru, Acyr Locatelli, Robert Kirk,				
847	Tim Rocktäschel, Edward Grefenstette,				
848	and Max Bartolo. 2024. Procedural knowl-				
849	edge in pretraining drives reasoning in				
850	large language models. <i>arXiv preprint</i>				
851	<i>arXiv:2411.12580</i> .				
852	Cansu Sen, Bingyang Ye, Javed Aslam,				
853	and Amir Tahmasebi. 2021. From extreme				
854	multi-label to multi-class: A hierarchical				
855	approach for automated icd-10 coding using				
856	phrase-level attention. <i>arXiv preprint</i>				
857	<i>arXiv:2102.09136</i> .				
858	Noah Shinn, Federico Cassano, Ash-				
859	win Gopinath, Karthik Narasimhan, and				
860	Shunyu Yao. 2023. <a href="#">Reflexion: Language</a>				
861	<a href="#">agents with verbal reinforcement learning</a> .				
862	In <i>Advances in Neural Information Pro-</i>				
863	<i>cessing Systems 36 (NeurIPS 2023)</i> . Cur-				
864	ran Associates Inc. ArXiv:2303.11366.				
865	Oscar Skean, Md Rifat Arefin, Dan Zhao,				
866	Niket Patel, Jalal Naghiyev, Yann Le-				
867	Cun, and Ravid Shwartz-Ziv. 2025. Layer				
868	by layer: Uncovering hidden representa-				
869	tions in language models. <i>arXiv preprint</i>				
870	<i>arXiv:2502.02013</i> .				
871	Taylor Sorensen, Benjamin Newman,				
872	Jared Moore, Chan Park, Jillian Fisher,				
873	Niloofer Miresghallah, Liwei Jiang, and				
874	Yejin Choi. 2025. Spectrum tuning: Post-				
875	training for distributional coverage and				
876	in-context steerability. <i>arXiv preprint</i>				
877	<i>arXiv:2510.06084</i> .				
878	Qwen Team and 1 others. 2024.				
879	Qwen2 technical report. <i>arXiv preprint</i>				
880	<i>arXiv:2407.10671</i> , 2:3.				
881	Guan Wang, Jin Li, Yuhao Sun, Xing				
882	Chen, Changling Liu, Yue Wu, Meng				
883	Lu, Sen Song, and Yasin Abbasi Yad-				
884	kori. 2025a. Hierarchical reasoning model.				
885	<i>arXiv preprint arXiv:2506.21734</i> .				
886	Haozhe Wang, Qixin Xu, Che Liu, Jun-				
887	hong Wu, Fangzhen Lin, and Wenhui Chen.				
888	2025b. Emergent hierarchical reasoning in				
889	llms through reinforcement learning. <i>arXiv</i>				
890	<i>preprint arXiv:2509.03646</i> .				
891	Zengzhi Wang, Fan Zhou, Xuefeng Li,				
892	and Pengfei Liu. 2025c. Octothinker:				
893	Mid-training incentivizes reinforcement				
894	learning scaling. <i>arXiv preprint</i>				
895	<i>arXiv:2506.20512</i> .				
896	Jason Wei, Xuezhi Wang, Dale Schuur-				
897	mans, Maarten Bosma, Brian Ichter, Fei				
	Xia, Ed Chi, Quoc Le, and Denny Zhou.				898
	2022. <a href="#">Chain-of-thought prompting elicits</a>				899
	<a href="#">reasoning in large language models</a> . In <i>Ad-</i>				900
	<i>vances in Neural Information Processing</i>				901
	<i>Systems</i> , volume 35, pages 24824–24837.				902
	Curran Associates, Inc.				903
	Fang Wu, Weihao Xuan, Ximing Lu,				904
	Mingjie Liu, Yi Dong, Zaid Harchaoui,				905
	and Yejin Choi. 2025. The invisible leash:				906
	Why rlvr may or may not escape its origin.				907
	<i>arXiv preprint arXiv:2507.14843</i> .				908
	An Yang, Anfeng Li, Baosong Yang,				909
	Beichen Zhang, Binyuan Hui, Bo Zheng,				910
	Bowen Yu, Chang Gao, Chengen Huang,				911
	Chenxu Lv, and 1 others. 2025a.				912
	Qwen3 technical report. <i>arXiv preprint</i>				913
	<i>arXiv:2505.09388</i> .				914
	Ling Yang, Zhaochen Yu, Bin Cui, and				915
	Mengdi Wang. 2025b. <a href="#">Reasonflux: Hier-</a>				916
	<a href="#">archical llm reasoning via scaling thought</a>				917
	<a href="#">templates</a> . <i>ArXiv</i> , abs/2502.06772.				918
	Yufan Ye, Ting Zhang, Wenbin Jiang, and				919
	Hua Huang. 2025. Process-supervised re-				920
	inforcement learning for code generation.				921
	<i>arXiv preprint arXiv:2502.01715</i> .				922
	Qiyang Yu, Zheng Zhang, Ruofei Zhu,				923
	Yufeng Yuan, Xiaochen Zuo, Yu Yue,				924
	Weinan Dai, Tiantian Fan, Gaohong Liu,				925
	Lingjun Liu, and 1 others. 2025. Dapo:				926
	An open-source llm reinforcement learn-				927
	ing system at scale. <i>arXiv preprint</i>				928
	<i>arXiv:2503.14476</i> .				929
	Jiaqing Yuan, Lin Pan, Chung-Wei Hang,				930
	Jiang Guo, Jiarong Jiang, Bonan Min,				931
	Patrick Ng, and Zhiguo Wang. 2024. To-				932
	wards a holistic evaluation of llms on				933
	factual knowledge recall. <i>arXiv preprint</i>				934
	<i>arXiv:2404.16164</i> .				935
	Eric Zelikman, Georges Harik, Yijia Shao,				936
	Varuna Jayasiri, Nick Haber, and Noah D.				937
	Goodman. 2024. <a href="#">Quiet-STaR: Language</a>				938
	<a href="#">models can teach themselves to think be-</a>				939
	<a href="#">fore speaking</a> . <i>Preprint</i> , arXiv:2403.09629.				940
	Denny Zhou, Nathanael Schärli, Le Hou,				941
	Jason Wei, Nathan Scales, Xuezhi Wang,				942
	Dale Schuurmans, Claire Cui, Olivier				943
	Bousquet, Quoc Le, and Ed Chi. 2023.				944
	<a href="#">Least-to-most prompting enables complex</a>				945
	<a href="#">reasoning in large language models</a> . In				946
	<i>The Eleventh International Conference on</i>				947
	<i>Learning Representations</i> .				948
	Noah Ziemis, Dilara Soylu, Lakshya A				949
	Agrawal, Isaac Miller, Liheng Lai, Chen				950
	Qian, Kaiqiang Song, Meng Jiang, Dan				951
	Klein, Matei Zaharia, and 1 others. 2025.				952
	Multi-module grpo: Composing policy gra-				953
	dients and prompt optimization for lan-				954
	guage model programs. <i>arXiv preprint</i>				955
	<i>arXiv:2508.04660</i> .				956

**Q: What is medical code 57.95 in ICD-9-CM?**

A) Other procedures for relief of elevated intraocular pressure **B) Replacement of indwelling urinary catheter** C) Arthrodesis of other joint D) Biopsy of heart

**Non-reasoning model (V3):  
Direct Information Recall**

"ICD-9-CM code 57.95 refers to  
**Other procedures for relief of  
elevated intraocular pressure**, hence  
A is the answer".

Answer: A ❌



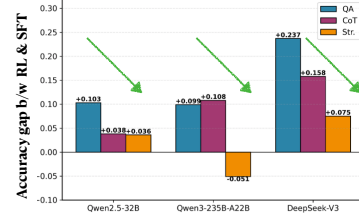
**Reasoning model (R1): Hierarchical Navigation**

1. **Define Problem Category:** "57.95 is ICD-9-CM Volume 5, Chapter 11, codes 57.0-57.99."  
2. **Identifies Main Procedure:** "This is a subset of 57.0 ... procedures relating to clearance of bladder ..."  
3. **Eliminate other options:** "Other procedures' in option A is not applicable ..."

... So 57.95 must be Replacement of unary catheter ...

Answer: B ✅

MedConceptsQA



Although RLed models show superior accuracy, structured prompting nearly closes this gap without training hinting that the **knowledge exists in SFT models but needs navigation**.

Figure 3: (Left) Overview of our main observation: When querying structured medical codes, non-reasoning models (DeepSeek-V3) rely on direct memorization attempts, often selecting incorrect answers (here choosing A). In contrast, reasoning-enhanced RL models (DeepSeek-R1) employ systematic hierarchical navigation—first categorizing the problem domain, then identifying relevant procedures, and finally interpreting ambiguous terminology—to successfully retrieve the correct answer (B). (Right) Reasoning models consistently outperform their instruction-tuned counterparts when prompted with conventional QA templates. This gap decreases when we optimize the prompt and is minimized with our hand-crafted structured prompt, hinting that the necessary knowledge exists in the instruct models.

## A Related Work

### A.1 The Alignment Tax and Factual Degradation

The trade-off between alignment and factual accuracy has been extensively explored. Lin et al. (2024) introduced the concept of “alignment tax”, showing systematic degradation on factual benchmarks as RLHF reward strength increases. Achiam et al. (2023) similarly reported that RLHF “does not improve exam performance (without active effort, it actually degrades it)” and can reduce calibration. Mechanistic analyses in Ghosh et al. (2024) reveal that instruction tuning primarily adjusts style rather than new knowledge, with responses generated from pre-trained knowledge consistently outperforming those from models learning new knowledge through instruction tuning. Both Li et al. (2025) and Kirk et al. (2023) show that base models’ parametric knowledge originates from pre-training while aligned models learn how to express it—training directly from base models mitigates knowledge forgetting and alignment tax incurred by SFT-based distillation. Recent work by Phan et al. (2025) reveals that optimizing for narrow verifiable rewards in reasoning-focused RL leads to regression in general capabilities, with models exhibiting increased hallucinations despite improved reasoning.

While these studies document factual

degradation from alignment, our work reveals a contrasting phenomenon: RL-enhanced models *outperform* their base counterparts on structured knowledge recall tasks. This apparent contradiction suggests that alignment tax may not uniformly affect all forms of parametric knowledge retrieval—particularly when retrieval demands systematic navigation through hierarchical structures rather than direct factual recall.

### A.2 Reasoning Enhancement Through RL

RL is commonly viewed as a means of amplifying reasoning ability. Process supervision and reward-driven methods (Lightman et al., 2023; Ye et al., 2025) demonstrate clear improvements on reasoning tasks, with process-supervised models solving substantially more problems than outcome-supervised variants. However, recent work hints at a more nuanced picture. Zelikman et al. (2024) introduce Quiet-STaR, showing that training models to generate internal rationales improves downstream reasoning by teaching systematic exploration of solution spaces—essentially navigation skills that achieve zero-shot improvements from 5.9% to 10.9% on GSM8K. Shinn et al. (2023) demonstrate that reinforcement learning primarily helps models learn from feedback to refine their search through problem spaces, rather than acquiring new problem-solving rules. Most strikingly, Guo et al.

(2025b) show that DeepSeek-R1 develops self-reflection, verification, and dynamic strategy adaptation through RL alone, without human-labeled reasoning trajectories, increasing pass@1 scores on AIME 2024 from 15.6% to 71.0%. Recent work on search-augmented reasoning (Jin et al., 2025; Chen et al., 2025b) demonstrates models learning to autonomously generate search queries and self-correct, with behaviors like pausing when detecting knowledge gaps emerging naturally.

These findings align with our hypothesis that RL enhances navigation of existing knowledge structures. However, while prior work focuses on mathematical and algorithmic reasoning, we examine whether these navigation improvements extend to retrieval from structured factual hierarchies, providing complementary evidence that RL’s benefits stem from improved access patterns rather than new knowledge acquisition.

### A.3 Hierarchical Reasoning and Structured Navigation

Hierarchical reasoning frameworks further support our knowledge navigation hypothesis. Wang et al. (2025a) present the Hierarchical Reasoning Model (HRM), a brain-inspired recurrent architecture that achieves near-perfect performance on complex tasks with only 27 million parameters trained on 1000 samples, without pre-training or chain-of-thought data. HRM’s architecture features interdependent modules for high-level abstract planning and low-level detailed computation, achieving 40.3% on ARC-AGI—precisely the type of structured traversal we hypothesize enables medical code lookup. Yang et al. (2025b) show that hierarchical reinforcement learning on template sequences rather than long chain-of-thought data achieves 91.2% on MATH, outperforming models trained on detailed reasoning traces. Wang et al. (2025b) reveal RL training induces emergent separation between high-level strategic planning and low-level procedural execution, with two-phase learning of procedural consolidation followed by strategic exploration.

In the medical domain, structured approaches demonstrate substantial gains. Liao et al. (2025) report that EHR-R1 achieves over 30 percentage points improvement on MIMIC-Bench (F1 of 0.6744

vs 0.3155 for GPT-4o) through graph-driven structured medical reasoning that converts raw EHR records into thinking graphs encoding temporal relations and causal hypotheses. Work on ICD code classification (Liu et al., 2022; Sen et al., 2021) shows that leveraging hierarchical structure through label-wise attention and multi-class reformulation improves classification, particularly at higher hierarchy levels.

While these works demonstrate that hierarchical architectures and structured representations improve reasoning, they typically attribute gains to enhanced reasoning capabilities. Our work provides an alternative interpretation: these improvements may stem from better *navigation* of knowledge hierarchies already encoded during pretraining, rather than acquiring new reasoning abilities. We test this by showing that structured prompting—which explicitly guides traversal without modifying model parameters—recovers most performance gaps between base and RL models.

### A.4 Prompting as an Alternative to RL

The possibility of achieving RL-like benefits through prompting has gained increasing attention. (Agrawal et al., 2024) demonstrate that Genetic-Evolution Prompt Alignment (GEPA) can outperform Group Relative Policy Optimization by up to 20% while using 35× fewer computational resources. They argue that “the interpretable nature of language provides a richer learning medium than sparse scalar rewards.” (Wei et al., 2022) show that chain-of-thought prompting can match fine-tuned performance on reasoning tasks, while (Zhou et al., 2023) demonstrate that optimized prompts can exceed supervised fine-tuning. The “Invisible Leash” phenomenon (Wu et al., 2025) reveals that much of RLHF’s apparent benefit comes from teaching models to follow implicit formatting patterns—effects reproducible through prompting.

### A.5 Knowledge Storage versus Knowledge Access

The distinction between knowledge acquisition and knowledge retrieval is crucial to our thesis. (Ovadia et al., 2024) show that models fine-tuned on new knowledge often

1136 “hallucinate” by incorrectly combining  
1137 existing knowledge rather than storing new  
1138 information. (Ruis et al., 2024) provide key  
1139 insights with their finding that models rely  
1140 on procedural knowledge extracted from  
1141 documents involving similar reasoning pro-  
1142 cesses rather than memorizing new facts.  
1143 This aligns with our hypothesis that RL  
1144 enhances navigation strategies rather than  
1145 expanding knowledge. (Berglund et al.,  
1146 2024) further support this through their  
1147 “Reversal Curse” findings—models trained  
1148 on “A is B” cannot infer “B is A,” sug-  
1149 gesting that training affects access patterns  
1150 rather than creating bidirectional knowl-  
1151 edge representations.

## 1152 A.6 Retrieval Complexity in 1153 Knowledge-Intensive Tasks

1154 Recent work has begun to to exam-  
1155 ine the relationship between retrieval  
1156 complexity and model performance in  
1157 knowledge-intensive tasks. (Gabburo  
1158 et al., 2024) show that retrieval complex-  
1159 ity extend beyond simple multi-hop rea-  
1160 soning—including temporal (15%), com-  
1161 parative (10%), and aggregate (16%) ques-  
1162 tions—suggesting that different types of  
1163 knowledge organization require distinct re-  
1164 trieval strategies. (Min et al., 2023)  
1165 demonstrate that in long-form generation,  
1166 factual accuracy in biographies drops as  
1167 entity rarity increases, suggesting that re-  
1168 trieval difficulty directly impacts knowl-  
1169 edge accessibility.

## 1170 B Technical Appendices and 1171 Supplementary Material

### 1172 B.1 Zero-Shot Prompt 1173 Templates

1174 We present three prompt templates used in  
1175 MedConceptsQA and IPC, which are de-  
1176 signed to elicit specific responses from lan-  
1177 guage models. These templates request:

- 1178 • Direct answers, both with and with-  
1179 out explanations.
- 1180 • Structural recall of codes and a step-  
1181 wise elimination of incorrect options.

#### Prompt Template 1: MCQ with Final Answer Only

Answer only A,B,C,D according to  
the answer to this multiple choice  
question.

[... Insert Question Text Here ...]

**Answer (only the letter of your  
choice (A, B, C, or D)):**

1182

#### Prompt Template 2: MCQ with Explanation

You are a medical research assis-  
tant. Read the following multiple-  
choice question carefully. Your task  
is to:

1. Answer each question with  
one of A/B/C/D, which cor-  
responds to the four options.
2. For my convenience, please  
give me a list of ANSWERS for  
the given instances in the for-  
mat 'Answer: ...', with addi-  
tional explanation for each an-  
swer in the format 'Explana-  
tion: ...'.

Respond in the following format:

**Answer:** <A/B/C/D>

**Explanation:** <your  
explanation here>

[... Insert Question Text Here ...]

**Answer:**

**Explanation:**

1183

**Prompt Template 3: MCQ with Stepwise Reasoning**

You are a medical classification expert. For each option, first **recall the general category and structure breakdown of the medical code**, then explain **why it might be wrong**. Finally pick the correct one.

\_\_\_\_\_

[... Insert Question Text Here ...]

**Steps to follow:**

1. Recall the general category and structural break down of the code.
2. Evaluate each option (A–D) briefly.
3. Choose the best option and justify.

**Answer format:**

Step 1: ...  
 Step 2A: ...  
 Step 2B: ...  
 Step 2C: ...  
 Step 2D: ...  
 Final Answer:  
 [A/B/C/D] because ...

\_\_\_\_\_

**C Layer-wise Representation Analysis**

**C.1 Question-Answer Pairwise Probing**

This section provides supplementary results for the layer-wise representation divergence analysis presented in Figure 2, extending the comparison across additional MedConceptsQA vocabularies for two model families.

**C.1.1 Qwen2.5 Series**

Figure 5 presents the analysis for the Qwen2.5-32B base model compared against its instruction-tuned (-Instruct), distilled (DeepSeek-R1-Distill-), and reasoning-enhanced (QwQ-32B) variants across the ATC, ICD10PROC, ICD9CM, and ICD10CM vocabularies.

**C.1.2 Mistral-Small-24B Series** 1203

Figure 4 shows the corresponding analysis for the Mistral-Small model family, comparing the base (-Base-2503), instruction-tuned (-Instruct-2503), and reasoning-enhanced (Magistral-Small-2507) variants across all five MedConceptsQA vocabularies (ATC, ICD9PROC, ICD9CM, ICD10CM, ICD10PROC). 1204-1211

**C.2 CoT Prompt Stepwise Probing** 1212-1213

To analyze model representations under chain-of-thought (CoT) prompting, we construct a series of hierarchical prompts. For example, for the question “What is the description of the medical code 743.63 in ICD9CM?”, the CoT series builds incrementally: 1214-1220

- “hmm let me think. 001-999.99 refers to diseases and injuries” 1221-1222
- “hmm let me think. 001-999.99 refers to diseases and injuries, and 740-759.99 refers to congenital anomalies” 1223-1224
- ... 1225-1226
- “hmm let me think. ...and 743.63 refers to other specified congenital anomalies of eyelid” 1227-1229

For each prompt in this series, we extract the activations from each layer of the model and group them by their corresponding vocabularies. 1230-1233

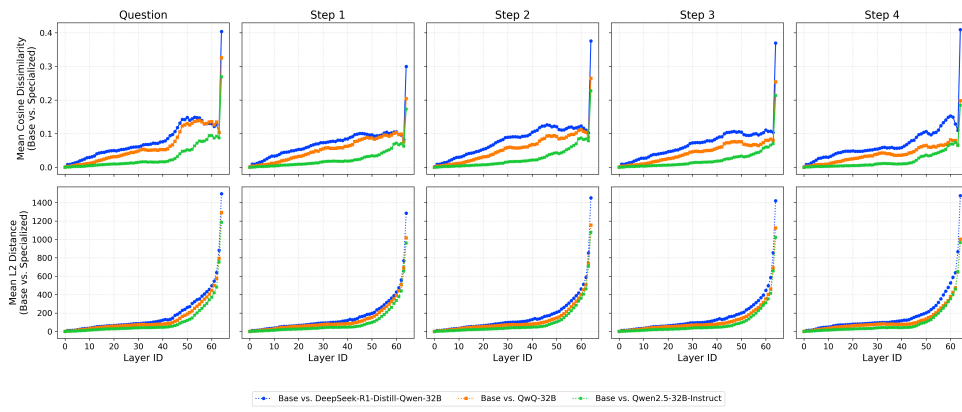
Additionally, we use **L2 distance** captures both directional and magnitude differences: 1234-1236

$$d_{L2}^{(a,b)}(\ell) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{h}_\ell^{(a)}(i) - \mathbf{h}_\ell^{(b)}(i)\|_2. \quad (2) \quad 1237$$

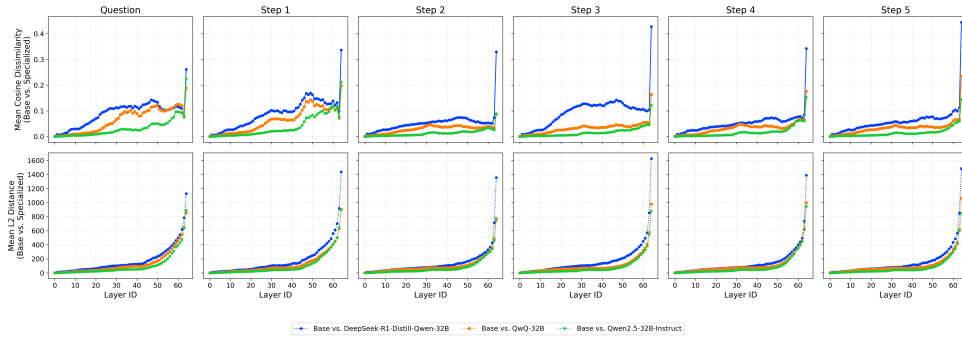
The number of CoT steps varies across vocabularies. To standardize this, we predefine all CoT sequences to be 5 steps long, with the exception of ICD10PROC, which uses 6 steps due to its more deeply embedded code structure (e.g., 0Q894Z). After grouping the activations by vocabulary for each layer, we compute the layerwise cosine similarity and L2 norm between the base and specialized models, following the methodology in Section 2.3. 1238-1247

1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202

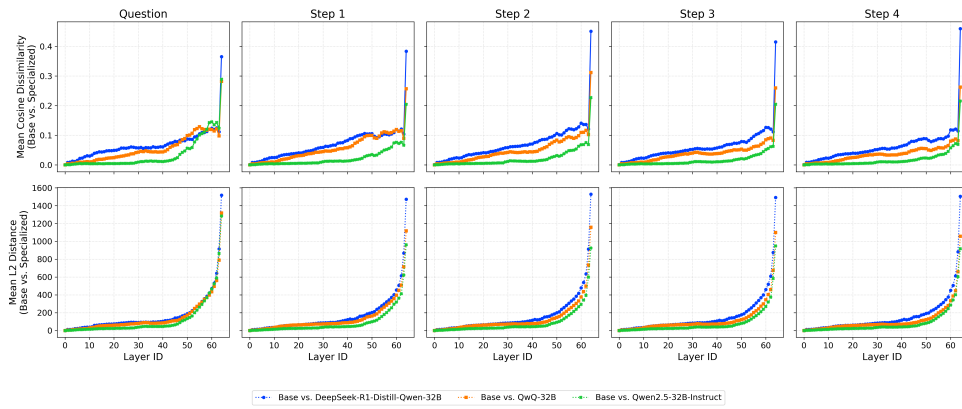
Layer-wise Similarity per COT Step: ATC Vocabulary



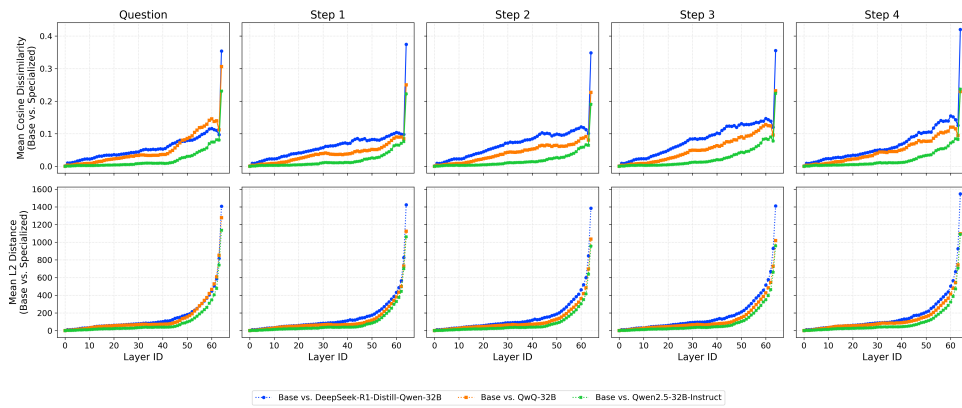
Layer-wise Similarity per COT Step: ICD10PROC Vocabulary

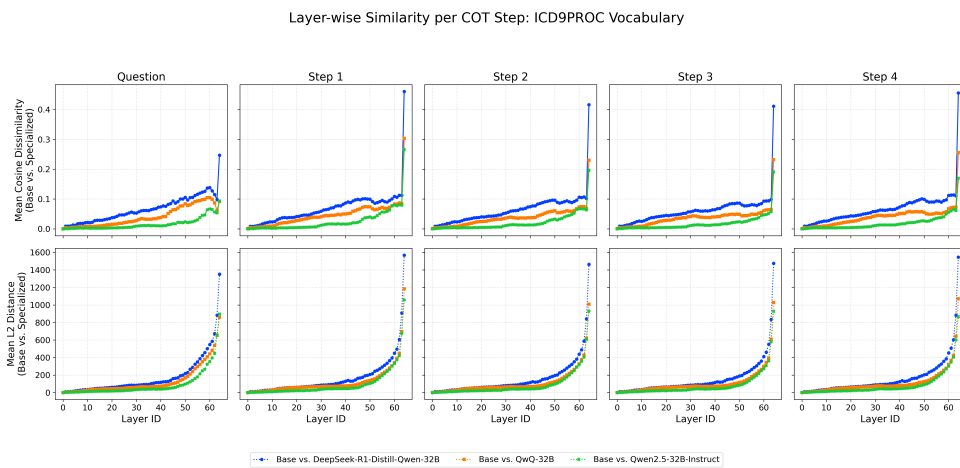


Layer-wise Similarity per COT Step: ICD9CM Vocabulary



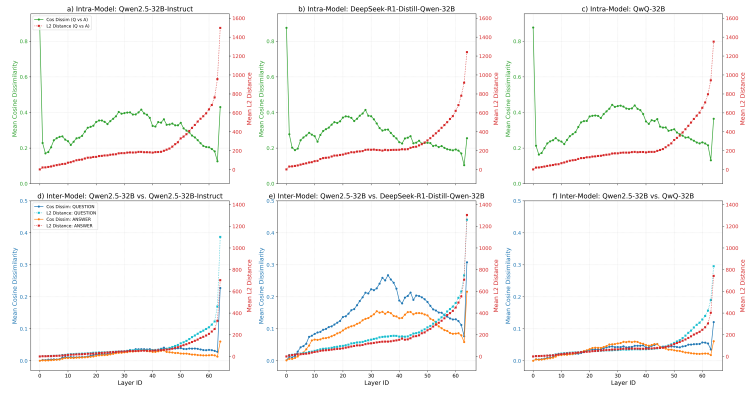
Layer-wise Similarity per COT Step: ICD10CM Vocabulary





**Figure 4: Layer-wise Representation Divergence Across CoT Steps for All MedConceptsQA Vocabularies.** This figure shows the divergence analysis results for the ATC, ICD9PROC, ICD10PROC, ICD9CM, and ICD10CM vocabularies. The top and bottom rows correspond to mean cosine similarity and L2 distance, respectively. Each column represents a distinct step in the Chain-of-Thought (CoT) process, from Step 0 (the original question) to the final step (the original question plus the complete hierarchical traversal to the correct answer).

### ATC



### ICD10PROC

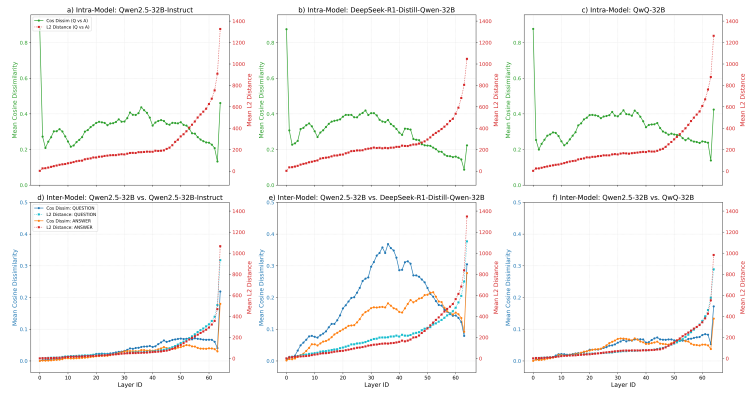
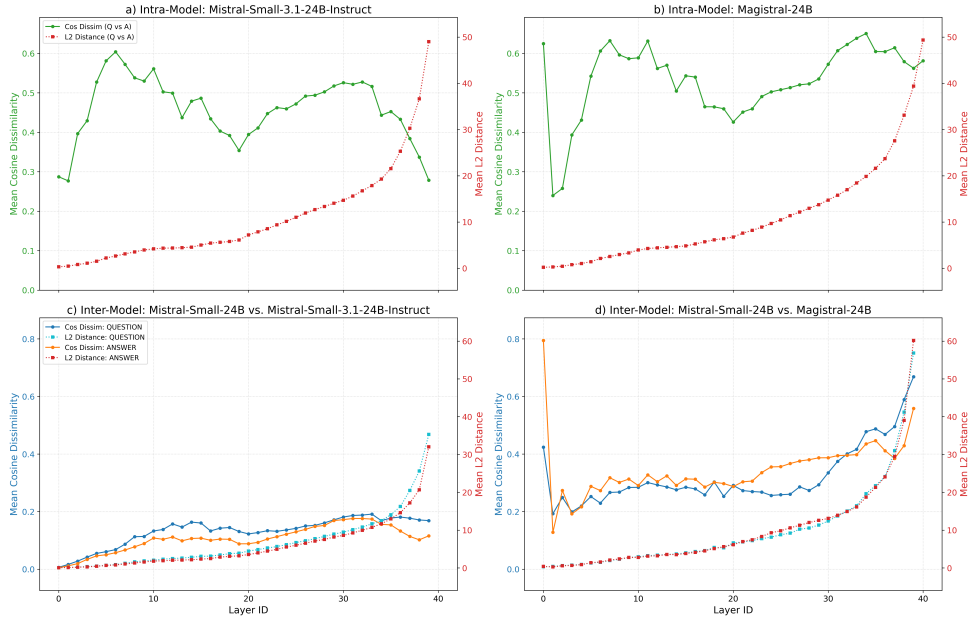


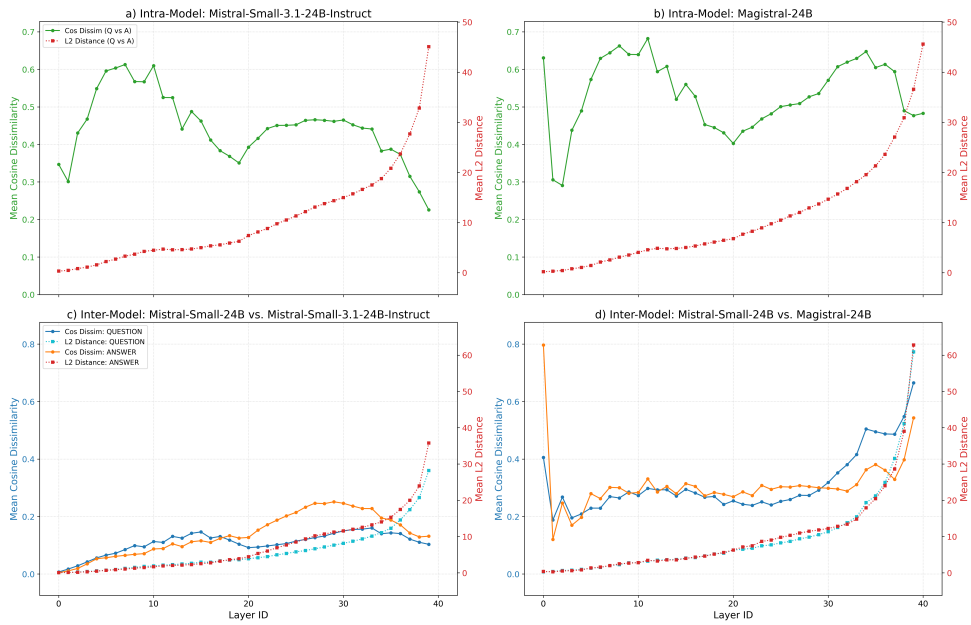


Figure 5: **Layer-wise Representation Divergence Across Remaining MedConceptsQA Vocabularies.** Same visualization format as Figure 2, showing results for ATC, ICD10PROC, ICD9CM, and ICD10CM vocabularies. Top and bottom rows correspond to intra- and inter-model divergence, respectively.

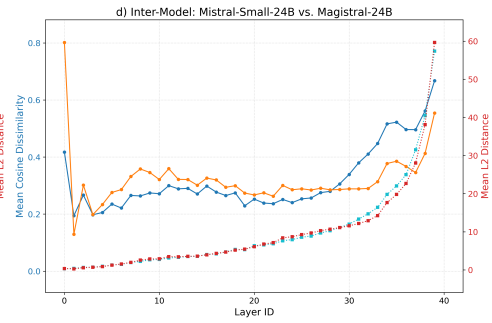
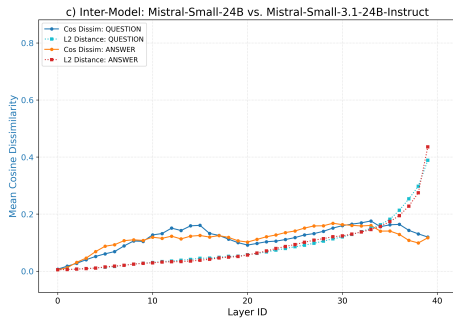
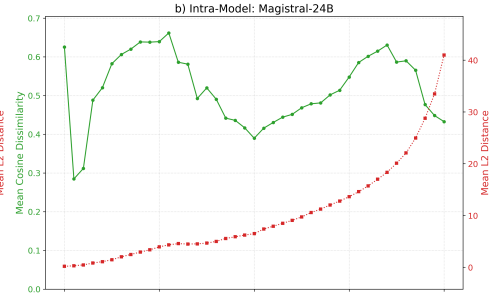
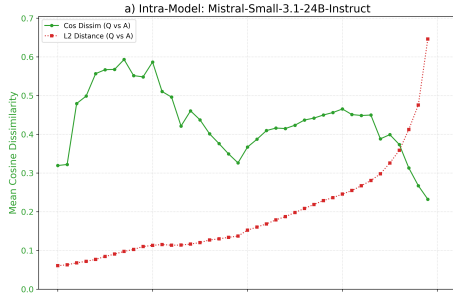
## ATC



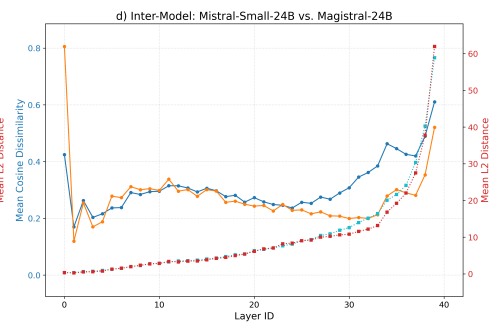
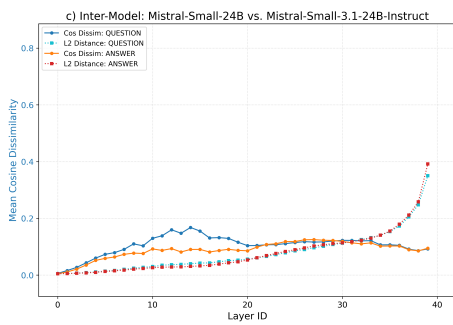
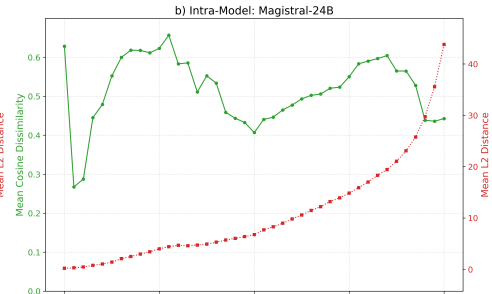
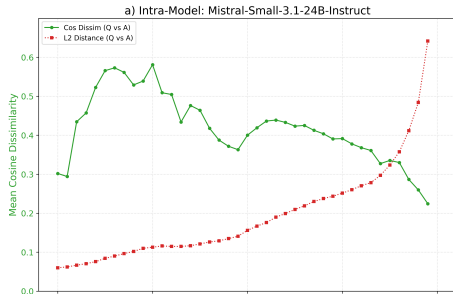
## ICD9PROC



## ICD10PROC



## ICD9CM



### ICD10CM

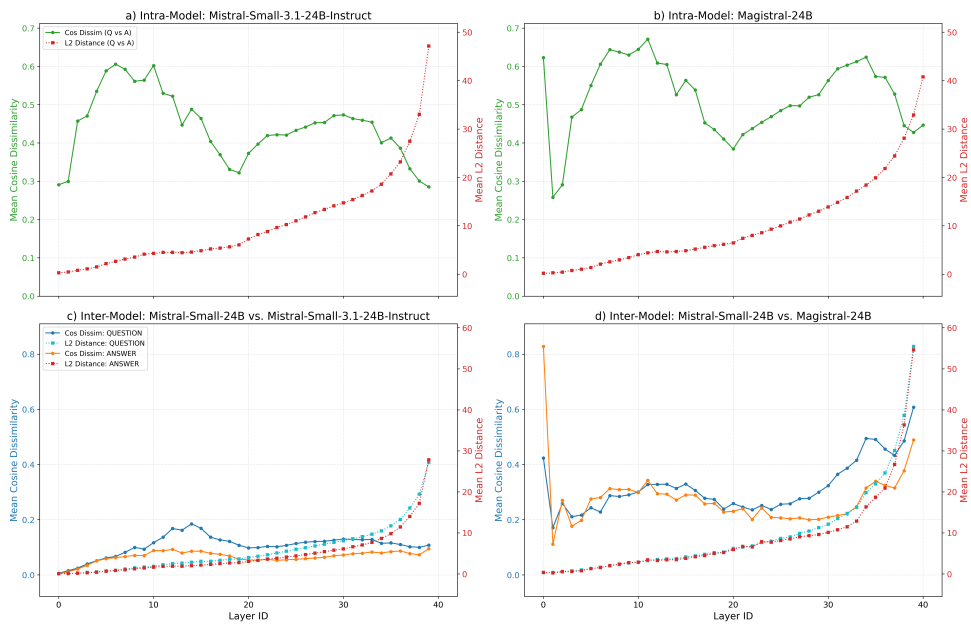


Figure 6: **Layer-wise Representation Divergence Across Remaining MedConceptsQA Vocabularies.** Same visualization format as Figure 2, showing results for ATC, ICD10PROC, ICD9CM, and ICD10CM vocabularies. Top and bottom rows correspond to intra- and inter-model divergence, respectively.