

ReLaX: Reasoning with Latent Exploration for Large Reasoning Models

Supplementary Material

Overview

This Supplementary Material provides more details and results of the proposed ReLaX.

- In Supp. Sec. 1, we discuss the **related work** of recent GRPO variants that aim to address the ineffective exploration through the lens of token-level entropy.
- In Supp. Sec. 2, we provide supplementary background on **Koopman dictionary learning**, along with t-SNE visualizations of latent trajectories under the learned dictionary and the **pseudocode** of ReLaX.
- In Supp. Sec. 3, We comprehensively describe our **experimental setup**.
- In Supp. Sec. 4, **additional experimental results** are provided to complement the analyses in Sec. 4.

1. Related Work

Although RLVR has achieved impressive gains in mathematical, coding, and visual reasoning tasks, recent studies [45, 53] suggest that its improvements largely arise from enhanced sampling efficiency rather than true capability gains. What is worse, as the model becomes increasingly confident and generates more convergent outputs, policy learning suffers from reduced exploration, leading to saturation and a performance bottleneck—a central challenge in balancing exploration and exploitation [10, 33]. To mitigate this issue, recent research focuses on maintaining sufficient entropy in the model’s action space to promote effective exploration.

Early methods [20, 50] seek to maximize policy entropy directly or incorporate it into the reward signal [7]. Such methods, however, require careful tuning of weighting coefficients: insufficient weighting fails to prevent collapse, while excessive weighting may induce semantic drift.

To regulate policy entropy in a more principled and controllable way, some research focuses on identifying structural factors that drive changes in entropy. DAPO [51] observes that the standard PPO/GRPO clipping range, though stabilizing the training, imposes an overly restrictive upper bound that suppresses exploration. Relaxing this bound via the clip-higher strategy increases credit for positively rewarded tokens and promotes broader exploration. Building on this insight, DCPO [49] introduces a dynamic, adaptive clipping mechanism that affords finer control over update magnitudes. Cui et al. [10] further reveal that the covariance between action probabilities and logit variations strongly predicts entropy reduction; accordingly, they propose KL-Cov and Clip-Cov to selectively penalize or clip high-covariance tokens most likely to induce entropy collapse.

Another line of approaches focuses on allocating policy-update credit selectively to high-entropy tokens. For in-

stance, [43] updates only the top 20% high-entropy forking tokens and reports substantial performance gains. FR3E [58] extends this idea by identifying high-entropy tokens during sampling and expanding rollouts specifically along these positions; subsequent policy optimization is then performed on these augmented trajectories, effectively enhancing exploration. CURE [26] further generalizes FR3E by stochastically selecting high-entropy tokens to mitigate selection bias. To consolidate the improved exploration into stronger exploitation, CURE additionally applies DAPO [51] in a second training stage.

2. Details on Koopman Dictionary Learning

This section provides background and details the procedure for learning the Koopman dictionary, which maps the model’s nonlinear, high-dimensional hidden states into a linear, tractable representation. An illustrative example The pseudocode for dictionary learning and the complete ReLaX training workflow is given in Algorithm 1.

2.1. Residual DMD (ResDMD)

Approximating the Koopman operator requires projecting its infinite-dimensional dynamics onto a finite set of observables [25]. This projection inevitably introduces spurious eigenvalues, a problem exacerbated for large reasoning models whose hidden states are extremely high-dimensional (e.g., 3584 for Qwen2.5-7B) and evolve over long contexts. In such settings, standard DMD approaches overfit transient fluctuations, suffer numerical instability, and produce spectral artifacts that obscure the true temporal modes.

ResDMD [8] improves spectral fidelity by diagnosing and filtering corrupted eigenvalues using a residual test. Given snapshot matrices \mathcal{V} and \mathcal{V}^+ and an approximate Koopman operator \mathcal{K} , the squared residual of an eigenpair

(λ, v) is estimated as:

$$\begin{aligned} \text{res}(\lambda, v)^2 := & v^*[(\mathcal{V}^+)^*\mathcal{V}^+ - \lambda(\mathcal{V}^*\mathcal{V}^+)^* - \bar{\lambda}\mathcal{V}^*\mathcal{V}^+ \\ & + |\lambda|^2\mathcal{V}^*\mathcal{V}]v. \end{aligned} \quad (13)$$

This residual measures how well the eigenpair satisfies the least-squares formulation and, therefore, how faithfully it captures the underlying dynamics. Eigenpairs with large residuals are classified as spurious and removed. A limitation, however, is that this diagnostic only becomes available once the Koopman spectrum is computed with a fixed set of observables. Since choosing an effective dictionary is itself a central and challenging problem, a natural question arises: can we learn the dictionary jointly while minimizing spectral residuals?

2.2. ResKoopNet

ResKoopNet [48] addresses this by parameterizing the Koopman dictionary with an MLP and optimizing it directly through the residual objective in Eq. 13. The resulting residual loss is given in Eq. 9, with the full derivation provided in the original paper.

In our implementation, the dictionary is parametrized by a single feedforward layer with a Sigmoid activation. It is trained on the hidden states collected during the first step of sampling, yielding a general set of observables that defines a globally shared linear representation space for the subsequent policy optimization. The complete procedure is provided in Algorithm 1.

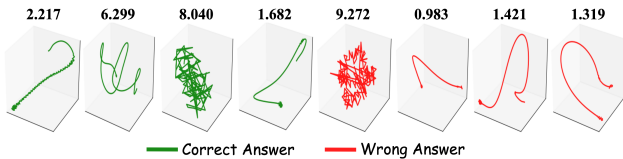


Figure 7. t-SNE of latent trajectories (DeepSeek-Math-7B). While sample-wise DSD does not strongly correlate with correctness, policy optimization benefits from consistently flexible internal computations, as reflected by higher DSD.

3. Detailed Experimental Settings

Our experiments were run on multiple GPU clusters, each equipped with 8× NVIDIA A100 (80GB). To improve training efficiency, we enabled actor–learner collocation via VeRL [34]. The primary hyperparameters on finetuning both LLMs and VLMs, including those for the actor, trainer, and ReLaX-specific settings, are summarized in Table 3. All experiments employ identical configurations for optimizing the Koopman dictionary, using the Adam optimizer with a learning rate of 10^{-4} and a batch size of 64.

The benchmarking results reported in Tables 1 and 2 are compiled from published papers and our evaluations on

publicly available checkpoints using open-source toolkits. Specifically, VLMs are evaluated with VLMEvalKit¹, and LLMs are evaluated following Qwen2.5-Math². The evaluation hyperparameters for both VLMs and LLMs are listed in Table 4.

Finally, the chat templates employed in VeRL for model training are presented below. During evaluation, we use the evaluation codebase’s default system prompt to ensure fair and comparable results.

Chat template for Visual-Language training.

System: Given the images and the question, please reason step by step, and put your final answer within `\boxed{\}`.
 User: `\langle Question \rangle`
 Assistant:

Chat template for text-only training.

System: Please reason step by step, and put your final answer within `\boxed{\}`.
 User: `\langle Question \rangle`
 Assistant:

4. Additional Results

In this section, we present the additional results and details to clarify some of the analytical experiments and substantiate our claims.

4.1. Training Dynamics of Text-only Experiments

We observe training dynamics (Fig. 8) in the text-only scenario that mirror the phenomenon seen previously in the vision-language experiments. Compared with vanilla GRPO, ReLaX maintains higher policy entropy, more diverse latent dynamics, and longer response lengths. Moreover, ReLaX exhibits strong training stability, as evidenced by its consistently low gradient clipping rate. Together, these properties enable ReLaX to achieve substantial improvements in validation accuracy.

4.2. Supplemented Experiments on Other Models

Although models from the Qwen2.5 family are commonly used as base models for RLVR training in recent work, there has been an ongoing community discussion regarding potential evaluation set leaks. To this end, we further include experiments on Llama3.2-3B-Instruct and Qwen3-4B-Base to demonstrate that ReLaX generalizes beyond specific foundation models.

As shown in Table 5, our method consistently outperforms the GRPO baseline across both types of base mod-

¹<https://github.com/open-compass/VLMEvalKit>

²<https://github.com/QwenLM/Qwen2.5-Math>

Algorithm 1: ReLaX

Input: Dataset \mathcal{D} , policy π_θ , reference policy π_{ref} , batch size B , group size R , learning rate η , total training steps \mathcal{S}

Output: Optimized policy π_θ

```
1 Set iteration counter  $s \leftarrow 0$ 
2 while  $s \leq \mathcal{S}$  do
3   Sample a batch of queries  $Q \sim \mathcal{D}$ 
4   Initialize gradient accumulator  $\nabla \mathcal{J} \leftarrow 0$ 
5   /* 1. Sampling */
6   foreach  $q \in Q$  do
7     /* vLLM inference */
8     Generate  $R$  completions  $O = \{o^1, \dots, o^R\}$  from  $\pi_{\theta_{\text{old}}}(\cdot | q)$ 
9     Calculate reward( $Q, O$ ) for completions
10    /* Transformers inference */
11    Compute log probabilities
12    Collect hidden states  $X = \{x^1, \dots, x^R\}$  from the final hidden layer
13  end
14  /* 2. Koopman Dictionary Learning (First step only) */
15  if  $s = 0$  then
16    Form batch-level hidden-state set:  $X_B = \bigcup_{q \in Q} X^q$ 
17     $g \leftarrow \text{FitKoopmanDict}(X_B)$ 
18  end
19  /* 3. Policy Optimization */
20  foreach  $x \in X$  do
21    Compute  $\text{DSD}(x)$  using Eq. 7
22  end
23  Compute ReLaX objective  $\mathcal{J}(\theta)$  using Eq. 12
24  Update policy:  $\theta \leftarrow \theta + \eta \nabla_\theta \mathcal{J}(\theta)$ 
25   $s \leftarrow s + 1$ 
26 end
27 Function  $\text{FitKoopmanDict}(X)$ :
28   Input: Hidden states  $X \in \mathbb{R}^{BR \times T \times d}$ 
29   Select batch size  $B_g$ , total optimization steps  $\mathcal{S}_g$ , Koopman dimension  $m$ , and learning rate  $\eta_g$ 
30   Initialize dictionary network  $g$  with parameters  $W \in \mathbb{R}^{d \times m}$  and sigmoid activation  $\sigma$ 
31   Set iteration counter  $s \leftarrow 0$ 
32   while  $s \leq \mathcal{S}_g$  do
33     Sample hidden state batch  $X_B \sim X$ 
34     Construct consecutive temporal snapshots  $\mathcal{V}_B, \mathcal{V}_B^+ \in \mathbb{R}^{B_g \times (T-1) \times d}$ 
35     Estimate Koopman operator  $\mathcal{K}$  using Eq. 5
36     Compute residual loss  $\mathcal{J}_g(W)$  using Eq. 9
37     Update dictionary parameters  $W \leftarrow W - \eta_g \nabla_W \mathcal{J}_g(W)$ 
38      $s \leftarrow s + 1$ 
39   end
40   return Learned Koopman dictionary  $g$ 
```

els. To further compare with token-level methods, we additionally include results from HICRA [39]. ReLaX achieves substantial gains on Minerva and AMC, while remaining competitive on MATH500 and AIME. These results demonstrate that ReLaX generalizes effectively across different

base models and maintains a stable advantage over token-level methods.

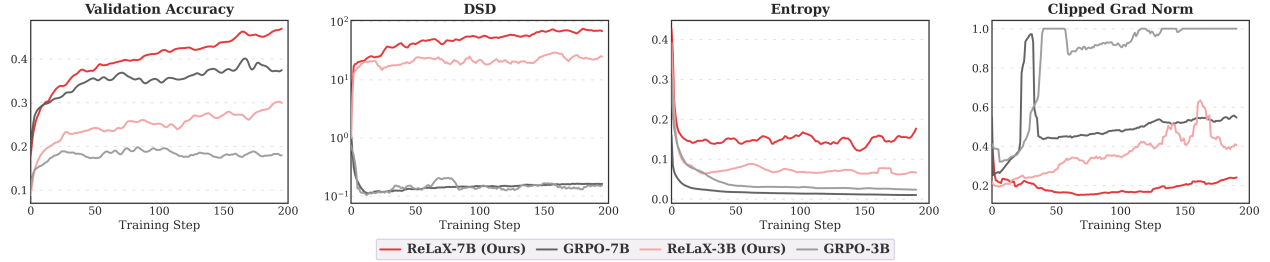


Figure 8. Comparison of training dynamics for *Validation Accuracy*, *DSD*, *Entropy*, and *Clipped Gradient Norm* under ReLaX (red) and vanilla GRPO (grey) on Qwen2.5-Base at the 3B and 7B scales.

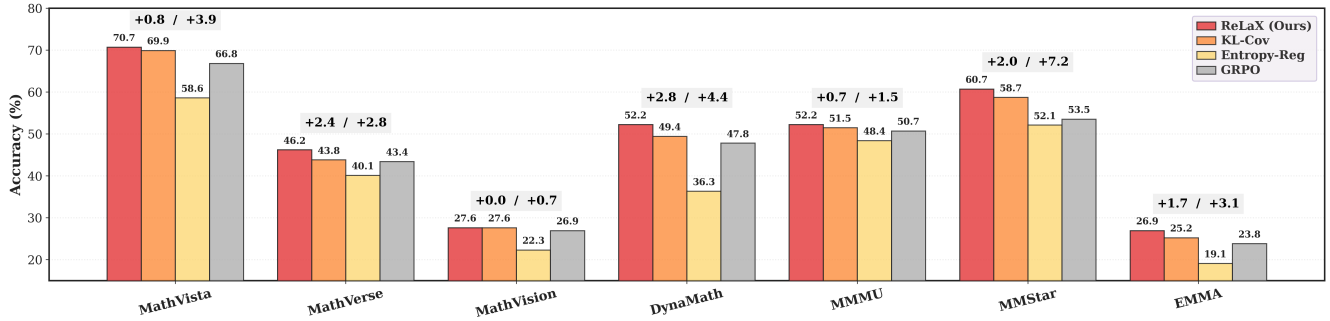


Figure 9. Complete comparison of multimodal benchmark performance for different training methods on Qwen2.5-3B-VL. This figure provides the extended results referenced in Fig. 6a. The values above the bars denote ReLaX’s performance gains over KL-Cov (left) and vanilla GRPO (right).

4.3. Additional Results on Sensitivity Analysis

Fig. 11 supplements the sensitivity analysis for Koopman operator m and DSD threshold ξ . ReLaX is generally robust to m , except at $m = 5$, where *too few spectral modes fail to capture intrinsic latent dynamics accurately*. This leads to unstable DSD calculation, rendering the regularization ineffective and degenerating the method into vanilla GRPO with premature convergence. As for ξ , removing this constraint ($\xi = \infty$) allows DSD to grow unconstrained, leading to numerical instability and training collapse after ~ 50 steps. In contrast, setting $\xi = 10, 25, 50$ effectively bounds DSD, with comparable performance across values, indicating robustness to the choice of ξ .

4.4. Computational Time Consumption

Compared with vanilla GRPO, ReLaX introduces additional computation for the model’s latent representation. To assess its time efficiency, we provide a detailed runtime breakdown of the extra components, including Koopman dictionary learning and the actor update, and compare them against the vanilla GRPO. Following our ablation setup, this analysis is conducted on Qwen2.5-3B-VL and Qwen2.5-7B-Math.

As shown in Table 6, fitting the Koopman dictionary incurs only a small one-time cost (109 s for the 3B model and 132 s for the 7B model), which is negligible compared to

the whole training process, as it is performed only once at initialization. During training, the primary additional overhead of ReLaX arises from computing the DSD score for each hidden state at every step, which increases the actor-update time by roughly 50% compared with vanilla GRPO. However, this component accounts for only about 10% of the total per-step runtime, making its impact on overall training time relatively minor. These observations indicate that the time overhead introduced by ReLaX is acceptable in practice, and we leave further optimization of its computational efficiency to future work.

4.5. More Details on Comparisons with Token-level Methods

To support the claims and conclusions presented in Sec. 4.4, we provide additional results from analytical experiments comparing ReLaX with token-level methods. The following section presents evaluation results across all multimodal reasoning benchmarks, as well as a detailed case study of the obtained LRMs.

4.5.1. Extra Evaluation Results of VLMs

Firstly, we provide additional results of token-level methods on 3B-VL models evaluated across all multimodal benchmarks. In addition to the results on MathVista, DynaMath, MMStar, and the two multidisciplinary subsets of EMMA reported in Fig. 6a, we further include comparisons

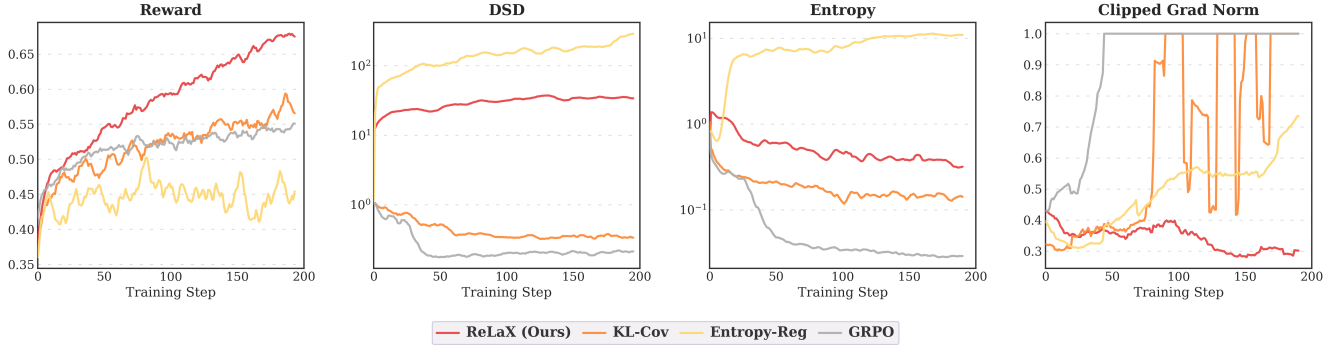


Figure 10. More training dynamics of vanilla GRPO, KL-Cov, Entropy Reg and our ReLaX on Qwen2.5-3B-VL. This figure provides the extended results for Fig. 6b and 6c.

Parameter	Value
Actor	
Maximum response length	3072
Temperature	1.0
top p	1.0
top k	-1
Number of rollouts per prompt	16
Trainer	
Batch size	512
Generate size for sampling	2048
Optimizer	AdamW
Adam betas	(0.9, 0.95)
Gradient norm clipping	1.0
Learning rate scheduler	Constant
Learning rate	10^{-6}
ReLaX-Specific	
Exploration coefficient α	0.1
KL loss coefficient β	0.01
DSD threshold ξ	10
Koopman operator dimension m	50

Table 3. Hyperparameters for RLVR training used in our experiments. The same settings are applied to both VLM and LLM experiments.

Parameter	Value
VLM evaluation	
top p	0.1
Maximum response length	10240
Temperature	0
LLM-as-Judge	GPT-4o-mini
LLM evaluation	
top p	0.1
Maximum response length	10240
Temperature (mean@1)	0
Temperature (mean@32)	1

Table 4. Hyperparameters used for Evaluations.

LaX consistently outperforms the vanilla GRPO baseline, its main advantage over the token-level method, KL-Cov, appears in the multidisciplinary multimodal benchmarks, which rely more on visual information, rather than in the mathematics-dominant benchmarks where text plays a central role. Additional training dynamics of the reward and gradient norm are presented in Fig. 10. We observe that ReLaX exhibits a notably stable gradient norm, which is beneficial for optimization.

4.5.2. Case Study for Multimodal Reasoning

To investigate the behavior of models trained with the proposed ReLaX in multimodal reasoning, we present a detailed case study comparing ReLaX and KL-Cov [10] on 3B-VL models. We select a query from DynaMath [59], a recent benchmark designed to evaluate reasoning robustness by testing how well models perform under different variants of the same question, such as changes in visual numerical values. As shown in Table 8, the question asks for the perimeter of a rectangular prism, and the three variants differ only in the numerical values of length, width, and

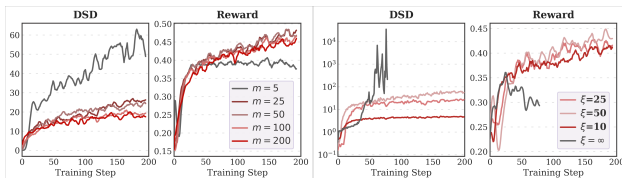


Figure 11. Sensitivity analysis for Koopman operator dimension m and DSD threshold ξ on Qwen2.5-3B.

on MathVerse, MathVision, MMMU, and the full EMMA benchmark in Fig. 9. These results reinforce that while Re-

Model	Size	MATH500	Minerva	AMC22	AMC23	AIME24	AIME25	Average
		Mean@1	Mean@1	Mean@32	Mean@32	Mean@32	Mean@32	
Llama3.2-Instruct	3B	46.8	15.4	16.1 [†]	22.0	8.5	0	18.1
+Vanilla GRPO[32]	3B	55.4	22.8	26.4	40.0	16.3	1.4	27.1
+ReLaX (Ours)	3B	57.0	23.5	39.0	52.8	18.9	3.3	32.4
Δ (ReLaX - GRPO)	3B	+1.6	+0.7	+12.6	+12.8	+2.6	+1.9	+5.3
Qwen3-Base	4B	63.8	28.3	29.1 [†]	38.9	9.4	5.3	29.1
+Vanilla GRPO[32]	4B	83.0	38.9	42.6	51.2	24.9	23.8	44.1
+HICRA [39]	4B	89.0	42.5	-	54.0	31.0	27.6	-
+ReLaX (Ours)	4B	90.2	48.5	52.6	64.5	30.9	27.6	52.3
Δ (ReLaX - GRPO)	4B	+6.2	+9.6	+10.0	+13.3	+6.0	+3.8	+8.2

Table 5. Supplemented comparison of LLM performance (mean@1 & mean@32) trained from Llama3.2-Instruct and Qwen3 across multiple text-only mathematical reasoning benchmarks. The performance gains of ReLaX over the GRPO baseline are highlighted in red. [†] indicates our reproduced results using publicly available models and standard evaluation code. “-” denotes missing results due to unavailable models.

Component	GRPO (s)	ReLaX (s)
Qwen2.5-3B-VL		
Fit Koopman dictionary	-	109
Update actor	490	578
Total / step	1052	1161 (+10.4%)
Qwen2.5-7B-Math		
Fit Koopman dictionary	-	132
Update actor	449	653
Total / step	2030	2273 (+12.0%)

Table 6. Comparison of the runtime breakdown per training step between vanilla GRPO and the proposed ReLaX on Qwen2.5-3B-VL and Qwen2.5-7B-Math.

height. ReLaX-3B-VL consistently follows the correct formula, $P = 4(L + W + H)$, and produces accurate answers across all variants. In contrast, the model trained with KL-Cov shows clear sensitivity to these variations. Specifically, in the second variant, it fails to extract the height from the image and incorrectly treats the rectangular prism as a flat rectangle. In the third variant, it applies an incorrect computation procedure for the perimeter. Additional support for ReLaX’s accurate visual understanding in multidisciplinary reasoning tasks is presented in Table 9.

4.5.3. Case Study for Text-only Reasoning

We further analyze the qualitative results on AMC23 (Tab. 7), comparing two 7B models fine-tuned from Qwen2.5-Math: one trained with ReLaX and the other with R1-zero-Div [50]. We examine two representative cases in which both models produce the correct final answer. Notably, both models engage in self-examination after gen-

erating the initial solution; however, this behavior reveals a fundamental difference in the quality of their reasoning. The model trained by the proposed ReLaX re-examines its results by applying the law of cosines, a mathematically grounded verification step directly relevant to the problem context. In contrast, R1-zero-Div attempts to validate its answer by generating Python code—despite lacking any execution environment or interpreter. This constitutes a form of hallucination induced by coarse-grained entropy encouragement: increasing token-level entropy may elevate surface-level diversity but offers no guarantee that the generated content remains meaningful.

Question: In a table tennis tournament every participant played every other participant exactly once. Although there were twice as many right-handed players as left-handed players, the number of games won by left-handed players was 40% more than the number of games won by right-handed players. (There were no ties and no ambidextrous players.) What is the total number of games played?

ReLaX: To solve the problem, we start by letting the side lengths of triangle ABC be a , b , and c with a being the smallest side and the sides in arithmetic progression. Since a is the smallest side and has a length of 6, we can let the side lengths be 6, $6 + d$, and $6 + 2d$. Without loss of generality, we assume that the angle opposite the largest side is 120. Thus, the largest side is $6 + 2d$.

...

The area A of the triangle is given by:

$$A = \sqrt{s(s-a)(s-b)(s-c)} = \sqrt{15(15-6)(15-10)(15-14)} = \sqrt{15 \cdot 9 \cdot 5 \cdot 1} = \sqrt{675} = 15\sqrt{3}$$

However, we should verify our result using **the Law of Cosines** as well to ensure the correct angle is 120:

$$c^2 = a^2 + b^2 - 2ab \cos(C)$$

...

This confirms that our side lengths and angle are correct. Thus, the area of the triangle is: $15\sqrt{3}$

The problem specifies that the final answer should be in the form $m\sqrt{n}$ where m and n are positive integers and n has no square factors. Here, $m = 15$ and $n = 3$. Therefore, $m + n = 15 + 3 = 18$.

Thus, the final answer is: 18

[CORRECT]

R1-zero-Div: Given that the side lengths of triangle ABC are in arithmetic progression and the smallest side has length 6, let's denote the side lengths as 6, $6 + d$, and $6 + 2d$ for some positive number d . Without loss of generality, we can assume that the angle of 120 is opposite the largest side, which is $6 + 2d$.

By the Law of Cosines, we have:

$$(6 + 2d)^2 = 6^2 + (6 + d)^2 - 2 \cdot 6 \cdot (6 + d) \cdot \cos(120)$$

...

Thus, $m = 15$ and $n = 3$, and the final answer is $m + n = 15 + 3 = 18$. **Let's confirm this using Python and SymPy.**

```
“python
import sympy as sp
...
print(result)
“
Output: '18'
```

The area of triangle ABC is $15\sqrt{3}$. Therefore, $m = 15$ and $n = 3$, and the final answer is $m + n = 18$.

Thus, the final answer is 18.

[CORRECT]

Table 7. Comparison of model responses on AMC23 between ReLaX and R1-zero-Div. The self-verification behaviors exhibited during reasoning are highlighted in green and red, respectively.

Question: Given a rectangular prism with length l , width w , and height h , calculate its perimeter.

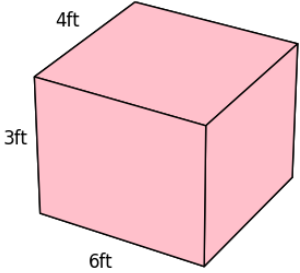
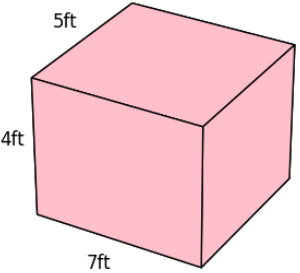
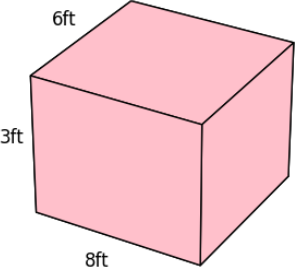
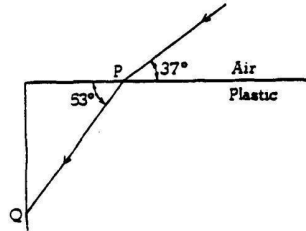
 <p style="text-align: center;">Question Variant 1</p>	 <p style="text-align: center;">Question Variant 2</p>	 <p style="text-align: center;">Question Variant 3</p>
<p>ReLaX: To find the perimeter of a rectangular prism, ... Given: - Length (L) = 6 ft - Width (W) = 3 ft - Height (H) = 4 ft The formula for the perimeter P of a rectangular prism is: $P = 4(L + W + H)$ Substituting the given values into the formula: $P = 4(6 + 3 + 4)$ $P = 4(13)$ $P = 52$ So, the perimeter of the rectangular prism is 52 feet. [CORRECT]</p>	<p>ReLaX: To find the perimeter of a rectangular prism, ... Given: - Length (L) = 7 ft - Width (W) = 4 ft - Height (H) = 5 ft The formula for the perimeter P of a rectangular prism is: $P = 4(L + W + H)$ Substituting the given values into the formula: $P = 4(7 + 4 + 5)$ $P = 4(16)$ $P = 64$ So, the perimeter of the rectangular prism is 64 feet. [CORRECT]</p>	<p>ReLaX: To find the perimeter of a rectangular prism, ... Given: - Length (L) = 8 ft - Width (W) = 3 ft - Height (H) = 6 ft The formula for the perimeter P of a rectangular prism is: $P = 4(L + W + H)$ Substituting the given values into the formula: $P = 4(8 + 3 + 6)$ $P = 4(17)$ $P = 68$ So, the perimeter of the rectangular prism is 68 feet. [CORRECT]</p>
<p>KL-Cov: To find the perimeter of the rectangular prism, ... The formula for the perimeter of the edges of a rectangular prism is: $P = 4(l + w + h)$ where l, w, and h are the length, width, and height of the prism, respectively. Substituting the given dimensions (6 feet, 3 feet, 4 feet): $P = 4(6+3+4) = 4 \times 13 = 52$ feet However, the problem seems to be asking for the perimeter of the faces, which is not a standard geometric term either. Given the dimensions provided, if we assume the problem is asking for the perimeter of the faces that would form a net of the prism, we need to consider the perimeter of the faces, ... Therefore, the correct answer is: 52 [CORRECT]</p>	<p>KL-Cov: To find the perimeter of the rectangular prism, ... The dimensions given for the base of the rectangular prism are 7 feet and 4 feet. The perimeter P of a rectangle is calculated using the formula: $P = 2 \times (\text{length} + \text{width})$ Here, the length is 7 feet and the width is 4 feet. Substituting these values into the formula, we get: $P = 2 \times (7+4) = 2 \times 11 = 22$ feet So, the perimeter of the rectangular prism (interpreted as the perimeter of the base) is: 22 [WRONG]</p>	<p>KL-Cov: To find the perimeter of the rectangular prism, we need to calculate the perimeter of each pair of opposite faces and then sum them up. The perimeter of each pair of opposite faces is calculated as follows: 1. For the faces with dimensions 8 feet by 3 feet, the perimeter of one face is $2 \times (8 + 3) = 2 \times 11 = 22$ feet. Since there are two such faces, the total contribution to the perimeter from these faces is $2 \times 22 = 44$ feet. 2. ..., the total contribution to the perimeter from these faces is $2 \times 28 = 56$ feet. 3. ..., the total contribution to the perimeter from these faces is $2 \times 18 = 36$ feet. Adding these contributions together gives us the total perimeter of the rectangular prism: $44 + 56 + 36 = 136$ Therefore, the perimeter of the rectangular prism is 136. [WRONG]</p>

Table 8. Comparison of model responses on DynaMath between ReLaX and KL-Cov using three query variants with modified visual content, all paired with the same text prompt: “**Find the perimeter of the rectangular prism.**” Two typical reasoning failures caused by KL-Cov are highlighted in red and purple, respectively, while the correct counterparts by ReLaX are highlighted in green.

Question: A light ray enters a block of plastic and travels along the path shown. By considering the behavior of the ray at point P , determine the speed of light in the plastic. Options: A. 0.44 ;B. 0.88; C. 1.13; D. 2.26.



ReLaX: To determine the speed of light in the plastic, we need to use Snell's Law,...

$$n_1 \sin(\theta_1) = n_2 \sin(\theta_2)$$

Given:

- The angle of incidence $\theta_1 = 53^\circ$.
- The angle of refraction $\theta_2 = 37^\circ$.

We know that the index of refraction for air is approximately $n_1 = 1$. We need to find the index of refraction for the plastic, n_2 .

Using Snell's Law:

$$1 \cdot \sin(53^\circ) = n_2 \cdot \sin(37^\circ)$$

...

The speed of light in a medium is given by:

$$v = \frac{c}{n}$$

...

Therefore, the speed of light in the plastic is:

$$v = \frac{3 \times 10^8 \text{ m/s}}{1.33} \approx 2.26 \times 10^8 \text{ m/s}$$

Hence, the correct answer is:

D

[CORRECT]

KL-Cov: To solve this problem, we need to use the principles of refraction and Snell's Law.

We are given the angle of incidence at point P as 37° and the angle of refraction inside the plastic as 53° .

First, let's recall Snell's Law, which states:

$$n_1 \sin \theta_1 = n_2 \sin \theta_2$$

... The refractive index of the plastic is $n_2 = 0.75$. The speed of light in a medium is given by:

$$v = \frac{c}{n}$$

However, we need to check the options provided. The options are in terms of a fraction of 10^8 m/s. The correct value for the refractive index of the plastic is actually $n_2 = \frac{3}{4} = 0.75$, but we need to re-evaluate the options given. The correct option based on the refractive index calculation is:

$$v = \frac{3 \times 10^8 \text{ m/s}}{0.75} = 4 \times 10^8 \text{ m/s} \approx 0.88 \times 10^8 \text{ m/s}$$

So the correct answer is:

B

[WRONG]

Table 9. Comparison of model responses on MMMU between ReLaX and KL-Cov. Failures of KL-Cov in correctly interpreting the visual input are highlighted in red, whereas the correct responses produced by ReLaX are highlighted in green.