

Beyond ‘Aha!’: Toward Systematic Meta-Abilities Alignment in Large Reasoning Models

Anonymous ACL submission

Abstract

Large reasoning models (LRMs) already possess a latent capacity for long chain-of-thought reasoning. Prior work has shown that outcome-based reinforcement learning (RL) can incidentally elicit advanced reasoning behaviors such as self-correction, backtracking, and verification—phenomena often referred to as the model’s “aha moment”. However, the timing and consistency of these emergent behaviors remain unpredictable and uncontrollable, limiting the scalability and reliability of LRMs’ reasoning capabilities. To address these limitations, we move beyond reliance on prompts and unpredictable “aha moments”. Instead, we explicitly align models with three meta-abilities: **deduction, induction, and abduction**, using automatically generated, self-verifiable tasks. Our three-stage pipeline (individual alignment, parameter-space merging, domain-specific reinforcement learning) boosts performance by over 10% relative to instruction-tuned baselines. Furthermore, domain-specific RL from the aligned checkpoint yields an additional gain in performance ceiling for both 7B and 32B models across math, coding, and science benchmarks, showing that explicit meta-ability alignment offers a scalable and dependable foundation for reasoning. Code and data can be found in Software and Data part in submission page.

1 Introduction

Large reasoning models, including OpenAI-o1 (Jaech et al., 2024), o3 (OpenAI, 2025), DeepSeek-R1 (Guo et al., 2025), Grok 3.5 (xAI, 2025), and Gemini 2.5 Pro (DeepMind, 2024), have demonstrated remarkable capabilities. These models excel at generating long Chain-of-Thought (CoT) (Wei et al., 2022) responses when tackling complex tasks and exhibit advanced, reflection-like reasoning behaviors. Recently, DeepSeek-R1 has shown that, starting from pretrained base or instruction-tuned models, pure reinforcement learning (RL) with

rule-based rewards can spontaneously lead to the emergence of long CoT reasoning, self-correction, self-reflection, and other advanced behaviors, collectively referred to as the “aha moment”. An “aha moment” is usually defined as after a few steps where the model exhibits a sudden performance surge and emergency of reasoning behaviors. Other open-source works, such as SimpleRL-Zoo (Zeng et al., 2025), tinyzero (Pan et al., 2025), and Logic-RL (Xie et al., 2025), which attempt to reproduce R1’s performance and technical details, have also observed similar aha moments. These behaviors, such as self-correction, self-verification, and backtracking, signal the model’s internal experience of strong reasoning ability.

However, relying solely on emergent behaviors is inherently unreliable and difficult to control. Models may fail to consistently manifest these advanced reasoning schemes, which limits both the predictability and scalability of LLM-based reasoning. To overcome this, we propose to explicitly align LLMs with three domain-general reasoning meta-abilities: deduction, induction, and abduction, drawn from Peirce’s classical inference triad (Peirce, 1931).

As Figure 1 demonstrated, deduction infers specific outcomes from general rules and hypotheses ($H + R \rightarrow O$), enabling rigorous prediction and verification. Induction abstracts rules from repeated co-occurrences ($H + O \rightarrow R$), facilitating pattern discovery and generalization. Abduction infers the most plausible explanation for surprising observations ($O + R \rightarrow H$), promoting creative and backward reasoning. Together, they form a closed inferential loop for hypothesis generation, testing, and revision, mirroring scientific method and supporting robust and interpretable reasoning.

To operationalize these meta-abilities, we construct a task suite with programmatically generated instances and automatic verifiability. Each task targets one core reasoning mode: Deduction: Proposi-

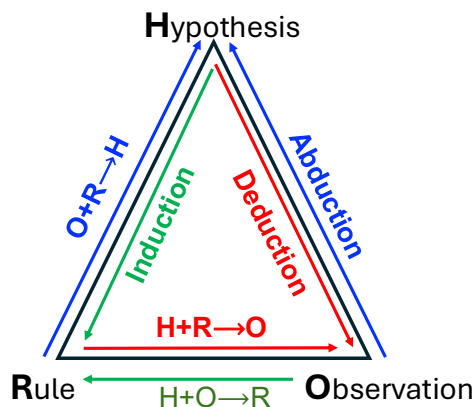


Figure 1: A closed H-R-O inferential loop

tional satisfiability tasks use rule sets R and candidate hypotheses H to test if all premises entail the observation O . Induction: Masked-sequence completion requires models to infer latent rules R from partial inputs H, O . Abduction: Inverse rule-graph search backchains from observed consequences O through a rule graph R to infer the minimal explanatory H . These tasks are constructed from synthetic distributions that lie out-of-distribution relative to common pretraining corpora, ensuring that gains reflect genuine meta-ability acquisition rather than memorization or shortcut exploitation.

We observe that models aligned to individual meta-abilities make complementary errors. Aggregating their predictions raises overall accuracy by more than 10% relative to a vanilla instruction-tuned baseline. To incorporate the three competencies into a single network, Rather than training the model on a mixed task corpus, we utilize parameter-space model merging to integrate these meta-abilities. Parameter-space merging lifts average accuracy by around 2% at 7B and about 4% at 32B over the instruct baseline, and by 16.3% over the 7B Base model, showing strong generalization of merged meta-abilities. Additionally, Our approach triggers a performance surge and emergent reasoning in around 25 iterations (400 examples), offering more control and efficiency than other methods. More details will be elaborated in Section 5.1 and 5.3.

Furthermore, to evaluate whether meta-ability alignment offers a stronger foundation for subsequent learning, we resumed domain-specific RL training from a checkpoint that had already been aligned and compared it with the same procedure applied to an instruction-tuned model. Starting from the meta-ability checkpoint raises the attainable performance ceiling: after identical continual domain-specific RL training, the model achieves an

average gain of about 2% over its instruction-only counterpart. Our key contributions are as follows:

- **Task suite for meta-abilities.** We introduce a novel RL task suite aligned with three classical meta-abilities: deduction, induction, and abduction, each constructed to train and validate domain-general reasoning skills in LLMs.
- **Recipe for reasoning mastery.** We propose a three-stage recipe: (1) independently align models to each meta-ability; (2) merge them via parameter-space integration; (3) fine-tune with domain-specific RL. This leads to improved generalization and downstream task accuracy.
- **Upper-bound boost and scalability.** We empirically demonstrate that meta-ability alignment raises the performance ceiling: our 7B and 32B models show consistent gains over instruction-tuned and base model baselines, across math, coding, and science benchmarks.

2 Related Work

RL-Driven Emergence of Reasoning Abilities.

Recent studies show that direct RL post-training can unlock long chain-of-thought reasoning beyond what supervised fine-tuning achieves (Zeng et al., 2025; Xiong et al., 2025; Guo et al., 2025). *SimpleRL-Zoo* (Zeng et al., 2025) proposes a zero-RL recipe using rule-based rewards, boosting math reasoning accuracy and inducing cognitive behaviors like self-verification across diverse base models. *DeepSeek-RL* (Guo et al., 2025) extends this idea to large-scale training. Its public replications, *Light-RL* (Wen et al., 2025), *Open-RL* (Face, 2025), and the minimal-cost *TinyZero* (Pan et al., 2025), confirm that curriculum schedules, DPO warm-up, and carefully shaped length rewards together yield stronger logical accuracy while keeping compute affordable. Complementary to these general pipelines, *Logic-RL* (Xie et al., 2025) applies rule-conditioned reinforcement learning to synthetic Knights-and-Knaves puzzles, Enabling transferable logical reasoning for math tasks. Together, these works establish RL as a viable path to large reasoning models.

Advanced reasoning ability. In addition to enhancing long-chain reasoning through RL, recent work investigates specific reasoning skills such as self-correction, counterfactual inference, self-verification and others. Chen et al. (2024) enhances

🗨️ DEDUCTION	🔍 INDUCTION	🔍 ABDUCTION
<p>Q: Determine whether a single truth assignment exists to satisfy **all** of these formulas:</p> <ol style="list-style-type: none"> $((\neg(\neg B \leftrightarrow \neg D) \leftrightarrow (\neg L \wedge \neg G) \leftrightarrow (B \oplus H)) \rightarrow (\neg(\neg L \oplus G) \oplus L \oplus \neg K)) \vee \dots$ $\neg \neg(((B \leftrightarrow K) \rightarrow (\neg G \oplus \neg L)) \wedge \dots)$ 	<p>Q: Given the following sequence, [0, 0, 0, 1, 2, 0, ?]. What is the value at the question mark?</p>	<p>Q: Premises (excerpt), $(P \wedge A) \rightarrow F, (E \wedge R) \rightarrow X, (F \wedge Q) \rightarrow X, (Q \wedge H) \rightarrow E, (Q \wedge X) \rightarrow H \dots$ (7 more). Observed facts X, H; Goal P. Analysis reachability of goal.</p>
<p>Thinking: Start with every atom = False. Φ_1 evaluates to False, so flip $K := \mathbf{True}$ to fix the consequent (hypothesis proposal). Re-check: Φ_1 is now True, but Φ_2 fails because $(B \leftrightarrow K)$ is False. Flip $B := \mathbf{True}$; contradiction resolved (self-correction).</p>	<p>Thinking: The sequence 0, 0, 0, 1, 2, 0, -2, ? yields differences 0, 0, 1, 1, -2, -2, -2. I spot that each delta repeats twice in a row (pattern abstraction manifested). From that I generalize “every difference comes in pairs,” so the next Δ should be -2 again, giving $-4 + (-2) = -6$ (generalization).</p>	<p>Thinking: H is observed, and the only rule for H is “$Q \wedge X \rightarrow H$.” X is already true, so I reverse-chain to seek Q (reverse thinking). Q itself comes from “$R \wedge B \rightarrow Q$,” neither known, so I assume $R = \mathbf{True}$ and $B = \mathbf{True}$ and propagate back through the rules (backward reasoning).</p>

Figure 2: Meta-Ability Alignment Task Example and Corresponding Thought Process

overall reasoning performance by training models to reason both forward and backward, demonstrating that reverse-thinking objectives can improve forward reasoning as well. Kumar et al. (2024) proposes SCoRe, an online RL method where a model iteratively critiques and improves its own answers, bridging the offline self-correction gap. Several recent studies also focus on equipping LLMs with self-verification and self-correction abilities, including ProCo (Wu et al., 2024), S²R (Ma et al., 2025), and SETS (Chen et al., 2025).

Investigation of Aha-moment. RL pipelines often show sudden accuracy jumps. Some concurrent analyses aim to uncover the internal “aha” moments that precede them. Gandhi et al. (2025) introduces four reasoning behaviors as diagnostic tools to explain and engineer models’ capacity for self-improvement under reinforcement learning. Yang et al. (2025) shows that “aha moments” emerge through anthropomorphic language, uncertainty adjustment, and latent-space shifts, helping models avoid reasoning collapse and adapt to problem difficulty. Additionally, Zhou et al. (2025) shows that similar emergence occurs in a 2B vision–language model without supervised warm-up.

3 Methodology

3.1 Task Design for Meta-Abilities Alignment

We design three reasoning tasks, each of which involves inferring one element: hypothesis (H), rules (R), or observation (O), given the other two. In *deduction* ($H + R \Rightarrow O$), the model is given a set of logical rules R and a candidate truth assignment H as hypothesis, and must verify whether the overall observation O (i.e., all formulas be-

ing true) follows, formulated as a **propositional satisfiability** task. In *induction* ($H + O \Rightarrow R$), the model is provided with observable items O and incomplete inputs H (e.g., masked tokens or implied guesses), and must abstract the underlying generative rule R to correctly complete the sequence, framed as a **masked-sequence completion** task. In *abduction* ($O + R \Rightarrow H$), the model is given observations O and a rule graph R , and must trace backward to recover the minimal set of hidden assumptions H that can logically explain the conclusion, posed as a **reverse rule-graph search** task. This design follows a strict two-known-one-infer schema, clearly ensuring a clean separation of reasoning types, and reformulates all tasks into a unified (H,R,O) triplet format. This enables consistent, comparable, and complementary training signals, systematically equipping the model with a full range of meta-reasoning capabilities. As illustrated in Figure 3, each instance is produced by an automated *Generator* and screened by a *Verifier*, yielding large-scale, self-checked training data entirely free of manual annotation. The pseudo code for data synthesis and additional task examples are provided in Appendix A.

Deduction. As the example shown in Figure 2, we pose task that present a concise cluster of nested propositional clauses involving the standard Boolean operators NOT, AND, OR, IMPLIES, IFF, and XOR; the model must either return a satisfying truth assignment or report that the clauses are unsatisfiable. The task poses a combinatorial explosion of possible variable assignments, with tightly coupled logical formulas such that the value of one variable may indirectly constrain or deter-

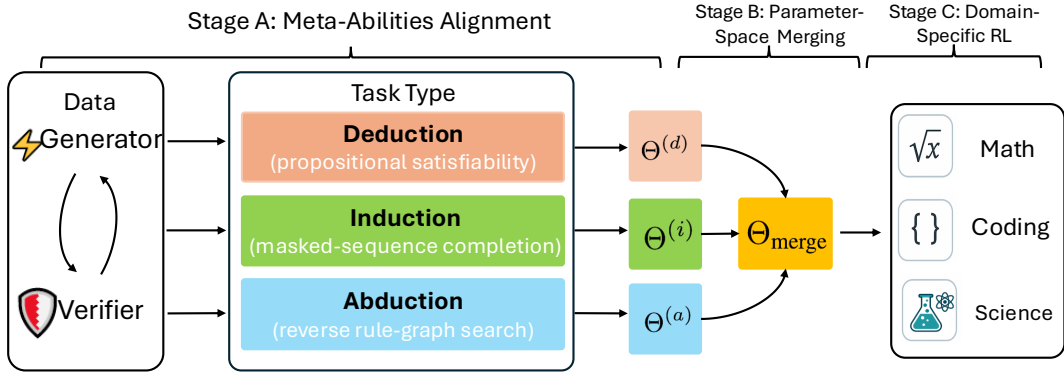


Figure 3: Overview of the three-stage pipeline: align deduction, induction, and abduction specialists, merge them in parameter space, and continually RL-adapt the unified model to downstream domains.

mine the values of many others through chains of logical dependencies. This interdependence renders purely enumerative or heuristically guessed assignments overwhelmingly unlikely to satisfy all constraints. The only tractable approach is to begin with a provisional assignment, treat each formula as a logical premise, systematically derive its consequences, identify any resulting contradictions, and revise the assignment accordingly. This iterative loop of hypothesis generation, logical consequence propagation, empirical inconsistency detection, and corrective refinement directly instantiates the core structure of deductive reasoning.

Induction. To develop inductive reasoning capabilities in models, we design a task for automatically-generated sequence with hidden terms. Each instance presents a series of elements following an undisclosed pattern (including numeric reasoning, symbolic patterns, and multi-step operation cycles) and requires identification of a missing element. This methodology specifically targets induction as the model must extract the underlying regularity governing the visible sequence and apply it to predict unseen values. Inductive learning through such structured sequences enhances the model’s fundamental capabilities in abstraction and generalization, which are essential components for robust reasoning across domains. We also provide the example in Figure 2 about induction task and ability behavior for better explanation.

Abduction. To cultivate backward reasoning ability, we introduce a reverse rule-graph search task in which forward inference is deliberately obstructed while backward inference remains efficient. Each instance is formulated as a directed rule graph, with atoms as nodes and implications encoded as hyperedges from premise sets to conclusions. Observed facts activate source nodes, while target hypotheses correspond to sink nodes with unknown truth

values. By inflating the branching factor in forward chaining, exhaustive exploration becomes computationally infeasible. In contrast, a backward strategy starts from a goal, hypothesizes minimal supporting premises, and verifies them against known facts. This approach can efficiently isolate relevant subgraphs. The design induces repeated cycles of goal-directed hypothesis formation, verification, and revision, thereby fostering the core mechanism of abductive reasoning.

3.2 Training Recipe for Reasoning

Figure 3 sketches how we transform the emergent “aha” moment into *controllable, composable* meta-abilities: we first carry out **Meta-Abilities Alignment**, independently training deduction, induction, and abduction specialists on synthetic diagnostics; we then fuse these specialists through **Parameter-Space Merging** to obtain a single checkpoint that retains their complementary strengths. Finally, **Domain-Specific Reinforcement Learning Training** refines the merged model on domain-specific data such as math, coding, and social dialogue.

3.2.1 Stage A: Meta-Abilities Alignment (RL)

We curate three synthetic but diagnostic datasets: *Deduction* (propositional satisfiability), *Induction* (masked-sequence completion), and *Abduction* (reverse rule-graph search). For a policy π_θ , we adopt the Group Relative Policy Optimization (GRPO) objective (Shao et al., 2024), where $r_{i,t}$ is the per-token weight and π_{ref} is a fixed reference model.:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \frac{1}{\sum_g |o_g|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \min \left\{ r_{i,t} \hat{A}_{i,t}, \text{clip}(r_{i,t}; 1-\epsilon, 1+\epsilon) \hat{A}_{i,t} \right\} - \beta D_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}}) \quad (1)$$

Each reward r_i is computed via a rule-based scheme combining **Format Reward** and **Answer Reward**. The Format Reward checks structural compliance using regex-based rules: the model must place reasoning in `<think>` tags and the final answer in `<answer>` tags. A correct format yields +1, while any deviation gives -1. The Answer Reward evaluates correctness relative to the task-specific ground truth: a fully correct answer receives +2, and an unparseable or missing answer -2. Task-specific criteria guide this evaluation. A deduction (Propositional satisfiability) output is correct only if it satisfies all Boolean formulas; an induction (masked-sequence completion) prediction is valid if the predicted term fits the sequence pattern; and an abduction (inverse rule-graph search) answer is accepted when its premises form the minimal consistent causal path from evidence to target. The total reward is then normalized across the group to produce \hat{A}_i .

3.2.2 Stage B: Parameter-Space Merging for Meta-Ability Integration

To unify the strengths of models specialized in distinct meta-abilities, we adopt *parameter-space merging*, which enables: (i) a cost-efficient combination of complementary competencies without additional training, and (ii) a high-quality initialization for domain-specific fine-tuning in Stage C.

We denote the parameters of the RL-deduction-, RL-induction-, and RL-abduction-aligned specialists as $\Theta^{(d)}$, $\Theta^{(i)}$, and $\Theta^{(a)}$, respectively. These models, trained separately on their respective meta-abilities, demonstrate highly complementary behaviors, aggregating their predictions. We construct the merged model Θ_{merge} by linearly interpolating the weights of the three specialists:

$$\Theta_{\text{merge}} = \lambda_d \Theta^{(d)} + \lambda_i \Theta^{(i)} + \lambda_a \Theta^{(a)} \quad (2)$$

We control the contribution of each specialist model via scalar weights λ_d , λ_i , and λ_a . These coefficients determine the relative influence of each meta-ability in the merged model. Notably, uniform weighting is not assumed. Unequal allocation may better reflect the asymmetry in task difficulty or generalization benefit across reasoning modes. Optimal weights are selected empirically based on the performance. In Section 5.2, we compare with nonlinear interpolation methods and assess the robustness of linear interpolation with respect to weight selection.

3.2.3 Stage C: Domain-Specific Reinforcement Learning Training

To evaluate whether meta-ability alignment provides a stronger foundation for downstream learning, we apply reinforcement learning to the aligned checkpoints using domain-specific data, specifically math tasks. For a fair and controlled comparison with instruction-tuned baselines, we follow the experimental settings of SimpleRL-Zoo (Zeng et al., 2025). Specifically, we adopt a rule-based reward function that assigns +1 to correct completions and 0 otherwise, and also use GRPO. These choices match SimpleRL-Zoo and help isolate the impact of initialization, ensuring performance gains arise from meta-ability alignment rather than optimization differences.

4 Experimental Performance

4.1 Experimental Setup

Dataset and Models. For each meta-ability task we introduce a specific difficulty-controlling parameter. Thus, we generate multiple difficulty levels for every task and adopt the curriculum learning strategy that trains the model level by level from easy to hard (More details about difficulty control are in Appendix B.). With this schedule, the 7B model converges by Level 2, and its reward does not improve further at higher levels, so we restrict training to Levels 1–2. The 32B model occasionally benefits from Level 3 but shows unstable reward curves. Therefore, we also use only the first two levels for it. We sample 200 instances per task per level for the 7B model and 2000 instances per task per level for the 32B model. For further domain-specific RL training, we adopt the same dataset as SimpleRL-Zoo (Zeng et al., 2025). The base model and instruct model we used in this work are from Qwen-2.5 series (Team et al., 2024).

Evaluation Setup. To validate the generalization of these meta-ability, we select 7 benchmarks from math, coding and science domain. In math tasks, we utilize MATH-500 (Hendrycks et al., 2021), the full AIME 1983–2024 corpus (Veeraboina, 2023), the recent AMC 2023 (Mathematical Association of America, 2023) and AIME 2024 (Mathematical Association of America, 2024) sets and the Olympiad-level OmniMath subset (Gao et al., 2025) as the evaluation benchmark. LiveCodeBench (LCB) is designed for code generation (Jain et al., 2024), and GPQA (Rein et al., 2023) is aimed for graduate-level science QA. For most

benchmarks, we report pass@1 results using temperature 0.0 and top-p 1.0. While, for AIME 2024 and AMC 2023, containing fewer problems, we use average accuracy (avg@32), computed 32 samples per problem with temperature 1.0 and top-p 0.95. **Hyperparameter for Three Stage Training.** We utilize VERL (Sheng et al., 2024) for meta-abilities alignment and continual reinforcement learning, and adopt MERGEKIT (Goddard et al., 2024) for parameter-space merging to integrate distinct meta-abilities. The optimal weighting coefficients are set to $\lambda_d = 1.0$, $\lambda_i = 0.2$, and $\lambda_a = 0.1$. We provide the comprehensive training setup in Appendix C.

4.2 Out-of-Domain Generalization of Meta-Abilities

Table 1 shows that meta-ability alignment, trained solely on synthetic diagnostic tasks, already transfers to seven unseen benchmarks (covering five math evaluations: Math500, AIME, AIME24, AMC23, and Olympic, as well as the LiveCodeBench coding test and the GPQA science QA set). At the 7B scale, the RL-induction-aligned model provides the largest average improvement, lifting the mean score by 1.7% (from 38.8% to 39.8%), whereas the RL-deduction-aligned model yields the largest single math task gain with a 2.8% increase on Math500 (from 73.0% to 75.8%). Integrating the three meta-abilities in the Merged model further raises the overall average by 2.5% (to 37.8%), confirming that the abilities combine constructively, e.g., boosting LCB from 25.7% to 26.0% and GPQA from 27.2% to 33.5%. An Oracle Ensemble (correct if any aligned model succeeds) lifts Math average by 11.1% to 49.9%, highlighting untapped complementarity for better fusion. Additionally, our merged model on 7B size, trained without domain data, tops concurrent AbsoluteZero Reasoner (Zhao et al., 2025) by 2.7% over AZR (Base) and 2.9% over AZR (Coder). While AbsoluteZero scales 768 seeds to 16,384 triplets via 500×384 self-play, our method reaches higher accuracy with only 1,200 examples, underscoring the effectiveness of factorization-plus-interpolation.

Scaling to the 32B model amplifies the pattern: each aligned model surpasses the Qwen2.5-32B-Instruct baseline, yielding a mean gain of 3.1% on the Math-overall metric and 2.6% on the overall average. For reference, the baseline already scores 46.9% on Math-overall and 44.6% overall. Deduction alignment contributes a +4.3% jump on Math-overall (to 51.2%), with peaks on AIME (+5.7%)

and AMC23 (+5.8%), plus gains of +2.6% on LiveCodeBench and +0.6% on GPQA. RL-Induction and RL-abduction-aligned models follow with respective Math-overall gains of +2.7% and +2.4%, and overall average lifts of +2.3% and +2.0%. Additionally, the Merged checkpoint improves the overall average by 3.5%, with a standout 4.4% gain on the Math-average (from 46.9% to 51.3%) and strong single-task performance such as +6.8% on AMC23. A full Oracle Ensemble run shows an additional 10.8% lift in the math average (to 57.7%), indicating that the three reasoning modes remain highly complementary even at larger scale.

4.3 Scalable Gains from Meta-Abilities Alignment

Table 2 shows that embedding meta-ability alignment into the domain-specific RL pipeline lifts performance at both 7B and 32B. For 7B, the instruction baseline scores 38.8 in math and 35.3 overall. Domain-RL-Ins reaches 41.2 and 37.8. Starting RL from the meta-merged checkpoint, Domain-RL-Meta reaches 43.0 and 39.0. Gains are most visible on AIME and Olympic, while code (LCB) and science (GPQA) stay stable or improve slightly.

At 32B, the instruction baseline is 46.9 in math and 44.6 overall. Domain-RL-Ins rises to 50.3 and 47.4, and Domain-RL-Meta to 52.3 and 48.8. Math improves by about 7% with Domain-RL-Ins and about twelve percent with Domain-RL-Meta, indicating that teaching core reasoning before domain adaptation both raises the starting point and increases the return from task-specific tuning.

4.4 Additional Baselines and Ablations

Unified multi-task RL and SFT controls on synthetic tasks. To rule out gains from simply seeing more synthetic data, we add controls trained on the *same* task generators and verifiers: (i) a unified multi-task RL baseline optimizing one model on the union of Deduction/Induction/Abduction tasks (curriculum and mixed-per-batch), and (ii) SFT baselines on the same tasks (answer-only; trace+answer). These comparisons show our gains come from RL-based, factorized meta-ability alignment rather than data exposure alone (Appendix E). **Early vs. late merging and Domain-RL initialization.** Beyond Table 2, we ablate *when* to merge and *what* to initialize from. We compare *early merging* (merge specialists, then run domain RL) vs. *late merging* (run domain RL per specialist, then merge), and Domain-RL initialized from each

Size	Model	Math				Coding	Science	Average		
		Math500	AIME	AIME24 (Avg@32)	AMC23 (Avg@32)	Olympic	LCB	GPQA	Math	Overall
7B	Qwen2.5-7B-Instruct	73.0	22.4	10.7	50.8	37.3	25.7	27.2	38.8	35.3
	RL-Deduction-Aligned	75.8	22.6	10.2	51.4	39.3	25.8	31.5	39.9(+1.1)	36.7(+1.4)
	RL-Induction-Aligned	75.0	22.5	11.8	52.3	37.5	27.0	33.0	39.8(+1.0)	37.0(+1.7)
	RL-Abduction-Aligned	75.2	22.7	11.4	49.1	38.4	26.8	31.9	39.4(+0.8)	36.5(+1.2)
	Merged Model	77.8	22.9	11.5	52.3	40.4	26.0	33.5	41.0(+2.2)	37.8(+2.5)
	Oracle Ensemble	85.5	32.1	18.0	67.1	46.7	–	–	49.9(+11.1)	–
32B	Qwen2.5-32B-Instruct	79.8	31.2	15.3	62.7	45.6	39.5	38.0	46.9	44.6
	RL-Deduction-Aligned	83.8	36.9	19.4	68.5	47.4	42.1	38.6	51.2(+4.3)	48.1(+3.5)
	RL-Induction-Aligned	82.6	34.8	18.6	66.2	45.8	41.7	38.4	49.6(+2.7)	46.9(+2.3)
	RL-Abduction-Aligned	83.8	33.3	17.5	65.9	46.2	41.1	38.6	49.3(+2.4)	46.6(+2.0)
	Merged Model	84.2	36.0	19.7	69.5	46.9	41.8	38.4	51.3(+4.4)	48.1(+3.5)
	Oracle Ensemble	88.1	43.2	27.3	76.8	53.0	–	–	57.7(+10.8)	–

Table 1: Performance of meta-ability-aligned models, merged ensembles, and oracle upper bounds on 7 benchmarks at both 7B and 32B parameter scales, illustrating consistent gains from scaling

Size	Model	Math				Coding	Science	Average		
		Math500	AIME	AIME24 (Avg@32)	AMC23 (Avg@32)	Olympic	LCB	GPQA	Math	Overall
7B	Qwen2.5-7B-Instruct	73.0	22.4	10.7	50.8	37.3	25.7	27.2	38.8	35.3
	Domain-RL-Ins	78.2	23.6	11.9	53.2	39.3	25.1	33.0	41.2(+2.4)	37.8(+2.5)
	Domain-RL-Meta	78.8	27.7	12.6	54.7	41.0	25.4	33.1	43.0(+4.2)	39.0(+3.7)
32B	Qwen2.5-32B-Instruct	79.8	31.2	15.3	62.7	45.6	37.5	38.0	46.9	44.6
	Domain-RL-Ins	83.0	36.5	18.6	67.5	46.1	41.8	38.2	50.3(+3.4)	47.4(+2.8)
	Domain-RL-Meta	84.6	38.2	19.8	70.4	48.7	41.6	38.6	52.3(+5.4)	48.8(+4.2)

Table 2: Comparison of 7B- and 32B-scale baseline instruction models and our domain-specific RL variants across math, code, and science benchmarks. *Domain-RL-Ins* denotes continual domain-specific RL starting from instruction model; *Domain-RL-Meta* applies the same RL schedule but from a meta-ability-merged initialization, yielding a higher attainable performance ceiling.

single specialist vs. the *merged* checkpoint. Results show merging is most effective *before* domain RL, as domain RL can homogenize specialists and reduce complementarity (Appendix F).

4.5 Robust gains from base initialization

To verify that our pipeline is not contingent on instruction-tuned checkpoints, we start from a *base* model and train meta-ability specialists, merge them, and optionally apply domain RL. As shown in Table 3, the merged model already delivers large gains over the base, and a light round of domain RL further lifts the ceiling, exceeding a strong RL-style baseline (SimpleRL-Zoo) on most metrics.

5 Analysis

5.1 Dose LRM really learn meta-abilities?

To further validate whether meta-ability alignment genuinely enhances the expected advanced reasoning abilities, we adopt the cognitive behavior framework proposed by Gandhi et al. (2025) and utilize OpenAI o4-mini (Hurst et al., 2024) to identify reasoning-related behaviors. A detailed correspondence between the cognitive behaviors associated with the three meta-abilities, their representative statements and evaluation prompts are provided in

Appendix D. We report the proportion of model responses that exhibit these cognitive behaviors across the three meta-abilities.

Relative to the instruction model, the deduction behavior ratio rises from 24.3% to 29.4% in the RL-Deduction-Aligned specialist and remains high at 28.3% in the merged model. Induction increases from 10.2% to 13.1% in the RL-Induction-Aligned specialist, and further to 14.0% in the merged model. Because abductive behaviors are sparse, we report behaviour counts rather than percentages. from 9 in instruct model to 26 in RL-Abduction-Aligned specialist and 18 in the merged model. These trends show that meta-ability alignment reliably amplifies targeted behaviors, while merged model preserves (and for induction, slightly improves) these gains, consistent with our deduction-heavy merge choice.

5.2 Which merging strategy is optimal for meta-ability integration?

Empirical comparison of merge methods. Beyond linear interpolation, we also compare three nonlinear merge families (Wortsman et al., 2022), TIES (task-informed sparsified delta add-back), DARE-TIES (density-adaptive, sign-consistent

Model	Math					Coding	Science	Average	
	Math500	AIME	AIME*24 (Avg@32)	AMC*23 (Avg@32)	Olympiad	LCB (Coding)	GPQA (Science)	Math Avg.	Overall Avg.
Qwen2.5-7B-Base	44.2	8.8	4.3	26.8	27.7	17.5	23.7	22.4	21.9
Base-Deduction-Aligned	70.0	18.0	9.8	46.3	36.5	22.8	30.4	36.1(+13.7)	33.4(+11.5)
Base-Induction-Aligned	68.5	17.0	9.5	45.5	35.0	21.8	29.5	35.1(+12.7)	32.4(+10.5)
Base-Abduction-Aligned	67.8	16.8	9.3	45.0	34.6	21.5	29.2	34.7(+12.3)	32.0(+10.1)
Base-Merged Model	73.5	19.2	10.1	47.8	38.0	23.0	31.0	37.7(+15.3)	34.7(+12.8)
SimpleRL-Zoo	75.6	26.5	11.7	53.2	38.4	23.3	29.8	41.1(+18.7)	36.9(+15.0)
Base Domain-RL-Meta	77.4	26.8	11.9	54.2	40.0	23.6	33.2	42.1(+19.7)	38.2(+16.3)

Table 3: Effectiveness from *base* initialization (7B). Only the average columns report absolute improvements over the base model in parentheses.

TIES), and Segmented SLERP (layer-wise spherical interpolation with a V-shaped mix). Across Math500, AIME, AIME24, and AMC23, linear interpolation consistently achieves the best scores while being the simplest and most compute-efficient, outperforming the more complex alternatives we tested.

Setting	Math500	AIME	AIME24	AMC23
TIES	67.0	19.1	9.9	47.9
DARE-TIES	73.0	20.7	10.5	50.2
SS	72.0	20.5	10.1	49.8
LI (ours)	77.8	22.9	11.5	52.3

Table 4: Merge method comparison. Linear interpolation (LI below) outperforms nonlinear methods while requiring no extra computing. SS is Segmented SLERP.

Sensitivity of linear interpolation and choice of best weights. To probe the robustness of linear interpolation, we sweep interpolation weights $(\lambda_d, \lambda_i, \lambda_a)$ over deduction, induction, and abduction specialists. Starting from a neutral average $(0.33, 0.33, 0.33)$, we observe that tilting toward a deduction-heavy mixture reliably improves all metrics; rotating dominance to induction or abduction degrades accuracy. Fixing the deduction expert at full weight $(\lambda_d = 1.0)$ and lightly “sprinkling” induction/abduction yielded a broad, stable optimum around $(1.0, 0.20, 0.10)$. This setting matches our qualitative analysis in which deduction contributes the highest-frequency reasoning micro-skills, with induction/abduction adding complementary but lower-frequency behaviors. Importantly, performance varies smoothly around this point, showing stability over knife-edge sensitivity.

5.3 Empirical Validation of “Aha Moment”

An “aha moment” is typically characterized by a sharp surge in accuracy (such as greater than 5–10%) over just a few training steps on benchmarks, coinciding with the emergence of advanced behaviors such as long chain-of-thought reasoning, self-correction, and backtracking. This phe-

$(\lambda_d, \lambda_i, \lambda_a)$	Math500	AIME	AIME24	AMC23
Qwen2.5-7B-Ins	73.0	22.4	10.7	50.8
$(0.33, 0.33, 0.33)$	74.0	20.8	10.7	50.1
$(0.70, 0.15, 0.15)$	76.2	22.5	11.4	52.0
$(0.15, 0.70, 0.15)$	73.8	21.4	11.3	49.1
$(0.15, 0.15, 0.70)$	73.8	22.9	10.5	49.4
$(1.00, 0.10, 0.10)$	76.8	23.7	11.2	52.7
$(1.00, 0.20, 0.10)$	77.8	22.9	11.5	52.3

Table 5: Linear-merge weight sweep. A deduction-heavy mixture is best. $(1.0, 0.20, 0.10)$ is a simple, robust choice. Qwen2.5-7B-Ins is the Instruct model.

nomenon is also noted by previous work (Gandhi et al., 2025), as well as in prior observations from DeepSeek R1, SimpleRL-Zoo, and Logic-RL. Empirically, compared with SimpleRL-Zoo’s 40–80 iterations and 10k–20k examples (batch = 256) to hit its accuracy knee, we elicit the inflection and those behaviors in around 25 iterations with batch = 16 (400 examples), yet still achieve robust gains 5.7% on a 7B model and 9.4% on a 32B model. This demonstrates a faster, data-efficient, and more predictable route to the “aha moment”.

6 Conclusion

In this work, we avoid unpredictable “aha moments” by aligning deduction, induction, and abduction with self-verifiable synthetic tasks, training specialists, and merging them into one model. The merged checkpoint outperforms the instruction baseline by > 10% on diagnostics and up to 2% on seven math/code/science benchmarks. Starting domain-specific RL from this meta-ability-aligned model adds around 4% and the gains widen from 7B to 32B, indicating scalable benefits. This pipeline also work robustly from base initialization, lifts average accuracy by by 16.3% over the 7B Base model, showing strong generalization of merged meta-abilities. Our method needs no human labels, supports rapid iteration with control, and improves interpretability by isolating meta-abilities. Moving from emergent “aha” to modular, self-verifying training makes reasoning predictable and reliable.

623 Limitations

624 Our work still leaves several open questions. First,
625 the training tasks we use are intentionally sym-
626 bolic—Boolean formulae, masked sequences, and
627 rule graphs—to guarantee automatic verification,
628 but this controlled design inevitably omits the rich-
629 ness and noise of natural-language arguments, dia-
630 grams, or multimodal evidence chains. Extending
631 the suite in those directions would test whether the
632 same alignment signals transfer. Second, although
633 we evaluate on seven established math, coding, and
634 science benchmarks, we have not probed areas such
635 as commonsense dialogue, legal reasoning, or cre-
636 ative writing, which demand different blends of
637 deduction, induction, and abduction; adding such
638 tasks would give a more rounded view of gener-
639 alization. Finally, we observe that deduction cur-
640 rently contributes the largest share of the overall
641 gain, while abduction lags behind, suggesting ei-
642 ther that our backward-reasoning task could be
643 refined or that its reward needs stronger shaping.
644 Addressing these points will not overturn the main
645 findings, but they offer concrete avenues for mak-
646 ing explicit meta-ability alignment more robust,
647 broadly applicable, and theoretically grounded.

648 References

649 Jiefeng Chen, Jie Ren, Xinyun Chen, Chengrun Yang,
650 Ruoxi Sun, and SercanÖ. Arik. 2025. [Sets: Leverag-
651 ing self-verification and self-correction for improved
652 test-time scaling](#). *arXiv preprint arXiv:2501.19306*.

653 Justin Chih-Yao Chen, Zifeng Wang, Hamid Palangi,
654 Rujun Han, Sayna Ebrahimi, Long Le, Vincent Perot,
655 Swaroop Mishra, Mohit Bansal, Chen-Yu Lee, and 1
656 others. 2024. Reverse thinking makes llms stronger
657 reasoners. *arXiv preprint arXiv:2411.19865*.

658 Google DeepMind. 2024. Gemini 2.5 pro: The lat-
659 est gemini multimodal model. [https://deepmind.
660 google/technologies/gemini/](https://deepmind.google/technologies/gemini/). Accessed: 2025-
661 05-15.

662 Hugging Face. 2025. [Open r1: A fully open reproduc-
663 tion of deepseek-r1](#).

664 Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh,
665 Nathan Lile, and Noah D Goodman. 2025. Cognitive
666 behaviors that enable self-improving reasoners, or,
667 four habits of highly effective stars. *arXiv preprint
668 arXiv:2503.01307*.

669 Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo
670 Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang
671 Chen, Runxin Xu, Zhengyang Tang, Benyou Wang,
672 Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei
673 Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu,

and Baobao Chang. 2025. Omnimath: A universal
olympiad level mathematical benchmark for large
language models. In *Proc. of ICLR*.

Charles Goddard, Shamane Siriwardhana, Malikeh
Ehghaghi, Luke Meyers, Vladimir Karpukhin, Brian
Benedict, Mark McQuade, and Jacob Solawetz. 2024.
[Arcee’s MergeKit: A toolkit for merging large lan-
guage models](#). In *Proceedings of the 2024 Confer-
ence on Empirical Methods in Natural Language
Processing: Industry Track*, pages 477–485, Miami,
Florida, US. Association for Computational Linguis-
tics.

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

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul
Arora, Steven Basart, Eric Tang, Dawn Song, and Ja-
cob Steinhardt. 2021. Measuring mathematical prob-
lem solving with the math dataset. *arXiv preprint
arXiv:2103.03874*.

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. *arXiv preprint
arXiv:2410.21276*.

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard-
son, Ahmed El-Kishky, Aiden Low, Alec Helyar,
Aleksander Madry, Alex Beutel, Alex Carney, and 1
others. 2024. Openai o1 system card. *arXiv preprint
arXiv:2412.16720*.

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fan-
jia Yan, Tianjun Zhang, Sida I. Wang, Armando
Solar-Lezama, Koushik Sen, and Ion Stoica. 2024.
Livecodebench: Holistic and contamination free eval-
uation of large language models for code. *arXiv
preprint arXiv:2403.07974*.

Aviral Kumar, Vincent Zhuang, Rishabh Agarwal,
Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli,
Shariq Iqbal, Colton Bishop, Rebecca Roelofs, and
1 others. 2024. Training language models to self-
correct via reinforcement learning. *arXiv preprint
arXiv:2409.12917*.

Ruotian Ma, Peisong Wang, Cheng Liu, Xingyan Liu,
Jiaqi Chen, Bang Zhang, Xin Zhou, Nan Du, and Jia
Li. 2025. [S²r: Teaching llms to self-verify and self-
correct via reinforcement learning](#). *arXiv preprint
arXiv:2502.12853*.

Mathematical Association of America. 2023. American
Mathematics Competitions (AMC) 2023.

Mathematical Association of America. 2024. American
Invitational Mathematics Examination (AIME) 2024.

OpenAI. 2025. Openai o3 and o4-mini system card.

674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728

729	Jiayi Pan, Junjie Zhang, Xingyao Wang, Lifan Yuan, Hao Peng, and Alane Suhr. 2025. Tinyzero. https://github.com/Jiayi-Pan/TinyZero . Accessed: 2025-01-24.	785
730		786
731		787
732		788
733	Charles Sanders Peirce. 1931. Collected papers of charles sanders peirce, volume 2: Elements of logic. pages Chapter 7, especially paragraphs 619–645. Harvard University Press. Peirce’s formulation of deduction, induction, and abduction as distinct forms of inference.	789
734		
735		
736		
737		
738		
739	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. GPQA: A graduate-level google-proof q&a benchmark. <i>arXiv preprint arXiv:2311.12022</i> .	790
740		791
741		792
742		793
743		794
744	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, and 1 others. 2024. Deepseek-math: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .	795
745		796
746		797
747		798
748		799
749		
750	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. Hybridflow: A flexible and efficient rlhf framework. <i>arXiv preprint arXiv:2409.19256</i> .	800
751		801
752		802
753		803
754		804
755	Qwen Team and 1 others. 2024. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> , 2:3.	805
756		806
757	Hemish Veeraboina. 2023. Aime problem set 1983-2024 .	807
758		808
759	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	809
760		
761		
762		
763		
764		
765	Liang Wen, Yunke Cai, Fenrui Xiao, Xin He, Qi An, Zhenyu Duan, Yimin Du, Junchen Liu, Lifu Tang, Xiaowei Lv, and 1 others. 2025. Light-r1: Curriculum sft, dpo and rl for long cot from scratch and beyond. <i>arXiv preprint arXiv:2503.10460</i> .	810
766		811
767		812
768		813
769		
770	Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and 1 others. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In <i>International conference on machine learning</i> , pages 23965–23998. PMLR.	
771		
772		
773		
774		
775		
776		
777		
778	Zhenyu Wu, Qingkai Zeng, Zhihan Zhang, Zhaoxuan Tan, Chao Shen, and Meng Jiang. 2024. Large language models can self-correct with key condition verification . <i>arXiv preprint arXiv:2405.14092</i> . EMNLP 2024.	
779		
780		
781		
782		
783	xAI. 2025. Grok 3.5: Advanced reasoning ai model by xai. https://grok.x.ai/ . Accessed: 2025-05-15.	
784		

A Data-Synthesis Pseudocode and Additional Task Examples

A.1 Deduction

Algorithm 1: Propositional Satisfiability

Input: Clause set Φ drawn with hyper-parameters (n_ℓ, f_ℓ, d_ℓ) as defined in Appendix B.

Output: A satisfying assignment σ or UNSAT.

```

1 Function DPLL( $\Phi, \sigma$ ):
2   if  $\forall C \in \Phi : C$  satisfied by  $\sigma$  then
3     return  $\sigma$ 
4   if  $\exists C \in \Phi : C$  falsified by  $\sigma$  then
5     return UNSAT
6   while  $\exists$  unit clause  $l$  do
7      $\sigma \leftarrow \sigma \cup \{l\}$ ;
8      $\Phi \leftarrow \text{SIMPLIFY}(\Phi, l)$ 
9   foreach literal  $p$  that appears with only
10    one polarity do
11      $\sigma \leftarrow \sigma \cup \{p\}$ ;
12      $\Phi \leftarrow \text{SIMPLIFY}(\Phi, p)$ 
13   pick the first unassigned variable  $x$ 
14   for  $b \in \{\text{TRUE}, \text{FALSE}\}$  do
15      $\sigma' \leftarrow \sigma \cup \{x=b\}$ 
16      $\Phi' \leftarrow \text{SIMPLIFY}(\Phi, x=b)$ 
17     result  $\leftarrow$  DPLL( $\Phi', \sigma'$ )
18     if result  $\neq$  UNSAT then return
19       result
20   return UNSAT

```

17 **Main:**

18 **return** DPLL(Φ, \emptyset)

Figure 4 presents examples of propositional satisfiability tasks at difficulty levels 1 through 3, illustrating how clause structures become increasingly complex. While exhaustive truth-assignment trials quickly become intractable, Algorithm 1 begins with unit propagation and pure-literal elimination to assign all forced literals. It then branches on unassigned variables, backtracking when conflicts arise. This process instantiates the hypothesis–consequence–contradiction–refinement cycle characteristic of deductive reasoning, while preserving an exact satisfaction oracle.

A.2 Induction

Algorithm 2: Masked-Sequence Completion

Input: Observed prefix

$S = [s_1, \dots, s_{n-1}, ?]$; cycle length

$k \leq 7$; operator alphabet

$\Sigma \subseteq \{+, -, \times\} \times \{1, \dots, 4\}$.

Output: Predicted missing value

```

1  $\mathcal{B} \leftarrow \{(\langle \rangle, \text{SCORE} = 0)\}$ 
2 for  $t \leftarrow 1$  to  $k$  do
3   if  $\mathcal{B} = \emptyset$  then
4     return FAIL
5   end
6    $\mathcal{B}' \leftarrow \emptyset$ 
7   foreach candidate  $(\vec{\sigma}, \text{SCORE}) \in \mathcal{B}$  do
8     foreach  $op \in \Sigma$  do
9        $\vec{\sigma}' \leftarrow \vec{\sigma} \parallel op$ 
10      if CONSISTENT( $\vec{\sigma}', S$ ) then
11        SCORE'  $\leftarrow$  MDL( $\vec{\sigma}'$ )
12        add  $(\vec{\sigma}', \text{SCORE}')$  to  $\mathcal{B}'$ 
13      end
14    end
15  end
16  keep top  $B$  elements of  $\mathcal{B}'$  by SCORE
17   $\mathcal{B} \leftarrow \mathcal{B}'$ 
18 end
19 return value produced by best  $\vec{\sigma}$  when
20    applied to  $S$ 

```

Figure 5 shows examples of masked-sequence completion tasks at difficulty levels 1 through 3, each featuring a hidden cycle of arithmetic operations of length k that induces a super-exponential hypothesis space. Algorithm 2 grows candidate operation sequences up to k , prunes those inconsistent with the observed prefix, scores survivors by minimal-description-length, and iteratively refines its beam until the optimal pattern is found, thereby realising the hypothesise–test–refine loop central to inductive reasoning.

A.3 Abduction

Figure 6 shows examples of reverse rule-graph structures at difficulty levels 1 through 3; Algorithm 3 begins at the goal g and recursively traverses candidate edges in reverse, memoizing results and pruning cycles until it either reaches a fact in F or exceeds the depth limit. The task therefore asks: given facts F and a goal g , determine reachability and, if successful, return the minimal explanation tree \mathcal{T} comprising only the utilized

DEDUCTION – Level 1	DEDUCTION – Level 2	DEDUCTION – Level 3
<p>Q: Below are some nested formulas:</p> <ul style="list-style-type: none"> - $(A \vee \neg C)$ - $(C \oplus A)$ - $(\neg C \oplus \neg D)$ - $(\neg A \leftrightarrow \neg B)$ <p>Please list the truth value of each variable</p>	<p>Q: Below are some nested formulas:</p> <ul style="list-style-type: none"> - $((\neg C \oplus \neg F) \leftrightarrow (E \rightarrow \neg C))$ - $((\neg E \oplus \neg C) \rightarrow (\neg D \leftrightarrow E))$ - $((\neg E \rightarrow B) \oplus \neg(\neg B))$ - $((\neg C \vee A) \wedge (\neg A \vee B))$ <p>Please list the truth value of each variable</p>	<p>Q: Below are some nested formulas:</p> <ul style="list-style-type: none"> - $\neg(((\neg F \vee H) \wedge (B \rightarrow \neg H)))$ - $((((G \vee \neg B) \rightarrow (A \vee F)) \oplus (\neg(\neg D) \leftrightarrow \neg(\neg B))))$ - $((((G \wedge D) \rightarrow (B \leftrightarrow \neg C)) \oplus (\neg(\neg A) \rightarrow (D \rightarrow \neg G))))$ - $((((A \oplus B) \leftrightarrow (F \rightarrow \neg B)) \leftrightarrow (\neg(\neg H) \oplus (\neg H \rightarrow \neg D))))$ - $((((A \vee \neg C) \wedge \neg(G)) \oplus ((D \vee \neg C) \vee (\neg D \vee A))))$ - $((\neg(\neg B) \oplus \neg(H)) \oplus (\neg(B) \wedge (F \leftrightarrow E)))$ <p>Please list the truth value of each variable</p>

Figure 4: Examples of propositional satisfiability problems with difficulty levels ranging from 1 to 3.

hyper-edges. This goal-first, premise-verification loop captures abductive reasoning, while graph depth, branching factor, and distractor rules enable us to scale difficulty yet preserve an exact oracle and dense edge-level supervision for RL agents.

B Controlling Task Difficulty in Synthesis

We expose *one to four orthogonal hyper-parameters per task* and let the generator sample from progressively wider intervals as the curriculum advances. Each hyper-parameter has a clear algorithmic interpretation, so the reader can verify that higher levels provably enlarge the effective search space.

• Deduction – Propositional Satisfiability.

Only the nested-formula mode is active in the released code. A level ℓ tuple $\langle n_\ell, f_\ell, d_\ell \rangle$ controls:

- n_ℓ : number of propositional variables (`n_vars` in the script).
- f_ℓ : number of independent formulas generated per instance.
- d_ℓ : maximum parenthesis nesting depth (`max_depth`).

The Boolean assignment space grows as 2^{n_ℓ} , while deeper nesting couples far-apart variables via implications, making local heuristics ineffective. We empirically observe that moving from $(n, d, f) = (6, 2, 3)$ to $(10, 4, 6)$ increases the search time of a basic backtracking solver, which systematically tries possible truth assignments, by two orders of magnitude.

• Induction – Masked-Sequence Completion.

For each level we sample *cycle_length* (*pattern length*, k ($1-7$)) $k \in \{1, \dots, 7\}$, which is the number of distinct operations that repeat, the operator alphabet $\Sigma \subseteq \{+, -, \times\} \times \{1, \dots, 4\}$, and the number of masked positions m ($1-2$).

- A longer cycle k means the model must retain a larger latent state to describe the full rule, e.g. $\langle +2, \times 2, -3, +4, \times 2, -5 \rangle$.
- Mixing symbolic with numeric operators further enlarges the hypothesis set to $|\Sigma|^k$; at $k = 7$ this already exceeds one million templates.
- Multiple blanks m break the symmetry that a single missing suffix would have, requiring the agent to interpolate rather than extrapolate.

Hence the difficulty rises super-exponentially in k and m , emphasizing abstraction and generalization over brute-force enumeration.

• Abduction – Reverse Rule-Graph Search.

The generator draws *chain_depth* d , *num_goals* g , *distractor_count* h and *cycle_prob* γ from level-specific ranges, e.g. $d=3-4, g=2-3, h=5-7, \gamma=0.10-0.25$ for level 3.

- Larger d lengthens the *minimal backward proof*, forcing deeper recursion.
- Each extra goal g multiplies the number of sink nodes the agent must explain in one episode.

916	– Distractors h are additional hyper-edges	on low-entropy regimes before encountering the	964
917	whose premises share symbols with real	full combinatorial explosion.	965
918	rules; they explode the forward branching		
919	factor, so blind forward chaining be-		
920	comes exponentially slower while a goal-		
921	directed agent is unaffected.		
922	– A non-zero cycle probability γ rewires		
923	some edges to form backward loops.		
924	Proof search must therefore keep a vis-		
925	ited set and handle NOT-YETS—a hall-		
926	mark of abductive reasoning.		
927	Jointly increasing $\langle d, g, h, \gamma \rangle$ produces a com-		
928	pounded state space whose size we approxi-		
929	mate as $\mathcal{O}((b+h)^d)$ for forward search, where		
930	b is the intrinsic out-degree of the knowledge		
931	graph.		
932	• Abduction – Reverse Rule-Graph Search.		
933	The generator draws <code>chain_depth</code> d (max-		
934	imal backward-proof depth), <code>num_goals</code> g		
935	(number of sink goals), <code>distractor_count</code>		
936	h (spurious hyper-edges), and <code>cycle_prob</code>		
937	γ (edge-rewiring probability) from level-		
938	specific ranges, e.g. $d=3-4$, $g=2-3$, $h=5-7$,		
939	$\gamma=0.10-0.25$ for level 3.		
940	– Larger d lengthens the <i>minimal back-</i>		
941	<i>ward proof</i> , forcing deeper recursion.		
942	– Each extra goal g multiplies the number		
943	of sink nodes the agent must explain in		
944	one episode.		
945	– Distractors h share symbols with true		
946	rules, exploding the forward branching		
947	factor—blind forward chaining slows ex-		
948	ponentially, while a goal-directed agent		
949	is unaffected.		
950	– A non-zero cycle probability γ rewires		
951	some edges to form backward loops;		
952	proof search must therefore maintain a		
953	visited set and handle NOT-YETS, a hall-		
954	mark of abductive reasoning.		
955	Jointly increasing $\langle d, g, h, \gamma \rangle$ produces a com-		
956	pounded state space whose size we approxi-		
957	mate as $\mathcal{O}((b+h)^d)$ for forward search, where		
958	b is the intrinsic out-degree of the knowledge		
959	graph.		
960	During curriculum training we feed the model in-		
961	stances in ascending level order (e.g. $d=2 \rightarrow 7$ for		
962	abduction, $n=6 \rightarrow 10$ for deduction, $k=1 \rightarrow 7$		
963	for induction), so it first acquires stable heuristics		
		C Comprehensive Multi-Stage Training Setup	966
			967
		Stage A: Meta-Abilities Alignment We warm-	968
		start from a frozen instruction model and train	969
		on three synthetic diagnostic corpora—Deduction	970
		(propositional SAT), Induction (masked-sequence	971
		completion) and Abduction (inverse rule-graph	972
		search), using the GRPO. Rollouts are generated	973
		with vLLM (temperature 0.7, $n = 4$, tensor-	974
		parallel size 2, max tokens $\approx 9\,000$), and optimized	975
		with AdamW (LR 5×10^{-7}) over 20 epochs. We	976
		use a training batch size of 32, PPO mini-batch	977
		size of 256, and PPO micro-batch size of 32. Re-	978
		wards are normalized per group to produce \hat{A}_i , and	979
		we regularize via a KL penalty ($\beta = 0.001$, low-	980
		variance estimator) to the frozen reference model.	981
		Stage C: Domain-Specific Reinforcement Learn-	982
		ing on Math Starting from the Stage A check-	983
		point (after MergeKit parameter merging), we fine-	984
		tune on the dataset from SimpleRL Zoo using	985
		GRPO objective to match SimpleRL-Zoo settings.	986
		Inputs are capped at $1\,024 + 3\,072$ tokens; train/val	987
		batch size 16; PPO mini-batch 128, micro per-GPU	988
		2; rollouts with $n = 8$, temperature 1.0, TP size	989
		2. Rewards are rule-based (+1 for correct, 0 oth-	990
		erwise) without format terms to isolate transfer	991
		effects. We use AdamW (LR 5×10^{-7}), clip ra-	992
		tio 0.2, entropy bonus 0.001, and KL to reference	993
		($\beta = 0.001$).	994
		D Quantifying Meta-Ability Behaviors	995
		through Targeted Prompts	996
		Below we present the three carefully crafted	997
		sentence-extraction prompts. One each for Deduc-	998
		tion, Induction and Abduction, that we feed to the	999
		model to capture its core reasoning behaviors. Each	1000
		prompt defines the target micro-skills, constrains	1001
		the extraction to only the relevant sentences, and	1002
		enforces a strict JSON output format so that we	1003
		can automatically compute the frequency of each	1004
		reasoning mode. These prompts form the basis for	1005
		quantifying how often the model engages in each	1006
		type of advanced inference.	1007

Deduction Prompt

Here is a thought record:

<thinking text>

We define ``Deduction'' to consist of exactly two sub-skills:

1. HypothesisProposing: sentences that explicitly introduce a testable hypothesis or scenario, such as:
 - * ``Assume <\dots>, then <\dots>.''
 - * ``Case <number>: <\dots>.''
 - * ``Suppose <\dots>, then <\dots>.''
2. SelfVerification/
SelfCorrection: sentences that explicitly validate or correct a prior result, such as:
 - * ``After checking, <\dots>, so <\dots>.''
 - * ``We see that <\dots> was wrong, so <\dots>.''
 - * ``This shows an error; we must <\dots>.''
 - * ``Correcting that, <\dots>.''

Extraction requirements:

- Extract **all and only** the sentences that clearly match one of the above patterns or equivalent wording that directly reflects the two sub-skills.
- Do **not** extract any explanatory, descriptive, or purely procedural sentences.
- Return a JSON array; each element must be exactly:

```
{
  "category": "
    HypothesisProposing"|
    SelfVerification/
    SelfCorrection",
  "sentence": "<the exact
    extracted sentence>"
}
```
- After the array, output one more field:

```
"total_count": <the total number
  of extracted sentences>
```

Please **strictly** follow this format without any additional commentary.

Induction Prompt

You are an expert reasoning analyst

TASK: Induction Sentence Extraction
(Wide-Net, English Only)

-
1. Input text is wrapped between THINKING: ... END THINKING.
 2. Scan every sentence (and its immediately preceding sentence) and extract only those that state or imply a general rule, trend, or pattern derived from multiple cases or examples.

Key English triggers include words or phrases such as:

- * therefore, thus, hence, it follows that
- * we observe that, observations show, we see that
- * typically, usually, in general, most, majority
- * suggests that, implies that, indicates that
- * pattern, trend, rule, conjecture, generalize
- * for all n, for any k, if ... then ...

Extraction is OK for probabilistic statements (e.g. \ ``Most A are B'') and recursive claims (e.g. \ ``If P(k) holds then P(k+1) holds, so P(n) for all n'').

3. Do NOT extract purely numerical computations, single-step deductions, or abduction sentences.

OUTPUT FORMAT (nothing else):

First, a JSON array of the exact extracted sentences as strings. Then, on the next line:

```
"total_count": <number_of_sentences
>
```

Example:

```
[
  "For any  $n \equiv 0 \pmod{5}$ ,
  the position is losing.",
  "Thus the sequence converges to
  zero for all starting values
  ."
]
"total_count": 2
```

THINKING:
<thinking text>
END THINKING

Remember: output valid JSON and the total_count line, and nothing else.

Abduction Prompt

You are an expert reasoning analyst

Task

1. I will give you a block of text labelled
THINKING: ... END THINKING.
2. Scan every sentence and extract **only** those that exhibit **abduction** (backward reasoning).
3. Exclude any sentences that use forward reasoning (deduction or induction) or are purely descriptive/mathematical.
4. Return a JSON array; each element must be exactly:

```
{
  "sentence": "<the exact
  extracted sentence>",
  "match_type": "template" | "
  close_paraphrase" | "
  clear_abduction_non_template
",
  "matched_pattern": "<template
  code or empty>"
}
```
5. After the array, output a final field on a **new line**:
"total_count": <the total number of extracted sentences>

Abduction vs. Forward Reasoning

- **Abduction** (backward reasoning) starts from an observation or surprising fact and proposes or tests a hypothesis/explanation.
- **Forward reasoning** (deduction / induction) starts from premises or rules and derives consequences --- **do not extract** those.

Abduction Templates (code: exact pattern)

- T1 Given-that -> ``Given that <observed fact>, it must be that <hypothesis>.''
- T2 Since-result -> ``Since <result> holds, one explanation is <hypothesis>.''
- T3 To-account-for -> ``To account for <outcome>, suppose <hypothesis>.''
- T4 Observing-suggests -> ``Observing <phenomenon> suggests <hypothesis>.''
- T5 Because-infer -> ``Because <result>, we infer <hypothesis>.''
- T6 Hypothesis-test-revise -> ``Testing that hypothesis, we find <contradiction>, so <revision>.''

- T7 Bad-consequence-fails -> ``However, this would imply <bad consequence>, thus the explanation fails.''
- T8 Closer-inspection -> ``On closer inspection, <hypothesis> does not hold, so <revision>.''

Close Paraphrase Examples

- ``One possible cause of X is Y.''
- ``A plausible explanation for X could be Y.''
- ``If X happened, then Y might be responsible.''
- ``X seems best explained by Y.''
- ``Assuming Y, we can explain X.''
- ``Y predicts X; therefore, Y is likely.''
- ``Y would predict Z, yet we observe not-Z; hence Y is unlikely.''

Strict Extraction Rules

- * Extract **all and only** sentences performing abduction as defined above.
- * Exclude any forward-reasoning or purely procedural/mathematical sentences.
- * If a sentence matches a template **exactly**, set `"match_type": "template"` and `"matched_pattern"` to its code (e.g. `"T1"`).
- * If it is a tight paraphrase (≤ 3 word changes per clause), set `"match_type": "close_paraphrase"`.
- * Otherwise, if it clearly does abduction but is not a template or close paraphrase, set `"match_type": "clear_abduction_non_template"` and leave `"matched_pattern"` blank.
- * Output must be valid JSON with **no** extra keys or explanatory text.

THINKING:
<thinking text>
END THINKING

Please **strictly** follow the specified output format and include nothing else.

1012

E Additional Baselines and Ablations

1013

This appendix consolidates additional baselines and design-justification ablations that were introduced in the rebuttal, including (i) unified multi-task RL and SFT baselines on the same synthetic meta-ability tasks, and (ii) early-vs-late merging and domain RL initialization from single specialists

1014

1015

1016

1017

1018

1019

1020	versus from the merged model.		
1021	E.1 Unified Multi-Task RL on the Synthetic		
1022	Meta-Ability Tasks		
1023	A natural baseline is to train a single model directly		
1024	on the union of Deduction, Induction, and Abduc-		
1025	tion tasks, instead of training three specialists and		
1026	merging. We evaluate two unified-training variants		
1027	under the same RL objective (GRPO; Eq. 1) and		
1028	the same synthetic task generators/verifiers:		
1029	Curriculum unified RL (Deduction → Induction		
1030	→ Abduction). We train on Deduction first and		
1031	progressively switch to Induction and then Abduc-		
1032	tion. After switching to Abduction, we observed		
1033	that optimization became unstable, with frequent		
1034	gradient spikes and no further improvement on ei-		
1035	ther the synthetic diagnostics or downstream bench-		
1036	marks.		
1037	Mixed unified RL (joint sampling within each		
1038	mini-batch). We mix Deduction/Induction/Ab-		
1039	duction examples in every mini-batch. In this set-		
1040	ting, the reward curve rises slowly at the beginning		
1041	and then quickly flattens, yielding substantially		
1042	weaker performance than the specialist-then-merge		
1043	pipeline.		
1044	Discussion. Our working hypothesis is that		
1045	jointly optimizing the three meta-abilities induces		
1046	conflicting update directions, making unified train-		
1047	ing brittle and difficult to tune without extensive		
1048	task-balancing or curriculum sweeps. We therefore		
1049	adopt the separate-then-merge strategy, which is		
1050	empirically more stable and compute-efficient in		
1051	our setting.		
1052	E.2 SFT Baselines on the Same Synthetic		
1053	Tasks		
1054	To control for the effect of “seeing more synthetic		
1055	data”, we also evaluate supervised fine-tuning		
1056	(SFT) baselines on the same synthetic meta-ability		
1057	datasets, using the same inputs as the RL setting.		
1058	Answer-only SFT. We perform standard SFT		
1059	where the synthetic question is the input and only		
1060	the final answer is the target. This variant can over-		
1061	fit to synthetic patterns and may hurt generalization		
1062	on downstream tasks.		
1063	Trace+answer SFT. We augment the target with		
1064	reasoning traces (chain-of-thought) from correct		
1065	solutions, and fine-tune the model to generate both		
		the trace and the final answer. This improves down-	1066
		stream performance relative to answer-only SFT,	1067
		but still lags behind RL-aligned specialists and the	1068
		merged model, consistent with prior observations	1069
		that outcome-based RL more directly refines rea-	1070
		soning behaviors.	1071
		Takeaway. These baselines indicate that our	1072
		gains are not explained by “training on extra syn-	1073
		thetic data” alone: naive SFT can degrade gen-	1074
		eralization, and even trace-augmented SFT re-	1075
		mains weaker than the full RL-based meta-ability	1076
		pipeline.	1077
		F Merging and Domain-RL Ablations	1078
		This section reports ablations that justify parameter-	1079
		space merging and quantify its interaction with	1080
		downstream domain RL.	1081
		F.1 Early vs. Late Merging under Domain RL	1082
		We compare early merging (merge specialists first,	1083
		then apply domain-specific RL; Domain-RL-Meta)	1084
		versus late merging (apply domain RL to each spe-	1085
		cialist first, then merge the resulting checkpoints;	1086
		Merged-Post-Domain-RL). We additionally report	1087
		domain RL from single specialists (Deduction/In-	1088
		duction/Abduction init) to isolate how much merg-	1089
		ing helps beyond the strongest single initialization.	1090
		Interpretation. The strongest single specialist	1091
		initialization is Deduction-Init-Domain-RL, while	1092
		Induction/Abduction initializations are comparable.	1093
		Late merging still improves over the best single	1094
		specialist, but the gain is smaller than the early-	1095
		merge recipe. A plausible explanation is that do-	1096
		main RL drives specialists toward a more similar	1097
		local optimum (“homogenization”), reducing com-	1098
		plementarity and thus diminishing the benefit of	1099
		late-stage merging.	1100
		F.2 Connection to the Main Domain-RL	1101
		Comparison	1102
		Table 2 in the main paper compares Domain-RL-	1103
		Ins (domain RL from the instruction-tuned model)	1104
		versus Domain-RL-Meta (domain RL from the	1105
		meta-ability merged initialization) under the same	1106
		SimpleRL-Zoo-style setup. Together with Table 6,	1107
		these results show that (i) merging provides a	1108
		stronger initialization than any single specialist,	1109
		and (ii) performing the merge <i>before</i> domain RL is	1110
		most beneficial.	1111

Algorithm 3: Reverse Rule-Graph Search

Input: Rule-graph $G = (V, E)$; facts

$F \subseteq V$; goal $g \in V$;

chain_depth d ; **distractor_count** h ;

cycle_prob γ (values sampled as in Appendix B).

Output: TRUE/FALSE and (optional) explanation tree \mathcal{T} .

```
1 Caches: memo Cache, recursion stack
   Path, edge set  $\mathcal{E}$ .
2 Function Prove( $q$ ,  $depth$ ):
3   if  $q \in F$  then
4      $\lfloor$  return TRUE
5   if  $depth > d \vee q \in \text{Path}$  then
6      $\lfloor$  return FALSE
7   if Cache[ $q$ ]  $\neq$  UNKNOWN then
8      $\lfloor$  return Cache[ $q$ ]
9   Path  $\leftarrow$  Path  $\cup$  { $q$ }
10  foreach  $e : (P \rightarrow q) \in E$  do
    // candidate hyper-edges
11    if  $\forall p \in P : \text{Prove}(p, depth+1)$ 
    then
12       $\mathcal{E} \leftarrow \mathcal{E} \cup \{e\}$ 
13      Cache[ $q$ ]  $\leftarrow$  TRUE;
      Path  $\leftarrow$  Path  $\setminus$  { $q$ }; return
      TRUE
14  Cache[ $q$ ]  $\leftarrow$  FALSE;
    Path  $\leftarrow$  Path  $\setminus$  { $q$ }; return FALSE
15 Main:
16 initialise caches;
17  $success \leftarrow$  Prove( $g, 0$ );
18 if  $success$  then
19    $\lfloor$  prune  $\mathcal{E}$  to minimal tree  $\mathcal{T}$ 
20 else
21    $\lfloor$   $\mathcal{T} \leftarrow \emptyset$ 
22 return  $success, \mathcal{T}$ 
```

Model (7B)	Math500	AIME	AIME24	AMC23
Deduction-Init-Domain-RL	76.8	25.1	12.3	54.2
Induction-Init-Domain-RL	76.4	24.8	11.1	53.9
Abduction-Init-Domain-RL	77.5	24.9	11.9	54.0
Merged-Post-Domain-RL	78.0	25.8	11.9	54.1
Domain-RL-Meta (Early-Merge)	78.8	27.7	12.6	54.7

Table 6: Early-vs-late merging and domain RL initialization from single specialists versus from the merged model (7B). Early merging followed by domain RL yields the best performance.

INDUCTION - Level 1

Q: Given the following sequence

['1', '-1', '2', '0', '?']

What is the value at the question mark?

INDUCTION - Level 2

Q: Given the following sequence

['1', '5', '7', '11', '13', '17', '19', '23', '?']

What is the value at the question mark?

INDUCTION - Level 3

Q: Given the following sequence

['3', '1', '5', '2', '0', '4', '1', '-1', '3', '0', '?']

What is the value at the question mark?

Figure 5: Examples of masked-sequence completion with difficulty levels ranging from 1 to 3.

ABDUCTION - Level 1

Q: premises: ['((KF AND A) => A', '((A OR KF) => F', '(F => N', '((N OR N) => KF', '((N OR C) => C', '((NOT F) OR N) => F']

known_atoms: ['A', 'N', 'F', 'C', 'KF']

goals: ['HE', 'OH', 'KF']

For each goal, analyze its reachability

ABDUCTION - Level 2

Q: premises: ['(C => I', '((A OR (NOT I)) OR (L OR NO) => I', '((N OR C) => O', '(NOT I) => A', '((N OR N) OR (NOT A) => NO', '((N OR I) OR A) => I', '((NO AND (I OR (NOT O))) => EN', '(I => N']

known_atoms: ['N', 'NO', 'L', 'O', 'I', 'C', 'A']

goals: ['FI', 'KB', 'NO']

For each goal, analyze its reachability

ABDUCTION - Level 3

Q: premises: ['((NOT D) => J', '((I OR J) => M', '((D AND M) OR B) => M', '((C OR (NOT B)) AND ((NOT D) AND (NOT C))) => B', '((I OR J) => HJ', '(((NOT M) AND (NOT M)) OR (I AND M) => I', '(((NOT I) AND (NOT H)) AND (ME OR (NOT D))) => C', '(((NOT M) OR M) OR ((NOT M) AND J)) => BC', '((M AND (NOT J)) => J', '(((M AND M) AND (NOT J)) => ME', '(ME => C', '((NOT C) => O', '(M => M']

known_atoms: ['B', 'J', 'ME', 'M', 'D', 'I', 'C']

goals: ['ME', 'H', 'DJ']

For each goal, analyze its reachability

Figure 6: Examples of reverse rule-graph search with difficulty levels ranging from 1 to 3.