

# Metaphor Reasoning is Meta-reasoning

Anonymous ACL submission

## Abstract

Metaphor reasoning is an essential cognitive ability that maps knowledge from familiar domains to more abstract domains. This ability functions as a meta-ability underlying many types of reasoning. However, existing work rarely investigates how metaphor reasoning affects other reasoning abilities. To bridge this gap, we systematically study how metaphor reasoning, particularly through metaphorical riddles, can enhance broader reasoning abilities in large language models. We propose METAR, an automated system for synthesizing metaphorical riddles that satisfy five quality dimensions: *diverse*, *balanced*, *reasoning-oriented*, *challenging*, and *verifiable*. Leveraging that answer categories determine riddle categories, we employ a hierarchical answer taxonomy for the former three criteria and a multi-agent refinement framework for the latter two, generating a high-quality dataset. Training with reinforcement learning on verifiable rewards using only thousands of metaphorical riddles, we demonstrate improvements across six out-of-distribution reasoning domains. Analysis reveals transfer effectiveness depends on model scale and pattern-target domain alignment. We will release our code and data.

## 1 Introduction

Metaphor reasoning is an essential cognitive ability (Lakoff, 1993). Through metaphors, humans are able to understand abstract concepts by mapping knowledge from familiar domains (i.e., source domains) to more abstract domains (i.e., target domains) (Lakoff and Johnson, 2024). As illustrated in Figure 1, comparing *ducks* that handle pests in rice fields to *wardens* conveys their role with only one word. Hence, metaphor reasoning can be considered as a meta-ability underlying many types of reasoning, such as abstract and analogical reasoning (Thibodeau and Boroditsky, 2011; Khatin-Zadeh et al., 2022).



Figure 1: An example of a *metaphorical riddle*.

However, existing work rarely investigates metaphor reasoning’s impact on reasoning abilities. Current research on metaphor reasoning focuses on detection (Tian et al., 2024; Chen et al., 2024), interpretation (Tong et al., 2024; Sanchez-Bayona and Agerri, 2025), creative generation (Chakrabarty et al., 2021; He et al., 2023a,b), or downstream tasks such as sentiment analysis (Li et al., 2022a), translation (Wang et al., 2024), jail-breaking (Yan et al., 2025), and word-games (Xu and Zhong, 2025). While recent work demonstrates that metaphor reasoning can enhance reasoning abilities (Kramer, 2025), this work lacks evaluation on public reasoning datasets. Therefore, metaphor reasoning’s impact on broader reasoning domains remains under-explored.

To investigate how metaphor reasoning affects reasoning abilities in broader domains, we first need to identify an appropriate task form. *Riddles* serve as a typical task form for metaphor reasoning (Panagiotopoulos et al., 2025; Le et al., 2025):

*For the essence of a riddle is to express true facts under impossible combinations. Now this cannot be done by any arrangement of ordinary words, but by the use of metaphor it can.*

—Poetics, Chapter XXII, by Aristotle

Specifically, riddles use ingenious and enigmatic language to describe a hidden answer, requiring

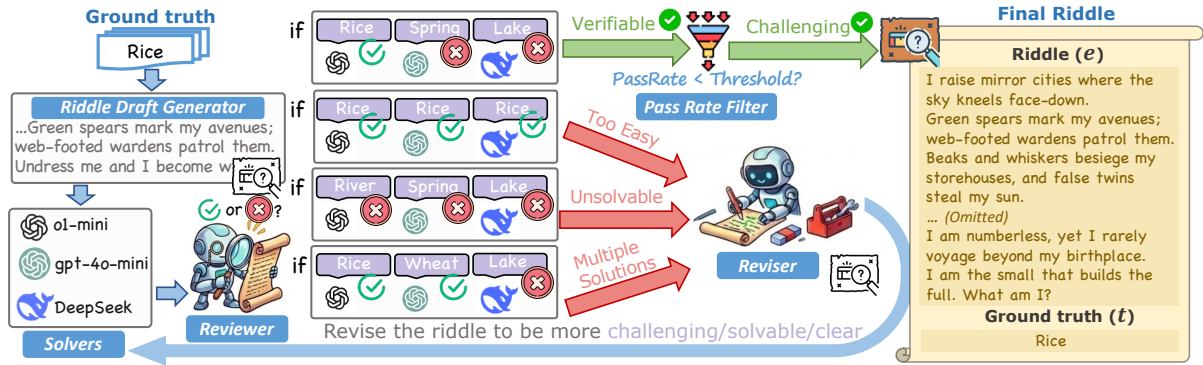


Figure 2: Multi-agent iterative refinement framework for riddle generation: (1) *Riddle Draft Generation* constructs initial riddles from metaphors given a ground truth answer; (2) *Solver Answer Collection* collects candidate answers from diverse solvers  $\mathcal{M}$ ; (3) *Reviewer Assessment* verifies *Verifiability* and *Challengingness* through quality assessment; (4) *Reviser Refinement* iteratively improves problematic riddles based on diagnostic outcomes; (5) *Pass Rate Filter* applies a final quality gate to ensure challengingness requirement.

solvers to infer the answer through deliberate reasoning. On one hand, *figurative language is common in riddles*, with 37.5% of riddles containing figurative language, of which 24% contain metaphors (Zhang and Wan, 2022). In our paper, we refer to riddles that contain metaphors as *metaphorical riddles*. On the other hand, *metaphorical riddle reasoning inherently involves multiple traditional metaphor reasoning tasks*, such as component identification (Li et al., 2022b) and metaphor interpretation (Tong et al., 2024). As shown in Figure 1, this metaphorical riddle contains multiple metaphorical clues that require mapping through metaphor reasoning to arrive at the final answer. Hence, we adopt metaphorical riddles as our task form to study how metaphor reasoning affects reasoning abilities in broader domains.

However, using metaphorical riddles to enhance models’ reasoning abilities faces several challenges: (1) *Scalable dataset construction*: Previous work typically collects metaphorical riddles through web scraping (Lin et al., 2021; Zhang and Wan, 2022), which suffers from limited scalability. (2) *Effective training recipe*: Existing methods either rely on prompting without improving model capabilities (Panagiotopoulos et al., 2025) or fine-tune small models on Question-Answer pairs, but fail to evaluate generalization to broader reasoning domains (Lin et al., 2021).

To bridge the gap, we systematically study how metaphor reasoning affects reasoning abilities in other domains using riddles as our task form. We propose METAR, an automated system for synthesizing metaphorical riddles with five key quality dimensions: *diverse, balanced, reasoning-oriented,*

*challenging, and verifiable*. Since each riddle describes a specific answer, the answer’s semantic category determines the riddle’s thematic category, enabling us to achieve these criteria through two mechanisms: a *riddle answer taxonomy* (3 hierarchical levels, 5,466 answers) ensuring the first three dimensions via category selection and popularity-based entity selection; and a *multi-agent iterative refinement framework* (Figure 2) ensuring the latter two dimensions. We adopt reinforcement learning with verifiable rewards (RLVR) (Guo et al., 2025; Yu et al., 2025) to enhance models’ metaphor reasoning abilities. Experiments show that training on only 3,444 metaphorical riddles improves reasoning in six out-of-distribution domains, demonstrating metaphor reasoning as a meta-reasoning ability and revealing the importance of model scale and domain alignment.

Overall, our contributions are as follows: (1) We are the first to systematically study whether metaphor reasoning can enhance models’ broader reasoning abilities. (2) We propose a scalable automated system for synthesizing metaphorical riddles, enabling large-scale generation of high-quality training data. (3) Through extensive experiments, we demonstrate that metaphor reasoning functions as a meta-ability and provide interpretable analysis of the underlying mechanisms.

## 2 Related Work

### 2.1 Metaphor Reasoning

Research on metaphor reasoning encompasses detection (Li et al., 2024; Tian et al., 2024; Chen et al., 2024), generation (Chakrabarty et al., 2021; He

et al., 2023a,b), and interpretation (He et al., 2022; Tong et al., 2024; Sanchez-Bayona and Agerri, 2025), with applications in sentiment analysis (Li et al., 2022a), translation (Wang et al., 2024), jail-breaking (Yan et al., 2025), and word games (Xu and Zhong, 2025). Metaphor reasoning may enhance general reasoning capabilities, as metaphors form the foundation of abstract reasoning (Lakoff and Johnson, 2024; Thibodeau and Boroditsky, 2011; Khatin-Zadeh et al., 2022). While conceptual metaphor-based prompting has shown promise to improve reasoning capabilities (Kramer, 2025), no prior work has investigated whether metaphor reasoning enhances model performance on established reasoning benchmarks.

## 2.2 Metaphor and Riddle Reasoning

Riddle reasoning datasets are typically collected from online sources and reformatted as question-answer (QA) questions (Lin et al., 2021; Jiang et al., 2024). Researchers have improved model performance through fine-tuning on QA datasets (Lin et al., 2021) and prompting techniques (Panagiotopoulos et al., 2025), but have not examined whether riddle reasoning training enhances reasoning capabilities in other domains.

Figurative language is common in riddles (Zhang and Wan, 2022). Panagiotopoulos et al. (2025) highlights metaphors when reconstructing in-context riddles, improving riddle-solving performance. Le et al. (2025) incorporates novel metaphors in riddle generation to increase difficulty. However, no prior work has used riddles as a medium to study how metaphor reasoning affects other reasoning capabilities.

## 2.3 Reinforcement Learning with Verifiable Rewards

Reinforcement Learning with Verifiable Rewards (RLVR) leverages objective, verifiable metrics as reward signals, enabling models to exhibit advanced deliberative thinking capabilities (Guo et al., 2025; Yu et al., 2025; Hu et al., 2025), such as planning and reflection, through extended reasoning chains. This approach has proven effective for complex reasoning tasks (Seed et al., 2025; Team et al., 2025; Chen et al., 2025). However, no prior work has applied RLVR to metaphor reasoning.

## 3 Task Formulation

We formulate metaphorical riddle reasoning as a question-answering (QA) task. Formally, let

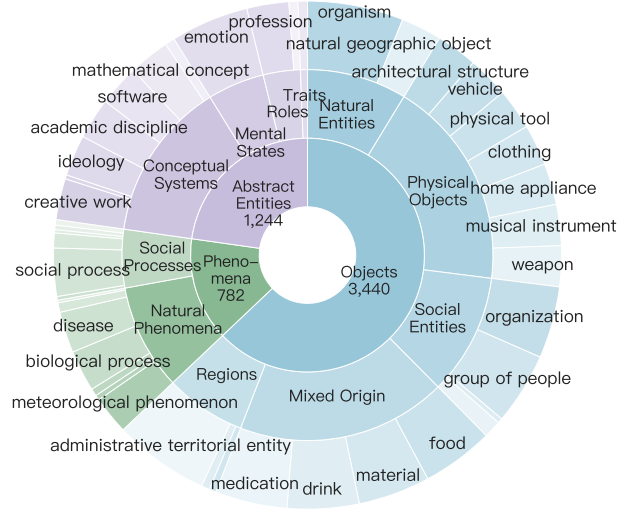


Figure 3: The statistics of the Riddle Answer Taxonomy.

$\mathcal{D} = \{(e_i, t_i)\}_{i=1}^{|\mathcal{D}|}$  denote a dataset of metaphorical riddles, where each pair  $(e, t)$  consists of a riddle  $e$  (which serves as the question) and its corresponding ground truth answer  $t$ . Given a riddle  $e$ , the task is to predict the answer  $\hat{t}$  that matches the ground truth  $t$ . An example is shown in Figure 2.

## 4 METAR

We propose METAR, an automated system for synthesizing metaphorical riddles. We establish design criteria for high-quality riddle datasets and present mechanisms to achieve them.

### 4.1 Design Criteria

Our riddle dataset satisfies five quality dimensions: (1) *diverse*, ensuring model generalization across different domains; (2) *balanced*, avoiding bias toward specific entity types; (3) *reasoning-oriented*, requiring metaphorical reasoning rather than knowledge retrieval; (4) *challenging*, incentivizing models' elaborate reasoning processes; (5) *verifiable*, providing reliable supervision signals for subsequent training.

Following §3, the *answer*  $t$  is the target entity, and the *question* (riddle)  $e$  is its metaphorical description. Since each riddle describes a specific answer, the answer's category determines the riddle's category. Hence, we achieve these criteria via two mechanisms: the first three through answer taxonomy design with category and popularity-based entity selection; the latter two through a multi-agent iterative refinement framework.

| $H_1$             | $H_2$              | $H_3$                     | Example Answers                              |
|-------------------|--------------------|---------------------------|--|
| Objects           | Mixed Origin       | drink                     | coffee, tea, milk, beer                      |
| Objects           | Natural Entities   | natural geographic object | mountain, ocean, valley, volcano             |
| Phenomena         | Natural Phenomena  | astronomical phenomenon   | eclipse, meteor, sunset                      |
| Phenomena         | Social Processes   | historical event          | battle, armistice, ceasefire                 |
| Abstract Entities | Conceptual Systems | ideology                  | nationalism, conservatism, fascism, feminism |
| Abstract Entities | Roles              | social status             | citizen, prisoner, slave, reputation         |

Table 1: Examples of riddle answers  $t$  organized by taxonomic hierarchy.

## 4.2 Riddle Answer Taxonomy

To achieve the criteria of *diverse*, *balanced*, and *reasoning-oriented*, we construct a three-level answer taxonomy ( $H_i$  for  $i \in \{1, 2, 3\}$ ), with statistics in Figure 3 and examples in Table 1<sup>1</sup>.

**Diverse Coverage.** To partition entities into fundamental categories at the top level  $H_1$ , we adopt UFO (Guizzardi et al., 2022)’s foundational categories: *Material Object*, *Phenomenon*, *Abstract Entities*. Since manually enumerating all relevant domain-specific subcategories would be impractical, we use GPT-5 (OpenAI, 2025b)<sup>2</sup> to generate second-level  $H_2$  categories. To ensure comprehensive coverage, we leverage Wikidata’s *subclass-of* relationship: for each  $h_2 \in H_2$ , we generate seeds  $\mathcal{S}_{h_2} = \{s_1, s_2, \dots, s_{n_{h_2}}\}$  where each seed  $s$  corresponds to a Wikidata Q-id  $Q_s$ .

**Common Answers.** To ensure *reasoning-oriented* criterion, we select *common* well-known entities, avoiding obscure answers that test knowledge rather than metaphor reasoning. For each seed  $s$  with Q-id  $Q_s$ , we employ two-stage sorting: (1) retrieve candidates via *subclass-of*, sort by *sitelinks* (cross-language Wikipedia links), select top candidates as  $C_s$ ; (2) then sort  $C_s$  by *popularity* using QRANK<sup>3</sup>(Arora et al., 2024), ranking entities by Wikimedia page views.

**Quality Control and Taxonomic Balance.** To ensure answer quality, we apply two filters: (1)

<sup>1</sup>Please refer to Appendix A.1.4 for more examples.

<sup>2</sup>Specifically, the version of the adopted model is high version of GPT-5-2025-08-17.

<sup>3</sup><https://qrank.toolforge.org/>

*Conciseness*: remove entities with labels exceeding 2 words, as overly complex answers make it difficult to satisfy the *verifiable* criterion; (2) *Quality threshold*: filter entities with sitelinks below threshold  $\tau$ , as they tend to be obscure. An entity  $x$  is valid,  $\text{Valid}(x)$ , if its label  $\leq 2$  words and sitelinks  $> \tau$ . To ensure *balanced* criterion, we allocate fixed quota  $K$  to each  $h_2 \in H_2$ , uniformly distributed among seeds  $\mathcal{S}_{h_2}$ . For  $h_2$  with  $n_{h_2}$  seeds, quota  $k_s$  per seed  $s \in \mathcal{S}_{h_2}$  is:

$$k_s = \min \left( \left\lfloor \frac{K}{n_{h_2}} \right\rfloor, |\{x \in C_s : \text{Valid}(x)\}| \right) \quad (1)$$

## 4.3 Riddle Generation

We design a multi-agent iterative refinement framework to achieve *challenging* and *verifiable* criteria. Figure 2 illustrates the framework. It takes ground truth answers  $t \in \mathcal{T}$  from the Riddle Answer Taxonomy and refines riddles through five stages. The pseudocode can refer to Algorithm 1 in Appendix A.1, and examples of generated riddles are provided in Appendix A.1.4.

**Stage 1: Riddle Draft Generation.** To obtain initial drafts, we require: (1) *metaphorical* language obscuring the answer; (2) sufficient *challenging* through multi-layered metaphors. Our prompt template (Table 4) includes: (1) example  $e_t$  demonstrating metaphorical reasoning patterns; (2) Wikidata descriptions  $d_t$  for  $t$ , providing contextual information. Formally, for  $t \in \mathcal{T}$ , we generate draft  $e^{(0)}$  using model  $M_g$ :  $e^{(0)} = M_g(t, e_t, d_t)$ .

**Stage 2: Solvers’ Answer Collection.** To prepare for subsequent assessment stages, we collect answers from a diverse set of solvers. We employ solver set  $\mathcal{M} = \{M_s^{(1)}, M_s^{(2)}, \dots, M_s^{(k)}\}$ , where each  $M_s^{(i)} \in \mathcal{M}$  represents different *architectural origins* (model families, training paradigms) and *capability levels*. To prevent bias, solvers must be *distinct* from generator, reviewer, and reviser models. Each  $M_s^{(i)} \in \mathcal{M}$  independently generates answer  $a_i \leftarrow M_s^{(i)}. \text{Solve}(e^{(r)})$  for draft  $e^{(r)}$  at round  $r$ , forming  $\mathcal{A}^{(r)} = \{a_1, a_2, \dots, a_k\}$ .

**Stage 3: Reviewer Assessment.** To assess riddle quality along *verifiability* and *challengingness*, we employ two steps: (1) reviewer independently evaluates each solver’s answer based solely on the riddle, without ground truth access; (2) then, independent judgments are cross-validated against ground truth to diagnose issues (excessive simplicity, unsolvability, multiple solutions).

**Reviewer Independent Assessment.** Reviewer  $M_r$  receives riddle  $e^{(r)}$  and answer set  $\mathcal{A}^{(r)} = \{a_1, a_2, \dots, a_k\}$  from Stage 2. For each  $a_i \in \mathcal{A}^{(r)}$ , reviewer evaluates correctness given only  $e^{(r)}$ , producing  $c_i = M_r(a_i, e^{(r)}) \in \{0, 1\}$ , yielding  $\mathcal{C}^{(r)} = \{c_1, c_2, \dots, c_k\}$ .

**Reviewer Cross-Validation.** A riddle is *solvable* if at least one answer is judged correct and matches ground truth:

$$\text{Solvable}(e^{(r)}, t) \iff \exists i : (c_i = 1) \wedge (a_i = t) \quad (2)$$

For *Pass (PASS)* status, a riddle must satisfy: (1) *Verifiable* and (2) *Challenging*:

$$\begin{cases} \text{Verifiable}(e^{(r)}, t) \iff \text{Solvable}(e^{(r)}, t) \wedge \\ \quad \neg(\exists i : (c_i = 1) \wedge (a_i \neq t)) \\ \text{Challenging}(e^{(r)}, t) \iff \text{PassRate}(e^{(r)}, t) \geq \theta_p \end{cases} \quad (3)$$

where  $\text{PassRate}(e^{(r)}, t)$  is the rate of answers both judged correct and matching ground truth:

$$\text{PassRate}(e^{(r)}, t) = \frac{|\{i : (c_i = 1) \wedge (a_i = t)\}|}{k} \quad (4)$$

Cross-validation yields four mutually exclusive outcomes: (1) *Multiple Solutions (MULTI)*: at least one answer judged correct but  $\neq t$ , indicating ambiguity; (2) *Too Easy (EASY)*: all answers =  $t$ , indicating insufficient challenge; (3) *Unsolvable (UNSOLV)*: no answer judged correct or no correct answer matches  $t$ , indicating excessive obscurity; (4) *Pass (PASS)*: solvable and challenging, indicating single unambiguous answer requiring genuine reasoning. Formally:

$$\begin{cases} \text{MULTI}(e^{(r)}, t) \iff \exists i : (c_i = 1) \wedge (a_i \neq t) \\ \text{EASY}(e^{(r)}, t) \iff \forall i : a_i = t \\ \text{UNSOLV}(e^{(r)}, t) \iff \neg \text{Solvable}(e^{(r)}, t) \\ \text{PASS}(e^{(r)}, t) \iff \text{Verifiable}(e^{(r)}, t) \wedge \\ \quad \text{Challenging}(e^{(r)}, t) \end{cases} \quad (5)$$

These four outcomes provide clear signals for the next refinement stage.

**Stage 4: Reviser Refinement.** Reviser  $M_v$  refines problematic riddles based on Stage 3 outcomes, addressing ambiguity, insufficient challenge, or unsolvability. Formally, given  $e^{(r)}$  and outcome  $d^{(r)} \in \{\text{MULTI}, \text{EASY}, \text{UNSOLV}\}$ , reviser generates  $e^{(r+1)} = M_v(e^{(r)}, d^{(r)})$ . Then, the refined riddle  $e^{(r+1)}$  returns to Stage 2, establishing an iterative loop until PASS status or exceeding  $R_{\max}$  rounds.

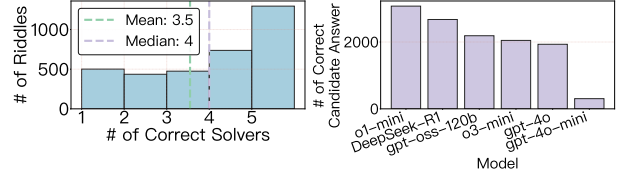


Figure 4: The statistics of valid generated riddles. The distribution of correct solver answer counts per riddle (left) and the correctness of each riddle across different solvers (right).

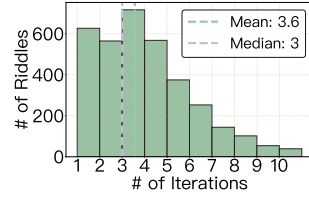


Figure 5: The distribution of iteration counts per riddle during generation.

**Stage 5: Pass Rate Filter.** We apply a final filter to ensure challengingness: only riddles with  $\text{PassRate}(e^{(r)}, t) \geq \theta_p$  (Equation 4) are retained; others are rejected.

Riddles passing all five stages are *valid riddles*, forming well-verified reasoning tasks. We adopt GPT-5 as generator, reviewer, and reviser<sup>4</sup>. With regard to the solver set<sup>5</sup>, we adopt o1-mini (OpenAI, 2024), Deepseek-R1 (Guo et al., 2025), GPT-oss-120b (OpenAI, 2025a), GPT-4o (Hurst et al., 2024), GPT-4o-mini (Hurst et al., 2024), and o3-mini (OpenAI, 2025c). Several parameters control riddle difficulty: (1) *Solver Set*: stronger solvers produce harder riddles; (2) *Pass-rate Threshold*: lower thresholds produce harder riddles. Detailed prompts and hyperparameter settings for each agent can refer to Appendix A.1.

#### 4.4 Riddle Statistics

**Riddle Answer Taxonomy Statistics.** Figure 3 presents taxonomy statistics: 5,466 riddle answers demonstrating (1) *diversity*: 3  $H_1$ , 11  $H_2$ , 44  $H_3$  categories; (2) *balance*: balanced  $H_3$  categories under each  $H_2$ .

**Riddle Generation Statistics.** Using six solvers, we generated 3,444 valid riddles satisfying all five

<sup>4</sup>Specifically, the version of the adopted models are medium version of GPT-5-2025-08-17.

<sup>5</sup>Specifically, the version of the adopted solvers are o1-mini-2024-09-12, GPT-4o-2024-11-20, GPT-4o-mini-2024-07-18

quality dimensions<sup>6</sup>. Figure 4 (left) shows the distribution of correct solver answer counts per riddle, indicating varying difficulty; median/mean suggest moderate overall difficulty. Figure 4 (right) shows solver performance: o1-mini perform the best while GPT-4o-mini perform the worst. Finally, Figure 5 shows iteration count distribution.

## 5 Training Recipe

We adopt the DAPO algorithm (Yu et al., 2025) for reinforcement learning with verifiable rewards (RLVR), which has demonstrated significant improvements in reasoning (Guo et al., 2025; Yu et al., 2025). The loss function is:

$$\mathcal{J}_{\text{DAPO}}(\theta) = \mathbb{E}_{(e,t) \sim \mathcal{D}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|e)} \left[ \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{j=1}^{|o_i|} \min \left( r_{i,j}(\theta) \hat{A}_{i,j}, \text{clip} \left( r_{i,j}(\theta), 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}} \right) \hat{A}_{i,j} \right) \right] \quad (6)$$

where  $\theta$  denotes model parameters,  $\mathcal{D}$  is the training dataset,  $(e, t)$  is a riddle-answer pair,  $G$  is the group size,  $\pi_{\theta_{\text{old}}}$  is the old policy for importance sampling,  $o_i$  is the  $i$ -th output sequence of length  $|o_i|$ ,  $r_{i,j}(\theta) = \frac{\pi_{\theta}(o_{i,j}|e, o_{i,<j})}{\pi_{\theta_{\text{old}}}(o_{i,j}|e, o_{i,<j})}$  is the importance sampling ratio at position  $j$ ,  $\hat{A}_{i,j}$  is the advantage function, and  $\varepsilon_{\text{low}}$  and  $\varepsilon_{\text{high}}$  are clipping parameters. DAPO’s *dynamic sampling* filters out groups with uniformly zero or maximum rewards, enabling effective and stable learning signals.

The advantage function uses group normalization:  $\hat{A}_{i,j} = R_i - \bar{R}$ , where  $R_i \in \{0, 1\}$  is the reward for output  $i$ , and  $\bar{R} = \frac{1}{G} \sum_{i=1}^G R_i$  is the baseline. Since generated riddles are verifiable, rewards are computed via exact match (1 if extracted answer matches ground truth, 0 otherwise).

## 6 Experiments

### 6.1 Experiment Setup

**Benchmark.** To evaluate the impact of Metaphor Reasoning on reasoning abilities in other domains, we employ six categories of Out-of-Distribution (OOD) reasoning across 7 benchmarks: Logical Reasoning (Ma et al.; Chen et al., 2025); Common-sense Reasoning (Talmor et al., 2019); Natural Language Inference (Zellers et al., 2019); Math Rea-

soning<sup>7</sup>; Science, Technology, Engineering, and Mathematics (STEM) Reasoning (Rein et al.); Out-of-Distribution (OOD) Riddle Reasoning (Lin et al., 2021). Detailed descriptions of each benchmark are provided in Appendix A.2.1.

**Backbones.** We employ reasoning models as the backbones for reinforcement learning (RL) training, as these models have been pre-trained with explicit reasoning capabilities that provide a solid foundation for further enhancement through metaphor reasoning (Guo et al., 2025). To comprehensively evaluate the generalization of our approach, we investigate models across different *parameter scales* and *architectural generations*, including Qwen3-8B (Yang et al., 2025), Qwen3-14B (Yang et al., 2025), and QwQ-32B (Team, 2025), which is a reasoning model training on Qwen2.5-32B-instruct (Qwen et al., 2025).

The implementation details of training and evaluation are provided in Appendix A.2.2.

### 6.2 Experiment Results

**Metaphor reasoning = Meta-reasoning?.** According to Table 2, Metaphor reasoning can broadly enhance the reasoning abilities of models across diverse reasoning domains. Remarkably, training on only 3,444 metaphorical riddles yields substantial improvements across six out-of-distribution reasoning domains, demonstrating the efficiency and meta-reasoning nature of metaphor reasoning. This small-scale training dataset highlights that metaphor reasoning functions as a meta-ability that can be effectively transferred with minimal data. First, QwQ-32B exhibits remarkable improvements across all six reasoning categories relative to its backbone. Also, models of different scales and versions also demonstrate improvements, further validating the generalization of our approach. With regard to overall performance, QwQ-32B<sub>MetaR</sub> surpasses significantly larger models such as DeepSeek-R1 (671B parameters), as well as closed-source models such as GPT-4o. Notably, QwQ-32B<sub>MetaR</sub> outperforms GPT-oss-120B on four benchmarks across reasoning domains.

**Model Scale.** Generalization depends critically on model scale. QwQ-32B shows consistent improvements across all categories, while Qwen3-14B exhibits selective gains with slight declines in

<sup>6</sup>The reduction from 5,466 answers to 3,444 riddles occurs because some answers fail to generate qualifying riddles even after  $R_{\text{max}}$  refinement rounds and are thus discarded.

<sup>7</sup>[https://artofproblemsolving.com/wiki/index.php/AIME\\_Problems\\_and\\_Solutions](https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions)

| Model                       | Math Reasoning        | STEM Reasoning        | Commonsense Reasoning | NL Inference           | Logical Reasoning      |                       | OOD Riddle Reasoning  | Overall               |
|-----------------------------|-----------------------|-----------------------|-----------------------|------------------------|------------------------|-----------------------|-----------------------|-----------------------|
|                             | AIME2024              | GPQA                  | CommonQA              | HellaSwag              | Enigmata               | KORBench              | RiddleSense           |                       |
| <b>Closed-Source Models</b> |                       |                       |                       |                        |                        |                       |                       |                       |
| GPT-4o-mini                 | 9.8                   | 42.4                  | 83.4                  | 85.0                   | 10.6                   | 43.2                  | 82.1                  | 50.9                  |
| GPT-4o                      | 11.9                  | 46.5                  | 84.3                  | 90.8                   | 21.8                   | 51.6                  | 88.7                  | 56.5                  |
| o1-mini                     | 87.7                  | 75.6                  | 85.3                  | 88.5                   | 89.0                   | 61.7                  | 94.4                  | 83.2                  |
| o3                          | 89.5                  | 83.1                  | 87.2                  | 93.7                   | 92.7                   | 64.2                  | 95.3                  | 86.5                  |
| GPT-5                       | 93.7                  | 89.4                  | 88.0                  | 92.9                   | 89.4                   | 67.2                  | 94.1                  | 87.8                  |
| <b>Open-Source Models</b>   |                       |                       |                       |                        |                        |                       |                       |                       |
| Qwen3-8B                    | 64.2                  | 59.2                  | 83.0                  | 79.3                   | 43.2                   | 51.4                  | 81.3                  | 65.9                  |
| QwQ-32B                     | 68.2                  | 60.0                  | 85.0                  | 78.5                   | 30.0                   | 66.0                  | 84.6                  | 67.5                  |
| Qwen3-14B                   | 69.0                  | 63.3                  | 84.2                  | 86.7                   | 44.5                   | 56.2                  | 85.2                  | 69.9                  |
| DeepSeek-R1                 | 77.4                  | 69.2                  | 69.9                  | 85.1                   | 49.8                   | 62.8                  | 79.2                  | 70.3                  |
| GPT-oss-120B                | 78.9                  | 71.2                  | 84.4                  | 83.1                   | 87.5                   | 63.6                  | 90.6                  | 79.9                  |
| <b>Ours</b>                 |                       |                       |                       |                        |                        |                       |                       |                       |
| Qwen3-8B <sub>MetaR</sub>   | 64.0 <sup>-0.2%</sup> | 56.0 <sup>-3.2%</sup> | 84.2 <sup>+1.2%</sup> | 81.2 <sup>+1.9%</sup>  | 41.3 <sup>-1.9%</sup>  | 51.4 <sup>+0.0%</sup> | 84.3 <sup>+3.0%</sup> | 66.1 <sup>+0.3%</sup> |
| Qwen3-14B <sub>MetaR</sub>  | 68.9 <sup>-0.1%</sup> | 62.9 <sup>-0.4%</sup> | 84.6 <sup>+0.4%</sup> | 87.7 <sup>+1.0%</sup>  | 51.8 <sup>+7.3%</sup>  | 57.6 <sup>+1.4%</sup> | 88.3 <sup>+3.1%</sup> | 71.7 <sup>+2.6%</sup> |
| QwQ-32B <sub>MetaR</sub>    | 74.4 <sup>+6.2%</sup> | 63.6 <sup>+3.6%</sup> | 86.1 <sup>+1.1%</sup> | 88.5 <sup>+10.0%</sup> | 40.3 <sup>+10.3%</sup> | 70.4 <sup>+4.4%</sup> | 91.4 <sup>+6.8%</sup> | 73.5 <sup>+8.9%</sup> |

Table 2: Main experimental results comparing baseline models and MetaR-enhanced models across different reasoning tasks. The superscript differences represent changes relative to the backbone model.

| Model                         | AIME2024              | GPQA                  | CommonQA              | HellaSwag             | Enigmata              | KORBench              | RiddleSense           | Overall               |
|-------------------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| Qwen3-14B <sub>MetaR</sub>    | 68.9                  | 62.9                  | 84.6                  | 87.7                  | 51.8                  | 57.6                  | 88.3                  | 71.7                  |
| Qwen3-14B <sub>MetaR-DS</sub> | 67.2 <sup>-1.7%</sup> | 63.0 <sup>+0.1%</sup> | 84.2 <sup>-0.4%</sup> | 86.2 <sup>-1.5%</sup> | 50.9 <sup>-0.9%</sup> | 55.6 <sup>-2.0%</sup> | 86.1 <sup>-2.2%</sup> | 69.8 <sup>-1.9%</sup> |
| Qwen3-8B <sub>MetaR</sub>     | 64.0                  | 56.0                  | 84.2                  | 81.2                  | 41.3                  | 51.4                  | 84.3                  | 66.1                  |
| Qwen3-8B <sub>MetaR-DS</sub>  | 58.5 <sup>-5.5%</sup> | 53.2 <sup>-2.8%</sup> | 83.6 <sup>-0.6%</sup> | 80.4 <sup>-0.8%</sup> | 42.3 <sup>+1.0%</sup> | 52.8 <sup>+1.4%</sup> | 83.7 <sup>-0.6%</sup> | 63.7 <sup>-2.4%</sup> |

Table 3: Dynamic ablation experimental results comparing MetaR models with and without dynamic reasoning.

Math and STEM. The 8B model shows limited generalization, improving only in OOD riddle reasoning, commonsense reasoning, and natural language inference. This scale-dependent pattern indicates that larger models better internalize and transfer metaphor reasoning strategies.

**Reasoning Domain.** Metaphor reasoning training exhibits differential effectiveness across reasoning domains. As shown in Table 2, OOD riddle reasoning shows the strongest improvements across all model scales, while commonsense reasoning and natural language inference also show consistent improvements. Math and STEM reasoning show gains only in larger models, suggesting greater reliance on domain-specific knowledge.

**Training Recipe.** Dynamic sampling is crucial for RL training. According to Table 3, removing dynamic sampling leads to performance degradation across different domains. We analyze the underlying reasons as follows. First, as shown in Figure 6, dynamic sampling yields smoother reward curves. Second, dynamic sampling increases entropy, ensuring greater sampling diversity (Cui

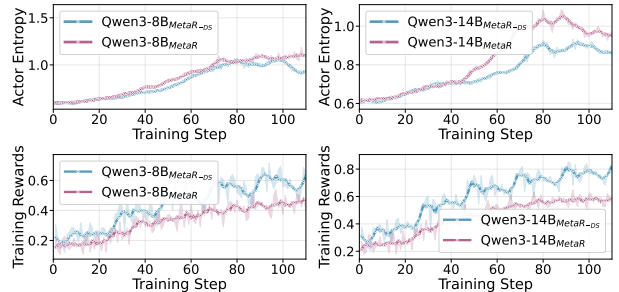


Figure 6: The training dynamics of models with and without dynamic sampling (indicated by “-DS”).

et al., 2025). Moreover, we found that dynamic sampling produces shorter responses, improving token efficiency (Luo et al., 2025) (Figure 7).

### 6.3 Analysis

We analyze the underlying mechanisms of metaphor reasoning training on general reasoning capabilities from two perspectives.

**Reasoning Domain Similarity.** To analyze why generalization results vary across domains, we take Qwen3-8B as an example and examine the outputs

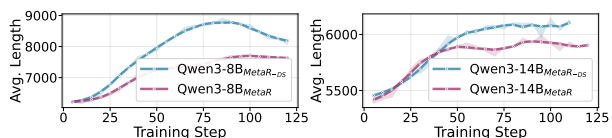


Figure 7: The average response length of models on the test set across different training steps.

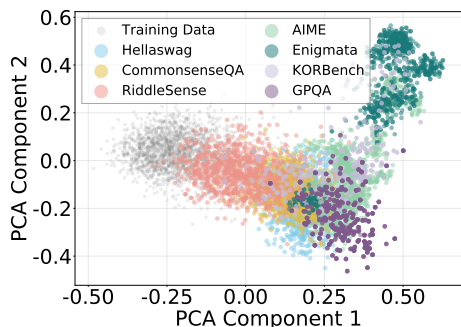


Figure 8: Response similarity analysis across different reasoning domains.

of the final model checkpoint on the training set and test sets from both embedding and vocabulary overlap perspectives. First, we perform semantic similarity analysis using Principal Component Analysis (PCA) visualization (Abdi and Williams, 2010)<sup>8</sup>. According to Figure 8, RiddleSense responses are closest to the training set in the semantic embedding space, while AIME and GPQA responses are farthest, indicating that RiddleSense shares more similar reasoning patterns with the metaphor reasoning training data. Also, we compute vocabulary overlap using Jaccard similarity (Jaccard, 1912) and Dice Coefficient (Dice, 1945), and according to Figure 9, the lexical patterns further validate this alignment. In conclusion, domains with aligned patterns benefit directly from metaphor reasoning training, while those requiring different approaches need larger model capacity to adapt.

**Reasoning Pattern Evolution.** We analyze the top-10 words with the largest frequency increases after metaphor reasoning training (Figure 10). These words fall into three key categories: *reflection* (e.g., “check”, “not”), *perspective switching* (e.g., “let’s”, “maybe”), and *careful deliberation* (e.g., “think”, “so”). This lexical shift demonstrates that metaphor reasoning training stimulates the model’s deep thinking capabilities.

<sup>8</sup>Sentence-transformers embeddings (Reimers and Gurevych, 2020) are used for semantic similarity analysis. Component 1 and Component 2 capture the largest variance in response embeddings for dimensionality reduction.

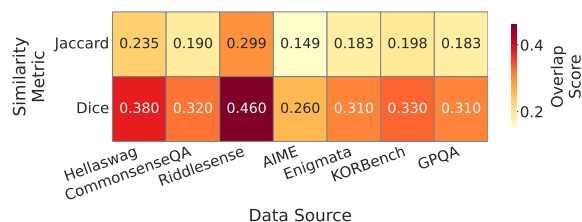


Figure 9: Vocabulary overlap analysis across different reasoning domains. The overlap score measures lexical similarity. Higher scores indicate greater shared vocabulary between domains.

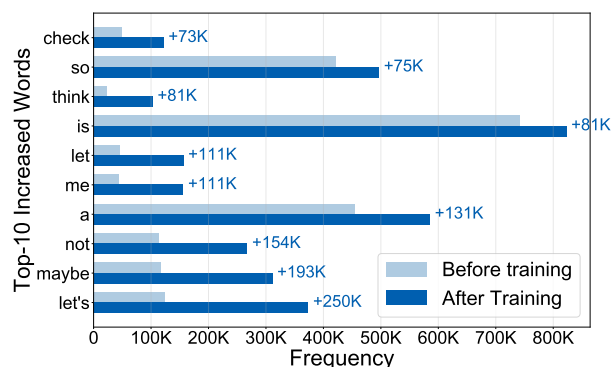


Figure 10: Top-10 increased words after metaphor reasoning training. K represents thousands. The counts are based on the number of words separated by spaces.

## 7 Conclusion

We present the first systematic study of metaphor reasoning’s impact on broader reasoning capabilities in LLMs. We propose METAR, an automated framework synthesizing high-quality metaphorical riddles via answer taxonomy and multi-agent refinement. Training on 3,444 metaphorical riddles improves performance across six out-of-distribution reasoning domains.

Despite using only riddle solving with thousands of examples, our experiments demonstrate that metaphor reasoning functions as a meta-reasoning ability that transfers effectively with minimal data. Our work offers an important perspective: rather than developing task-specific solutions, we should identify core meta-abilities that generalize across diverse scenarios. Once properly trained, such meta-abilities serve as powerful transferable skills, reducing computational costs and data annotation requirements. Future work will explore additional task formulations beyond riddles to further investigate metaphor reasoning’s transfer potential.

## 527 Limitations

528 While our work demonstrates the effectiveness  
529 of metaphor reasoning as a meta-reasoning abil-  
530 ity, several limitations should be acknowledged.  
531 First, despite our multi-agent iterative refinement  
532 framework, guaranteeing completely unique an-  
533 swers for every riddle remains challenging due  
534 to natural language ambiguity. However, our  
535 reviewer agent employs independent assessment  
536 and cross-validation (§4.3), ensuring only riddles  
537 passing the *Verifiable* criterion are included, with  
538 experimental results demonstrating effectiveness  
539 (Table 2). Second, while our evaluation cov-  
540 ers six reasoning categories across 7 benchmarks,  
541 important domains like causal or temporal rea-  
542 soning are excluded. However, our evaluation  
543 breadth—from abstract logical to concrete mathe-  
544 matical reasoning—provides strong evidence for  
545 the meta-reasoning nature. Finally, experiments  
546 are conducted exclusively on the Qwen family  
547 (Qwen3-8B, Qwen3-14B, QwQ-32B), though other  
548 architectures like Llama and DeepSeek (Guo et al.,  
549 2025) exist. However, comprehensive validation  
550 across different versions, parameter scales (8B,  
551 14B, 32B), and training paradigms provides robust  
552 evidence for generalization.

## 553 Acknowledgments

554 We acknowledge the use of Cursor ([https://](https://github.com/cursor/cursor)  
555 [github.com/cursor/cursor](https://github.com/cursor/cursor)) as an AI-assisted  
556 writing tool in the preparation of this manuscript.  
557 Specifically, Cursor assisted with writing polish-  
558 ment for the initial draft of this paper. Additionally,  
559 for the experiments section, Cursor helped gener-  
560 ate code for data visualization and figure genera-  
561 tion. During the development of the riddle answer  
562 taxonomy in §4 and the multi-agent framework,  
563 Cursor assisted with debugging and code develop-  
564 ment. The core ideas, framework design, and ex-  
565 perimental design of this work were independently  
566 conceived and developed by the authors.

## 567 References

568 Hervé Abdi and Lynne J Williams. 2010. Principal  
569 component analysis. *Wiley interdisciplinary reviews:*  
570 *computational statistics*, 2(4):433–459.

571 Abhishek Arora, Emily Silcock, Melissa Dell, and Le-  
572 ander Heldring. 2024. Contrastive entity coreference  
573 and disambiguation for historical texts. In *Proceed-*  
574 *ings of the 2024 Conference on Empirical Methods*  
575 *in Natural Language Processing*, pages 6174–6186.

Tuhin Chakrabarty, Xurui Zhang, Smaranda Muresan, 576  
and Nanyun Peng. 2021. *MERMAID: Metaphor gen-* 577  
*eration with symbolism and discriminative decoding.* 578  
In *Proceedings of the 2021 Conference of the North* 579  
*American Chapter of the Association for Computa-* 580  
*tional Linguistics: Human Language Technologies,* 581  
pages 4250–4261, Online. Association for Computa- 582  
tional Linguistics. 583

Jiangjie Chen, Qianyu He, Siyu Yuan, Aili Chen, 584  
Zhicheng Cai, Weinan Dai, Hongli Yu, Qiyang Yu, 585  
Xuefeng Li, Jiase Chen, and 1 others. 2025. Enig- 586  
mata: Scaling logical reasoning in large language 587  
models with synthetic verifiable puzzles. In *The* 588  
*Thirty-Ninth Annual Conference on Neural Informa-* 589  
*tion Processing Systems.* 590

Puli Chen, Cheng Yang, and Qingbao Huang. 2024. 591  
*Merely judging metaphor is not enough: Research* 592  
*on reasonable metaphor detection.* In *Findings of the* 593  
*Association for Computational Linguistics: EMNLP* 594  
*2024*, pages 5850–5860, Miami, Florida, USA. Asso- 595  
ciation for Computational Linguistics. 596

Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, 597  
Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, 598  
Huayu Chen, Weize Chen, Zhiyuan Liu, Hao Peng, 599  
Lei Bai, Wanli Ouyang, Yu Cheng, Bowen Zhou, and 600  
Ning Ding. 2025. *The entropy mechanism of rein-* 601  
*forcement learning for reasoning language models.* 602  
*CoRR*, abs/2505.22617. 603

Lee R Dice. 1945. Measures of the amount of ecologic 604  
association between species. *Ecology*, 26(3):297– 605  
302. 606

Giancarlo Guizzardi, Alessandro Botti Benevides, Clau- 607  
denir M. Fonseca, Daniele Porello, João Paulo A. 608  
Almeida, and Tiago Prince Sales. 2022. *UFO: unified* 609  
*foundational ontology.* *Appl. Ontology*, 17(1):167– 610  
210. 611

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao 612  
Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi- 613  
rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. 614  
Deepseek-r1: Incentivizing reasoning capability in 615  
llms via reinforcement learning. *arXiv preprint* 616  
*arXiv:2501.12948.* 617

Qianyu He, Sijie Cheng, Zhixu Li, Rui Xie, and 618  
Yanghua Xiao. 2022. Can pre-trained language mod- 619  
els interpret similes as smart as human? In *Proceed-* 620  
*ings of the 60th Annual Meeting of the Association for* 621  
*Computational Linguistics (Volume 1: Long Papers),* 622  
pages 7875–7887. 623

Qianyu He, Xintao Wang, Jiaqing Liang, and Yanghua 624  
Xiao. 2023a. Maps-kb: A million-scale probabilistic 625  
simile knowledge base. In *Proceedings of the AAAI* 626  
*Conference on Artificial Intelligence*, volume 37, 627  
pages 6398–6406. 628

Qianyu He, Yikai Zhang, Jiaqing Liang, Yuncheng 629  
Huang, Yanghua Xiao, and Yunwen Chen. 2023b. 630  
*HAUSER: Towards holistic and automatic evaluation* 631

|     |  |     |
|-----|--|-----|
| 632 | of simile generation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12557–12572, Toronto, Canada. Association for Computational Linguistics.   |     |
| 633 |  |     |
| 634 |  |     |
| 635 |  |     |
| 636 |  |     |
| 637 | Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum. 2025. Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base model. In <i>The Thirty-Ninth Annual Conference on Neural Information Processing Systems</i> .  |     |
| 638 |  |     |
| 639 |  |     |
| 640 |  |     |
| 641 |  |     |
| 642 |  |     |
| 643 | Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. <i>arXiv preprint arXiv:2410.21276</i> .   |     |
| 644 |  |     |
| 645 |  |     |
| 646 |  |     |
| 647 |  |     |
| 648 | Paul Jaccard. 1912. The distribution of the flora in the alpine zone. <i>New phytologist</i> , 11(2):37–50.  |     |
| 649 |  |     |
| 650 | Yifan Jiang, Filip Ilievski, and Kaixin Ma. 2024. SemEval-2024 task 9: BRAINTEASER: A novel task defying common sense. In <i>Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)</i> , pages 1994–2008, Mexico City, Mexico. Association for Computational Linguistics.                                     |     |
| 651 |  |     |
| 652 |  |     |
| 653 |  |     |
| 654 |  |     |
| 655 |  |     |
| 656 | Omid Khatin-Zadeh, Hassan Banaruee, and Babak Yazdani-Fazlabadi. 2022. A cognitive perspective on basic generic metaphors and their specific-level realizations. <i>Polish Psychological Bulletin</i> , pages 60–65.   |     |
| 657 |  |     |
| 658 |  |     |
| 659 |  |     |
| 660 |  |     |
| 661 | Oliver Kramer. 2025. Conceptual metaphor theory as a prompting paradigm for large language models. <i>arXiv preprint arXiv:2502.01901</i> .  |     |
| 662 |  |     |
| 663 |  |     |
| 664 | George Lakoff. 1993. The contemporary theory of metaphor.  |     |
| 665 |  |     |
| 666 | George Lakoff and Mark Johnson. 2024. <i>Metaphors we live by</i> . University of Chicago press.   |     |
| 667 |  |     |
| 668 | Duy Le, Kent Ziti, Evan Girard-Sun, Sean O’Brien, Vasu Sharma, and Kevin Zhu. 2025. Filtering for creativity: Adaptive prompting for multilingual riddle generation in llms. <i>arXiv e-prints</i> , pages arXiv–2508.   |     |
| 669 |  |     |
| 670 |  |     |
| 671 |  |     |
| 672 | Yu Xi Li, Bo Peng, Yu-Yin Hsu, and Chu-Ren Huang. 2024. EmbodiedBERT: Cognitively informed metaphor detection incorporating sensorimotor information. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 16868–16876, Miami, Florida, USA. Association for Computational Linguistics.                   |     |
| 673 |  |     |
| 674 |  |     |
| 675 |  |     |
| 676 |  |     |
| 677 |  |     |
| 678 |  |     |
| 679 | Yucheng Li, Frank Guerin, and Chenghua Lin. 2022a. The secret of metaphor on expressing stronger emotion. In <i>Proceedings of the 3rd Workshop on Figurative Language Processing (FLP)</i> , pages 39–43, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.  |     |
| 680 |  |     |
| 681 |  |     |
| 682 |  |     |
| 683 |  |     |
| 684 |  |     |
|     | Yucheng Li, Chenghua Lin, and Frank Guerin. 2022b. CM-gen: A neural framework for Chinese metaphor generation with explicit context modelling. In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 6468–6479, Gyeongju, Republic of Korea. International Committee on Computational Linguistics. | 685 |
|     |  | 686 |
|     |  | 687 |
|     |  | 688 |
|     |  | 689 |
|     |  | 690 |
|     |  | 691 |
|     | Bill Yuchen Lin, Ziyi Wu, Yichi Yang, Dong-Ho Lee, and Xiang Ren. 2021. Riddlesense: Reasoning about riddle questions featuring linguistic creativity and commonsense knowledge. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 1504–1515.   | 692 |
|     |  | 693 |
|     |  | 694 |
|     |  | 695 |
|     |  | 696 |
|     |  | 697 |
|     | Haotian Luo, Li Shen, Haiying He, Yibo Wang, Shiwei Liu, Wei Li, Naiqiang Tan, Xiaochun Cao, and Dacheng Tao. 2025. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning. <i>CoRR</i> , abs/2501.12570.   | 698 |
|     |  | 699 |
|     |  | 700 |
|     |  | 701 |
|     |  | 702 |
|     | Kaijing Ma, Xeron Du, Yunran Wang, Haoran Zhang, Xingwei Qu, Jian Yang, Jiaheng Liu, Xiang Yue, Wenhao Huang, Ge Zhang, and 1 others. Kor-bench: Benchmarking language models on knowledge-orthogonal reasoning tasks. In <i>The Thirteenth International Conference on Learning Representations</i> .                                       | 703 |
|     |  | 704 |
|     |  | 705 |
|     |  | 706 |
|     |  | 707 |
|     |  | 708 |
|     | OpenAI. 2024. Learning to reason with llms.  | 709 |
|     | OpenAI. 2025a. gpt-oss-120b & gpt-oss-20b model card.  | 710 |
|     |  | 711 |
|     | OpenAI. 2025b. Introducing gpt-5.  | 712 |
|     | OpenAI. 2025c. Introducing o3 and o4 mini.   | 713 |
|     | Ioannis Panagiotopoulos, George Filandrianos, Maria Lymperaiou, and Giorgos Stamou. 2025. Riscore: Enhancing in-context riddle solving in language models through context-reconstructed example augmentation. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 9431–9455.                     | 714 |
|     |  | 715 |
|     |  | 716 |
|     |  | 717 |
|     |  | 718 |
|     |  | 719 |
|     |  | 720 |
|     | Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. Qwen2.5 technical report. <i>Preprint</i> , arXiv:2412.15115.                               | 721 |
|     |  | 722 |
|     |  | 723 |
|     |  | 724 |
|     |  | 725 |
|     |  | 726 |
|     |  | 727 |
|     | Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.   | 728 |
|     |  | 729 |
|     |  | 730 |
|     |  | 731 |
|     |  | 732 |
|     |  | 733 |
|     | David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In <i>First Conference on Language Modeling</i> .  | 734 |
|     |  | 735 |
|     |  | 736 |
|     |  | 737 |
|     |  | 738 |

|     |  |   |     |
|-----|--|---|-----|
| 739 | Elisa Sanchez-Bayona and Rodrigo Agerri. 2025.                   | <i>Natural Language Processing</i> , pages 11343–11358,         | 795 |
| 740 | Metaphor and large language models: When surface                 | Miami, Florida, USA. Association for Computational              | 796 |
| 741 | features matter more than deep understanding. In                 | Linguistics.  | 797 |
| 742 | <i>Findings of the Association for Computational Lin-</i>        |   |     |
| 743 | <i>guistics: ACL 2025</i> , pages 17462–17477.                   |   |     |
| 744 | ByteDance Seed, Jiaze Chen, Tiantian Fan, Xin Liu,               | Shuhang Xu and Fangwei Zhong. 2025. <b>CoMet:</b>               | 798 |
| 745 | Lingjun Liu, Zhiqi Lin, Mingxuan Wang, Chengyi                   | <b>Metaphor-driven covert communication for multi-</b>          | 799 |
| 746 | Wang, Xiangpeng Wei, Wenyuan Xu, and 1 others.                   | <b>agent language games</b> . In <i>Proceedings of the 63rd</i> | 800 |
| 747 | 2025. Seed1. 5-thinking: Advancing superb reason-                | <i>Annual Meeting of the Association for Computational</i>      | 801 |
| 748 | ing models with reinforcement learning. <i>arXiv</i>             | <i>Linguistics (Volume 1: Long Papers)</i> , pages 7892–        | 802 |
| 749 | <i>preprint arXiv:2504.13914</i> .                               | 7917, Vienna, Austria. Association for Computa-                 | 803 |
|     |  | tional Linguistics.   | 804 |
| 750 | Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,                  | Yu Yan, Sheng Sun, Zenghao Duan, Teli Liu, Min Liu,             | 805 |
| 751 | Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu,                  | Zhiyi Yin, LeiJingyu LeiJingyu, and Qi Li. 2025.                | 806 |
| 752 | and Daya Guo. 2024. <b>Deepseekmath: Pushing the</b>             | from benign import toxic: Jailbreaking the language             | 807 |
| 753 | <b>limits of mathematical reasoning in open language</b>         | model via adversarial metaphors. In <i>Proceedings</i>          | 808 |
| 754 | <b>models</b> . <i>CoRR</i> , abs/2402.03300.                    | <i>of the 63rd Annual Meeting of the Association for</i>        | 809 |
|     |  | <i>Computational Linguistics (Volume 1: Long Papers)</i> ,      | 810 |
| 755 | Alon Talmor, Jonathan Herzig, Nicholas Lourie, and               | pages 4785–4817.  | 811 |
| 756 | Jonathan Berant. 2019. Commonsenseqa: A question                 |   |     |
| 757 | answering challenge targeting commonsense knowl-                 | An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,                | 812 |
| 758 | edge. In <i>Proceedings of the 2019 Conference of</i>            | Binyuan Hui, Bo Zheng, Bowen Yu, Chang                          | 813 |
| 759 | <i>the North American Chapter of the Association for</i>         | Gao, Chengen Huang, Chenxu Lv, and 1 others.                    | 814 |
| 760 | <i>Computational Linguistics: Human Language Tech-</i>           | 2025. Qwen3 technical report. <i>arXiv preprint</i>             | 815 |
| 761 | <i>nologies, Volume 1 (Long and Short Papers)</i> , pages        | <i>arXiv:2505.09388</i> .                                       | 816 |
| 762 | 4149–4158.   |   |     |
| 763 | Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen,                  | Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,                | 817 |
| 764 | Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru                    | Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan,                 | 818 |
| 765 | Chen, Yuankun Chen, Yutian Chen, and 1 others.                   | Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo:             | 819 |
| 766 | 2025. Kimi k2: Open agentic intelligence. <i>arXiv</i>           | An open-source llm reinforcement learning system                | 820 |
| 767 | <i>preprint arXiv:2507.20534</i> .                               | at scale. In <i>The Thirty-Ninth Annual Conference on</i>       | 821 |
|     |  | <i>Neural Information Processing Systems</i> .                  | 822 |
| 768 | Qwen Team. 2025. <b>Qwq-32b: Embracing the power of</b>          | Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali                  | 823 |
| 769 | <b>reinforcement learning</b> .                                  | Farhadi, and Yejin Choi. 2019. Hellaswag: Can a                 | 824 |
|     |  | machine really finish your sentence? In <i>Proceedings</i>      | 825 |
| 770 | Paul H Thibodeau and Lera Boroditsky. 2011.                      | <i>of the 57th Annual Meeting of the Association for</i>        | 826 |
| 771 | Metaphors we think with: The role of metaphor in                 | <i>Computational Linguistics</i> , pages 4791–4800.             | 827 |
| 772 | reasoning. <i>PLoS one</i> , 6(2):e16782.                        |   |     |
| 773 | Yuan Tian, Nan Xu, and Wenji Mao. 2024. <b>A theory</b>          | Yunxiang Zhang and Xiaojun Wan. 2022. Birdqa: A                 | 828 |
| 774 | <b>guided scaffolding instruction framework for LLM-</b>         | bilingual dataset for question answering on tricky rid-         | 829 |
| 775 | <b>enabled metaphor reasoning</b> . In <i>Proceedings of the</i> | dles. In <i>Proceedings of the AAAI Conference on Arti-</i>     | 830 |
| 776 | <i>2024 Conference of the North American Chapter of</i>          | <i>ficial Intelligence</i> , volume 36, pages 11748–11756.      | 831 |
| 777 | <i>the Association for Computational Linguistics: Hu-</i>        |   |     |
| 778 | <i>man Language Technologies (Volume 1: Long Pa-</i>             |   |     |
| 779 | <i>pers)</i> , pages 7738–7755, Mexico City, Mexico. Asso-       |   |     |
| 780 | ciation for Computational Linguistics.                           |   |     |
| 781 | Xiaoyu Tong, Rochelle Choenni, Martha Lewis, and                 |   |     |
| 782 | Ekaterina Shutova. 2024. <b>Metaphor understanding</b>           |   |     |
| 783 | <b>challenge dataset for LLMs</b> . In <i>Proceedings of the</i> |   |     |
| 784 | <i>62nd Annual Meeting of the Association for Com-</i>           |   |     |
| 785 | <i>putational Linguistics (Volume 1: Long Papers)</i> ,          |   |     |
| 786 | pages 3517–3536, Bangkok, Thailand. Association                  |   |     |
| 787 | for Computational Linguistics.                                   |   |     |
| 788 | volcengine. 2025. <b>verl: Volcano engine reinforcement</b>      |   |     |
| 789 | <b>learning for llms</b> .                                       |   |     |
| 790 | Shun Wang, Ge Zhang, Han Wu, Tyler Loakman, Wen-                 |   |     |
| 791 | hao Huang, and Chenghua Lin. 2024. <b>MMTE: Cor-</b>             |   |     |
| 792 | <b>pus and metrics for evaluating machine translation</b>        |   |     |
| 793 | <b>quality of metaphorical language</b> . In <i>Proceedings</i>  |   |     |
| 794 | <i>of the 2024 Conference on Empirical Methods in</i>            |   |     |

## 832 A Appendix

### 833 A.1 Riddle Generation

#### 834 A.1.1 Riddle Generation Algorithm

835 We formalize our multi-agent iterative refinement  
836 framework for riddle generation as Algorithm 1.

#### 837 A.1.2 Prompt Templates for Riddle 838 Generation

839 This section presents the prompt templates used  
840 in our multi-agent iterative refinement framework  
841 for riddle generation. Each stage of the framework  
842 employs a specialized prompt template to guide the  
843 language models in their respective roles.

844 Table 4 shows the prompt template for Stage 1  
845 (Riddle Draft Generation), which guides the gener-  
846 ator to create metaphor-based riddles using the  
847 ground truth answer and its descriptions. The tem-  
848 plate variables are populated as follows (referring  
849 to Algorithm 1):

- 850 1. `{label}`: The ground truth answer  $t \in \mathcal{T}$   
851 from the Riddle Answer Taxonomy, provided  
852 as input to the algorithm (line 11 in Algo-  
853 rithm 1).
- 854 2. `{descriptions}`: The Wikidata descriptions  
855  $d_t$  for the ground truth answer, which pro-  
856 vide rich contextual information to enable the  
857 model to craft sophisticated metaphorical con-  
858 nections (line 11 in Algorithm 1).

859 Table 5 shows the prompt template for Stage 2  
860 (Solver Answer Collection), which guides language  
861 models to solve riddles by analyzing metaphors  
862 and logical clues from multiple perspectives. The  
863 template variables are populated as follows:

- 864 1. `{riddle}`: The current riddle draft  $e^{(r)}$  at revi-  
865 sion round  $r$ , where  $e^{(0)}$  is the initial draft gener-  
866 ated in Stage 1 (line 17 in Algorithm 1), and  
867  $e^{(r)}$  for  $r > 0$  is the refined riddle from previ-  
868 ous revision rounds (line 65 in Algorithm 1).

869 Table 6 shows the prompt template for Stage 3  
870 (Reviewer Assessment), which guides the reviewer  
871 to independently evaluate each solver’s answer  
872 based solely on the riddle itself, without access  
873 to the ground truth, thus avoiding bias caused by  
874 exposing the ground truth. The template variables  
875 are populated as follows:

- 876 1. `{riddle}`: The current riddle draft  $e^{(r)}$  being  
877 evaluated (line 31 in Algorithm 1).

2. `{solver_summary}`: A formatted summary  
878 of all solver answers from Stage 2, con-  
879 structed from the answer set  $\mathcal{A}^{(r)} =$   
880  $\{a_1, a_2, \dots, a_k\}$  collected in lines 22-26 in  
881 Algorithm 1, where each  $a_i$  is generated by  
882 solver  $M_s^{(i)} \in \mathcal{M}$ .  
883

884 Table 7 shows the prompt template for Stage  
885 4 (Reviser Refinement), which guides the reviser  
886 to refine problematic riddles based on diagnos-  
887 tic outcomes from Stage 3, addressing ambiguity  
888 (MULTI), insufficient challenge (EASY), or unsolv-  
889 ability (UNSOLV). The template variables are popu-  
890 lated as follows:

- 891 1. `{riddle}`: The previous riddle draft  $e^{(r-1)}$   
892 that needs refinement (line 65 in Algorithm 1).
- 893 2. `{groundtruth_answer}`: The ground truth  
894 answer  $t$  from the algorithm input (line 11 in  
895 Algorithm 1).
- 896 3. `{reviewer_feedback}`: The diagnostic out-  
897 come  $d^{(r-1)}$  determined in Stage 3, which can  
898 be MULTI (line 42 in Algorithm 1), EASY (line  
899 44 in Algorithm 1), or UNSOLV (line 46 in Al-  
900 gorithm 1), indicating the specific issue that  
901 needs to be addressed.
- 902 4. `{solver_feedback}`: Information derived  
903 from the solver answer set  $\mathcal{A}^{(r-1)}$  and cor-  
904 rectness judgments  $\mathcal{C}^{(r-1)}$  collected in previ-  
905 ous stages, providing additional context about  
906 how different solvers interpreted the riddle  
907 (line 48-52 in Algorithm 1).

#### 908 A.1.3 Hyperparameters

909 We set the following hyperparameters in our riddle  
910 generation pipeline: (1) *Quality threshold*  $\tau =$   
911 20 for filtering entities with sitelinks below the  
912 minimum threshold; (2) *Taxonomic quota*  $K =$   
913 1000 allocated to each second-level category; (3)  
914 *Maximum revision rounds*  $R_{\max} = 10$ ; (4) *Pass-*  
915 *rate threshold*  $\theta_p = 0.9$ .

#### 916 A.1.4 Examples

917 **Riddle Answer Examples.** Table 8 presents ex-  
918 amples of riddle answers organized by the hierar-  
919 chical taxonomy structure (H1, H2, H3), demon-  
920 strating the diversity of concepts, objects, and phe-  
921 nomena covered in our dataset.

922 **Riddle Examples.** Tables 9 and 10 present ex-  
923 ample riddles from different taxonomic categories,  
924 demonstrating the diversity and quality of riddles  
925 generated by our framework.

## A.2 Experimental Setup

### A.2.1 Benchmark Details

To evaluate the impact of Metaphor Reasoning on reasoning abilities in other domains, we employ six categories of Out-of-Distribution (OOD) reasoning across 7 benchmarks: (1) *Logical Reasoning*: KORBench (Ma et al.), a knowledge-orthogonal benchmark, and Enigmata (Chen et al., 2025), containing 36 puzzle reasoning tasks across seven categories; (2) *Commonsense Reasoning*: CommonsenseQA (Talmor et al., 2019), which tests the ability to reason about everyday situations and world knowledge; (3) *Natural Language Inference*: HellaSwag (Zellers et al., 2019), which requires inferring the most plausible next situation based on context; (4) *Math Reasoning*: American Invitational Mathematics Examination (AIME) 2024<sup>9</sup>, containing challenging problems from mathematical competitions; (5) *Science, Technology, Engineering, and Mathematics (STEM) Reasoning*: GPQA Diamond (Rein et al.), a graduate-level science test; (6) *Out-of-Distribution (OOD) Riddle Reasoning*: RiddleSense (Lin et al., 2021), a collection of manually curated riddles in multiple-choice format.

### A.2.2 Implementation Details

**Training Details.** We implement DAPO training using the verl framework (volcengine, 2025), a scalable reinforcement learning system for large language models. DAPO (Yu et al., 2025) is a variant of GRPO (Shao et al., 2024) that employs group-based advantage normalization and asymmetric policy clipping for stable reinforcement learning with verifiable rewards. Unlike traditional PPO-based methods, DAPO does not require a critic network and computes advantages directly from group-normalized rewards within each group of generated outputs. All hyperparameters used in our experiments are summarized in Table 11.

**Evaluation Details.** To ensure sufficient training signal and stable learning, we employ dataset repetition during training. Each dataset is repeated multiple times to balance the representation of different reasoning domains. Specifically, the repetition counts for each dataset are as follows: AIME is repeated 32 times, GPQA\_RuleVerifier is repeated 16 times, and all other datasets are repeated 3 times. All results reported in this paper are statistically significant with  $p < 0.001$ .

<sup>9</sup>[https://artofproblemsolving.com/wiki/index.php/AIME\\_Problems\\_and\\_Solutions](https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions)

---

### Algorithm 1: Riddle Generation via Solvers-Reviewer-Reviser Workflow

---

```
Input: Ground truth answer  $t \in \mathcal{T}$  from Riddle Answer Taxonomy  $\mathcal{T}$ ; Solver set  $\mathcal{M} = \{M_s^{(1)}, M_s^{(2)}, \dots, M_s^{(k)}\}$ ; Reviewer  $M_r$ ; Reviser  $M_v$ ; Generator  $M_g$ ; Example  $e_t$ ; Descriptions  $d_t$ ; Maximum revision rounds  $R_{\max}$ ; Pass-rate threshold  $\theta_p$ .  
Output: Final riddle  $e^{(r)}$  or REJECTED.  
// Stage 1: Riddle Draft Generation  
 $e^{(0)} \leftarrow M_g(t, e_t, d_t)$   
 $r \leftarrow 0$   
while  $r < R_{\max}$  do  
  // Stage 2: Solver Answer Collection  
   $\mathcal{A}^{(r)} \leftarrow \{\}$   
  for  $M_s^{(i)} \in \mathcal{M}$  do  
     $a_i \leftarrow M_s^{(i)}.Solve(e^{(r)})$   
     $\mathcal{A}^{(r)}.add(a_i)$   
  // Stage 3: Reviewer Assessment  
   $\mathcal{C}^{(r)} \leftarrow \{\}$   
  for  $a_i \in \mathcal{A}^{(r)}$  do  
     $c_i \leftarrow M_r(a_i, e^{(r)}) \in \{0, 1\}$   
     $\mathcal{C}^{(r)}.add(c_i)$   
   $pass\_rate \leftarrow \frac{|\{i: (c_i=1) \wedge (a_i=t)\}|}{k}$   
   $solvable \leftarrow \exists i: (c_i = 1) \wedge (a_i = t)$   
   $verifiable \leftarrow solvable \wedge \neg(\exists i: (c_i = 1) \wedge (a_i \neq t))$   
   $challenging \leftarrow (pass\_rate \geq \theta_p)$   
  // Determine diagnostic outcome  
  if  $\exists i: (c_i = 1) \wedge (a_i \neq t)$  then  
     $d^{(r)} \leftarrow MULTI$   
  else if  $\forall i: a_i = t$  then  
     $d^{(r)} \leftarrow EASY$   
  else if  $\neg solvable$  then  
     $d^{(r)} \leftarrow UNSOLV$   
  else if  $verifiable \wedge challenging$  then  
     $d^{(r)} \leftarrow PASS$   
  // Stage 5: Pass Rate Filter  
  if  $d^{(r)} = PASS$  then  
    return  $e^{(r)}$   
  // Stage 4: Reviser Refinement  
  if  $d^{(r)} = MULTI$  then  
     $revision\_type \leftarrow$  "clarify ambiguity"  
  else if  $d^{(r)} = EASY$  then  
     $revision\_type \leftarrow$  "increase challenge"  
  else if  $d^{(r)} = UNSOLV$  then  
     $revision\_type \leftarrow$  "address unsolvability"  
   $r \leftarrow r + 1$   
   $e^{(r)} \leftarrow M_v(e^{(r-1)}, d^{(r-1)})$   
return REJECTED
```

---

### Prompt for Riddle Draft Generation

You are an expert riddle creator specializing in crafting challenging, metaphor-based puzzles that test advanced reasoning abilities.

#### Background

Riddles achieve their mystique by using metaphors to obscure the answer (A) through comparisons to seemingly unrelated objects (B, C, D), creating an atmosphere of mystery and intrigue.

#### Example:

- **Riddle:** “I am a fortress with no doors, a treasure chest with golden stores.”
- **Answer:** An egg
- **Analysis:** The “fortress” and “treasure chest” serve as metaphors for the egg. The shell resembles a sturdy “fortress,” while the yolk represents the golden “treasure.”

#### Task

Please generate a riddle using the following label as the answer. I have provided extensive information that you can use to design your riddle. Ensure there is no obvious information leakage—the riddle should be as cryptic as possible through the use of metaphors.

#### The answer of the riddle:

{label}

#### Description:

{descriptions}

#### Requirements:

- Use metaphorical language to obscure the answer
- Avoid direct information leakage
- Create a challenging, multi-layered puzzle
- Consider nested metaphors for increased difficulty
- Ensure the riddle tests genuine reasoning abilities

#### Output Format

Please wrap your riddle in the following XML tags:

```
<riddle>Your riddle here</riddle>
```

Table 4: The prompt template used for riddle draft generation in Stage 1. The template guides the language model to create metaphor-based riddles using the ground truth answer and its descriptions.

### Prompt for Solver Answer Collection

You are an expert riddle solver with deep knowledge and strong logical reasoning skills. You have the ability to analyze complex riddles from multiple perspectives and angles.

When solving the given riddle, please:

- Think broadly and explore various possible answers from different dimensions
- Identify and analyze key words, metaphors, and logical clues in the riddle
- Synthesize all possibilities and select the most fitting and elegant answer
- Present only your final answer as the conclusion

Please enclose your answer within <answer></answer> tags.

**Riddle:**  
{riddle}

Table 5: The prompt template used for solver answer collection in Stage 2. The template guides language models to solve riddles by analyzing metaphors and logical clues from multiple perspectives.

### Prompt for Reviewer Assessment

You are a senior puzzle editor evaluating whether proposed answers correctly solve a riddle. Judge each answer based solely on the riddle text, independent of any groundtruth answer.

**Riddle:**  
{riddle}

**Solver responses:**  
{solver\_summary}

**Instructions:**

For each answer provided, analyze and judge whether it correctly solves the riddle:

1. Explain your reasoning: Evaluate how well the answer fits the riddle constraints, clues, and requirements.
2. Determine correctness: Based on your reasoning, decide if the answer is correct or incorrect.
3. Note that answers can be: all correct, all incorrect, or a mix of both.
4. If the riddle has ambiguities or multiple valid interpretations, explicitly note them.

Return your judgments using ONLY the per-solver tagged format below. Preserve the solver order from the summary and emit no extra narration.

```
<solver solver_id="solver_precise">  
<answer>the solver's final answer</answer>  
<analysis>concise single-line justification</analysis>  
<correctness>true</correctness>  
</solver>
```

**Formatting requirements:**

- Use lowercase true or false inside <correctness>.
- Keep each <analysis> on one line; replace newlines with spaces.
- Include exactly one <solver> block for every solver shown in the summary.
- Do not add any other tags or text outside the solver blocks.

Table 6: The prompt template used for reviewer assessment in Stage 3. The template guides the reviewer to independently evaluate each solver's answer based solely on the riddle itself, without access to the ground truth, thus avoiding bias caused by exposing the ground truth.

### Prompt for Reviser Refinement

You are an expert riddle editor. Given an existing riddle and editorial feedback, produce a revised riddle that preserves the original mystery while ensuring only the groundtruth answer fits.

**Old Riddle**

{riddle}

**Groundtruth Answer**

{groundtruth\_answer}

**Reviewer Feedback**

{reviewer\_feedback}

**Solver Feedback**

{solver\_feedback}

**Task**

Write a brand-new riddle that resolves the reviewer feedback, avoids leaking the groundtruth answer directly, and keeps a comparable difficulty level.

Return only the revised riddle wrapped exactly in <riddle></riddle> tags with no extra commentary.

Table 7: The prompt template used for reviser refinement in Stage 4. The template guides the reviser to refine problematic riddles based on diagnostic outcomes from Stage 3, addressing ambiguity (MULTI), insufficient challenge (EASY), or unsolvability (UNSOLV).

| $H_1$             | $H_2$              | $H_3$                             | <b>Example of Riddle Answers</b>                              |
|-------------------|--------------------|-----------------------------------|---|
| Abstract Entities | Conceptual Systems | ideology                          | nationalism, conservatism, fascism, feminism                  |
| Abstract Entities | Conceptual Systems | mathematical concept              | area, volume, length, height                                  |
| Abstract Entities | Roles              | social status                     | citizen, prisoner, slave, reputation                          |
| Abstract Entities | Traits             | personality trait                 | curiosity, geek, cardinal virtues                             |
| Objects           | Mixed Origin       | drink                             | coffee, tea, milk, beer                                       |
| Objects           | Mixed Origin       | food                              | bread, rice, apple, pizza                                     |
| Objects           | Natural Entities   | natural geographic object         | mountain, ocean, valley, volcano                              |
| Objects           | Physical Objects   | architectural structure           | house, building, bridge, road                                 |
| Objects           | Physical Objects   | clothing                          | dress, hat, trousers, jeans                                   |
| Objects           | Physical Objects   | home appliance                    | refrigerator, washing machine, mobile phone, desktop computer |
| Objects           | Physical Objects   | musical instrument                | guitar, violin, drum, flute                                   |
| Objects           | Physical Objects   | weapon                            | sword, rifle, pistol, firearm                                 |
| Objects           | Regions            | administrative territorial entity | country, state, province, county                              |
| Objects           | Regions            | city                              | metropolis, big city, capital city, global city               |
| Objects           | Regions            | country                           | empire, sovereign state, nation state, developing country     |
| Objects           | Social Entities    | educational institution           | school, university, college, kindergarten                     |
| Objects           | Social Entities    | organization                      | company, government, army, religion                           |
| Objects           | Social Entities    | political party                   | communist party, green party, big tent                        |
| Phenomena         | Natural Phenomena  | astronomical phenomenon           | eclipse, meteor, sunset, solar eclipse                        |
| Phenomena         | Natural Phenomena  | disease                           | cancer, diabetes, asthma, influenza                           |
| Phenomena         | Natural Phenomena  | meteorological phenomenon         | rain, snow, storm, thunder                                    |
| Phenomena         | Natural Phenomena  | natural disaster                  | earthquake, flood, tsunami, avalanche                         |
| Phenomena         | Natural Phenomena  | season                            | spring, summer, autumn, winter                                |
| Phenomena         | Social Processes   | historical event                  | battle, armistice, ceasefire                                  |
| Phenomena         | Social Processes   | public election                   | general election, presidential election                       |
| Phenomena         | Social Processes   | social process                    | education, migration, revolution, reform                      |

Table 8: Example of Riddle Answers organized by taxonomic hierarchy.

| H1        | H2                | H3           | Riddle Answer | Riddle  |
|-----------|-------------------|--------------|---------------|---|
| Objects   | Mixed Origin      | Food         | Banana        | <p>I am a crescent wrought by a moonsmith, sheathed in a cloak with three long seams. My choir hangs like a chandelier near the ceiling of a giant that is grass yet mimics a tree.</p> <p>As days pass, night-dust freckles my garment; the bread in my heart turns slowly into honey.</p> <p>By the scholars' ledger I am called a berry, though my beads are ghosts and I'm born without a wedding.</p> <p>Some of my kin must learn the language of fire before their song becomes sweet.</p> <p>What am I?</p>   |
| Objects   | Social Entities   | Organization | College       | <p>A learning village rings its hours by a tower; paths braid toward a central green. Porches don Greek in the swift season of choosing; whispers and handclasps bind those lettered doors.</p> <p>Three wardens tend my thresholds: one keeps the purse, one the rolls and times, one the yes-or-no.</p> <p>At a bell's strike, blue books bloom; hush gathers while credits accrue by the hour. Some houses send you off in two winters with short initials; others keep you four for the longer pair.</p> <p>My kin answer to surnames like Community, Junior, and Liberal Arts; crisp Saturdays drum my name across the field.</p> <p>Between red margin and blue lines, the ruling bears my title; a numbered plan—five-two-nine—saves toward me by name.</p> <p>Here, they say they're going to me; across the sea, they go to "uni."</p> <p>When tassels turn, AA, AS, BA, or BS trail your name; taller hoods are seldom cut beneath my roof.</p> <p>And when grit must do, they urge the old try that carries my name.</p> <p>What am I?</p> |
| Objects   | Physical Objects  | Clothing     | Jeans         | <p>I am a night-skinned map that brightens where the world keeps touching me. Two hollow roads run my length—pilgrims stand inside me to meet the ground. My edges wear paired scars, and the mouths of my caves are pinned by tiny copper suns that do not set.</p> <p>Within one cavern, a smaller echo hides—fit for a spark, a coin, a whisper. My front can grin with tempered teeth, then hush with a single cold kiss. Five mute moons circle my crown to leash a tamed serpent.</p> <p>Born to wrestle grit and gravel, I learned the manners of parlors, yet I never forgot the language of dust.</p> <p>Name me.</p>  |
| Phenomena | Natural Phenomena | Disease      | Vaccine       | <p>I was christened with a pastoral echo, though I graze nowhere. I knock at the cistern before the siege; I mend no wounds—I drill the watch. I come disguised: a ghost of the raider, a splinter, or a rumor written in letters. I travel with winter as escort; through a small hill of flesh I pin my notice. At times I hid in sweetness; at others I rode a breeze to the gate. I hire no mercenaries; I teach your forgemasters to cast their own steel and keep the plans.</p> <p>Rehearsals ring by appointment; the red river's numbers speak of my lesson. My stamped trail opens ports and halls; batch marks and dates betray my path. When true banners rise, the city answers at once, already sure.</p>   |

Table 9: Example riddles from different taxonomic categories (Part 1 of 2).

| H1                | H2               | H3             | Riddle Answer | Riddle   |
|-------------------|------------------|----------------|---------------|--|
| Phenomena         | Social Processes | Social Process | Middle Ages   | <p>My gate is the hinge between marble courts and untested mirrors.<br/> I woke when the western eagle's colors bled from its robe.<br/> Deep stone forests rose; praying boughs caught wind and brewed saints in light.<br/> Debt sworn on steel parceled fields; bread bowed to blade and altar.<br/> Lamplight tillers of parchment sowed letters;<br/> iron tongues weighed the hours.<br/> Empty fairs were counted by a pale collector;<br/> earth, silk, and incense ate together.<br/> Across seals and hands, towns bought voices;<br/> shells and crosses stitched roads to hungry shrines.<br/> Grain of thunder humbled glittering mail;<br/> a lodestone taught keels to speak straight.<br/> Etched metal whispered more words than throats;<br/> halls knotted three above four beneath contentious roofs.<br/> Speak me plain—ten letters and a gap; no herald walks before my name.</p>   |
| Abstract Entities | Roles            | Social Status  | Reputation    | <p>I am the foreword they read before your mouth writes a line.<br/> I am woven from footprints on days now dry;<br/> the weavers are not on your payroll.<br/> Parliament and notaries may declare me anew, yet their stamps pass through me.<br/> Scales quarrel with digits and calm at a sigh; light air edits my measure.<br/> I unlatch iron paths and make velvet doors grow thorns<br/> with no hand upon a bolt.<br/> I leave a colorless patina; baths and polish do nothing.<br/> Your mirror will not find me; their memory will;<br/> a shaped tale can bruise me while your truth stands.<br/> A crown hoards me in silence; a fool spends me with a grin and wakes to poverty.<br/> A syllable can split me; a decade can stitch me; I refuse pockets.<br/> I may be born of roofs you never raised<br/> and die for nights you never lived.<br/> Marks, brands, and charters pretend to leash me;<br/> merchants hire keepers to herd my smoke.<br/> No ledger numbers me, yet prices bow when I swell or sink.<br/> In rooms you enter, the air has already chosen its lean;<br/> I set the chairs.</p> |
| Abstract Entities | Traits           | Personal Trait | Geek          | <p>In crowded rooms I dodge the clink of shallow glasses,<br/> choosing one tick, one topic, to tune until it sings.<br/> I string stray facts like LEDs on a midnight cable,<br/> small sparks others swept from the floor.<br/> The name they stuck on me first pricked like a burr—<br/> I hammered it flat into a badge I wear at cons.<br/> Under sawdust lights a lifetime ago it meant a sideshow,<br/> a tent-oddity with a chicken and a crowd.<br/> Now it powers squads that come to mend your screens,<br/> and turns verb when joy makes me explain.<br/> Not Greek, though many write me so;<br/> my middle peers through twin lenses.<br/> What word am I?</p>  |

Table 10: Example riddles from different taxonomic categories (Part 2 of 2).

| <b>Hyperparameter</b>   | <b>Value</b>            |
|---|-------------------------|
| Training Framework  | verl (volcengine, 2025) |
| Train Batch Size  | 256                     |
| Mini Batch Size   | 256                     |
| Group Size ( $G$ , num_bon)                                     | 16                      |
| Input Length  | 2K                      |
| Response Length   | 32K                     |
| Training Steps  |                         |
| QwQ-32B   | 125                     |
| Qwen3-14B   | 130                     |
| Qwen3-8B  | 120                     |
| Actor Learning Rate   | $1 \times 10^{-6}$      |
| Clip Ratio ( $\epsilon_{\text{low}} / \epsilon_{\text{high}}$ ) | 0.28 / 0.2              |
| KL Coefficient  | 0.0                     |
| Number of GPUs  |                         |
| QwQ-32B   | 256 A100                |
| Qwen3-14B   | 64 A100                 |
| Qwen3-8B  | 64 A100                 |

Table 11: DAPO training hyperparameters for different model scales. Note that K denotes thousands of tokens.