

From Static Snapshots to Dynamic Trajectories: Evaluating and Enhancing the Learning Pathways of Multimodal Large Language Models

Supplementary Material

A. Additional Experiments

A.1. Detailed Analysis of Human Performance

To contextualize the models’ performance and quantify the gap in learning gains, we conducted a human study. We recruited 16 testers of varying skill levels and tasked them with a challenging subset of our benchmark, containing 132 questions (2 per knowledge point). For practical considerations of test duration and to prevent participant fatigue, we limited the number of examples to two and designed the experiment to follow the specific learning pathway shown in Table 4 of the main paper.

As shown in Table 4 of the main paper, human testers exhibit a clear and monotonic learning trajectory. Then we quantify the benefits of each instructional type for humans in Table 5 of the main paper. The results show that humans derive substantial and consistently positive gains from all materials. Strategic guidance provides a large initial boost, and examples also yield significant, positive gains. This data confirms that human learning is highly reliable and efficient.

As shown in Table 6 of the main paper, we present the gap between the learning benefits of various models and humans. The larger the number, the greater the gap between the learning benefits of the model and humans, indicating poorer performance of the model. The large, positive $\Delta\bar{G}$ values reveal that models fail to leverage abstract guidance as effectively as humans. This gap widens when learning from examples, where some models’ performance even degrades. This contrast highlights that the models’ learning processes are fundamentally fragile and unreliable.

A.2. Sensitivity to Prompt Structure

Beyond the content of instructional materials, we investigated whether the structure of the prompt itself affects learning. We tested two formats: the default Question-Last format, where the test question appears after all instructional materials, and an alternative Question-First format. The results for the open-source models, shown in Table 8, reveal a profound sensitivity to this structural change.

The position of the question dramatically and unpredictably alters the learning gains. For Qwen2.5-VL, using the Question-First format exacerbates the negative learning from examples in multimodal tasks (MR-M), with the gain from three examples plummeting from -3.12% to -10.50%. This suggests that seeing the question first makes it even harder for the model to correctly utilize the subsequent mul-

Algorithm 1 Adaptive Reasoning Pathways (ARP)

Require: New problem $P_{\text{new}} = \{x_j\}_{j=1}^L$, MLLM model \mathcal{M} , Top-K parameter K , Adaptation steps T , Experience Pool $\mathcal{E} = \{\pi_1, \dots, \pi_n\}$.

- 1: Initialize the experience pool $\mathcal{E} = \mathcal{E} \cup \{\pi_0 = \mathbf{0}\}$
- 2: **for** each π_i in \mathcal{E} **do**
- 3: Compute NLL loss $\mathcal{L}(\pi_i|P_{\text{new}})$ as Eq.(8)
- 4: **end for**
- 5: Select K vectors with the lowest loss for voting:
 $\Pi_K = \text{top-K}_{\pi_i \in \mathcal{E}} - \mathcal{L}(\pi_i|P_{\text{new}})$
- 6: **for** each $\pi_k \in \Pi_K$ **do**
- 7: $\pi_k^{(0)} \leftarrow \pi_k$
- 8: **for** $t = 0$ to $T - 1$ **do**
- 9: $\pi_k^{(t+1)} \leftarrow \text{OptimizerStep}(\pi_k^{(t)}, \nabla_{\pi_k} \mathcal{L})$
- 10: **end for**
- 11: **end for**
- 12: **for** each updated policy $\pi_k^{(T)}$ **do**
- 13: $H'_{k,t} = H_t + \pi_k^{(T)}$
- 14: $R_k = \mathcal{M}(H'_{k,t})$
- 15: **end for**
- 16: **for** each R_k **do**
- 17: Identify the ultimate choice:
 $c_k = \text{ExtractChoice}(R_k)$
- 18: **end for**
- 19: Vote for final answer:
 $A_{\text{final}} = \underset{c \in \mathcal{C}}{\text{argmax}} \sum_{k=1}^K \mathbb{I}(c_k = c)$
- 20: Select the winner policy with minimum loss:
 $\pi_{\text{new}} = \underset{c_k = A_{\text{final}}}{\text{argmin}} \mathcal{L}(\pi_k^{(T)}|P_{\text{new}})$
- 21: Update the experience pool: $\mathcal{E} = \mathcal{E} \cup \pi_{\text{new}}$
- 22: **return** A_{final} .

Experience Retrieval

Self-Adaptation and Parallel Guidance

Ensemble Voting and Memory Consolidation

timodal examples.

Conversely, for InternVL3-78B, the same Question-First structure mitigates the catastrophic interference in text-only tasks (MR-T). The negative gain from three examples improves from -10.70% to -4.53%. This counter-intuitive result may indicate that seeing the question first helps the model anchor its reasoning, making it slightly more resilient to confusing examples. These findings demonstrate that prompt structure is another critical factor that reveals the fragility of the models’ learning mechanisms.

A.3. Resource and Cost Analysis

To ensure transparency and reproducibility, we report the computational and financial costs associated with this study. The extensive evaluation of model learning trajectories involved traversing multiple instructional dimensions (Strategy, Examples, Video) across models.

Table 8. Impact of prompt structure on Learning Gains (%). This table compares the default Question-Last format with the Question-First format.

	Mode		MR-T					MR-M				
	Question-Last	Question-First	\bar{G}_s	$\bar{G}_{e@1}$	$\bar{G}_{e@2}$	$\bar{G}_{e@3}$	\bar{G}_v	\bar{G}_s	$\bar{G}_{e@1}$	$\bar{G}_{e@2}$	$\bar{G}_{e@3}$	\bar{G}_v
Qwen2.5-VL-72B-Instruct	✓		1.20	2.37	2.68	2.32	-0.02	2.96	-2.27	-1.65	-3.12	0.53
		✓	0.14	-0.68	0.32	0.25	-1.27	0.74	-4.80	-9.35	-10.50	1.20
	Mode		MR					FR				
	Question-Last	Question-First	\bar{G}_s	$\bar{G}_{e@1}$	$\bar{G}_{e@2}$	$\bar{G}_{e@3}$	\bar{G}_v	\bar{G}_s	$\bar{G}_{e@1}$	$\bar{G}_{e@2}$	$\bar{G}_{e@3}$	\bar{G}_v
✓		1.53	1.50	1.87	1.30	0.08	0.29	3.95	1.80	1.25	0.47	
	✓	0.26	-1.38	-1.37	-1.57	-0.75	1.83	2.50	2.68	2.88	0.27	
InternVL3-78B	Mode		MR-T					MR-M				
	Question-Last	Question-First	\bar{G}_s	$\bar{G}_{e@1}$	$\bar{G}_{e@2}$	$\bar{G}_{e@3}$	\bar{G}_v	\bar{G}_s	$\bar{G}_{e@1}$	$\bar{G}_{e@2}$	$\bar{G}_{e@3}$	\bar{G}_v
	✓		-0.06	-5.95	-8.90	-10.70	-15.87	-0.27	-8.12	-6.70	-8.15	2.05
		✓	1.51	-2.18	-3.15	-4.53	-5.80	5.76	0.43	-1.27	-1.02	0.52
InternVL3-78B	Mode		MR					FR				
	Question-Last	Question-First	\bar{G}_s	$\bar{G}_{e@1}$	$\bar{G}_{e@2}$	$\bar{G}_{e@3}$	\bar{G}_v	\bar{G}_s	$\bar{G}_{e@1}$	$\bar{G}_{e@2}$	$\bar{G}_{e@3}$	\bar{G}_v
	✓		-0.10	-6.36	-8.49	-10.22	-12.50	0.04	-0.98	-2.00	-1.75	0.82
		✓	2.30	-1.82	-2.95	-4.00	-4.62	-0.54	1.95	2.00	1.40	0.53

- **Model Inference Cost:** The total cost for API calls to proprietary models (including GPT-4o, o1, Gemini-2.5-Pro, etc.) was approximately **\$9,000 USD**.
- **Human Study Cost:** For the human performance baseline, we compensated the 16 participants for their time and effort, totaling approximately **\$600 USD**.

B. ARP Method Algorithm

The whole procedure of ARP method is outlined in Algorithm 1.

C. HieraLogic and LTBench Details

C.1. Data Source and Statistics

All reasoning problems in HieraLogic were harvested from publicly available educational websites and resources on the internet. The dataset is rigorously cleaned and anonymized.

We utilize a balanced subset of HieraLogic to construct the LTBench. The specific distribution of LTBench samples across different tasks is as follows:

- **Mathematical Reasoning (MR):** This domain comprises a total of 820 questions, split into two modalities:
 - **Text-only (MR-T):** 666 questions.
 - **Multimodal (MR-M):** 154 questions.
- **Figure Reasoning (FR):** This domain is entirely multimodal and consists of 500 questions.

In total, LTBench contains 1,320 evaluation samples covering 66 distinct knowledge points (20 samples per point).

C.2. Three-level Hierarchical Structure

As shown in Table 9, this section provides a detailed overview of the hierarchical taxonomy structure, which is divided into two domains, 30 categories, and 66 knowledge points.

C.3. HieraLogic Structure

HieraLogic is a hierarchical logical reasoning dataset with the following unified fields:

- **question:** the task question, MR-M/FR include an image.
- **options:** candidate answers (e.g., A/B/C/D).
- **answer:** the ground-truth option label.
- **analysis:** textual rationale.
- **video:** filename identifier of the explanation video.
- **domain:** level-1 taxonomy, e.g., *Mathematical Reasoning*, *Figure Reasoning*.
- **category:** level-2 taxonomy (problem type within a domain), e.g., *Probability Problems* under *Mathematical Reasoning*, *Spatial* under *Figure Reasoning*.
- **knowledge point:** level-3 taxonomy (fine-grained concept), e.g., *Given Situation*, *Find Probability*, *Spatial-3D Assembly*.

HieraLogic is organized by task and modality into: *Mathematical Reasoning—Text (MR-T)*, *Mathematical Reasoning—Multimodal (MR-M)*, and *Figure Reasoning (FR, multimodal)*.

C.4. HieraLogic Examples

C.4.1. MR-T (Mathematical Reasoning Text)

Question: A and B participated in a math, physics, and chemistry competition held in a certain city. The probabilities of A passing the preliminary round in the three subjects are 80%, 50%, and 60%, respectively, while the probabilities for B are 60%, 60%, and 50%, respectively. The correct statement among the following is:

Options:

A: The probability that A passes at least 2 subjects in the preliminary round exceeds 80%.

B: The probability that B passes at least 2 subjects in the preliminary round exceeds 70%.

C: The probability that both A and B have passed only one subject in the preliminary competition exceeds 10%.

D: The probability that at least one of A or B passes all three subjects in the preliminary round is less than 40%.

Answer: D

Analysis: Let $P_A(k)$ denote A's probability of passing exactly k subjects and $P_B(k)$ for B. Assuming independence,

$$\begin{aligned} P_A(1) &= 0.8 \times (1 - 0.5) \times (1 - 0.6) \\ &\quad + (1 - 0.8) \times 0.5 \times (1 - 0.6) \\ &\quad + (1 - 0.8) \times (1 - 0.5) \times 0.6 = 0.26, \end{aligned}$$

$$\begin{aligned} P_A(2) &= 0.8 \times 0.5 \times (1 - 0.6) + 0.8 \times (1 - 0.5) \times 0.6 \\ &\quad + (1 - 0.8) \times 0.5 \times 0.6 = 0.46, \end{aligned}$$

$$P_A(3) = 0.8 \times 0.5 \times 0.6 = 0.24;$$

$$\begin{aligned} P_B(1) &= 0.6 \times (1 - 0.6) \times (1 - 0.5) \\ &\quad + (1 - 0.6) \times 0.6 \times (1 - 0.5) \\ &\quad + (1 - 0.6) \times (1 - 0.6) \times 0.5 = 0.32, \end{aligned}$$

$$\begin{aligned} P_B(2) &= 0.6 \times 0.6 \times (1 - 0.5) \\ &\quad + 0.6 \times (1 - 0.6) \times 0.5 \\ &\quad + (1 - 0.6) \times 0.6 \times 0.5 = 0.42, \end{aligned}$$

$$P_B(3) = 0.6 \times 0.6 \times 0.5 = 0.18.$$

Hence, A's probability of "at least two passes" is $P_A(2) + P_A(3) = 0.70$ (so A is false), B's is $0.42 + 0.18 = 0.60$ (B is false), and "both pass exactly one" is $P_A(1)P_B(1) = 0.26 \times 0.32 < 0.1$ (C is false). The probability that "at least one passes all three" is

$$\begin{aligned} &1 - (1 - P_A(3))(1 - P_B(3)) \\ &= 1 - (1 - 0.24)(1 - 0.18) \\ &= 1 - 0.76 \times 0.82 < 0.4, \end{aligned}$$

thus D is correct.

Video: MRT_video.mp4

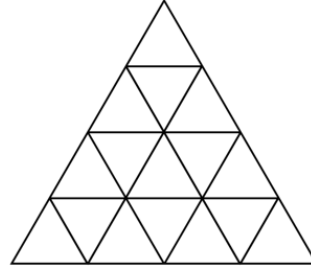
Domain: Mathematical Reasoning

Category: Probability Problems

Knowledge points: Given Situation, Find Probability

C.4.2. MR-M (Mathematical Reasoning Multimodal)

Question: As shown in the image, there are 16 triangular boxes. If two identical balls are randomly placed into two different boxes, the probability that these two boxes are not adjacent (meaning they do not share a common edge) is:



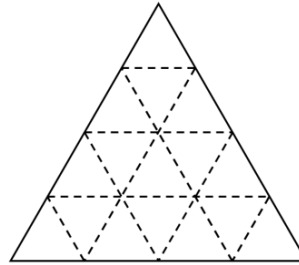
Options:

A: $\frac{17}{20}$ B: $\frac{3}{20}$ C: $\frac{2}{15}$ D: $\frac{13}{15}$

Answer: A

Analysis: Count the complement. "Adjacent" pairs are determined by a shared edge: there are 18 adjacent pairs in total, so the probability of adjacency is $\frac{\binom{18}{1}}{\binom{16}{2}} = \frac{18}{120} = \frac{3}{20}$.

Therefore, "not adjacent" is $1 - \frac{3}{20} = \frac{17}{20}$.



Video: MRM_video.mp4

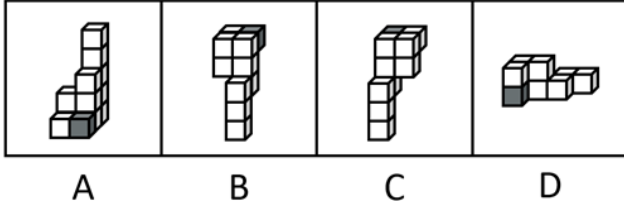
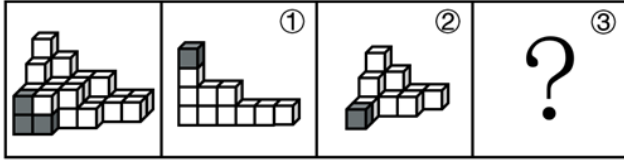
Domain: Mathematical Reasoning

Category: Probability Problems

Knowledge points: Given Situation, Find Probability

C.4.3. FR (Figure Reasoning)

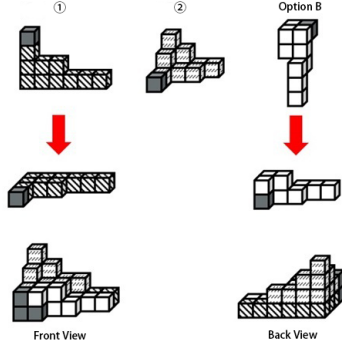
Question: The left image shows a polyhedron composed of 3 equally sized gray cubes and 31 white cubes, which can be divided into three smaller polyhedra (1) (2) and (3) The polyhedron represented by (3) could be:



Options: A/B/C/D

Answer: B

Analysis: This task assesses 3D assembly. First flatten shape (1) by pushing it from back to front, then translate option B from right to left and flip it; combine it with shape (2) to obtain the composite shown on the left.



Video: FR.video.mp4

Domain: Figure Reasoning

Category: Spatial

Knowledge points: Spatial-3D Assembly

C.5. LTBench Structure

LTBench leverages HieraLogic to assemble trajectory tasks. We formalize an LTBench instance by the tuple

$$\langle \text{Task}, \text{SG}, \text{ED}(k), \text{VA} \rangle,$$

where:

- **Task** (question/options from HieraLogic):
 - question: the original question.
 - options: candidate answers (e.g., A/B/C/D).
- **SG** (Strategic Guidance): an optional high-level hint targeting the current knowledge point.
- **ED** (Example Depth k): a set of k worked HieraLogic examples with solutions targeting the current knowledge point.
- **VA** (Video Augmentation): an optional instructional video accompanying each example when $k > 0$.

For evaluation, a configuration is denoted by (s^σ, e_k, v^τ) with $\sigma, \tau \in \{+, -\}$, indicating the presence/absence of SG and VA, and example depth k .

C.6. LTBench Examples

C.6.1. Example A: Zero-shot without SG

Configuration: (s^-, e_0) (VA is irrelevant when $k = 0$).

Task-Question: A and B participated in a math, physics, and chemistry competition held in a certain city. The probabilities of A passing the preliminary round in the three subjects are 80%, 50%, and 60%, respectively, while the probabilities for B are 60%, 60%, and 50%, respectively. The correct statement among the following is:

Task-Options:

A: The probability that A passes at least 2 subjects in the preliminary round exceeds 80%.

B: The probability that B passes at least 2 subjects in the preliminary round exceeds 70%.

C: The probability that both A and B have passed only one subject in the preliminary competition exceeds 10%.

D: The probability that at least one of A or B passes all three subjects in the preliminary round is less than 40%.

SG: None (s^-)

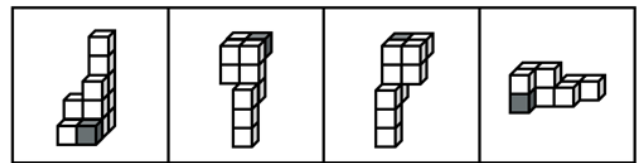
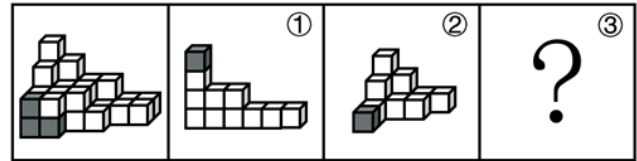
Examples: None ($k = 0$)

VA: N/A

C.6.2. Example B: With SG, one example, with video

Configuration: (s^+, e_1, v^+) .

Task-Question: The left image shows a polyhedron composed of 3 equally sized gray cubes and 31 white cubes, which can be divided into three smaller polyhedra (1) (2) and (3) The polyhedron represented by (3) could be:



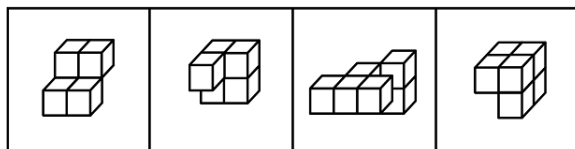
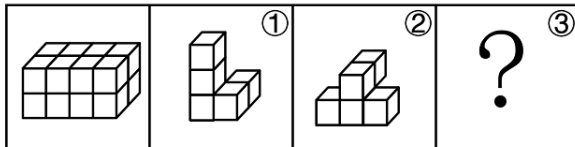
Task-Options: A/B/C/D

SG: For non-cube blocks, it is important to pay attention to matching the concave and convex parts, ensuring that the pieces fit together by their indentations and protrusions. For cube-shaped blocks, first, if the number of blocks in the options differs, prioritize counting the pieces to ensure that the number of blocks corresponds to the puzzle's structure. Next, start assembling from the most complex and distinct shapes, focusing on pieces with the most blocks or the most unique shapes, especially those that clearly match specific parts of the 3D model. These unique shapes will make the assembly process more manageable as they fit into obvi-

ous parts of the overall structure. After placing the blocks that can be definitively positioned, focus on identifying the remaining pieces by layer. Examine each layer of the structure to determine which specific blocks are needed, either by shape or by number, to complete the remaining parts.

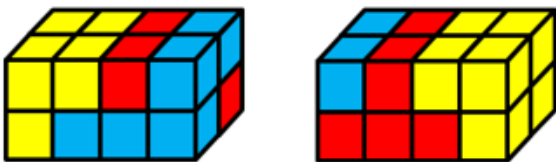
Examples ($k = 1$):

- **ex_question:** The image on the left shows a cuboid formed by stacking identical cubes. This cuboid can be composed of three polyhedra: (1), (2), and (3). Which of the following can fill the place of the question mark?



A B C D

- **ex_options:** A/B/C/D
- **ex_answer:** D
- **ex_analysis:** Observing the cuboid given in the question, we can see that the figure is divided into 2 layers, with a total of 16 small cubes. Figures (1) and (2) have a total of 10 small cubes, so there should be 6 small cubes at the ? position, eliminating options A and B. Rotate Figure (1) counterclockwise 90° and combine it with Figure (2) according to the principle of consistent concavity and convexity. The remaining shape only matches option D. The specific combination method is shown in the figure below: Therefore, the correct answer is D.



- **ex_video:** EX_FR.mp4

VA (v⁺): Enabled; the example includes an instructional video.

Domain	Category	Knowledge Point
Mathematical Reasoning	Work-Rate Problems	Given Completion Time Given Efficiency Ratio Given Concrete Units Work-Rate (Other)
	Extremum Problems	Atypical Extremum Problems Constructed Sequences Worst-Case Construction Multi-Set Reverse Construction
	Age Problems	Age Problems
	Sum-Difference-Multiple-Ratio Problems	Sum-Difference-Multiple-Ratio Problems
	Periodicity Problems	Periodic Meeting Problems Periodic Remainder Problems Periodicity (Other)
	Sequence Problems	Sequence Problems
	Motion Problems	Train Crossing Bridge Average Speed Standard Motion Meeting and Pursuit Boat in Current Motion Problems (Other)
	Geometry Problems	Geometry Formulas - Plane Figures Geometry Formulas - Solid Figures Geometric Theorems - Triangle Related
	Inclusion-Exclusion Problems	Two Sets Three Sets
	Permutations and Combinations	Basic Permutations and Combinations Adjacency Problems Non-Adjacent Problems Distribution of Identical Items Circular Permutation Problems Derangements Permutations and Combinations (Other)
	Probability Problems	Given Situation, Find Probability Given Probability, Find Probability Probability (Other)
	Profit Problems	Profit Problems
	Planning and Scheduling Problems	Planning and Scheduling Problems
	Day/Date Problems	Day/Date Problems
	LCM and GCD Problems	LCM and GCD Problems
Piecewise Calculation Problems	Piecewise Calculation Problems	
Function Extrema Problems	Function Extrema Problems	
Figure Reasoning	Position Patterns	Position Patterns - Translation Position Patterns - Rotation Position Patterns - Mixed Position Patterns - Reflection
	Style Patterns	Style Patterns - Black/White Operations Style Patterns - Add/Sub Same/Diff Style Patterns - Enumeration
	Attribute Patterns	Attribute Patterns - Symmetry Attribute Patterns - Open/Closed Attribute Patterns - Curved/Straight Attribute Patterns - Composite
	Quantity Patterns	Quantity Patterns - Point Quantity Patterns - Line Quantity Patterns - Plane Quantity Patterns - Prime Quantity Patterns - Angle Quantity Patterns - Composite
	Special Patterns	Special Patterns - Inter-Figure Relations Special Patterns - Functional Elements
	Spatial	Spatial - 3D Assembly Spatial - Three-View Drawings Spatial - Cross-Section Spatial - Reconstruction - Hexahedron
	Text-Letter-Number	Text-Letter-Number
	Black-White Blocks	Black-White Blocks

Table 9. Three-level hierarchical taxonomy structure.