

# How and Why LLMs Generalize: A Fine-Grained Analysis of LLM Reasoning from Cognitive Behaviors to Low-Level Patterns

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) display strikingly different generalization behaviors: supervised fine-tuning (SFT) often narrows capability, whereas reinforcement-learning (RL) tuning tends to preserve it. The reasons behind this divergence remain unclear, as prior studies have largely relied on coarse accuracy metrics. We address this gap by introducing a novel benchmark that decomposes reasoning into atomic core skills such as calculation, fact retrieval, simulation, enumeration, and diagnostic, providing a concrete framework for addressing the fundamental question of what constitutes reasoning in LLMs. By isolating and measuring these core skills, the benchmark offers a more granular view of how specific cognitive abilities emerge, transfer, and sometimes collapse during post-training. Combined with analyses of low-level statistical patterns such as distributional divergence and parameter statistics, it enables a fine-grained study of how generalization evolves under SFT and RL across mathematical, scientific reasoning, and non-reasoning tasks. Our meta-probing framework tracks model behavior at different training stages and reveals that RL-tuned models maintain more stable behavioral profiles and resist collapse in reasoning skills, whereas SFT models exhibit sharper drift and overfit to surface patterns. This work provides new insights into the nature of reasoning in LLMs and points toward principles for designing training strategies that foster broad, robust generalization.

## 1 Introduction

LLMs fine-tuned with long Chain-of-Thought (CoT) reasoning, DeepSeek-R1 [17], OpenAI-o1 [39], Claude-Sonnet [49], achieve strong results on math and science benchmarks, yet their ability to *generalize* remains fragile. A notable pattern has emerged: models trained with supervised fine-tuning (SFT) often narrow their capability and overfit to surface patterns [48], while those tuned with

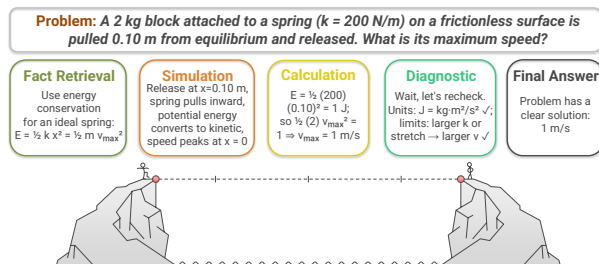


Figure 1: Decomposing reasoning into atomic cognitive skills. A spring-block example is solved step by step: fact retrieval, simulation, calculation, and diagnostic, checking each function as an isolated cognitive skill. When combined, these skills form a coherent reasoning trace that enables solving real-world problems, illustrating how complex reasoning emerges from the composition of simpler, specialized components.

reinforcement learning (RL) better preserve or even enhance generalization [22; 46]. The underlying reasons remain unclear, partly because prior studies rely on coarse metrics such as end-to-end accuracy or pass@k, which obscure the behavioral dynamics that drive reasoning success or failure [54; 38].

We argue that reasoning is not a monolithic property but rather an emergent composition of *atomic core skills*, calculation, simulation, fact retrieval, enumeration, and diagnostic self-checking, that interact with statistical patterns acquired during training [26; 36]. Accuracy on a final answer can mask weaknesses in these intermediate skills: a model may recall the correct formula yet incorrectly simulate a process, or arrive at the right solution through superficial pattern matching. Understanding how post-training reshapes these skills is crucial for explaining divergent generalization trends.

**Illustrative Example.** Figure 1 shows a classic spring block problem: a 2 kg block attached to a spring is pulled 0.10 m from equilibrium and released. Solving it requires sequential use of several atomic core skills. The solver first **retrieves** the

conservation of energy relation for a spring, then **simulates** the block’s motion as potential energy converts to kinetic, **calculates** the peak speed at the equilibrium point, and finally performs a **diagnostic** of units and boundary conditions. Each skill plays a distinct functional role; failure at any stage breaks the reasoning chain. While not needed in this particular example, a separate **enumeration** skill, systematically listing possible cases or configurations, is equally crucial for comprehensive reasoning, especially in combinatorial or case-based problem settings. This illustrates that what we call “reasoning” in practice is the coordinated composition of such simpler cognitive elements.

Existing benchmarks and analyses fall short of capturing this structure. Large mixed-domain corpora such as *Numina-Math* [30] and *Omni-Math* [16] obscure which skills drive success; controlled datasets like *GSM-Symbolic* [33] and *GSM-PLUS* [31] target narrow sub-skills; and conventional accuracy metrics hide shifts in intermediate behavior [13]. To fill this gap, we introduce a novel benchmark that explicitly decomposes reasoning into atomic core skills across mathematics, scientific reasoning, coding, and non-reasoning tasks, complemented by probes for low-level statistical patterns such as distributional divergence and shifts in word-frequency profiles. This design enables us to track how individual skills and statistical tendencies evolve through the course of post-training.

Our empirical results reveal several findings: (a) **RL preserves balanced cognitive skills.** RL-tuned models maintain rounded radar profiles across core skills: calculation, enumeration, simulation, fact retrieval, and diagnostic, indicating more stable and broad-spectrum reasoning. This balance persists across both in-domain (math) and out-of-domain (e.g., physics) settings. (b) **SFT induces overspecialization and drift.** SFT-tuned models display jagged radar shapes, with strong spikes in a narrow skill (often diagnostic or calculation) but dips below baseline in others, such as simulation and enumeration. This pattern signals over-fitting to surface heuristics and loss of transferable reasoning capacity. (c) **Training objective matters more than parameter magnitude.** The observed divergence arises despite SFT and RL modifying a comparable proportion of model parameters. The difference stems from their optimization objectives, which shape how cognitive skills are retained or distorted during post-training.

Our findings provide a behavioral account of

why post-training strategies diverge in their generalization effects and highlight the importance of strengthening atomic core skills as a foundation for robust and interpretable reasoning.

## 2 Related Works

**Post-training for Reasoning in LLMs.** Recent progress in large language models has underscored the importance of specialized post-training strategies, especially fine-tuning for reasoning [58; 48]. Chain-of-thought (CoT) prompting introduced by [54] encourages step-by-step explanation and significantly improves performance on symbolic and multi-step reasoning tasks [32; 27]. Recent extensions such as DeepSeek-R1 (Team, 2025a) combine CoT with reinforcement-learning-based optimization to further enhance reasoning and have achieved state-of-the-art results on math, logic, and competitive programming benchmarks [19; 17].

Post-training methods for reasoning typically fall into two categories: supervised fine-tuning (SFT) and reinforcement-learning-based tuning (RL) [58; 12]. SFT trains models to replicate explicit reasoning traces collected from annotated solutions [53; 54], whereas RL rewards models for accurate and logically coherent reasoning without requiring explicit intermediate supervision [12; 58]. This difference in optimization objective leads to distinct generalization patterns: SFT often narrows the behavioral diversity of models and overfits to surface heuristics [42], whereas RL tends to preserve or even enhance broader reasoning capacities but can introduce reward hacking and biases [40; 8]. Despite these insights, most prior analyses rely primarily on coarse outcome-based metrics such as final-task accuracy, leaving the process-level and representation-level dynamics of post-training effects underexplored.

**Generalization and cross-domain reasoning.** Generalization to tasks or domains outside the training distribution, remains a central challenge for LLMs. This challenge has been extensively studied in discriminative models, where systems often fail under covariate, diversity, or semantic shifts despite strong in-distribution performance [5; 7; 55; 4; 6]. While scaling laws highlight global performance trends across model size and data [25; 20], fine-tuning often induces qualitative shifts in reasoning robustness and error modes [56; 19]. Comparative studies suggest that SFT-dominated reasoning models frequently over-specialize, losing robustness to

new task formats or domains, whereas RL-based fine-tuning helps retain transferable skills [38]. For example, OpenAI’s o1 model excels in STEM reasoning but has raised concerns about versatility on other tasks [23]. Recent evaluations emphasize that reasoning-oriented fine-tuning can improve performance on target tasks but sometimes comes at the cost of cross-domain generalization [47; 11]. Our work complements these studies by moving beyond aggregate accuracy to analyze reasoning at a finer granularity. We decompose reasoning into core cognitive behaviors: calculation, simulation, enumeration, fact retrieval, and diagnostic, and reveal how SFT and RL tuning differentially shape these components, providing a clearer account of divergent generalization effects.

**Cognitive behaviors and representation-level shifts.** Accuracy alone can hide weaknesses in intermediate reasoning. We therefore decompose performance into five measurable behaviors, *calculation*, *simulation*, *fact retrieval*, *enumeration*, and *diagnostic checking*, making these skills explicit and comparable across domains and training stages. While BIG-bench, MATH, and related psychometric work touch on similar abilities [45; 19; 28], our framework operationalizes them with targeted prompts and metrics; see also domain-specific analyses such as *From Scores to Skills* in finance [26] and broader links to cognitive science [36].

Fine-tuning not only changes outward behavior but also the internal representations that support it. CoT can elicit more systematic reasoning [54], yet traces may be unfaithful to internal computations [51]. Sparse autoencoders and activation steering reveal interpretable subspaces tied to reasoning features [52; 50]. By jointly measuring behavioral skills and representational structure, we characterize how SFT and RL differentially shape both the internal and external facets of LLM reasoning.

### 3 Decompose “Reasoning” into Fine-grained Cognitive Behaviors

A central question of this study is how different post-training regimes, supervised fine-tuning (SFT) versus reinforcement-learning-based tuning (RL), reshape not only the overall accuracy of large language models (LLMs) but also the composition of their underlying reasoning skills. To answer this, we develop a controlled benchmark and a meta-analysis framework that reveals how individual cognitive behaviors evolve under these regimes.

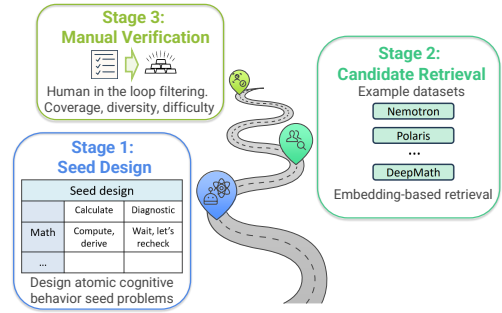


Figure 2: Overview of our three-stage benchmark construction pipeline. (1) **Seed Design** defines atomic seeds by domain and behavior; (2) **Candidate Retrieval** expands them via embedding search over public datasets; (3) **Manual Verification** filters for coverage, diversity, and difficulty, yielding a curated benchmark for fine-grained reasoning analysis.

Our approach proceeds in three stages (Figure 2). First, we identify representative model families and collect checkpoints at multiple stages of SFT and RL. Second, we construct a structured benchmark that spans four domains (Mathematics, Scientific Reasoning, Coding, and Non-Reasoning QA) and targets five core cognitive behaviors: *calculation*, *enumeration*, *simulation*, *fact retrieval*, and *diagnostic checking*. Third, we use behavior-focused probes to analyze how training dynamics affect the balance among these sub-skills and associated statistical patterns (e.g., distributional divergence). This controlled design moves beyond aggregate accuracy to trace how SFT and RL redistribute effort across fundamental behaviors, clarifying their distinct generalization profiles.

#### 3.1 Cognitive Grounding and Functional Completeness of Atomic Skills

We define “reasoning” at a functional level as a composition of elementary information-processing operations, rather than in neuroanatomical terms, following classic cognitive accounts of goal-directed problem solving. In the problem-space tradition, reasoning proceeds through a recurrent loop of (i) retrieving relevant knowledge, (ii) generating candidate states or hypotheses, (iii) applying operators that transform internal representations, and (iv) evaluating and monitoring intermediate results to guide control and termination [35; 29; 37].

This functional cycle provides a principled basis for our taxonomy. Fact retrieval instantiates access to declarative knowledge [3; 2]; enumeration corresponds to structured candidate gener-

Table 1: Example problem templates across four domains and five cognitive behaviors. Representative tasks illustrate the diversity of problem types used to evaluate model capabilities.

Category	Math Reasoning	Scientific Reasoning	Coding	Non-reasoning
<b>Calculation</b>	<p><b>Definition:</b> The execution of arithmetic, algebraic, or symbolic operations to derive quantitative or formal results from given inputs.</p> <p><b>Example:</b> Find the second largest prime factor of 49766.</p>	<p><b>Classical mechanics:</b> A proton enters a uniform magnetic field perpendicularly, compute the radius of its circular path.</p> <p><b>Magnetism:</b> In an RC circuit, compute capacitor’s charge at a given time after the circuit is closed.</p>	<p><b>Matrix primitivity:</b> Given a non-negative matrix <math>A \in \mathbb{R}_{\geq 0}^{n \times n}</math>, decide whether there exists an integer <math>k \geq 1</math> such that <math>A^k &gt; 0</math></p> <p><b>Unique array.</b> Unique <math>s</math>, construct <math>a, b \geq 0</math>, <math>a_i + b_i = s_i</math>, each unique after removing <math>\leq \lfloor n/2 \rfloor</math> entries.</p>	N/A
<b>Enumeration</b>	<p><b>Definition:</b> The systematic and exhaustive generation of all elements, cases, or options satisfying a set of explicit constraints.</p> <p><b>Example:</b> You have letters ‘q’: 4, ‘l’: 2, ‘b’: 2. Distribute them into 3 labeled boxes of capacities [4, 2, 2]. How many ways?</p>	<p><b>Spin Projection:</b> 8 non-interacting spin-<math>\frac{1}{2}</math> particles, find number of microstates with total spin projection <math>M_s = +1</math>.</p> <p><b>Atomic-Orbital Counting:</b> The number of microstates for a <math>d^3</math> electron configuration consistent with the Pauli principle.</p>	<p><b>Set Partitions:</b> All partitions of <math>\{1..n\}</math> in canonical order; consecutive partitions differ by moving one element between blocks.</p> <p><b>Primes in Ranges:</b> For <math>t</math> queries <math>[m, n]</math> (with <math>n - m \leq 10^5</math>), list all primes in each range.</p>	<p><b>Rewards program features enumeration:</b> Create a webpage section listing Eventsfy’s Sparks Awards program features.</p> <p><b>Article structure enumeration:</b> Create academic article on integrating ASHA workers into geriatric care and NCD clinic limitations.</p>
<b>Simulation</b>	<p><b>Definition:</b> The symbolic enactment of a process, system, or sequence of operations to predict or trace its behavior over time.</p> <p><b>Example:</b> Integers 1–120 on board. Each minute replace <math>n</math> by <math>d(n+3)</math>, divisor-count. Find total sum on board after 1 day.</p>	<p><b>Cooling Process:</b> An object cools from <math>90^\circ\text{C}</math> to ambient <math>20^\circ\text{C}</math> Newton’s law; compute <math>T</math> after 30 min with <math>k = 0.02 \text{ min}^{-1}</math>.</p> <p><b>Decay Chain:</b> Nuclide <math>A</math> decays to <math>B</math> (half-life 2 h), then <math>B</math> decays to stable <math>C</math> (half-life 8 h); find the time when <math>B</math> peaks.</p>	<p><b>Robot Path with One Extra Move.</b> Given a walk <math>s</math> over <math>\{W, A, S, D\}</math>, insert one move to minimize the bounding box area.</p> <p><b>Furlo/Rublo Coin Game.</b> Replace <math>x</math> by <math>y \in [\lfloor x/4 \rfloor, \lfloor x/2 \rfloor]</math>; winner given by xor of Grundy <math>g(x)</math>, which periodic in <math>\log x</math>.</p>	<p><b>Weather Simulator:</b> A program forecasts weather and suggests clothing based on temperature and rain conditions.</p> <p><b>Gaming-related simulation:</b> Creating a vodcast script analyzing how sound displays mood in the mobile game Pixel Gun 3D.</p>
<b>Fact Retrieval</b>	<p><b>Definition:</b> (1) Pose the problem naming the required theorem. (2) Pose it without the theorem and check if the model applies it correctly.</p> <p><b>Example:</b> Surface distance on radius <math>r</math> sphere between <math>(p_1, q_1)</math>, and <math>(p_2, q_2)</math>.</p>	<p><b>Cylindrical Capacitor:</b> Recall the capacitance formula for a long cylindrical capacitor; compute <math>C</math> for <math>a = 1 \text{ mm}</math>, <math>b = 5 \text{ mm}</math>, <math>L = 0.5 \text{ m}</math>.</p> <p><b>Bragg’s Law:</b> Recall <math>2d \sin \theta = n\lambda</math>; find the smallest <math>\theta</math> for <math>d = 0.2 \text{ nm}</math>, <math>\lambda = 0.15 \text{ nm}</math>.</p>	<p><b>Mixing Rule in the Minimum-Smoke:</b> When combining adjacent mixtures of colors <math>a</math> and <math>b</math>, what are the smoke produced and the resulting color?</p> <p><b>GCD-Constrained BST:</b> What gcd condition must each edge between adjacent vertices satisfy?</p>	<p><b>Cross-Cultural Information Retrieval:</b> Identify sources for an article on the Alcatraz occupation and Indigenous Peoples’ Day.</p> <p><b>Procedural Knowledge:</b> Prerequisites for AP Chemistry topics (stoichiometry, electrochemistry).</p>
<b>Diagnostic</b>	<p><b>Definition:</b> Diagnostic in LLMs often accompanied by explicit self-check (e.g., "wait", "let’s recheck").</p> <p><b>Example:</b> <math>f_n = \frac{nx_1 + n^2x_2}{nx}</math> is perturbed as <math>\frac{nx_1 - n^2x_2}{nx}</math>.</p>	<p><b>Pursuit of a fox by a hound:</b> Centripetal acceleration perturbed with false small-deviation expansion.</p> <p><b>Falling hinged rods:</b> Energy-conservation problem perturbed to suggest small-angle oscillation approach.</p>	<p><b>Uniqueness vs. Duplicates:</b> counting the same solution twice because it appears in different positions.</p> <p><b>Optimal vs. Sub-optimal:</b> Picking the shortest edge next in a path problem that needs full DP.</p>	<p><b>Counter-Factual:</b> "Is tree planting an eco-social activity?" with perturbed response claiming it is not.</p> <p><b>Self-contradictory:</b> Document identifies "Villainous" as Cartoon Network series, then labels incorrectly.</p>

ation and hypothesis search in problem spaces [35; 29]; calculation and simulation partition operator application into symbolic transformation versus model-based rollout of consequences, aligning with mental-model theories of reasoning [24]; and diagnostic checking operationalizes evaluation and metacognitive control, consistent with classical accounts of metacognition, conflict monitoring, and error detection [34; 14; 10; 21].

Under this characterization, the set of five behaviors is complete at our level of analysis: every step in a reasoning trace performs one of these functions, while higher-level reasoning modes such as analogy or abduction arise from structured compositions of these primitives rather than requiring additional atomic categories [29; 24; 34]. For a more detailed theoretical grounding, see Appendix H.

### 3.2 Benchmark Construction

**Principles.** Each item targets a single atomic behavior (calculation, enumeration, simulation, fact retrieval, diagnostic) and is curated for *skill isolation*, coverage, diversity, and calibrated difficulty. Implementation details (templates, keyword lists, perturbations, rubrics) appear in Appendix F.

**Pipeline.** Figure 2 summarizes the process.

1. **Stage 1: Seed Design** creates atomic seeds for each behavior–domain pair, avoiding multi-skill conflation.
2. **Stage 2: Candidate Retrieval** performs embedding-based nearest-neighbor search over large open repositories (e.g., Nemotron [9], Polaris [1], DeepMath [18]) to harvest surface variants that instantiate the same behavior in diverse contexts.
3. **Stage 3: Manual Verification** applies human-in-the-loop filtering for skill isolation, coverage, diversity, and difficulty. The result is a high-fidelity benchmark that decouples reasoning behaviors from surface task formats.

**Behavior Assembly.** *Calculation, Enumeration, Simulation:* template seeds plus keyword queries retrieve items with explicit single-behavior traces from open sources, followed by deduplication and removal of multi-skill bleed. *Fact Retrieval:* select questions whose solution hinges on a named person, theorem, or definition. We use two modes: *Guided* (theorem named in the prompt) and *Unguided* (theorem not named; verify its use in the solution). *Diagnostic:* sample questions and apply minimal perturbations to a valid trace (logical

contradiction, dropped condition, counterfactual). The prompt includes the question and the perturbed trace; the model must detect and correct the error (see example in Appendix F.2).

**Difficulty.** Items are bucketed as *easy*, *medium*, or *hard* using: (i) lightweight heuristics (operand scale, step count, constraint breadth, perturbation subtlety); (ii) reference-model success rates; and (iii) human adjustment prioritizing skill purity.

**Prompting (Appendix).** Prompts consist of a standardized instruction header, behavior-specific payload (e.g., theorem tag or perturbed trace), output format, and fixed decoding settings, which are summarized in Appendix F.1.

**Human Evaluation.** We check coverage across (behavior, domain, difficulty) cells, diversity of forms and entities, and the validity of difficulty labels; disagreements are adjudicated.

**Metrics.** The primary metric is **accuracy** under standardized decoding. For Calculation, Enumeration, and Simulation, we use an exact match with unit normalization. For Fact Retrieval, *Guided* requires a correct answer and theorem-consistent steps; *Unguided* requires a correct answer and verified theorem use. For Diagnostic, we require correct answer and explicit identification of self-check.

### 3.3 Behavior–Domain Grid

To support interpretable analysis, we organize the benchmark as a two-dimensional grid crossing five cognitive behaviors with four domains (Table 1). For each behavior–domain cell, the table provides a concise behavior definition and, in most cases, two representative examples. Due to space constraints, the *Math* column anchors the behavior definitions and includes only one representative example per behavior; additional mathematical examples are provided in Appendix J, Table G.

The five behaviors capture complementary facets of reasoning: (1) Calculation: quantitative manipulation of explicit formulas or equations; (2) Enumeration: systematic generation of combinatorial possibilities; (3) Simulation: mental or symbolic rollout of dynamics; (4) Fact Retrieval: access to stored knowledge such as definitions or constants; and (5) Diagnostic: identifying and correcting faulty reasoning or self-contradictions.

For example, as shown in Table 1, calculation tasks range from prime-factor queries in mathematics to computing capacitor charge in physics; enumeration spans classical combinatorics to set partitioning in code; simulation includes Newto-

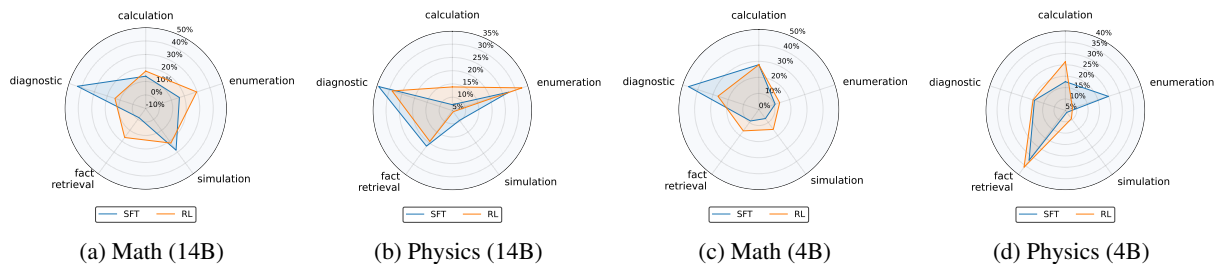


Figure 3: Radar plots of five cognitive behaviors (calculation, enumeration, simulation, fact retrieval, diagnostic checking). Each panel contrasts SFT (blue) vs. RL (orange) for (a) Math-14B, (b) Physics-14B, (c) Math-4B, and (d) Physics-4B. RL is more balanced; SFT often concentrates on a few skills.

nian cooling and other dynamical rollouts; fact-retrieval probes recall of theorems or scientific laws; and diagnostic checking tests recognition of flawed proofs or perturbed reasoning templates. This grid yields comprehensive yet disentangled coverage of reasoning behaviors, enabling us to track how SFT and RL affect each component skill within and across domains.

## 4 Experiments

**Experimental setup.** We evaluate Qwen3-14B-Base and its SFT- and RL-tuned variants, along with smaller Qwen3-4-Base and Qwen3-1.7B-Base counterparts. To analyze training dynamics, we evaluate intermediate checkpoints throughout training. Following models used in [22], the RL model uses the Verl framework [44] with a GRPO setup [43] for RL fine-tuning, optimizing for answer-correctness rewards. The SFT model uses LLaMA-Factory [57] to train on teacher-generated chain-of-thought traces through reject sampling. Evaluation is conducted on the above cognitive skills benchmark using accuracy. Detailed subcategory distributions for the proposed benchmark are provided in Appendix G. More details about training datasets, baseline models, and hyperparameters can be found in Appendix A.

### 4.1 Cognitive-skill Profiles

**RL yields balanced skills and SFT over-specializes.** A striking pattern in Figure 3 is the irregular, jagged shape of the SFT curves compared to the rounded RL curves. The SFT-tuned models consistently exhibit asymmetric profiles, often showing pronounced spikes in a single skill (most commonly diagnostic) while dipping well below baseline in others, such as enumeration and simulation. This shape reflects a tendency toward over-specialization or over-fitting, where the model

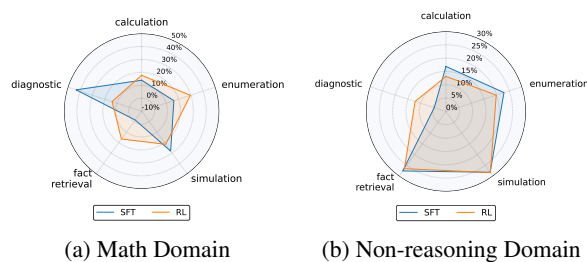
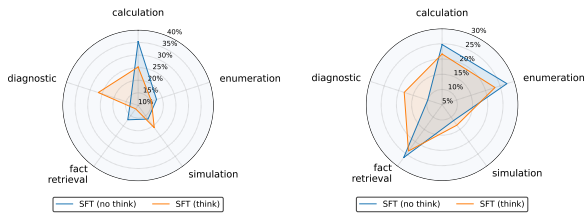


Figure 4: Radar plots of post-training effects on math (a) and non-reasoning (b) using math-trained models. Both transfer weakly to non-reasoning; RL preserves a more balanced skill profile.

exploits superficial patterns in a limited subset of reasoning behaviors rather than maintaining a balanced skill set. For example, in Math-4B, SFT shows a sharp spike on diagnostic, while enumeration and simulation fall below baseline; in Math-14B, it drops below baseline for calculation, simulation, and especially fact retrieval. In contrast, the RL-tuned profiles remain rounded and smoother in performance across all five cognitive behaviors, indicating that reinforcement-learning post-training promotes more even retention of diverse reasoning skills. This shape contrast is informative even when raw accuracy gaps between SFT and RL are modest, emphasizing that the training regime, not merely model scale or base capability, drives these differences in cognitive skill balance. To complement these visual patterns, Appendix B provides quantitative evidence using behavior-wise accuracies and smoothness metrics, confirming that RL yields significantly more balanced cognitive-skill profiles than SFT.

**Long CoT rebalances SFT toward systematic reasoning.** Prior work shows that SFT with teachers in thinking mode, producing long CoT data, generally outperforms SFT with short CoT [22]. Figure 5 compares SFT skill profiles under think vs. no-think modes for Math and Physics. The



(a) Math (think vs no think) (b) Physics (think vs no think)

Figure 5: Shape comparison of SFT profiles with thinking (orange) vs no thinking (blue) on (a) Math and (b) Physics. Compared to the think model, the no-think model is even more spiky and calculation-dominated. All models are derived from the Qwen3-14B base; SFT and RL share the same base.

no-think models spike on calculation but lag on simulation, diagnostic checking, and enumeration, whereas long-CoT models exhibit more rounded profiles, redistributing capacity toward complementary processes essential for multi-step reasoning. This effect is especially pronounced in the out-of-domain Physics setting, where long CoT reduces calculation dominance and strengthens other components. Overall, long-CoT training reshapes not only absolute performance but also the geometry of skill profiles, yielding more systematic, transferable reasoning behaviors.

## 4.2 Cross-Domain Generalization

**RL preserves a more balanced skill mix and degrades less than SFT under cross-domain transfer.** Prior work [22] shows that both RL and SFT degrade when transferring from math to non-reasoning domains, but RL typically degrades less and maintains a more balanced skill profile. As shown in Figure 4, RL-tuned models exhibit smoother, more uniform radar shapes in the non-reasoning setting, whereas SFT displays irregular profiles with sharp spikes and drops across skills. This contrast suggests that RL post-training acts as a regularizer that preserves general reasoning competencies, while SFT tends to overfit to math-specific patterns at the expense of cross-domain transfer. Although absolute gains remain limited outside math, RL consistently retains broader cognitive behaviors under domain shift.

## 4.3 Linking Behaviors to Sparse Latent Features

**Shallow layers are near-Gaussian; deeper layers become heavier-tailed across base, SFT, and RL.** We analyze joint skewness–kurtosis of SAE [15]

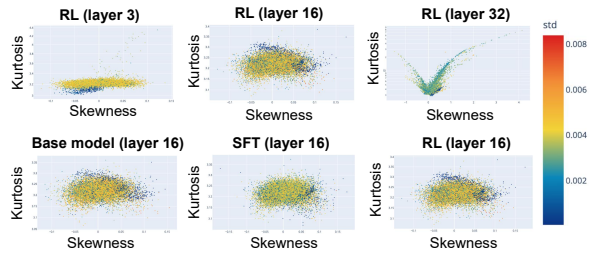
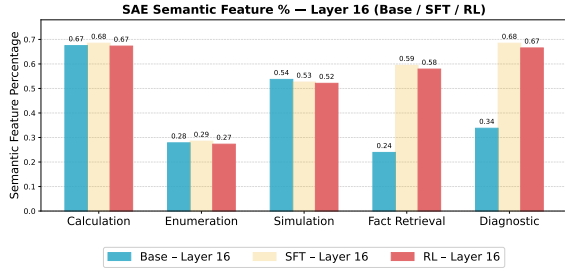


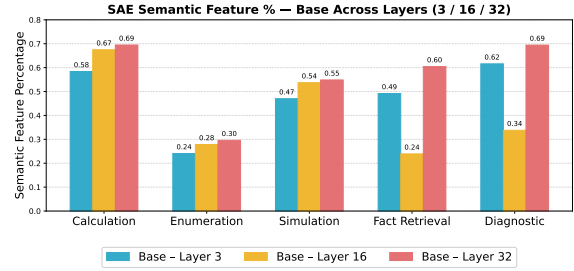
Figure 6: Skewness–kurtosis of logit weights across layers. RL shifts distributions from near-Gaussian (layer 3) to heavy-tailed, structured patterns (layer 32). At layer 16, base is compact and symmetric, SFT slightly broadens, while RL introduces more dispersed, structured variations.

latent weight distributions to characterize how training regimes reshape higher-order structure. In shallow layers, the base model is approximately symmetric and near-Gaussian, indicating relatively homogeneous activations, while SFT introduces mild increases in skewness and kurtosis consistent with modest sharpening of feature selectivity. With depth, all three models transition toward heavier-tailed, non-Gaussian regimes, reflecting growing feature specialization. RL slightly amplifies this trend, exhibiting broader skewness–kurtosis dispersion, but does not induce a qualitatively different latent regime. Overall, RL increases heterogeneity and sparsity while preserving the depth-dependent progression present in the base and SFT models.

**Interpretable subspaces tied to reasoning-related features.** Figure 7 uses sparse autoencoders (SAEs) [15] to analyze Qwen3-4B and its SFT/RL variants, revealing behavior-aligned latent subspaces. We focus on Layer 16, the model midpoint, where representations are most compositional and provide a stable basis for comparison. At this depth, both post-training strategies selectively reallocate capacity toward knowledge access (*Fact Retrieval*) and self-verification (*Diagnostic*), while *Calculation* and *Simulation* remain largely stable and *Enumeration* stays sparse. Across depth in the base model, *Calculation* and *Simulation* increase gradually, whereas *Fact Retrieval* and *Diagnostic* follow a U-shaped trajectory, prominent in early and late layers. Together, these trends suggest a functional pipeline, retrieval and checking in early and deep layers, transformation in mid layers, with post-training amplifying retrieval and diagnostic subspaces without disrupting core computational circuits. These patterns also generalize to physics-domain reasoning (Appendix E).



(a) Different Strategy



(b) Different Layers

Figure 7: SAE-based cognitive-feature composition across strategies and layers. We train sparse autoencoders (SAEs) on hidden states for different models (Base, SFT, RL) and layers, then compute the fraction of SAE features aligned with five cognitive behaviors: Calculation, Enumeration, Simulation, Fact Retrieval, and Diagnostic. Panel (a) compares post-training strategies at layer 16; panel (b) shows layer-wise trends for the Base model.

Metric	SFT (no-think)	SFT (think)	RL
<i>Global parameter shift vs. Qwen3-14B base</i>			
Changed parameters (#)	$1.44 \times 10^9$	$1.44 \times 10^9$	$1.44 \times 10^9$
Change percentage (%)	97.83	97.80	97.81
Total change magnitude	29 462.2	29 033.8	29 134.1
Average change magnitude	66.51	65.54	65.77
Maximum change magnitude	1 544.5	1 540.9	1 556.3
<i>Distribution of change across major components</i>			
Embed tokens change (%)	97.04	97.04	97.03
Layer-wise change (%)	97.87	97.84	97.94
Norm change (%)	55.23	55.23	55.23
LM-head change (%)	97.89	97.88	97.94

Table 2: Parameter-space shifts of Qwen3-14B under SFT (no-think/think) and RL relative to the base. All methods update a similar share of parameters with comparable magnitudes, and changes are distributed similarly across embeddings, layers, norms, and the LM head, indicating that behavioral differences stem from training objectives rather than update scope.

#### 4.4 Low-Level Statistics & Surface Patterns

##### Parameter-space shifts under post-training.

To quantify post-training effects, we measure the fraction of parameters updated and their shift magnitudes for Qwen3-14B variants. As Table 2 shows, roughly 98% of parameters change in every variant, updates are pervasive, not sparse. RL’s total change magnitude is slightly below SFT (no-think) and comparable to SFT (think), indicating reward-based tuning does not necessarily induce larger global perturbations.

**Component-level differences.** Decomposing changes by module (Table 2) shows that embeddings and transformer layers account for most shifts across regimes, while normalization layers change much less (55.23%). RL also exhibits a slightly larger maximum shift in the LM head than SFT (think), consistent with its objective of optimizing the output distribution via reward

signals. This supports the view that RL primarily nudges the decision boundary by reshaping the output layer, whereas SFT distributes adjustments more uniformly across the network.

## 5 Conclusion

We introduce a controlled benchmark that decomposes LLM reasoning into five core cognitive behaviors: calculation, enumeration, simulation, fact retrieval, and diagnostic checking, across mathematics, scientific reasoning, coding, and non-reasoning QA. Our three-stage pipeline enables reproducible and interpretable analysis of how post-training reshapes these skills. Our results highlight three main findings. First, RL-tuned models preserve a more balanced distribution of cognitive skills, supporting broader generalization. Second, SFT often induces over-specialization, strengthening narrow abilities such as diagnostic while weakening others like simulation. Third, these differences arise primarily from the training objective rather than model scale. A discussion of the complementary roles of SFT and RL in practice is provided in Appendix C.

These results suggest that future reasoning-oriented LLMs should emphasize not only aggregate accuracy but also balanced skill development. Promising directions include behavior-aware objectives to prevent skill collapse, curriculum strategies to reinforce both domain-specific and transferable skills, and representation-level methods such as sparse autoencoders and activation steering to monitor and guide reasoning-related subspaces. By revealing how post-training reshapes cognitive behaviors, our benchmark points toward training strategies that foster more robust, interpretable, and transferable reasoning capabilities.

## 544 Limitations

545 Our analysis leverages SAEs to link behavior-level  
546 signals to internal features; while informative, cur-  
547 rent SAEs techniques are compute-heavy and sen-  
548 sitive to design choices (e.g., hookpoints, sparsity  
549 targets). We therefore treat them as high-fidelity  
550 but non-definitive probes. A next step is to develop  
551 lightweight methods that still connect low-level  
552 statistics to hidden representations.

## 553 References

- 554 [1] Chenxin An, Zhihui Xie, Xiaonan Li, Lei Li, Jun  
555 Zhang, Shansan Gong, Ming Zhong, Jingjing Xu,  
556 Xipeng Qiu, Mingxuan Wang, and Lingpeng Kong.  
557 2025. *Polaris: A post-training recipe for scaling rein-*  
558 *forcement learning on advanced reasoning models.*
- 559 [2] John R Anderson. 2014. *Rules of the mind.* Psychol-  
560 ogy Press.
- 561 [3] John R Anderson and Christian J Lebiere. 2014. *The*  
562 *atomic components of thought.* Psychology Press.
- 563 [4] Haoyue Bai, Yifei Ming, Julian Katz-Samuels, and  
564 Yixuan Li. Hypo: Hyperspherical out-of-distribution  
565 generalization. In *The Twelfth International Confer-*  
566 *ence on Learning Representations.*
- 567 [5] Haoyue Bai, Rui Sun, Lanqing Hong, Fengwei Zhou,  
568 Nanyang Ye, Han-Jia Ye, S-H Gary Chan, and Zhen-  
569 guo Li. 2021. Decaug: Out-of-distribution general-  
570 ization via decomposed feature representation and  
571 semantic augmentation. In *Proceedings of the AAAI*  
572 *Conference on Artificial Intelligence*, volume 35,  
573 pages 6705–6713.
- 574 [6] Haoyue Bai, Jifan Zhang, and Robert Nowak. 2024.  
575 Aha: Human-assisted out-of-distribution generaliza-  
576 tion and detection. *Advances in Neural Information*  
577 *Processing Systems*, 37:33863–33890.
- 578 [7] Haoyue Bai, Fengwei Zhou, Lanqing Hong, Nanyang  
579 Ye, S-H Gary Chan, and Zhenguo Li. 2021. Nas-ood:  
580 Neural architecture search for out-of-distribution gen-  
581 eralization. In *Proceedings of the IEEE/CVF inter-*  
582 *national conference on computer vision*, pages 8320–  
583 8329.
- 584 [8] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, and  
585 1 others. 2022. Constitutional ai: Harmlessness from  
586 ai feedback. *arXiv preprint arXiv:2212.08073*.
- 587 [9] Akhiad Bercovich, Itay Levy, Izik Golan, Moham-  
588 mad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil,  
589 Zach Moshe, Tomer Ronen, Najeeb Nabwani, and 1  
590 others. 2025. Llama-nemotron: Efficient reasoning  
591 models. *arXiv preprint arXiv:2505.00949*.
- 592 [10] Matthew M Botvinick, Todd S Braver, Deanna M  
593 Barch, Cameron S Carter, and Jonathan D Cohen.  
594 2001. Conflict monitoring and cognitive control.  
595 *Psychological review*, 108(3):624.

- [11] Zhoujun Cheng, Shibo Hao, Tianyang Liu, Fan  
Zhou, Yutao Xie, Feng Yao, Yuexin Bian, Yong-  
hao Zhuang, Nilabjo Dey, Yuheng Zha, and 1 oth-  
ers. 2025. Revisiting reinforcement learning for llm  
reasoning from a cross-domain perspective. *arXiv*  
*preprint arXiv:2506.14965*. 596  
597  
598  
599  
600  
601
- [12] Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Sheng-  
bang Tong, Saining Xie, Dale Schuurmans, Quoc V  
Le, Sergey Levine, and Yi Ma. 2025. Sft mem-  
orizes, rl generalizes: A comparative study of  
foundation model post-training. *arXiv preprint*  
*arXiv:2501.17161*. 602  
603  
604  
605  
606  
607
- [13] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
Nakano, and 1 others. 2021. Training verifiers  
to solve math word problems. *arXiv preprint*  
*arXiv:2110.14168*. 608  
609  
610  
611  
612  
613
- [14] Jonathon D Crystal and Allison L Foote.  
2011. Evaluating information-seeking approaches  
to metacognition. *Current zoology*, 57(4):531–542. 614  
615  
616
- [15] Hoagy Cunningham, Aidan Ewart, Logan Riggs,  
Robert Huben, and Lee Sharkey. 2023. Sparse au-  
toencoders find highly interpretable features in lan-  
guage models. *arXiv preprint arXiv:2309.08600*. 617  
618  
619  
620
- [16] Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai,  
Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma,  
Liang Chen, Runxin Xu, and 1 others. 2024. Omni-  
math: A universal olympiad level mathematic bench-  
mark for large language models. *arXiv preprint*  
*arXiv:2410.07985*. 621  
622  
623  
624  
625  
626
- [17] 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*. 627  
628  
629  
630  
631  
632
- [18] Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu,  
Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu,  
Zhenwen Liang, Wenxuan Wang, and 1 others. 2025.  
Deepmath-103k: A large-scale, challenging, decon-  
taminated, and verifiable mathematical dataset for ad-  
vancing reasoning. *arXiv preprint arXiv:2504.11456*. 633  
634  
635  
636  
637  
638
- [19] Dan Hendrycks, Collin Burns, Steven Basart, and  
1 others. 2021. Measuring massive multitask lan-  
guage understanding. In *International Conference*  
*on Learning Representations (ICLR)*. 639  
640  
641  
642
- [20] Jordan Hoffmann, Sebastian Borgeaud, Arthur  
Mensch, and 1 others. 2022. Training compute-  
optimal large language models. *arXiv preprint*  
*arXiv:2203.15556*. 643  
644  
645  
646
- [21] Clay B Holroyd and Michael GH Coles. 2002. The  
neural basis of human error processing: reinforce-  
ment learning, dopamine, and the error-related nega-  
tivity. *Psychological review*, 109(4):679. 647  
648  
649  
650

651	[22] Maggie Huan, Yuetai Li, Tuney Zheng, Xiaoyu Xu,	Barret Zoph, Jason Wei, and 1 others. 2023. The flan	707
652	Seungone Kim, Minxin Du, Radha Poovendran, Gra-	collection: Designing data and methods for effective	708
653	ham Neubig, and Xiang Yue. 2025. Does math reason-	instruction tuning. In <i>International Conference on</i>	709
654	ing improve general llm capabilities? understanding	<i>Machine Learning</i> , pages 22631–22648. PMLR.	710
655	transferability of llm reasoning. <i>arXiv preprint</i>		
656	<i>arXiv:2507.00432</i> .		
657	[23] Aaron Jaech, Adam Kalai, Adam Lerer, Adam	[33] Iman Mirzadeh, Keivan Alizadeh, Hooman	711
658	Richardson, Ahmed El-Kishky, Aiden Low, Alec Hel-	Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad	712
659	lyar, Aleksander Madry, Alex Beutel, Alex Carney,	Farajtabar. 2024. Gsm-symbolic: Understanding the	713
660	and 1 others. 2024. Openai o1 system card. <i>arXiv</i>	limitations of mathematical reasoning in large lan-	714
661	<i>preprint arXiv:2412.16720</i> .	guage models. <i>arXiv preprint arXiv:2410.05229</i> .	715
662	[24] Philip N Johnson-Laird. 2010. Mental models	[34] Thomas O Nelson. 1990. Metamemory: A theoret-	716
663	and human reasoning. <i>Proceedings of the National</i>	ical framework and new findings. In <i>Psychology of</i>	717
664	<i>Academy of Sciences</i> , 107(43):18243–18250.	<i>learning and motivation</i> , volume 26, pages 125–173.	718
		Elsevier.	719
665	[25] Jared Kaplan, Sam McCandlish, Tom Henighan,	[35] Allen Newell, Herbert Alexander Simon, and 1 oth-	720
666	Tom B Brown, and 1 others. 2020. Scaling	ers. 1972. <i>Human problem solving</i> , volume 104.	721
667	laws for neural language models. <i>arXiv preprint</i>	Prentice-hall Englewood Cliffs, NJ.	722
668	<i>arXiv:2001.08361</i> .		
669	[26] Ziyang Kuang, Feiyu Zhu, Maowei Jiang, Yanzhao	[36] Qian Niu, Junyu Liu, Ziqian Bi, Pohsun Feng,	723
670	Lai, Zelin Wang, Zhitong Wang, Meikang Qiu, Jiajia	Benji Peng, Keyu Chen, Ming Li, Lawrence KQ	724
671	Huang, Min Peng, Qianqian Xie, and 1 others. 2025.	Yan, Yichao Zhang, Caitlyn Heqi Yin, and 1 others.	725
672	From scores to skills: A cognitive diagnosis frame-	2024. Large language models and cognitive science:	726
673	work for evaluating financial large language models.	A comprehensive review of similarities, differences,	727
674	<i>arXiv preprint arXiv:2508.13491</i> .	and challenges. <i>arXiv preprint arXiv:2409.02387</i> .	728
675	[27] Nathan Lambert, Jacob Morrison, Valentina Py-	[37] Stellan Ohlsson. 2012. The problems with problem	729
676	atkin, Shengyi Huang, Hamish Ivison, Faeze Brahma-	solving: Reflections on the rise, current status, and	730
677	man, Lester James V Miranda, Alisa Liu, Nouha	possible future of a cognitive research paradigm. <i>The</i>	731
678	Dziri, Shane Lyu, and 1 others. 2024. Tulu 3: Push-	<i>Journal of problem solving</i> , 5(1):7.	732
679	ing frontiers in open language model post-training.	[38] OpenAI. 2023. Gpt-4 technical report. <i>arXiv</i>	733
680	<i>arXiv preprint arXiv:2411.15124</i> .	<i>preprint arXiv:2303.08774</i> .	734
681	[28] Andrew K Lampinen, Ishita Dasgupta,	[39] OpenAI. 2024. <a href="#">Learning to reason with llms</a> .	735
682	Stephanie CY Chan, Kory Matthewson,		
683	Michael Henry Tessler, Antonia Creswell, James L	[40] Long Ouyang, Jeff Wu, Xu Jiang, and 1 oth-	736
684	McClelland, Jane X Wang, and Felix Hill. 2022. Can	ers. 2022. Training language models to follow in-	737
685	language models learn from explanations in context?	structions with human feedback. <i>arXiv preprint</i>	738
686	<i>arXiv preprint arXiv:2204.02329</i> .	<i>arXiv:2203.02155</i> .	739
687	[29] Pat Langley, Lorenzo Magnani, Christian Schunn,	[41] Gonalo Paulo, Alex Mallen, Caden Juang, and	740
688	and Paul Thagard. 2005. An extended theory of hu-	Nora Belrose. 2024. Automatically interpreting mil-	741
689	man problem solving. In <i>Proceedings of the annual</i>	lions of features in large language models. <i>arXiv</i>	742
690	<i>meeting of the cognitive science society</i> , volume 27.	<i>preprint arXiv:2410.13928</i> .	743
691	[30] Jia Li, Edward Beeching, Lewis Tunstall, Ben	[42] Ethan Perez, Douwe Kiela, and Kyunghyun Cho.	744
692	Lipkin, Roman Soletskyi, Shengyi Costa Huang,	2021. True few-shot learning with language mod-	745
693	Kashif Rasul, Longhui Yu, Albert Jiang, Ziju	els. In <i>Advances in Neural Information Processing</i>	746
694	Shen, Zihan Qin, Bin Dong, Li Zhou, Yann	<i>Systems (NeurIPS)</i> .	747
695	Fleureau, Guillaume Lample, and Stanislas Polu.	[43] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	748
696	2024. Numinamath. [ <a href="https://huggingface.co/AI-M0/NuminaMath-1.5">https://huggingface.co/</a>	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	749
697	<a href="https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf">AI-M0/NuminaMath-1.5</a> ]( <a href="https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf">https://github.com/</a>	Zhang, YK Li, Yang Wu, and 1 others. 2024.	750
698	<a href="https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf">project-numina/aimo-progress-prize/blob/</a>	Deepseekmath: Pushing the limits of mathematical	751
699	<a href="https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf">main/report/numina_dataset.pdf</a> ).	reasoning in open language models. <i>arXiv preprint</i>	752
700	[31] Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng	<i>arXiv:2402.03300</i> .	753
701	Kong, and Wei Bi. 2024. Gsm-plus: A comprehen-	[44] Guangming Sheng, Chi Zhang, Zilingfeng Ye,	754
702	sive benchmark for evaluating the robustness of llms	Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,	755
703	as mathematical problem solvers. <i>arXiv preprint</i>	Haibin Lin, and Chuan Wu. 2025. Hybridflow: A	756
704	<i>arXiv:2402.19255</i> .	flexible and efficient rlhf framework. In <i>Proceedings</i>	757
705	[32] Shayne Longpre, Le Hou, Tu Vu, Albert Webson,	<i>of the Twentieth European Conference on Computer</i>	758
706	Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le,	<i>Systems</i> , pages 1279–1297.	759

760	[45] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, and 1 others. 2022. Beyond the imitation game benchmark for measuring and extrapolating the capabilities of language models. <i>arXiv preprint arXiv:2206.04615</i> .	<i>International Conference on Learning Representations (ICLR)</i> .	815
761			816
762			
763			
764			
765	[46] Yiyu Sun, Yuhan Cao, Pohao Huang, Haoyue Bai, Hannaneh Hajishirzi, Nouha Dziri, and Dawn Song. 2025. Delta-code: How does rl unlock and transfer new programming algorithms in llms? <i>arXiv preprint arXiv:2509.21016</i> .	[57] Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. <i>arXiv preprint arXiv:2403.13372</i> .	817
766			818
767			819
768			820
769			821
770	[47] Yiyu Sun, Shawn Hu, Georgia Zhou, Ken Zheng, Hannaneh Hajishirzi, Nouha Dziri, and Dawn Song. 2025. Omega: Can llms reason outside the box in math? evaluating exploratory, compositional, and transformative generalization. <i>arXiv preprint arXiv:2506.18880</i> .	[58] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. <i>arXiv preprint arXiv:1909.08593</i> .	822
771			823
772			824
773			825
774			826
775			
776	[48] Yiyu Sun, Georgia Zhou, Haoyue Bai, Hao Wang, Dacheng Li, Nouha Dziri, and Dawn Song. 2025. Climbing the ladder of reasoning: What llms can and still can't-solve after sft? <i>arXiv preprint arXiv:2504.11741</i> .		
777			
778			
779			
780			
781	[49] Anthropic Team. <a href="#">The claude 3 model family: Opus, sonnet, haiku</a> .		
782			
783	[50] Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini, and Monte MacDiarmid. 2023. Activation addition: Steering language models without optimization. <i>arXiv e-prints</i> , pages arXiv–2308.		
784			
785			
786			
787			
788	[51] Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2023. Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting. <i>Advances in Neural Information Processing Systems</i> , 36:74952–74965.		
789			
790			
791			
792			
793	[52] Various Contributors. 2024. Hookedsaetransformer: Sparse autoencoder interpretability toolkit. GitHub repository.		
794			
795			
796	[53] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. <i>arXiv preprint arXiv:2203.11171</i> .		
797			
798			
799			
800			
801	[54] Jason Wei, Xuezhi Wang, Dale Schuurmans, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .		
802			
803			
804			
805	[55] Nanyang Ye, Kaican Li, Haoyue Bai, Runpeng Yu, Lanqing Hong, Fengwei Zhou, Zhenguo Li, and Jun Zhu. 2022. Ood-bench: Quantifying and understanding two dimensions of out-of-distribution generalization. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 7947–7958.		
806			
807			
808			
809			
810			
811			
812	[56] Chiyuan Zhang, Samy Bengio, Moritz Hardt, and 1 others. 2021. Understanding deep learning requires rethinking generalization under distribution shift. In		
813			
814			

## A Additional Experimental Details

### A.1 Setup for Base, SFT, and RL Model Training and Evaluation

**Models.** We adopt the recipe of [22]: we (i) use their released Qwen3-14B SFT and RL checkpoints, (ii) replicate the procedure at 4B to obtain Qwen3-4B SFT/RL, and (iii) train Qwen3-1.7B SFT/RL while saving intermediate checkpoints for dynamics visualization.

Using Verl with GRPO and answer-correctness rewards for RL, we set: lr =  $1 \times 10^{-6}$ ; global batch = 512; clip = 0.22–0.28; context = 16k tokens; 16 rollouts/prompt; mini-batch = 128; KL/entropy coeffs = 0; train 140 steps (and use the corresponding checkpoint). For SFT, we fine-tune with LLaMA-Factory on teacher CoT traces from Qwen3-32B-Instruct via reject sampling; lr =  $5 \times 10^{-5}$ ; global batch = 512; 1.5 epochs (mirroring the RL horizon). The dataset used is 47K curated math problems, with CoT generated by Qwen3-32B-Instruct. The baselines are Qwen3-14B-Base and Qwen3-14B-Instruct (evaluated in *think/no-think* modes).

**Generation settings.** Unless noted otherwise, we use nucleus sampling with temperature 0.6 and top- $p = 0.95$ . For reasoning tasks, we cap generation at 32k new tokens and apply stop sequences aligned with our standardized output headers. All decoding settings are held fixed across models for fair comparison.

**Prompt template.** We use a single minimal prompt for all tasks:

```
{problem}
```

```
Please reason step by step,  
and put your final answer within \boxed{...}.
```

Here, {problem} contains the task statement and, when applicable, any behavior-specific context (e.g., theorem tags for *Guided* fact-retrieval or a lightly perturbed trace for diagnostic). Decoding settings are identical across models; evaluation reads only the final boxed answer.

**Core metrics.** Primary metric is **accuracy** under fixed decoding. We also compute behavior profiles (per-behavior accuracy vectors), normalized skill indices relative to the base model, and

distributional drift of intermediate features when applicable.

**Evaluation harness.** We run batched inference with vLLM and a deterministic evaluation harness that records per-item seeds, prompts, responses, and scores. All artifacts (prompts, outputs, logs) are versioned and exported for audit.

### A.2 Setup for SAE Training and Interpretation

**Activation collection.** We extract activations from specified transformer layers (e.g., blocks 3, 16, 32) at hook\_mlp\_out with context length 2048 over a balanced mixture of behavior-labeled prompts. We follow a stratified sampler to equalize behaviors and difficulties.

**SAE objective & launch.** We train per-layer sparse autoencoders (SAEs) on transformer MLP outputs by minimizing  $\mathcal{L} = \|h - Dz\|_2^2 + \lambda \|z\|_1$ , where  $h$  is the hidden activation at blocks.  $< L >$ . hook\_mlp\_out,  $z$  the sparse code, and  $D$  the decoder. We use the open-source sparsify repo. Unless noted: -ctx\_len 2048 and -batch\_size 5.

Example command (single layer):

```
python -m sparsify <BASE_MODEL_NAME> <DATASET_OR_REPO> \  
--layers 16 --batch_size 5 --ctx_len 2048 \  
--data_preprocessing_num_proc 5 --k 250 \  
--run_name SAE-layer16-example
```

Replace <BASE\_MODEL\_NAME> (e.g., Qwen3-4B-Base) and <DATASET\_OR\_REPO> with model tag and dataset identifier or local path.

**SAE interpretation with Delphi.** For automated explanation of sparse autoencoder (SAE) and transcoder features, we use the open-source Delphi library [41]. Delphi programmatically *generates* and *scores* natural-language explanations for features, enabling large-scale interpretability; models can run locally or via OpenRouter. In our pipeline, for each learned feature we feed top-activating examples, obtain  $k$  candidate explanations, score them with DELPHI’s built-in scorer, and keep the top-ranked explanation subject to a quality threshold.

Example generation command:

```
python -m delphi <BASE_MODEL> <SAE_CHECKPOINT_DIR> \  
--hookpoints layers.16 \  
--scorers detection \  
--n_tokens 10_000_000 \  
--max_latents 10000 \  
--dataset_repo <DATASET_REPO> \  
--dataset_split 'train[:100%]' \  
--filter_bos \  
--name sae-layer16-explain
```

925	Replace placeholders with model tag, the directory containing the trained SAE for the target layer, and the dataset source. The flag <code>-hookpoints</code>	983
926		984
927	layers. <code>16</code> selects the MLP-out features of layer 16 (other layers analogous); <code>-filter_bos</code> removes BOS tokens; <code>-scorers</code> detection invokes the default explanation scorer.	985
928		
929		
930		
931		
932	<b>Scoring SAE features for cognitive behaviors.</b>	986
933	After obtaining per-feature natural-language explanations with DELPHI, we assign scores for the five behaviors (calculation, enumeration, simulation, fact retrieval, diagnostic) using a lightweight LLM grader. For each behavior, we supply a behavior-specific rubric (name, description, keywords, indicators) and batch a list of feature explanations; the model returns a JSON record per feature with a soft score, confidence, and a binary label. We decode <i>greedily</i> (no sampling) to ensure schema validity and reproducibility. A feature is counted as “related” to a behavior if score $\geq 0.5$ ; we also report confidence-weighted variants ( $\tilde{s} = \text{score} \times \text{confidence}$ ). Percentages are computed per layer and per model.	987
934		988
935		989
936		990
937		991
938		992
939		993
940		994
941		995
942		996
943		997
944		998
945		999
946		1000
947		1001
948		1002
949	<b>Behavior-parameterized prompt:</b>	
950	You are analyzing feature explanations to determine if they relate to <BEHAVIOR>.	
951		
952	Category: <NAME>	
953	Description: <DESCRIPTION>	
954	Keywords: <KEYWORDS>	
955	Indicators: <INDICATORS>	
956		
957	Analyze the following <N> features:	
958	Feature 1 (ID: <feature_id_1>): <explanation_1>	
959		
960	...	
961		
962	For each feature, determine:	
963	1) The extent to which it relates to <BEHAVIOR> (score 0.0-1.0).	
964		
965	2) Your confidence in this judgment (0.0-1.0).	
966	3) A binary classification (0 or 1) for whether it relates to <BEHAVIOR>.	
967		
968		
969	Respond with a JSON array (one element per feature) using EXACTLY:	
970		
971	[	
972	{"feature_id": "<feature_id_1>", "score	
973	": <0.0-1.0>,	
974	"confidence": <0.0-1.0>, "classification": <0 or	
975	1>},	
976	...	
977	]	
978	Only return the JSON array.	
979	<i>Example instantiation (diagnostic).</i> Replace <BEHAVIOR> with mathematical diagnostics and populate <NAME>, <DESCRIPTION>, <KEYWORDS>, <INDICATORS> from our diag-	
980		
981		
982		
	nostic rubric (Appendix F). The same template is used for the other four behaviors by swapping in the corresponding rubrics.	983
		984
		985
	<b>Interpretation.</b> We align SAE latents to the five behaviors with a two-stage procedure. (i) <i>Text-based alignment:</i> using DELPHI, we generate per-feature explanations and grade them with a behavior-parameterized LLM prompt (shown above), which returns a soft score $s \in [0, 1]$ , a confidence $\hat{c} \in [0, 1]$ , and a binary label; a feature is considered related if $s \geq 0.5$ , and we also report confidence-weighted scores $s \cdot \hat{c}$ . (ii) <i>Signal-based validation:</i> we fit linear probes and compute mutual information between latent activations and behavior labels; any discrepancies trigger manual auditing of top-activating tokens/spans. We further summarize higher-order structure via skewness-kurtosis of decoder weights and visualize feature selectivity with activation-triggered exemplars, reporting layer-wise percentages for each model.	986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000
		1001
		1002
	<b>A.3 Additional Details</b>	1003
	<b>Reproducibility and runtime.</b> We fix seeds for our experiments, and the evaluation harness, enable deterministic kernels where supported, and pin vLLM to a specific engine build. Experiments run on H100 (80 GB) GPUs; inference uses tensor parallelism with static max sequence lengths, and batch sizes are tuned to avoid OOM while keeping GPU utilization $\geq 80\%$ . We cache tokenizer outputs and KV states when permitted by the harness to ensure identical re-runs.	1004
		1005
		1006
		1007
		1008
		1009
		1010
		1011
		1012
		1013
	<b>AI Assistance in Research and Writing.</b> This work uses AI tools for sentence-level proofreading.	1014
		1015
	<b>B Quantifying Skill Balance Beyond Radar Plots</b>	1016
		1017
	<b>Quantifying skill-balance beyond radar visualizations.</b> In the main text, we argue that reinforcement learning (RL) tends to produce a more balanced cognitive-skill profile than supervised fine-tuning (SFT), as illustrated by radar plots. To provide quantitative support for this claim, we report both raw behavior-wise accuracies (Table 3) and three complementary smoothness metrics that directly capture how uniformly skills are distributed across cognitive dimensions: (i) the coefficient of variation (CV), measuring dispersion across skills; (ii) adjacent-difference smoothness (AdjDiff), measuring local jaggedness in the skill profile; and	1018
		1019
		1020
		1021
		1022
		1023
		1024
		1025
		1026
		1027
		1028
		1029
		1030

Setting	Strategy	Calculation	Enumeration	Simulation	Fact Retrieval	Diagnostic
Math (14B)	SFT	13.90	16.22	28.19	1.54	43.24
Math (14B)	RL	17.70	29.67	21.69	16.59	14.00
Physics (14B)	SFT	7.20	27.37	9.60	21.61	34.21
Physics (14B)	RL	13.89	32.60	5.56	19.44	28.89
Math (4B)	SFT	27.90	10.70	6.98	8.84	46.05
Math (4B)	RL	27.82	13.63	15.33	16.47	26.76
Physics (4B)	SFT	17.60	24.88	6.09	32.15	19.29
Physics (4B)	RL	26.39	8.27	9.52	35.91	19.90

Table 3: Behavior-wise accuracy improvements across cognitive skills. Accuracy is reported relative to the base model for five reasoning-related behaviors, Calculation, Enumeration, Simulation, Fact Retrieval, and Diagnostic, under SFT and RL across math and physics domains and two model scales (14B and 4B). These results provide the raw behavioral profiles that underlie the smoothness metrics reported in Table 4.

Setting	Strategy	CV ↓	AdjDiff ↓	Circularity ↑
Math (14B)	SFT	0.6851	22.396	0.0356
Math (14B)	RL	0.2741	6.268	0.4719
Physics (14B)	SFT	0.5152	17.912	0.2105
Physics (14B)	RL	0.4900	16.816	0.1706
Math (4B)	SFT	0.7453	15.628	0.1516
Math (4B)	RL	0.3014	5.676	0.4899
Physics (4B)	SFT	0.4308	13.336	0.1894
Physics (4B)	RL	0.5203	13.652	0.2303
<b>Average</b>	<b>SFT</b>	<b>0.5941</b>	<b>17.3180</b>	<b>0.1468</b>
<b>Average</b>	<b>RL</b>	<b>0.3965</b>	<b>10.6030</b>	<b>0.3407</b>

Table 4: Quantitative smoothness metrics for cognitive-skill profiles under SFT and RL. We report three complementary measures of skill balance: coefficient of variation (CV, lower is better), adjacent-difference smoothness (AdjDiff, lower is better), and circularity (higher is better). Across all settings, RL consistently yields smoother and more uniformly distributed skill profiles than SFT.

(iii) circularity, defined as the ratio between the minimum and maximum radius in the radar plot, capturing overall profile roundness.

Across all evaluated settings, RL consistently exhibits lower CV and lower AdjDiff than SFT, indicating substantially reduced variance and smoother transitions across cognitive skills. At the same time, RL achieves higher circularity, reflecting a more uniform distribution of behavioral competencies. Averaged across all models, RL reduces CV from 0.59 to 0.40 and AdjDiff from 17.3 to 10.6, while increasing circularity from 0.15 to 0.34, corresponding to roughly a 33% reduction in overall variance, a 39% reduction in local unevenness, and a 2.3× increase in profile uniformity. All numbers are averaged over five independent runs.

Importantly, the accuracies reported in Table 3 are measured relative to the base model, rather than

as absolute accuracies, ensuring a fair comparison of how different post-training objectives reshape the distribution of skills. Together, these quantitative results align closely with the radar visualizations and confirm that RL yields a smoother and more evenly balanced cognitive-skill profile.

## C Interpreting the Roles of SFT and RL

**Complementary roles of SFT and RL in post-training.** Our findings should not be interpreted as suggesting that supervised fine-tuning is unnecessary or that reinforcement learning alone is sufficient for post-training. In our experimental design, SFT and RL are applied independently to the same base model in order to enable a controlled comparison of how each objective reshapes the distribution of cognitive behaviors. This setup isolates their inductive effects on skill organization, but does not evaluate broader dimensions such as alignment, safety, or instruction-following, where SFT plays a central role in practice.

More broadly, the two objectives exhibit complementary strengths. In our experiments, RL primarily rebalances and smooths existing competencies, reducing over-specialization and promoting a more uniform distribution of skills across behaviors. By contrast, SFT often produces sharper peaks and valleys, and by reconstructing output distributions from supervised examples, can introduce or amplify capabilities that are not prominent in the base model. These differences suggest distinct but synergistic roles: RL is particularly effective at mitigating unevenness and stabilizing general reasoning performance, while SFT remains essential for alignment-oriented behaviors and for extending the usable knowledge of the base model.

Accordingly, our conclusions should be read as clarifying when and how each post-training objec-

1086 tive is most effective, rather than advocating an  
1087 RL-only pipeline.

## 1088 D Training Dynamic

1089 Under 1.7B RL, all five skills improve over training,  
1090 with fast early gains that settle into a stable plateau.  
1091 In math, simulation is consistently strongest, diag-  
1092 nostic rises steadily, and calculation shows clear  
1093 improvement; enumeration and fact retrieval make  
1094 smaller but persistent gains (see Fig. 8). In physics,  
1095 calculation leads throughout, fact retrieval strength-  
1096 ens over time, simulation improves to a moder-  
1097 ate level, and enumeration and diagnostic advance  
1098 more gradually. The shaded bands in Fig. 8 remain  
1099 tight, indicating stable progress without regressions  
1100 and balanced multi-skill benefits from RL.

## 1101 E Additional Results on Sparse Latent 1102 Features

1103 Figure 9 extends our SAE-based cognitive-feature  
1104 analysis to the physics domain. At the midpoint  
1105 layer (Layer 16), we observe a pattern closely  
1106 aligned with the math domain: both SFT and RL  
1107 selectively reweight semantic capacity toward *Fact*  
1108 *Retrieval* and *Diagnostic* features, while leaving  
1109 *Calculation* and *Simulation* largely stable and *Enu-*  
1110 *meration* consistently sparse. This suggests that  
1111 post-training enhances knowledge access and self-  
1112 verification behaviors in physics reasoning with-  
1113 out disrupting core computational and forward-  
1114 modeling circuits.

1115 Across depth in the base model, physics tasks ex-  
1116 hibit the same functional organization observed in  
1117 math. *Calculation* and *Simulation* strengthen with  
1118 depth, whereas *Fact Retrieval* and *Diagnostic* fol-  
1119 low a U-shaped trajectory, prominent in early and  
1120 late layers, with reduced emphasis in the mid stack.  
1121 Together, these results indicate that the pipeline  
1122 structure identified in the main text—retrieval and  
1123 checking in early and deep layers, transformation  
1124 and transport in mid layers, generalizes beyond  
1125 mathematics to scientific reasoning more broadly.

## 1126 F Cognitive Behaviors: Definitions, 1127 Keywords, and Indicators

1128 Definition of the five cognitive behaviors: These be-  
1129 haviors represent fundamental modes of reasoning  
1130 and problem solving—calculation, diagnosis, fact  
1131 retrieval, simulation, and enumeration—each cap-  
1132 turing a distinct way in which a system or learner  
1133 processes information, manipulates symbols, and

monitors or generates solutions across mathemati- 1134  
cal, scientific, and general domains. 1135

- **Calculation:** The execution of arithmetic, 1136  
algebraic, or symbolic operations to derive 1137  
quantitative or formal results from given in- 1138  
puts. Calculation emphasizes step-by-step ma- 1139  
nipulation under well-defined mathematical or 1140  
logical rules, and may involve simplification, 1141  
substitution, or transformation of expressions 1142  
to reach a final result. 1143

**Keywords:** compute, calculate, derive, sim- 1144  
plify, evaluate, solve, result, step-by-step, for- 1145  
mula, expression. 1146

**Indicators:** explicit arithmetic/algebraic op- 1147  
erations, transformations of equations, numer- 1148  
ical outputs. 1149

- **Enumeration:** The systematic and exhaustive 1150  
generation of all elements, cases, or options 1151  
satisfying a set of explicit constraints. Enu- 1152  
meration emphasizes coverage and complete- 1153  
ness, often under conditions where solution 1154  
space must be fully explored or cataloged for 1155  
correctness, diversity, or comprehensiveness. 1156

**Keywords:** list, all cases, options, possibili- 1157  
ties, exhaustive, generate, under constraints, 1158  
coverage. 1159

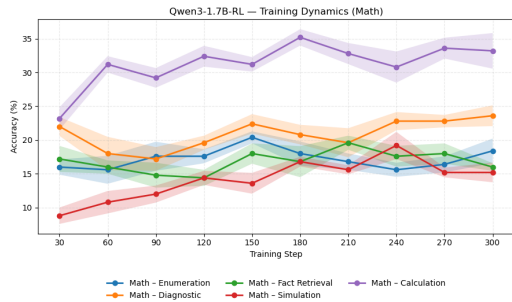
**Indicators:** systematic listing, covering 1160  
search space, combinatorial completeness. 1161

- **Simulation:** The mental or symbolic enact- 1162  
ment of a process, system, or sequence of 1163  
operations to predict or trace its behavior over 1164  
time. Simulation requires stepwise reason- 1165  
ing consistent with formal rules (e.g., alge- 1166  
braic manipulation, causal transitions, or state- 1167  
machine updates) and supports forecasting, 1168  
testing, or illustrating dynamic evolution. 1169

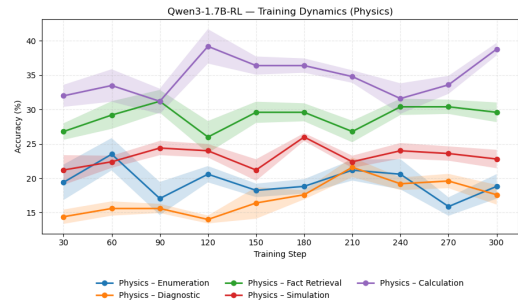
**Keywords:** simulate, step, next state, transi- 1170  
tion, apply rule, run through, dynamic, causal, 1171  
evolve, predict. 1172

**Indicators:** walking through process rules, 1173  
state transitions, mental enactment of system 1174  
behavior. 1175

- **Fact retrieval:** The recall or extraction of 1176  
canonical, domain-grounded knowledge units 1177  
(e.g., theorems, definitions, constants, statutes, 1178  
or factual attributes). Fact retrieval empha- 1179  
sizes accuracy, grounding, and direct mapping 1180

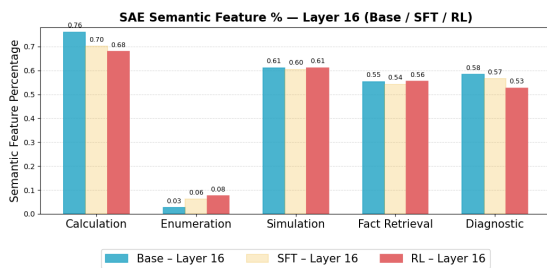


(a) Math dynamic

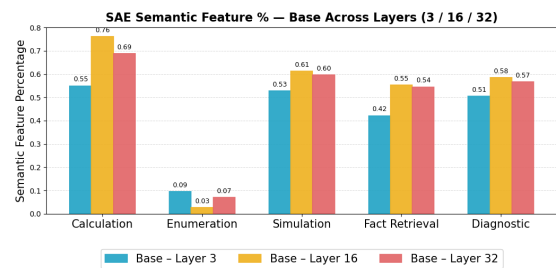


(b) Physics dynamic

Figure 8: Training dynamics (Qwen3-1.7B-RL). Accuracy vs. step for math (a) and physics (b), with shaded variability. Calculation dominates; enumeration stays weakest/most volatile. Fact-retrieval, diagnostic, and simulation rise steadily (simulation relatively stronger in physics), indicating balanced RL gains rather than specialization.



(a) Different Strategy



(b) Different Layers

Figure 9: SAE-based cognitive-feature composition in the physics domain across strategies and layers. We train sparse autoencoders (SAEs) on hidden states of Qwen3-4B, Qwen3-4B-SFT, and Qwen3-4B-RL for physics-domain tasks, and compute the fraction of SAE features aligned with five reasoning-related behaviors: Calculation, Enumeration, Simulation, Fact Retrieval, and Diagnostic. Panel (a) compares post-training strategies at Layer 16; panel (b) shows layer-wise trends (Layers 3/16/32) for the Base model. Overall, the physics domain exhibits patterns consistent with the math domain: post-training selectively amplifies retrieval and diagnostic subspaces, while depth organizes a pipeline from transformation-focused mid layers to retrieval and verification in early and late layers.

1181 to authoritative sources, rather than inference  
1182 or elaboration.

1183 **Keywords:** recall, definition, theorem, law,  
1184 constant, fact, property, known, canonical,  
1185 from memory.

1186 **Indicators:** citing established knowledge, di-  
1187 rect recall of named entities or rules without  
1188 derivation.

- 1189 • **Diagnostic:** The act of identifying, localizing,  
1190 and explaining the source of an error, inconsis-  
1191 tency, or divergence in a reasoning process,  
1192 system execution, or narrative account. In the  
1193 context of LLMs, diagnostic reasoning is of-  
1194 ten accompanied by explicit self-monitoring  
1195 or self-check language (e.g., “wait,” “let’s  
1196 recheck,” “that step seems wrong”), which  
1197 signals recognition of potential failure modes,  
1198 causal attribution, and corrective reasoning.

**Keywords:** error, mistake, wrong, recheck, 1199  
bug, contradiction, mismatch, correction, fail- 1200  
ure mode. 1201

**Indicators:** self-check language (“wait,” 1202  
“let’s recheck”), identifying divergence, ex- 1203  
plaining why an approach fails. 1204

## F.1 Template and prompt for the dataset 1205 construction 1206

**Candidate retrieval.** We retrieve candi- 1207  
dates via dense *embedding* search using 1208  
text-embedding-3-small. Each item is in- 1209  
dexed by the concatenation of its *question* and 1210  
*reasoning trace*. For every behavior–domain seed, 1211  
we compose a query from the seed summary, 1212  
behavior keywords, and a few template sentences, 1213  
then select the top-*k* nearest items by embedding 1214  
similarity (the model’s native similarity). We 1215  
de-duplicate with an embedding-similarity thresh- 1216

1217	old $\tau$ , retaining the highest-scoring item among	1267
1218	any pairs above $\tau$ .	1268
1219	<b>Retrieval query template:</b>	1269
1220	{SEED_SUMMARY}	1270
1221	Behavior: {BEHAVIOR}   Domain: {DOMAIN}	1271
1222	Keywords: {k1}, {k2}, {k3}	1272
1223	Template sentences: {TEMPLATE_SENTENCES}	1273
1224	(Find items semantically similar to this query	1274
1225	using embedding similarity in	1275
1226	text-embedding-3-small.)	1276
1227		1277
1228	<b>F.2 Diagnostic Example</b>	1278
1229	For the example in figure 1.5, starting from	1279
1230	$\sqrt{x+4} = x - 2$ , a first pass squares both sides	1280
1231	to obtain $x + 4 = (x - 2)^2$ , which simplifies	1281
1232	to $x^2 - 5x = 0$ and yields candidate solutions	
1233	$x \in \{0, 5\}$ . A diagnostic self-check then enforces	
1234	the original domain $\sqrt{x+4} \geq 0$ , giving $x - 2 \geq 0$	
1235	and thus $x \geq 2$ , which rules out $x = 0$ . Sub-	
1236	stituting into the original equation confirms the	
1237	extraneous nature of $x = 0$ since $\sqrt{4} = 2 \neq -2$ .	
1238	Therefore the only valid solution is $x = 5$ . This se-	
1239	quence—derivation, reflection, constraint checking,	
1240	and correction—illustrates diagnostic behavior by	
1241	identifying an error introduced by a non-invertible	
1242	operation (squaring) and revising the solution ac-	
1243	cordingly.	
1244		
1245	<b>F.3 Benchmark Scale and Extensibility.</b>	
1246	The benchmark spans $5 \times 4$ (behavior $\times$ domain)	
1247	cells, with each cell containing 50 to 200 carefully	
1248	curated, high quality examples vetted for skill iso-	
1249	lation, correctness, and calibrated difficulty. Seed	
1250	questions are templatable and can be faithfully ex-	
1251	panded, for example by renaming variables, rescal-	
1252	ing parameters, swapping units or entities, or con-	
1253	trolled paraphrasing, thereby increasing coverage	
1254	without altering the targeted skill.	
1255		
1256	<b>G Datasets Statistics</b>	
1257	<b>Dataset statistics.</b> Figures 10, 11, 12, and 13	
1258	present the distributions of subcategories across	
1259	our newly created mathematics, scientific reason-	
1260	ing (physics), coding, and non-reasoning datasets,	
1261	each further grouped by key behavioral facets:	
1262	Calculation, Enumeration, Fact Retrieval, Simu-	
1263	lation, and Diagnostic. The figures showcase a bal-	
1264	anced and comprehensive coverage of both domain-	
1265	specific and cross-domain skills, from core mathe-	
1266	matical topics such as prime factorization, combi-	
	inatorics, and binomial theorems, to scientific rea-	
	soning tasks in physics such as mechanics, thermo-	
	dynamics, and optics, to algorithmic and compu-	
	tational challenges including graph/tree problems,	1267
	dynamic programming, and simulation-based tasks,	1268
	and finally to non-reasoning tasks such as writing	1269
	and summarization, error detection, knowledge re-	1270
	trieval, and scenario-based simulations. This de-	1271
	liberate and systematic construction ensures that	1272
	the dataset spans diverse problem types and cogni-	1273
	tive demands, enabling robust evaluation of both	1274
	reasoning-intensive and non-reasoning capabilities	1275
	of large language models. The broad yet well-	1276
	curated distributions highlight the quality, diversity,	1277
	and representativeness of our benchmark, estab-	1278
	lishing it as a reliable resource for comprehensive	1279
	assessment of models across reasoning, problem-	1280
	solving, and real-world task performance.	1281
	<b>H Cognitive Grounding and Functional</b>	1282
	<b>Completeness of the Atomic Skill</b>	1283
	<b>Taxonomy</b>	1284
	We conceptualize “reasoning” at a functional	1285
	level as the coordinated operation of elementary	1286
	information-processing mechanisms, rather than in	1287
	neuroanatomical or implementation-specific terms.	1288
	This perspective follows a long tradition in cogni-	1289
	tive science that characterizes problem solving as	1290
	goal-directed search in a problem space, governed	1291
	by representations, operators, and control processes	1292
	[35; 29; 37].	1293
	In the classic problem-space framework, reason-	1294
	ing unfolds through a recurring control loop: (i) re-	1295
	trieving relevant knowledge from long-term mem-	1296
	ory, (ii) generating candidate states, hypotheses, or	1297
	partial solutions, (iii) applying operators that trans-	1298
	form internal representations, and (iv) evaluating	1299
	intermediate results to guide further search, revisi-	1300
	on, or termination [35; 29]. This loop provides an	1301
	abstract and operational description of how humans	1302
	and artificial systems organize complex reasoning	1303
	across domains.	1304
	Our five atomic skills map directly onto this func-	1305
	tional decomposition. Fact retrieval instantiates	1306
	access to declarative and semantic knowledge, con-	1307
	sistent with cognitive architectures that distinguish	1308
	retrieval of stored facts from procedural transforma-	1309
	tion [3; 2]. Enumeration corresponds to structured	1310
	candidate generation and hypothesis search, a core	1311
	component of problem-space exploration in both	1312
	symbolic and connectionist models [35; 29]. Cal-	1313
	culation and simulation jointly operationalize oper-	1314
	ator application, but capture two distinct modes: (a)	1315
	static symbolic or arithmetic transformation over	1316

well-defined representations (calculation), and (b) model-based rollout over time, counterfactuals, or hypothetical scenarios (simulation). This distinction aligns with mental-model theories of reasoning, which posit that reasoners construct internal models of situations and draw inferences by simulating their consequences [24]. Finally, diagnostic checking implements evaluation and metacognitive control, processes that monitor uncertainty, detect conflict or error, and trigger revision, additional search, or termination, consistent with classical accounts of metacognition, conflict monitoring, and error detection [34; 14; 10; 21].

Under this standard functional characterization, the proposed set of atomic skills is complete at our level of analysis. Any step in a reasoning trace can be classified as one of four fundamental operations: retrieving information, generating alternatives, transforming representations via calculation or simulation, or evaluating and controlling the process through diagnostic checking. Higher-level reasoning modes such as analogy, abduction, or causal inference do not require additional primitives; rather, they emerge from structured compositions of these same operations. For example, analogy combines retrieval of prior cases with simulation of relational mappings and diagnostic evaluation of fit, while abduction integrates enumeration of hypotheses with simulation of explanatory consequences and diagnostic selection [29; 24; 34].

This functional completeness is not only theoretically motivated but also empirically supported by our analysis. Across domains and training regimes, we find that observed reasoning behaviors can be consistently decomposed into these five skills, and that performance shifts under SFT and RL manifest as redistributions among them rather than as the emergence of qualitatively new categories. In this sense, our benchmark operationalizes a long-standing cognitive-science view of reasoning, translating it into a measurable and model-agnostic framework for analyzing generalization, robustness, and skill transfer in large language models.

## I Detailed Description of the Meta-Evaluation Benchmark

Table 1 provides a comprehensive overview of meta-evaluation benchmark, which systematically categorizes evaluation tasks across 20 distinct problem types (4 domains  $\times$  5 cognitive behaviors). This structured framework enables granular assess-

ment of LLM capabilities while maintaining consistency in evaluation methodology.

### I.1 Domain Categories

The benchmark spans four complementary domains that represent core areas of AI capability assessment:

**Math Reasoning** encompasses problems requiring mathematical logic, symbolic manipulation, and quantitative reasoning. Tasks in this domain range from elementary arithmetic operations to advanced mathematical proof construction, testing models’ ability to handle formal symbolic systems and rigorous logical deduction.

**Scientific Reasoning** evaluates understanding of physical principles, scientific methodologies, and domain-specific knowledge from fields such as physics, chemistry, and biology. These tasks assess whether models can apply scientific laws, perform calculations with physical units, and reason about real-world phenomena governed by natural principles.

**Coding** tests computational thinking, algorithm design, and programming proficiency. Problems in this category require models to understand data structures, optimize algorithms, debug code, and translate problem specifications into executable solutions across various programming paradigms.

**Non-reasoning** captures tasks that depend primarily on factual knowledge retrieval, procedural execution, or creative generation rather than complex logical inference. This category serves as a control to distinguish reasoning capabilities from memorization or pattern matching.

### I.2 Cognitive Behaviors

The five cognitive behaviors represent fundamental modes of thinking that cut across domains:

**Calculation** involves executing well-defined computational procedures to derive quantitative or symbolic results. These tasks test precision, attention to detail, and mastery of algorithmic processes. Examples include computing prime factorizations, determining molecular radii from physical constants, implementing matrix operations, and calculating environmental metrics from model equations. The complexity varies from single-step arithmetic to multi-stage computations requiring careful tracking of intermediate results.

**Enumeration** requires systematic generation of all elements satisfying specified constraints. This

1416	behavior tests completeness, organizational thinking, and the ability to explore solution spaces exhaustively without omission or duplication. Tasks include combinatorial problems (such as distributing items into labeled boxes), generating quantum mechanical states (atomic orbital configurations), partitioning data structures (set partitions with ordering constraints), and cataloging instances (program features or article structures). Success demands both algorithmic rigor and verification that all cases have been considered.	
1417		
1418		
1419		
1420		
1421		
1422		
1423		
1424		
1425		
1426		
1427		
1428	<b>Simulation</b> involves mental or symbolic enactment of dynamic processes over time or through state transitions. These tasks evaluate sequential reasoning, state tracking, and the ability to project forward through cause-and-effect chains. Examples include iterating mathematical functions on a board, modeling physical processes (cooling or radioactive decay), tracing algorithmic execution (robot pathfinding or coin replacement dynamics), and predicting outcomes (weather forecasting or game simulation). The challenge lies in maintaining consistency across multiple steps while managing increasing complexity.	
1429		
1430		
1431		
1432		
1433		
1434		
1435		
1436		
1437		
1438		
1439		
1440	<b>Fact Retrieval</b> tests access to declarative knowledge and the ability to recall, recognize, or apply learned information. Tasks range from theorem retrieval in mathematics, recalling physical laws and formulas in science, accessing API documentation in coding, to retrieving cultural or procedural knowledge in non-reasoning contexts. This behavior assesses both the breadth of a model's knowledge base and its ability to retrieve relevant information when needed. Unlike pure memorization, many tasks require applying retrieved facts to novel situations.	
1441		
1442		
1443		
1444		
1445		
1446		
1447		
1448		
1449		
1450		
1451		
1452	<b>Diagnostic</b> evaluates meta-cognitive abilities, including error detection, consistency checking, and critical evaluation of reasoning processes. This highest-level cognitive behavior requires models to step back from problem-solving to assess the validity of solutions, identify logical flaws, or recognize ambiguities. Examples include verifying mathematical proofs, identifying errors in physical reasoning (such as perturbation analysis), debugging code, and detecting logical inconsistencies in arguments. Diagnostic tasks are particularly challenging because they require models to reason about reasoning itself.	
1453		
1454		
1455		
1456		
1457		
1458		
1459		
1460		
1461		
1462		
1463		
1464		
	<b>I.3 Problem Design Principles</b>	1465
	Each problem template follows several key design principles. First, problems are <i>self-contained</i> , providing all necessary information within the problem statement to avoid ambiguity. Second, they have <i>verifiable solutions</i> , enabling objective evaluation without subjective judgment. Third, they exhibit <i>scalable difficulty</i> , allowing generation of instances ranging from simple to complex by adjusting parameters. Fourth, they are <i>domain-representative</i> , reflecting authentic tasks that practitioners in each field actually encounter. Finally, problems are designed to be <i>minimally ambiguous</i> , with clear success criteria that reduce evaluation uncertainty.	1466 1467 1468 1469 1470 1471 1472 1473 1474 1475 1476 1477 1478
	<b>I.4 Cross-Domain Patterns</b>	1479
	Several interesting patterns emerge from the taxonomy. Calculation tasks share common structure across domains—all involve applying defined procedures to inputs—but differ in the nature of operations (arithmetic vs. physical vs. algorithmic). Enumeration tasks universally require exhaustive search but vary in the combinatorial structure being explored. Simulation tasks all involve temporal or sequential progression but differ in whether the dynamics are deterministic or stochastic, discrete or continuous. Fact Retrieval varies most dramatically across domains, reflecting the distinct knowledge bases required. Diagnostic tasks converge on common meta-cognitive skills despite different surface manifestations.	1480 1481 1482 1483 1484 1485 1486 1487 1488 1489 1490 1491 1492 1493 1494
	<b>I.5 Use Cases</b>	1495
	This benchmark serves multiple purposes in AI evaluation research. It enables <i>comparative analysis</i> of different models across specific cognitive dimensions, revealing whether a model excels at calculation but struggles with enumeration, for instance. It facilitates <i>targeted improvement</i> by identifying precise capability gaps. It supports <i>meta-evaluation</i> by providing a diverse testbed for assessing evaluation methods themselves—different evaluation approaches may perform better or worse across different cells of the taxonomy. Finally, it enables <i>robustness testing</i> by generating multiple problem instances within each template to assess consistency of model performance.	1496 1497 1498 1499 1500 1501 1502 1503 1504 1505 1506 1507 1508 1509

### Math Diagnostic Example

**Problem:** Solve for all real  $x$ :

$$\sqrt{x+4} = x-2$$

**First Attempt:** The square root is always non-negative, so  $x-2 \geq 0$ , hence  $x \geq 2$ . Squaring both sides:

$$x+4 = (x-2)^2 \implies x+4 = x^2-4x+4$$

Simplify:

$$0 = x^2 - 5x \implies x(x-5) = 0$$

Therefore,  $x = 0$  or  $x = 5$ .

**Self-check/correction:** Wait, let's re-check:

Since  $x \geq 2$ , the candidate solution  $x = 0$  is invalid.

Also, substituting  $x = 0$  into the original equation gives

$$\sqrt{4} \neq -2$$

so it does not satisfy the equation.

Thus, the only valid solution is  $x = 5$ .

**Final Answer:**

$$x = 5$$

**Diagnostic:**

The diagnostic behavior is marked by the phrase “Wait, let's re-check” (or similar), followed by detecting a domain or consistency issue and revising the solution.

ative redistribution, and compound interest growth. *Fact retrieval* targets canonically stated results (Euler's circle geometry formulae, Chinese Remainder Theorem and totient, multiple-angle identities, inclusion-exclusion, spherical distance). *Diagnostic* provides templated failure modes—including statement drift, inconsistent traces, theorem misuse, domain or legality mistakes, and answer or format mismatch—so models must detect and correct errors rather than compute alone. Each template admits numeric instantiations with a unique short answer, enabling scalable generation while preserving skill isolation.

1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536

## J Additional Math Examples

Table G illustrates representative problem templates for each behavior in mathematics. *Calculation* covers core symbolic and numeric skills such as greatest common divisor, prime factors, eigenvalues and singular values of small matrices, and solving linear equations. *Enumeration* includes combinatorial distributions, grid path counting, lattice enumeration, seating permutations with constraints, symmetry counting on polyhedra, and expression rationalization. *Simulation* features stepwise processes such as number sequence updates, random walk boundary hitting, expected stopping time, iter-

Table 5: Example problem templates across different mathematical domains. Each row shows a problem category, specific problem type, a template example with its solution, illustrating the diversity of mathematical reasoning tasks from basic arithmetic to error diagnosis.

Category	Problem Name	Template Example	Answer
Calculation	Arithmetic gcd	Calculate the greatest common divisor of 702 and 86814.	234
	prime factors	Find the second-largest prime factor of 49766.	149
	arithmetic matrix eigenvalues	Determine the largest eigenvalue by absolute value of the matrix: $[-8, 4, 4], [4, -1, 6], [4, 6, 5]$	10
	arithmetic matrix svd	Calculate the rounded sum of all singular values from the SVD of: $[[3, -1, -3, 3], [-2, 2, 3, 4], [-2, 2, -3, 3], [-3, -4, 2, -1]]$	21
	algebra linear equation	Solve $-523x + 5q + 30 = -522x, 4x + 0x - 33q = 146$ for $x$ .	20
Enumeration	combinatory distribution	You have letters $\{ 'q' : 4, 'l' : 2, 'b' : 2 \}$ . Distribute them into 3 labeled boxes of capacities $[4, 2, 2]$ . How many ways?	26
	logic gridworld	In a $7 \times 6$ grid, count all monotone paths from $(0, 0)$ to $(6, 5)$ avoiding cells $(6, 0), (3, 5), (1, 1)$ .	183
	geometric lattice enumeration	Determine the largest eigenvalue by absolute value of the matrix: $[-8, 4, 4], [4, -1, 6], [4, 6, 5]$	1540
	combinatorial permutation constraint	Five couples sit in a row of 10 seats; no one sits next to their partner. How many seating arrangements are possible?	1263360
	geometric symmetry counting	How many distinct ways to color the faces of a cube with 6 distinct colors so that adjacent faces have different colors?	30
Simulation	Number-sequence simulation	The integers 1–120 are written on a board. Each minute replace each $n$ by $d(n+3)$ , the divisor-count function. After a day, what is the sum of the numbers on the board?	240
	Stochastic walk / boundary-hitting	A frog starts at $(1, 2)$ inside a $4 \times 4$ square, jumping randomly by unit steps along coordinate axes until it hits the boundary. What is the probability it stops on a vertical side?	$\frac{5}{8}$
	Expected-value stopping time	Joseph rolls a fair die repeatedly until he gets 3 identical consecutive rolls. What is the expected number of rolls?	43
	Iterated redistribution process	Five balls are placed in the first five boxes of a row of six boxes. If any box has $\geq 2$ balls, move one ball to the next box on the right. Repeat until all first five boxes have $\leq 1$ ball. In how many initial placements does the last box end up empty?	1296
	Financial iterative growth	Bao invests \$1000 at 10% annual compound interest. How much total interest is earned after 3 years?	331
Fact Retrieval	Euler's formula / circle geometry	Midpoint $H$ of chord $PQ$ in unit circle; chord length $2 \sin \frac{t}{2}$ . Find $H$ in terms of $t$ .	$(\cos \frac{t}{2} \cos t, \cos \frac{t}{2} \sin t)$
	Chinese Remainder Theorem / Euler's totient	Smallest $x \in \mathbb{N}$ with $x \equiv 3^{234} \pmod{700}$ .	169
	Binomial/Chebyshev cosine multiple-angle	If $\cos \theta = \frac{1}{4}$ , find $\cos 5\theta$ .	$\frac{61}{64}$
	Inclusion–Exclusion	Count integers $7 \leq n \leq 59$ relatively prime to 15.	29
	Great-circle distance (spherical law of cosines)	Surface distance on radius- $r$ sphere between $(p_1, q_1)$ and $(p_2, q_2)$ .	$r \arccos(\sin q_1 \sin q_2 \cos(p_2 - p_1) + \cos q_1 \cos q_2)$
Diagnostic	Problem-statement drift (sign/parameter flip)	Solver inadvertently changes the original problem (flips sign, replaces parameter), solving a different equation. <b>Example:</b> Asked to solve $x^2 - 3x + 2 = 0$ but solves $x^2 + 3x + 2 = 0$ .	NA
	Internal inconsistency in trace	Derivation steps contradict stated givens or earlier results (e.g., variable values change mid-solution). <b>Example:</b> Declares $f(0) = 1$ initially but later uses $f(0) = 0$ .	NA
	Tool / theorem mislabeling or misuse	Applies wrong theorem or uses one without meeting its conditions. <b>Example:</b> Uses Mean Value Theorem on a non-continuous or non-differentiable function.	NA
	Domain / legality mistakes	Ignores domain restrictions (division by zero, square roots of negatives, invalid trig ranges). <b>Example:</b> Cancels $(x - 1)$ without noting $x = 1$ is excluded.	NA
	Answer / format mismatch	Final answer doesn't match derived quantity or required format (units, boxed form, multiple-part). <b>Example:</b> Computes 12 but reports $3x^2 = 12$ , or omits required components.	NA

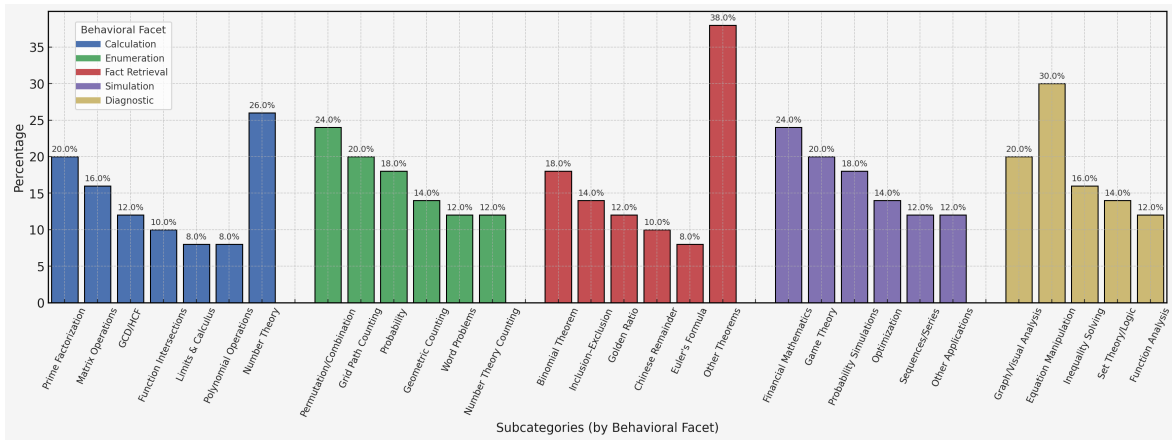


Figure 10: Distribution of math problems subcategories in our newly created dataset, grouped by behavioral facets: Calculation, Enumeration, Fact Retrieval, Simulation, and Diagnostic. Each bar shows the proportion (percentage) of problems belonging to a subcategory, with colored segments indicating the corresponding facet. This visualization highlights the relative prevalence of different cognitive skill types across the dataset.

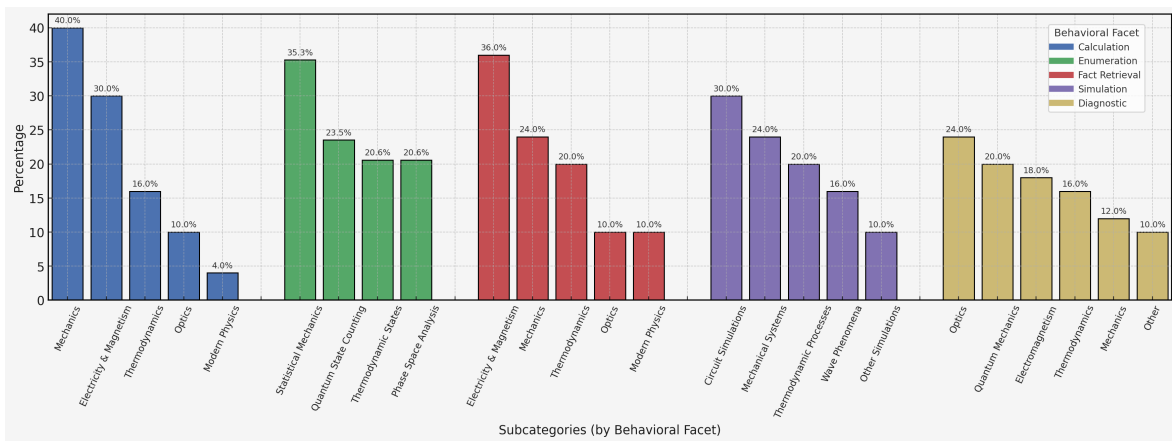


Figure 11: Distribution of physics subcategories in our newly created dataset, grouped by behavioral facets: Calculation, Enumeration, Fact Retrieval, Simulation, and Diagnostic. Each bar indicates the percentage of problems belonging to a specific subcategory, with colors corresponding to the facet type. The plot highlights the diverse coverage of physics domains such as Mechanics, Electromagnetism, Thermodynamics, Optics, and others across different cognitive skill requirements.

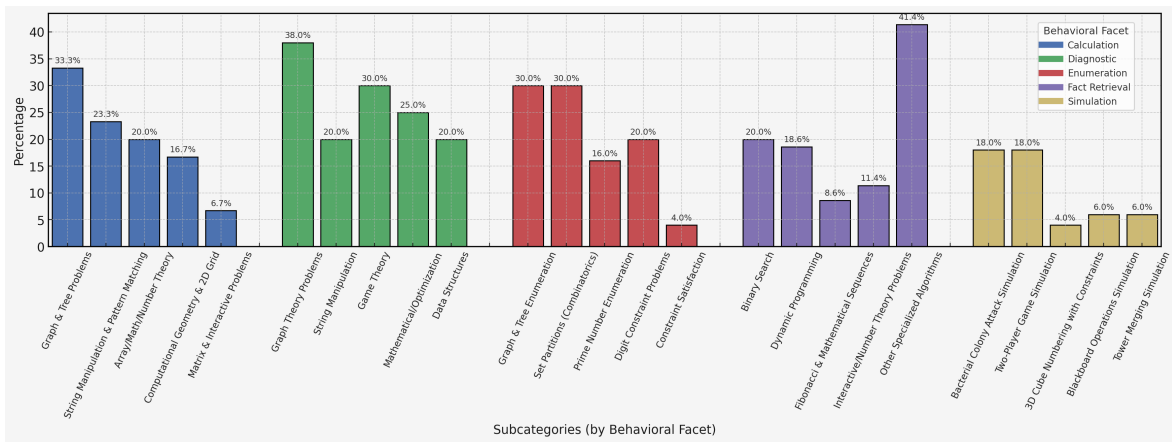


Figure 12: Distribution of code problem subcategories in our dataset, grouped by behavioral facets: Calculation, Diagnostic, Enumeration, Fact Retrieval, and Simulation. The figure highlights a diverse and well-balanced coverage of programming tasks including graph/tree problems, string manipulation, enumeration, dynamic programming, and simulation-based challenges, demonstrating the dataset’s quality and breadth for evaluating reasoning in algorithmic contexts.

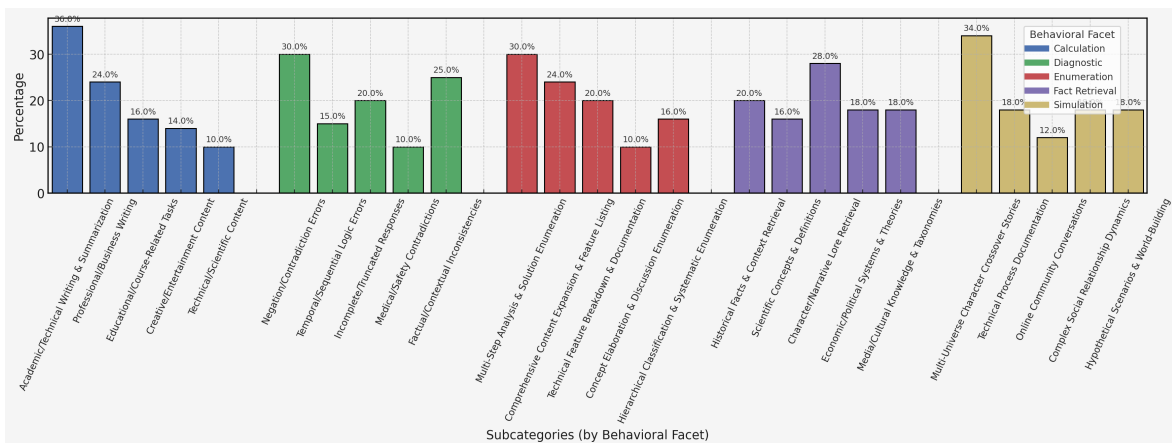


Figure 13: Distribution of non-reasoning subcategories in our dataset, grouped by behavioral facets: Calculation, Diagnostic, Enumeration, Fact Retrieval, and Simulation. The figure shows a diverse and balanced coverage of real-world tasks including writing and summarization, error detection, enumeration of content, knowledge retrieval, and interactive or scenario-based activities, underscoring the dataset’s breadth and quality for evaluating models beyond core reasoning skills.