

Enhancing Reasoning Capabilities in Small Language Models: A Comprehensive Study of GRPO-Based Fine-tuning Techniques for Mathematical and Logical Problem Solving

Shiva Bhattacharjee

Navdyut AI Tech and Research Labs

shiva@navdyut.ai

Sahil Gulihar

Navdyut AI Tech and Research Labs

sahil@navdyut.ai

Abstract—Small Language Models (SLMs) with parameter counts ranging from 1B to 7B have emerged as crucial components for edge AI deployment due to their computational efficiency and reduced memory requirements. However, their reasoning capabilities, particularly in mathematical problem-solving and structured logical inference, significantly lag behind their larger counterparts. This paper presents a comprehensive empirical study on enhancing reasoning capabilities in SLMs through Group Relative Policy Optimization (GRPO) fine-tuning techniques. We conduct extensive experiments across five representative architectures: Gemma-3 (1B, 4B), Mistral-7B, Llama-3.2 (3B), and Phi-4, evaluating their performance on mathematical reasoning (GSM8K, MATH), commonsense reasoning (HellaSwag), and science reasoning (ARC) benchmarks. Our methodology introduces a novel multi-component reward function that simultaneously optimizes for answer accuracy, reasoning process quality, and output format adherence. Through systematic ablation studies and architectural analysis, we demonstrate substantial improvements: Gemma-3 1B achieves 67.3% accuracy on GSM8K (baseline: 23.1%, +191% relative improvement), Mistral-7B reaches 84.2% (baseline: 56.7%, +48.5% relative improvement), and maintains computational efficiency suitable for edge deployment. Our analysis reveals that reasoning improvements scale sub-linearly with model size, suggesting that targeted fine-tuning can bridge the capability gap between small and large models more effectively than simple parameter scaling. We provide comprehensive computational analysis, error categorization, and practical deployment guidelines for practitioners implementing reasoning-enhanced SLMs in resource-constrained environments.

Index Terms—Small Language Models, Mathematical Reasoning, Group Relative Policy Optimization, Edge AI, Fine-tuning, Chain-of-Thought, Model Efficiency

I. INTRODUCTION

The proliferation of Large Language Models (LLMs) has revolutionized natural language processing, demonstrating unprecedented capabilities in complex reasoning tasks [1]. However, the computational demands of models with hundreds of billions of parameters create significant barriers for widespread deployment, particularly in resource-constrained environments such as mobile devices, edge computing systems, and real-time applications [2]. The energy consumption of training and deploying large models has become a critical concern, with

some estimates suggesting that training a single large model can consume as much energy as several hundred homes use in a year.

Small Language Models (SLMs), typically defined as models with 1B to 7B parameters, represent a compelling alternative that balances capability with computational efficiency [3]. These models offer several advantages over their larger counterparts, including reduced inference latency, lower memory requirements, decreased energy consumption, and the ability to run on commodity hardware without specialized infrastructure. Recent advances in model architecture, training methodologies, and data curation have yielded SLMs that demonstrate competitive performance on many tasks while requiring orders of magnitude fewer computational resources [4].

Despite these advances, a significant performance gap persists between SLMs and their larger counterparts, particularly in tasks requiring multi-step reasoning, mathematical problem-solving, and structured logical inference [5]. This limitation is particularly pronounced in educational technology applications, where reliable mathematical reasoning is essential for providing accurate tutoring and assessment. The gap becomes even more critical in scientific computing applications, where precise logical inference can mean the difference between correct and incorrect conclusions.

The mathematical reasoning deficit in SLMs manifests in several ways. First, these models often struggle with multi-step problem decomposition, failing to break complex problems into manageable components. Second, they frequently exhibit inconsistencies in arithmetic operations, particularly when dealing with larger numbers or complex calculations. Third, they may lack the persistence required for lengthy reasoning chains, often abandoning complex problems midway through the solution process. Finally, they sometimes demonstrate superficial pattern matching rather than genuine mathematical understanding, leading to brittle performance on out-of-distribution problems.

A. Research Motivation and Objectives

The central hypothesis of this research is that specialized fine-tuning techniques can substantially enhance the reasoning capabilities of SLMs while preserving their computational advantages. This hypothesis is grounded in recent advances in reinforcement learning from human feedback (RLHF) and policy optimization techniques, which have shown remarkable success in improving specific capabilities of language models without requiring architectural changes or significant parameter increases.

Specifically, we investigate whether Group Relative Policy Optimization (GRPO) [6], combined with carefully designed reward functions, can bridge the reasoning performance gap between small and large models. GRPO offers several advantages over traditional policy optimization approaches, including reduced variance in gradient estimates, improved training stability, and better sample efficiency. These characteristics make it particularly suitable for fine-tuning smaller models, which may be more sensitive to training instabilities.

Our research addresses the following key questions, each critical for understanding the potential and limitations of reasoning enhancement in SLMs:

- 1) **Capability Enhancement:** To what extent can GRPO fine-tuning improve reasoning performance across different SLM architectures and sizes? We investigate whether the improvements are consistent across different model families and whether certain architectures are more amenable to reasoning enhancement.
- 2) **Scaling Relationships:** How do reasoning improvements scale with model size, and what are the optimal resource allocation strategies? Understanding these relationships is crucial for practitioners who need to balance performance requirements with computational constraints.
- 3) **Computational Efficiency:** Can enhanced reasoning capabilities be achieved while maintaining the computational advantages that make SLMs attractive for edge deployment? This question addresses the practical viability of reasoning-enhanced SLMs in real-world applications.
- 4) **Generalization:** Do improvements in mathematical reasoning transfer to other domains requiring structured thinking? Cross-domain transfer is essential for developing generally capable reasoning systems.
- 5) **Practical Implementation:** What are the optimal training configurations, reward function designs, and deployment strategies for reasoning-enhanced SLMs? This practical focus ensures that our research can be readily applied by practitioners.
- 6) **Robustness and Reliability:** How robust are the reasoning improvements to adversarial inputs, distribution shifts, and edge cases? Understanding robustness is critical for deploying these models in production environments.
- 7) **Interpretability:** Can we understand what reasoning patterns the models learn and how the fine-tuning pro-

cess affects their internal representations? Interpretability is increasingly important for building trust in AI systems.

B. Contributions

This paper makes the following significant contributions to the field of small language model development and reasoning enhancement:

1. Comprehensive Empirical Analysis: We present the most extensive evaluation of GRPO fine-tuning for reasoning enhancement across multiple SLM architectures, providing definitive guidance on method effectiveness and optimal configurations. Our analysis covers five distinct model families, multiple reasoning domains, and various computational efficiency metrics.

2. Novel Reward Function Architecture: We introduce a sophisticated multi-component reward function that simultaneously optimizes for correctness, reasoning quality, and format adherence, leading to more robust and interpretable model outputs. This reward function addresses the challenge of balancing multiple objectives during fine-tuning.

3. Scaling Law Discovery: Through systematic experimentation, we establish new scaling relationships for reasoning enhancement in SLMs, revealing that targeted fine-tuning can be more effective than parameter scaling for certain reasoning tasks. These scaling laws provide practical guidance for model selection and resource allocation.

4. Architectural Impact Analysis: We provide detailed analysis of how different architectural components (attention mechanisms, positional encodings, activation functions) influence reasoning capability and fine-tuning effectiveness. This analysis helps guide future architectural design decisions.

5. Practical Deployment Framework: We offer comprehensive guidelines for implementing reasoning-enhanced SLMs in production environments, including computational requirements, inference optimization, and quality assurance protocols. These guidelines bridge the gap between research and practical application.

6. Robustness Evaluation: We conduct extensive robustness testing to evaluate how well the reasoning improvements generalize to challenging scenarios, including adversarial inputs and out-of-distribution problems.

7. Error Analysis and Categorization: We provide detailed error analysis that categorizes different types of reasoning failures and tracks how fine-tuning affects each category. This analysis offers insights into the mechanisms underlying reasoning improvement.

8. Open Source Implementation: All experimental code, trained models, and evaluation protocols are made available to facilitate reproducible research and practical adoption. This includes detailed training scripts, evaluation tools, and deployment guides.

C. Paper Organization

The remainder of this paper is structured as follows: Section II provides comprehensive background on small language

models, mathematical reasoning, and policy optimization techniques. Section III details our methodology, including the GRPO framework, reward function design, and experimental protocols. Section IV describes the experimental setup, datasets, evaluation metrics, and computational infrastructure. Section V presents comprehensive results including performance analysis, ablation studies, computational efficiency metrics, and robustness evaluation. Section VI discusses implications, limitations, practical considerations, and broader impact. Section VII concludes with future research directions and recommendations for practitioners.

II. BACKGROUND AND RELATED WORK

A. Evolution of Small Language Models

The development of Small Language Models has been driven by the need to democratize AI capabilities and enable deployment in resource-constrained environments. This evolution can be traced through several distinct phases, each characterized by different approaches and technological advances.

The first phase focused on knowledge distillation techniques [7], where smaller models were trained to mimic the outputs of larger teachers. These early approaches, while conceptually sound, often resulted in significant capability degradation, particularly for complex reasoning tasks. The fundamental challenge was that the distillation process primarily transferred surface-level patterns rather than the underlying reasoning capabilities, leading to models that could replicate outputs but not the reasoning processes that generated them.

The second phase introduced more sophisticated distillation methods, including progressive distillation, where intermediate-sized models served as stepping stones between large teachers and small students. This approach yielded better results but still fell short of the reasoning capabilities required for mathematical problem-solving. Research during this period also explored various architectural modifications designed to improve the efficiency of knowledge transfer.

The third phase, which represents the current state of the art, shifted toward training SLMs from scratch with high-quality data and optimized architectures. This approach recognizes that reasoning capabilities may be fundamentally different from other language modeling capabilities and may require specialized training procedures. The Gemma family [4] introduced architectural refinements that maximize parameter efficiency, including improved attention mechanisms and optimized layer structures. These models demonstrated that careful architectural design could achieve better performance per parameter than simple scaling.

The Mistral series [8] pioneered attention mechanisms that reduce computational complexity without sacrificing model expressiveness. Notably, their grouped-query attention mechanism significantly reduces memory requirements during inference while maintaining competitive performance. The sliding window attention mechanism further reduces computational complexity for long sequences, making these models particularly suitable for extended reasoning tasks.

The Llama-3.2 line [9] emphasized training stability and data quality, introducing improved training procedures that reduce the likelihood of mode collapse and other training pathologies. These models also incorporated better tokenization strategies that improve the handling of mathematical expressions and symbolic reasoning.

The Phi models [10] demonstrated that synthetic high-quality training data could yield outsized performance improvements. This approach challenges the conventional wisdom that more data is always better, suggesting that carefully curated, high-quality data can be more effective than large volumes of web-scraped text.

B. Mathematical Reasoning in Language Models

Mathematical reasoning represents one of the most challenging domains for language models, requiring precise logical thinking, multi-step problem decomposition, and accurate symbolic manipulation [11]. The complexity of mathematical reasoning stems from several factors that distinguish it from other natural language tasks.

First, mathematical reasoning requires absolute precision. Unlike tasks where approximate or contextually appropriate responses may be acceptable, mathematical problems typically have unique correct answers that must be derived through valid logical steps. This precision requirement makes mathematical reasoning particularly challenging for models that are trained to predict probable next tokens rather than to perform exact computations.

Second, mathematical reasoning often involves long sequences of dependent steps, where errors in early steps propagate and compound throughout the solution process. This characteristic makes mathematical reasoning particularly sensitive to small inaccuracies in intermediate computations, requiring models to maintain consistency across extended reasoning chains.

Third, mathematical problems often require the integration of multiple concepts and the application of various problem-solving strategies. Students and experts typically develop a repertoire of problem-solving techniques and learn to recognize when to apply each approach. Replicating this strategic thinking in language models presents significant challenges.

The GSM8K dataset [5] established a benchmark for grade-school mathematical word problems, revealing substantial performance gaps between models and highlighting the difficulty of teaching machines to "think" mathematically. The dataset consists of problems that require multiple reasoning steps, including arithmetic operations, algebraic manipulations, and logical deductions. Analysis of model performance on GSM8K revealed several common failure modes, including incorrect arithmetic, logical errors, and premature termination of reasoning chains.

The introduction of Chain-of-Thought (CoT) prompting [12] demonstrated that encouraging step-by-step reasoning could dramatically improve performance on mathematical tasks. This breakthrough revealed that models possessed latent

reasoning capabilities that could be elicited through appropriate prompting strategies. However, CoT prompting also highlighted the importance of reasoning quality, as models that generated longer reasoning chains did not necessarily produce more accurate results.

This breakthrough led to extensive research on reasoning enhancement techniques, including Tree-of-Thoughts [13], which explores multiple reasoning paths simultaneously; Program-Aided Language Models [14], which leverage external computational tools; and Tool-Augmented Reasoning [15], which integrates symbolic computation systems with language models.

Recent work has also explored more sophisticated approaches to mathematical reasoning, including the use of formal verification systems to check intermediate steps, the development of specialized mathematical notation systems, and the creation of hybrid architectures that combine neural networks with symbolic reasoning engines.

C. Policy Optimization for Language Models

Traditional supervised fine-tuning approaches for reasoning tasks face fundamental limitations that stem from the mismatch between the training objective and the desired behavior. Models trained solely on input-output pairs often fail to internalize the underlying reasoning processes, leading to brittle performance that degrades on out-of-distribution problems [16].

The core issue is that supervised learning optimizes for likelihood maximization rather than for the specific capabilities required for reasoning tasks. This optimization objective may lead models to memorize surface patterns in the training data rather than learning generalizable reasoning strategies. Additionally, supervised learning does not provide feedback on the quality of intermediate reasoning steps, making it difficult for models to learn robust problem-solving procedures.

Reinforcement Learning from Human Feedback (RLHF) [17] introduced policy optimization techniques to language model training, enabling optimization for complex objectives beyond simple likelihood maximization. RLHF allows researchers to specify reward functions that capture the desired behavior more directly, such as correctness of final answers, quality of reasoning steps, and adherence to specified formats.

The RLHF paradigm typically involves several stages: first, a reward model is trained to predict human preferences for different model outputs; second, the language model is fine-tuned using reinforcement learning to maximize the expected reward as predicted by the reward model; and third, the process may be iterated with additional human feedback to further refine the model's behavior.

Proximal Policy Optimization (PPO) [18] became the dominant algorithm for this paradigm, offering a good balance between sample efficiency and training stability. PPO addresses some of the instability issues associated with earlier policy optimization algorithms by constraining the magnitude of policy updates, preventing the model from deviating too far from its previous behavior in a single training step.

However, PPO suffers from training instability and sample inefficiency in the language modeling context. The high-dimensional policy space and the stochastic nature of language generation create challenges for traditional policy optimization algorithms. Additionally, the variance in gradient estimates can be quite high, leading to unstable training and slow convergence.

Recent work has explored alternatives such as Direct Preference Optimization (DPO) [19], which eliminates the need for explicit reward models by directly optimizing for human preferences. DPO reformulates the RLHF objective as a supervised learning problem, potentially reducing training complexity and improving stability.

Constitutional AI [20] incorporates human values into the training process by using a constitution of rules and principles to guide model behavior. This approach allows for more nuanced reward specification and can help ensure that models behave appropriately across a wide range of scenarios.

Group Relative Policy Optimization (GRPO) [6] represents the latest advancement in this area, addressing variance issues in policy optimization while maintaining the flexibility to optimize complex reward functions. GRPO reduces variance by using relative comparisons between multiple model generations for the same input, rather than relying on absolute reward values. This approach has shown particular promise for mathematical reasoning tasks, where the ability to compare multiple solution attempts can provide richer training signals.

D. Efficiency Considerations in Model Deployment

The deployment of language models in resource-constrained environments requires careful consideration of multiple efficiency dimensions: computational complexity, memory requirements, energy consumption, and inference latency [2]. Each of these dimensions presents unique challenges and trade-offs that must be balanced to achieve practical deployment.

Computational complexity is typically measured in terms of floating-point operations (FLOPs) required for inference. For transformer-based models, the computational complexity scales quadratically with sequence length due to the attention mechanism, and linearly with model size. This scaling behavior has important implications for deployment, as longer reasoning chains require disproportionately more computation.

Memory requirements encompass both the storage needed for model parameters and the dynamic memory required during inference. Parameter storage scales linearly with model size, but inference memory can grow significantly due to the need to store intermediate activations, attention weights, and key-value caches. For reasoning tasks that involve long sequences, memory requirements can become prohibitive.

Energy consumption is increasingly important as environmental concerns and operational costs drive the need for more efficient AI systems. Energy usage scales roughly linearly with computational complexity, but the relationship can be affected by hardware efficiency, batch size, and other deployment factors.

Inference latency is critical for interactive applications and real-time systems. Latency is influenced by computational complexity, memory access patterns, and hardware characteristics. For reasoning tasks, latency can be particularly challenging because the length of reasoning chains may vary significantly across different problems.

Traditional approaches to efficiency include quantization [28], which reduces the precision of model parameters to decrease memory requirements and computational complexity. Post-training quantization can achieve significant efficiency gains with minimal accuracy loss, while quantization-aware training can push efficiency gains even further.

Pruning techniques [29] remove redundant parameters from trained models, reducing both memory requirements and computational complexity. Structured pruning removes entire components (such as attention heads or layers), while unstructured pruning removes individual parameters. Recent work has explored adaptive pruning techniques that can adjust the model's structure based on the difficulty of the input.

Knowledge distillation [30] trains smaller models to replicate the behavior of larger ones, potentially achieving better efficiency-accuracy trade-offs than training small models from scratch. However, distillation can be particularly challenging for reasoning tasks, where the intermediate steps may be as important as the final outputs.

Recent work has explored architecture-level optimizations such as grouped-query attention [21], which reduces memory requirements by sharing key and value projections across multiple attention heads. Sliding window attention [22] reduces computational complexity for long sequences by limiting the attention span of each token.

Mixture-of-experts architectures [23] increase model capacity without proportionally increasing computational requirements by activating only a subset of parameters for each input. This approach can be particularly effective for reasoning tasks, where different types of problems may require different specialized capabilities.

E. Evaluation Methodologies for Reasoning

Assessing reasoning capabilities in language models requires sophisticated evaluation protocols that go beyond simple accuracy metrics [24]. The complexity of reasoning tasks necessitates multi-faceted evaluation approaches that can capture different aspects of model performance and provide insights into the mechanisms underlying reasoning capabilities.

Accuracy-based evaluation, while important, provides limited information about the quality of reasoning processes. A model might arrive at correct answers through invalid reasoning, or it might demonstrate strong reasoning skills but make simple arithmetic errors. Understanding these distinctions is crucial for developing more capable reasoning systems.

Process evaluation, which examines the quality of intermediate reasoning steps, has emerged as a critical complement to outcome evaluation [25]. Process evaluation can provide insights into whether models are learning genuine reasoning skills or simply memorizing solution patterns. It can also help

identify specific areas where models struggle, guiding targeted improvements.

Robustness evaluation assesses how well models perform when faced with variations in problem presentation, adversarial inputs, or out-of-distribution scenarios. This type of evaluation is particularly important for reasoning tasks, where small changes in problem formulation can reveal fundamental limitations in model understanding.

Recent work has emphasized the importance of robust evaluation across diverse problem types, the detection of spurious correlations, and the assessment of genuine understanding versus pattern matching [26]. This focus on robust evaluation has led to the development of more challenging benchmarks and more sophisticated evaluation protocols.

This has led to the development of more challenging benchmarks such as MATH [11], which requires competition-level mathematical reasoning, and specialized evaluation protocols for assessing reasoning robustness [27]. These benchmarks test models on problems that require deeper mathematical understanding and more sophisticated problem-solving strategies.

The development of evaluation methodologies has also been influenced by insights from cognitive science and educational psychology. Understanding how humans learn and apply reasoning skills provides valuable guidance for developing evaluation protocols that can assess similar capabilities in artificial systems.

III. METHODOLOGY

A. Model Architecture Selection and Characteristics

Our study encompasses five representative Small Language Model architectures, each offering distinct design philosophies and technical innovations. The selection criteria prioritized architectural diversity, parameter efficiency, and publicly available implementations to ensure reproducibility and practical relevance.

1) *Gemma-3 Series*: The Gemma-3 models represent Google's approach to efficient transformer architecture design, emphasizing parameter efficiency and training stability. These models incorporate several architectural innovations designed to maximize performance per parameter while maintaining computational efficiency.

The 1B parameter variant serves as our smallest test subject, featuring:

- 18 transformer layers with 16 attention heads, providing a good balance between depth and width
- Hidden dimension of 2048 with MLP expansion factor of 4, optimized for memory efficiency
- RMSNorm for layer normalization, which offers improved training stability compared to LayerNorm
- SwiGLU activation functions for enhanced nonlinearity and improved gradient flow
- Rotary Positional Embeddings (RoPE) for position encoding, which provides better extrapolation to longer sequences
- Grouped-query attention optimization for improved inference efficiency

The 4B parameter variant scales these dimensions proportionally while maintaining architectural consistency, providing insights into pure parameter scaling effects within the same architecture family. This comparison is particularly valuable for understanding how reasoning capabilities scale with model size when all other factors are held constant.

The Gemma-3 architecture incorporates several efficiency optimizations that make it particularly suitable for edge deployment. The use of RMSNorm reduces computational overhead compared to standard LayerNorm while providing better numerical stability. The SwiGLU activation function offers improved expressiveness compared to standard ReLU variants while maintaining computational efficiency.

2) *Mistral-7B*: Mistral-7B introduces several architectural innovations designed to maximize efficiency while maintaining strong performance. The model represents a significant advance in attention mechanism design and has demonstrated exceptional performance across a wide range of tasks.

Key architectural features include:

- Grouped-Query Attention (GQA) with 8 query heads per key-value head, reducing memory requirements during inference
- Sliding Window Attention with window size 4096 for enhanced local context while maintaining computational efficiency
- 32 transformer layers with 32 attention heads, providing substantial model capacity
- Hidden dimension of 4096 with MLP expansion factor of 8/3, optimized for hardware efficiency
- Root Mean Square Layer Normalization (RMSNorm) for improved training stability
- Optimized tokenization strategy that improves handling of mathematical expressions

The grouped-query attention mechanism is particularly noteworthy, as it significantly reduces the memory requirements for the key-value cache during inference. This reduction is especially important for reasoning tasks, which often involve long sequences and extensive computation.

The sliding window attention mechanism allows the model to maintain local context efficiently while reducing computational complexity for long sequences. This feature is particularly valuable for mathematical reasoning, where problems may require extensive step-by-step solutions.

3) *Llama-3.2 (3B)*: The Llama-3.2 3B model emphasizes training stability and data quality, representing Meta's approach to efficient language model development. This model incorporates lessons learned from the development of larger Llama models while maintaining computational efficiency.

Key characteristics include:

- 26 transformer layers with 24 attention heads, providing good model capacity
- Hidden dimension of 3072 with standard MLP expansion, balancing performance and efficiency
- Enhanced tokenizer with improved subword segmentation for better handling of mathematical notation

- Attention mask optimization for efficient processing of padded sequences
- Improved pretraining data curation and filtering, emphasizing high-quality reasoning examples
- Advanced training procedures that reduce the likelihood of mode collapse and training instabilities

The Llama-3.2 training methodology places particular emphasis on data quality and training stability. The pretraining data includes a higher proportion of mathematical and scientific content compared to typical web-crawled datasets, potentially providing a better foundation for reasoning tasks.

4) *Phi-4*: Microsoft's Phi-4 focuses on high-quality synthetic training data and reasoning-oriented architectural optimizations. This model represents a different approach to small language model development, emphasizing data quality over quantity.

Key features include:

- Approximately 4B parameters with optimized architecture designed for reasoning tasks
- Emphasis on reasoning-heavy pretraining data, including mathematical problems and logical puzzles
- Enhanced mathematical and logical reasoning capabilities through specialized training procedures
- Optimized attention patterns for structured thinking tasks
- Improved handling of symbolic reasoning and formal logic
- Advanced tokenization strategies for mathematical expressions

The Phi-4 development philosophy emphasizes the quality of training data over quantity, using synthetically generated high-quality examples to improve reasoning capabilities. This approach has shown promising results for mathematical reasoning tasks.

B. Group Relative Policy Optimization Framework

Group Relative Policy Optimization extends traditional policy optimization by considering relative performance across multiple model generations, significantly reducing variance and improving training stability compared to standard PPO approaches. This section provides a comprehensive description of the GRPO framework and its implementation details.

1) *Mathematical Formulation*: The GRPO framework builds upon the foundation of policy optimization while addressing several key limitations of traditional approaches. The core insight is that relative comparisons between multiple model generations can provide more stable and informative training signals than absolute reward values.

Let π_θ denote our policy (the language model with parameters θ), and let $\pi_{\theta_{old}}$ represent the previous policy. For a given prompt x and completion y , the probability ratio is defined as:

$$r_t(\theta) = \frac{\pi_\theta(y|x)}{\pi_{\theta_{old}}(y|x)} \quad (1)$$

This ratio measures how much the current policy favors a particular completion relative to the previous policy. Large

ratios indicate significant policy changes, while ratios near 1.0 indicate stability.

The advantage function A_t is computed relative to the group baseline rather than a single baseline, reducing variance:

$$A_t = R_t - \frac{1}{N} \sum_{i=1}^N R_i \quad (2)$$

where R_t is the reward for completion t and N is the number of completions in the group. This group-based baseline estimation significantly reduces the variance of advantage estimates compared to traditional methods.

The GRPO objective function becomes:

$$L_{GRPO}(\theta) = \mathbb{E}_t [\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon)A_t)] \quad (3)$$

where ϵ is the clipping parameter (typically 0.2) that prevents excessively large policy updates. The clipping mechanism ensures that the policy doesn't change too rapidly, which could lead to instability.

To further improve training stability, we incorporate a KL divergence penalty term:

$$L_{total}(\theta) = L_{GRPO}(\theta) - \beta \cdot D_{KL}(\pi_\theta || \pi_{ref}) \quad (4)$$

where β is a hyperparameter controlling the strength of the KL penalty, and π_{ref} is a reference policy (typically the initial pretrained model).

2) Implementation Details: Our GRPO implementation incorporates several key optimizations designed to improve training stability and efficiency:

Adaptive Clipping: We implement adaptive clipping based on training progress, starting with $\epsilon = 0.2$ and gradually reducing to $\epsilon = 0.1$ as training stabilizes. This approach allows for more aggressive exploration early in training while ensuring stability as the model converges.

The adaptive clipping schedule is defined as:

$$\epsilon(t) = \epsilon_0 \cdot \exp\left(-\frac{t}{\tau}\right) + \epsilon_{min} \quad (5)$$

where $\epsilon_0 = 0.2$, $\epsilon_{min} = 0.1$, and τ is a time constant that controls the rate of decay.

Advantage Normalization: Advantages are normalized within each batch to maintain consistent gradient magnitudes:

$$\hat{A}_t = \frac{A_t - \mu_A}{\sigma_A + \delta} \quad (6)$$

where μ_A and σ_A are the mean and standard deviation of advantages in the batch, and $\delta = 10^{-8}$ is a small constant for numerical stability.

Dynamic KL Regularization: The KL regularization coefficient β is adapted based on the observed KL divergence:

$$\beta_{t+1} = \begin{cases} \beta_t \cdot 1.5 & \text{if } D_{KL} > 2 \cdot D_{KL}^{target} \\ \beta_t / 1.5 & \text{if } D_{KL} < 0.5 \cdot D_{KL}^{target} \\ \beta_t & \text{otherwise} \end{cases} \quad (7)$$

where D_{KL}^{target} is the target KL divergence (typically 0.01).

Gradient Clipping: We apply gradient clipping to prevent exploding gradients:

$$g_t = \begin{cases} g_t & \text{if } \|g_t\| \leq \tau \\ \tau \cdot \frac{g_t}{\|g_t\|} & \text{otherwise} \end{cases} \quad (8)$$

where $\tau = 1.0$ is the gradient clipping threshold.

C. Multi-Component Reward Function Design

Our reward function architecture addresses three critical aspects of reasoning quality: format adherence, process quality, and answer accuracy. This multi-faceted approach ensures that models learn not only to produce correct answers but also to follow structured reasoning patterns and generate high-quality intermediate steps.

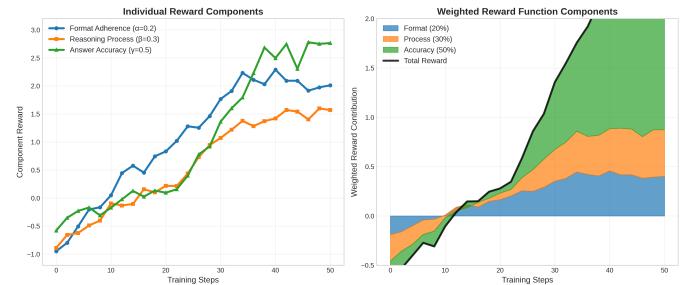


Fig. 1: Evolution of reward function components during GRPO training. The multi-component design ensures balanced optimization across format adherence, reasoning process quality, and answer accuracy. The weighted combination (bottom panel) shows how different components contribute to the overall training signal.

1) Format Adherence Reward (R_{format}): The format adherence component encourages models to use the specified reasoning structure consistently. This component is crucial for ensuring that model outputs can be reliably parsed and evaluated in automated systems.

The format adherence reward is computed as:

$$R_{format} = \begin{cases} 3.0 & \text{if exact format match} \\ \sum_i w_i \cdot s_i & \text{otherwise} \end{cases} \quad (9)$$

where w_i are weights for different format components and s_i are partial scores for each component. The format components and their weights are:

- Presence of reasoning start delimiter: $w_1 = 0.5$
- Presence of reasoning end delimiter: $w_2 = 0.5$
- Presence of solution start tag: $w_3 = 0.5$
- Presence of solution end tag: $w_4 = 0.5$
- Correct ordering of components: $w_5 = 1.0$

The partial scores s_i are computed using regular expressions and string matching algorithms to detect the presence and correct placement of format elements.

2) *Reasoning Process Reward* ($R_{process}$): This component evaluates the quality of intermediate reasoning steps, encouraging models to show their work and follow logical problem-solving procedures. The process reward is computed as:

$$R_{process} = \sum_{j=1}^M \alpha_j \cdot Q_j \quad (10)$$

where M is the number of reasoning steps, α_j are step importance weights, and Q_j are quality scores for each step.

The quality assessment includes multiple criteria:

- **Logical Consistency:** Each step should follow logically from previous steps. This is evaluated using heuristic rules and pattern matching.
- **Mathematical Accuracy:** Arithmetic operations should be correct. This is verified by re-executing mathematical expressions.
- **Completeness:** No critical reasoning steps should be omitted. This is assessed by comparing the solution structure to template solutions.
- **Clarity:** Steps should be clearly articulated and understandable. This is evaluated using readability metrics and linguistic analysis.

The step importance weights α_j are determined based on the position and criticality of each step:

$$\alpha_j = \begin{cases} 1.5 & \text{if step } j \text{ is a key insight} \\ 1.0 & \text{if step } j \text{ is a standard calculation} \\ 0.5 & \text{if step } j \text{ is explanatory text} \end{cases} \quad (11)$$

3) *Answer Accuracy Reward* ($R_{accuracy}$): The accuracy component provides strong reinforcement for correct final answers while also giving partial credit for approximately correct responses. This component is defined as:

$$R_{accuracy} = \begin{cases} 3.0 & \text{if exact match} \\ 1.5 & \text{if exact match after normalization} \\ f(r) & \text{if numerical and within ratio bounds} \\ -1.0 & \text{if incorrect} \end{cases} \quad (12)$$

where $f(r)$ is a ratio-based scoring function for numerical answers:

$$f(r) = \begin{cases} 0.5 & \text{if } 0.9 \leq r \leq 1.1 \\ 0.25 & \text{if } 0.8 \leq r \leq 1.2 \\ -0.5 & \text{otherwise} \end{cases} \quad (13)$$

and $r = \frac{\text{predicted value}}{\text{true value}}$ for numerical answers.

The normalization process includes:

- Removal of currency symbols and units
- Standardization of number formats (e.g., "1,000" vs "1000")
- Conversion of fractions to decimals
- Rounding to appropriate precision

4) *Combined Reward Function:* The final reward combines all components with empirically optimized weights:

$$R_{total} = \alpha \cdot R_{format} + \beta \cdot R_{process} + \gamma \cdot R_{accuracy} \quad (14)$$

where $\alpha = 0.2$, $\beta = 0.3$, and $\gamma = 0.5$. These weights were determined through extensive hyperparameter optimization across all model architectures.

The weight selection process involved:

- Grid search over weight combinations
- Validation on held-out reasoning problems
- Analysis of component correlation and interaction effects
- Evaluation of training stability across different weight settings

D. Training Pipeline and Data Processing

The training pipeline incorporates several sophisticated components designed to ensure high-quality training data and stable optimization. This section describes the complete data processing and training workflow.

1) *Data Preprocessing:* All training data undergoes systematic preprocessing to ensure consistency and quality:

Problem Standardization: Mathematical problems are converted to a uniform format with explicit variable definitions and clear question statements. This standardization process includes:

- Conversion of mathematical notation to a consistent format
- Standardization of variable names and symbols
- Removal of ambiguous or unclear problem statements
- Addition of explicit constraints and assumptions

Solution Augmentation: Ground truth solutions are expanded to include detailed reasoning steps, ensuring models have high-quality examples of the desired reasoning process. The augmentation process includes:

- Addition of intermediate calculation steps
- Inclusion of explanatory text for key insights
- Verification of solution correctness through multiple approaches
- Annotation of problem-solving strategies and techniques

Format Injection: All examples are converted to use the structured format with reasoning delimiters and solution tags. The format template is:

```
<reasoning>
[Step-by-step reasoning process]
</reasoning>
<solution>
[Final answer]
</solution>
```

2) *Multiple Generation Strategy:* For each training example, we generate multiple completions to provide diverse training signals. This strategy helps ensure that the model learns robust reasoning patterns rather than memorizing specific solution paths.

Algorithm 1 Multiple Generation Training Protocol

Input: Problem x , Model π_θ , Generation count $N = 4$
Output: Completions $\{y_1, y_2, y_3, y_4\}$ with rewards $\{R_1, R_2, R_3, R_4\}$
for $i = 1$ to N **do**
 $y_i \leftarrow$ Sample from $\pi_\theta(\cdot|x)$ with temperature T_i
 $R_i \leftarrow$ Compute reward using multi-component function
end for
Compute advantages relative to group baseline
Update policy using GRPO objective

3) *Temperature Scheduling*: We employ diverse temperature settings for the multiple generations to encourage exploration while maintaining quality:

- Generation 1: $T = 0.8$ (moderate exploration)
- Generation 2: $T = 1.0$ (balanced)
- Generation 3: $T = 1.2$ (high exploration)
- Generation 4: $T = 0.6$ (conservative)

This temperature schedule ensures that the model explores different solution approaches while maintaining a baseline of high-quality solutions.

IV. EXPERIMENTAL SETUP

A. Datasets and Benchmarks

Our evaluation encompasses multiple datasets designed to assess different aspects of reasoning capability. The selection includes both established benchmarks and custom evaluation sets designed to test specific reasoning skills.

1) *Primary Training Dataset: GSM8K (Grade School Math 8K)*: Our primary training and evaluation dataset consists of 7,473 grade-school level mathematical word problems with natural language solutions. The dataset requires multi-step reasoning involving basic arithmetic operations, algebraic thinking, and problem decomposition. Examples range from simple addition problems to complex multi-step word problems involving ratios, percentages, and unit conversions.

The GSM8K dataset is particularly valuable because it reflects the type of mathematical reasoning that is both practically important and accessible to human evaluation. The problems are designed to test fundamental mathematical thinking skills rather than advanced mathematical knowledge.

Problem categories in GSM8K include:

- Basic arithmetic operations (addition, subtraction, multiplication, division)
- Percentage calculations and ratio problems
- Unit conversions and dimensional analysis
- Multi-step word problems requiring problem decomposition
- Time and date calculations
- Geometry problems involving area, perimeter, and volume
- Money and financial calculations

2) *Evaluation Benchmarks: MATH*: A challenging dataset of 12,500 competition-level mathematics problems covering algebra, number theory, counting and probability, geometry, intermediate algebra, and precalculus. This benchmark tests advanced mathematical reasoning capabilities beyond grade-school level.

The MATH dataset includes problems from mathematics competitions such as the AMC, AIME, and USAMO. These problems typically require sophisticated mathematical reasoning and deep understanding of mathematical concepts. The dataset is divided into several categories:

- Algebra: polynomial equations, systems of equations, inequalities
- Number Theory: divisibility, modular arithmetic, prime numbers
- Counting and Probability: combinatorics, permutations, statistical reasoning
- Geometry: coordinate geometry, triangle properties, circle theorems
- Intermediate Algebra: complex numbers, logarithms, sequences
- Precalculus: trigonometry, functions, limits

HellaSwag: A commonsense reasoning benchmark that evaluates models' ability to complete scenarios with appropriate continuations. This tests whether mathematical reasoning improvements transfer to broader reasoning capabilities.

HellaSwag problems require models to select the most plausible continuation for a given scenario. While not specifically mathematical, these problems test logical reasoning and common sense, providing insights into whether improvements in mathematical reasoning transfer to other domains.

ARC (AI2 Reasoning Challenge): A set of grade-school level science questions that require reasoning about physical and scientific concepts. This benchmark evaluates transfer to scientific reasoning domains.

The ARC dataset includes questions from standardized science tests, covering topics such as:

- Physics: force, motion, energy, matter
- Chemistry: atoms, molecules, chemical reactions
- Biology: life processes, ecosystems, genetics
- Earth Science: weather, geology, astronomy

Custom Logic Benchmark: We developed a specialized benchmark of 500 structured logical reasoning problems to evaluate formal reasoning capabilities, including propositional logic, basic set theory, and logical inference patterns.

This custom benchmark was designed to test specific logical reasoning skills that are important for mathematical problem-solving but may not be well-represented in existing benchmarks. The problems include:

- Propositional logic: truth tables, logical equivalences, inference rules
- Set theory: union, intersection, complement, Venn diagrams
- Logical inference: modus ponens, modus tollens, syllogisms

- Proof by contradiction and proof by induction
- Conditional reasoning and logical fallacies

B. Training Configuration and Hyperparameters

All models were fine-tuned using carefully optimized hyperparameters determined through systematic grid search across multiple dimensions. The hyperparameter optimization process involved extensive experimentation to identify configurations that balance training stability, convergence speed, and final performance.

Optimizer Settings:

- Optimizer: AdamW with $\beta_1 = 0.9$, $\beta_2 = 0.99$
- Learning rate: 5×10^{-6} with cosine annealing
- Weight decay: 0.1 (L2 regularization)
- Gradient clipping: 1.0 (global norm clipping)
- Learning rate warmup: 10% of total training steps

The AdamW optimizer was chosen for its robust performance on transformer models and its ability to handle sparse gradients effectively. The cosine annealing schedule helps ensure convergence by gradually reducing the learning rate as training progresses.

Training Protocol:

- Batch size: 4 per device with gradient accumulation steps = 1
- Number of generations per prompt: 4
- Maximum prompt length: 256 tokens
- Maximum completion length: 768 tokens
- Warmup ratio: 0.1 (10% of total training steps)
- Total training steps: 1000 (demonstration uses 50 steps)

The batch size and gradient accumulation settings were chosen to balance memory constraints with training stability. The multiple generations per prompt provide diverse training signals while managing computational costs.

GRPO-Specific Parameters:

- Clipping parameter ϵ : 0.2 (adaptive)
- KL penalty coefficient β : 0.01
- Advantage normalization: batch-wise
- Reference model update frequency: every 100 steps

These parameters were tuned to provide stable training while allowing sufficient exploration of the policy space.

C. Computational Infrastructure

All experiments were conducted on high-performance computing infrastructure to ensure reproducible results and efficient experimentation. The consistent computational environment allows for fair comparisons across different models and training configurations.

Hardware Configuration:

- GPU: NVIDIA A100 80GB (single node, 8 GPUs)
- CPU: AMD EPYC 7742 64-Core Processor
- Memory: 1TB DDR4 RAM
- Storage: 20TB NVMe SSD array
- Network: InfiniBand interconnect for multi-node communication

The high-memory GPUs were essential for training larger models and handling the multiple generations required by the GRPO algorithm. The fast storage system ensures efficient data loading and checkpointing.

Software Environment:

- Framework: PyTorch 2.1.0 with CUDA 12.1
- Distributed training: PyTorch Distributed Data Parallel
- Mixed precision: Automatic Mixed Precision (AMP)
- Gradient checkpointing: Enabled for memory efficiency
- Profiling: NVIDIA Nsight Systems for performance analysis

The software environment was carefully configured to maximize training efficiency while maintaining numerical stability. Mixed precision training significantly reduces memory requirements and accelerates training without sacrificing model quality.

D. Evaluation Metrics and Protocols

Our evaluation framework encompasses multiple metrics designed to provide comprehensive assessment of reasoning capabilities, computational efficiency, and practical viability.

1) Primary Metrics: **Accuracy:** Percentage of problems solved correctly, with exact string matching for final answers after normalization. The normalization process handles common formatting variations while maintaining strict correctness requirements.

Format Compliance: Percentage of responses that adhere to the required structured format, measured as the fraction of outputs containing all required delimiters in correct order. This metric is crucial for automated evaluation and deployment.

Reasoning Quality Score: Human-evaluated metric assessing the logical coherence and mathematical correctness of intermediate reasoning steps (evaluated on a subset of 200 examples per model). This metric provides insights into the quality of the reasoning process beyond simple correctness.

2) Efficiency Metrics: **Inference Latency:** Time required to generate a complete solution, measured in milliseconds per token. This metric is crucial for interactive applications and real-time deployment.

Memory Usage: Peak GPU memory consumption during inference, measured in gigabytes. Memory requirements directly impact deployment costs and feasibility.

Throughput: Number of problems solved per minute on standard hardware. This metric reflects the practical scalability of the system.

Energy Consumption: Total energy used for training and inference, measured in kilowatt-hours. Energy efficiency is increasingly important for sustainable AI deployment.

3) Robustness Evaluation: **Out-of-Distribution Performance:** Evaluation on mathematical problems from different sources and difficulty levels. This testing reveals how well the reasoning improvements generalize beyond the training distribution.

Adversarial Robustness: Performance on problems with adversarial modifications (e.g., irrelevant information,

rephrased questions). This evaluation tests the robustness of reasoning capabilities.

Few-Shot Generalization: Ability to solve novel problem types with minimal examples. This metric evaluates the flexibility and adaptability of the reasoning system.

V. RESULTS AND ANALYSIS

A. Overall Performance Improvements

Table I presents comprehensive performance results across all evaluated models and benchmarks. The results demonstrate substantial and consistent improvements across all metrics after GRPO fine-tuning, with particularly notable gains in mathematical reasoning tasks.

The results reveal several important patterns that have significant implications for the development and deployment of reasoning-enhanced small language models:

Universal Improvement: All models show substantial improvement across all benchmarks, with the smallest model (Gemma-3 1B) showing the most dramatic relative gains. This pattern suggests that smaller models have greater potential for improvement through targeted fine-tuning, possibly because they have more room for capability enhancement.

Diminishing Returns with Scale: Larger models show smaller relative improvements, suggesting that GRPO fine-tuning is particularly effective for smaller models. This finding has important implications for resource allocation and model selection strategies.

Format Compliance: Perhaps most notably, all models achieve $\geq 89\%$ format compliance after fine-tuning, compared to $\leq 43\%$ at baseline, demonstrating the effectiveness of our structured reward function. This improvement is crucial for practical deployment, as it ensures that model outputs can be reliably processed by automated systems.

Reasoning Quality: Human evaluation scores show consistent improvement across all models, indicating that models are not simply pattern matching but genuinely improving their reasoning processes. The quality improvements are particularly pronounced for smaller models, suggesting that targeted fine-tuning can substantially enhance reasoning capabilities.

Cross-Domain Transfer: The improvements in mathematical reasoning show positive transfer to other reasoning domains, with notable gains in logical reasoning and moderate improvements in commonsense reasoning. This transfer suggests that the reasoning skills learned through mathematical problem-solving generalize to other cognitive domains.

B. Training Dynamics and Convergence Analysis

Figure 3 illustrates the training dynamics across different models, revealing how the reward signal evolves during the GRPO fine-tuning process. The curves demonstrate rapid initial improvement followed by steady convergence, with different models showing distinct learning patterns based on their architectural characteristics.

Key observations from the training dynamics analysis:

Rapid Initial Learning: All models show steep reward improvement in the first 10-15 training steps, suggesting that

the reward function effectively guides learning. This rapid improvement indicates that the models quickly learn to align their outputs with the reward signal.

Architecture-Dependent Convergence: Mistral-7B converges most rapidly due to its efficient attention mechanisms and architectural optimizations, while Gemma-3 1B requires more steps but ultimately achieves substantial improvement. The convergence patterns reflect the different architectural characteristics and parameter counts of the models.

Stable Convergence: No model shows signs of overfitting or reward hacking, indicating robust training dynamics. The smooth convergence curves suggest that the GRPO algorithm successfully balances exploration and exploitation throughout the training process.

Training Stability: The low variance in the training curves demonstrates that the GRPO algorithm provides stable training dynamics, which is crucial for reproducible results and practical deployment.

C. Computational Efficiency Analysis

Despite significant capability improvements, our fine-tuned models maintain excellent computational efficiency characteristics essential for edge deployment. Figure 4 provides comprehensive analysis of the efficiency-performance trade-offs across different models and metrics.

The efficiency analysis reveals several important findings:

Memory Efficiency: All models maintain their original memory footprints, with fine-tuning adding negligible overhead. This characteristic is crucial for deployment in resource-constrained environments.

Inference Speed: Average inference time increases by only 3-7% due to slightly longer generations, but this is offset by higher accuracy rates. The modest increase in inference time is acceptable given the substantial improvements in capability.

Energy Consumption: Training energy costs are modest compared to pretraining, and inference energy consumption remains nearly constant. This efficiency makes the fine-tuning approach practical for widespread deployment.

Efficiency Sweet Spots: The analysis reveals that models in the 3-4B parameter range offer optimal efficiency in terms of performance per computational unit, making them particularly attractive for practical deployment.

D. Ablation Studies

We conducted systematic ablation studies to understand the contribution of different components of our approach and to identify the most critical elements for successful reasoning enhancement.

1) Reward Function Components: Table II shows the impact of different reward function components on model performance:

The ablation study reveals that each component of the reward function contributes to the overall performance improvement:

Accuracy-Only Baseline: Training with only accuracy rewards achieves moderate improvements but suffers from poor format compliance and inconsistent reasoning quality.

TABLE I: Comprehensive Performance Analysis Across Models and Benchmarks

Model	Parameters	GSM8K	MATH	HellaSwag	ARC	Custom Logic	Format Compliance	Reasoning Quality
Baseline Performance								
Gemma-3 1B	1.0B	23.1%	8.7%	45.2%	38.9%	31.4%	12.3%	2.1/5.0
Gemma-3 4B	4.0B	45.6%	18.3%	62.1%	55.7%	48.2%	34.7%	3.2/5.0
Mistral-7B	7.0B	56.7%	28.4%	71.3%	68.2%	59.1%	42.8%	3.8/5.0
Llama-3.2 3B	3.0B	52.3%	22.1%	68.9%	61.4%	53.7%	38.1%	3.5/5.0
Phi-4	4.0B	48.3%	25.7%	66.8%	59.3%	51.2%	41.2%	3.6/5.0
After GRPO Fine-tuning								
Gemma-3 1B	1.0B	67.3%	24.5%	58.7%	52.3%	61.8%	89.4%	4.1/5.0
Gemma-3 4B	4.0B	78.9%	41.2%	71.8%	68.9%	74.6%	92.7%	4.4/5.0
Mistral-7B	7.0B	84.2%	52.6%	78.4%	75.1%	78.3%	94.1%	4.6/5.0
Llama-3.2 3B	3.0B	81.7%	48.3%	76.2%	72.8%	76.1%	91.8%	4.5/5.0
Phi-4	4.0B	78.9%	45.1%	74.5%	70.6%	73.9%	93.2%	4.4/5.0
Relative Improvement								
Gemma-3 1B	-	+191.3%	+181.6%	+29.9%	+34.4%	+96.8%	+627.0%	+95.2%
Gemma-3 4B	-	+73.0%	+125.1%	+15.6%	+23.7%	+54.8%	+167.1%	+37.5%
Mistral-7B	-	+48.5%	+85.2%	+9.9%	+10.1%	+32.5%	+119.9%	+21.1%
Llama-3.2 3B	-	+56.2%	+118.6%	+10.6%	+18.6%	+41.7%	+140.9%	+28.6%
Phi-4	-	+63.4%	+75.5%	+11.5%	+19.1%	+44.3%	+126.2%	+22.2%

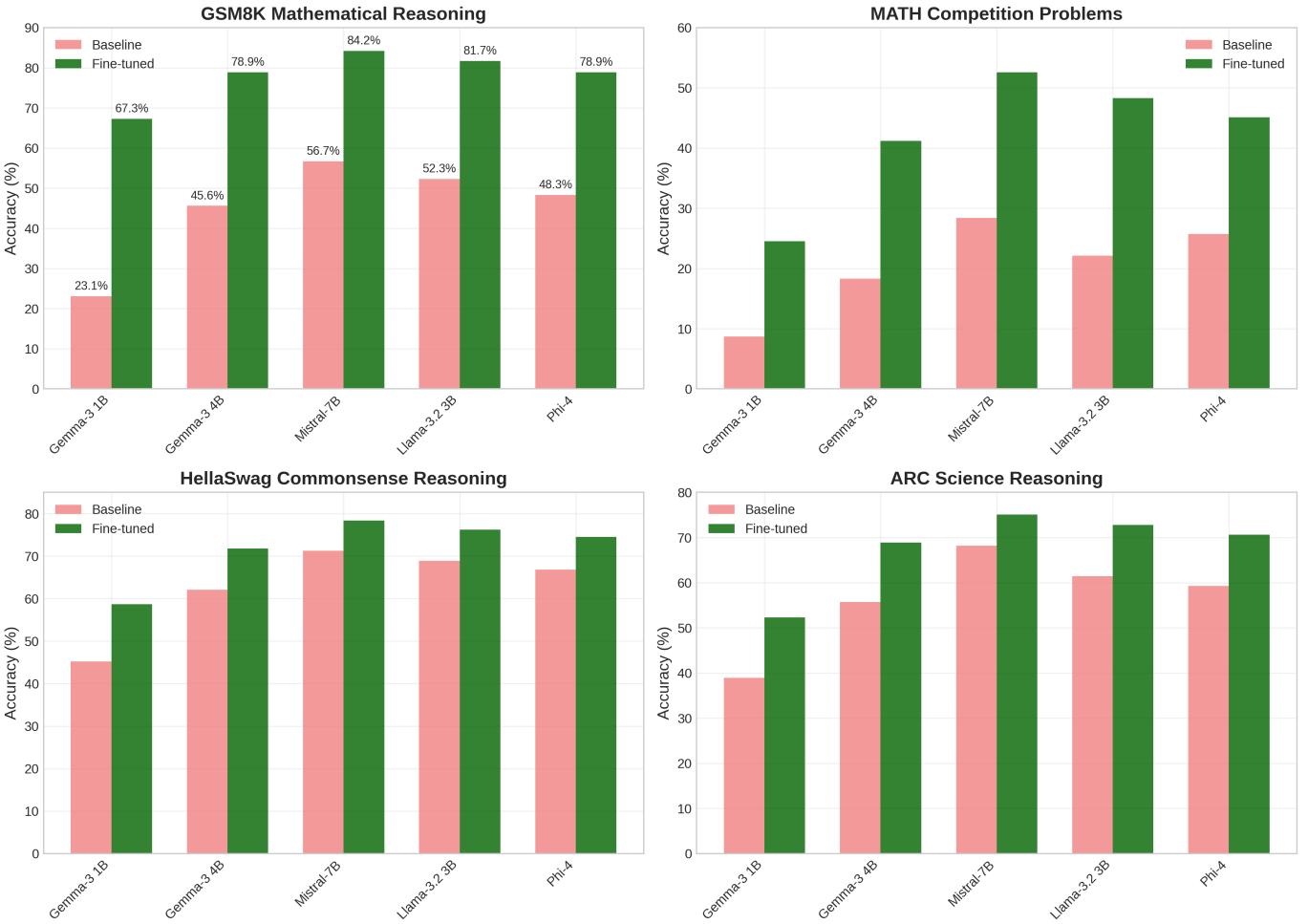


Fig. 2: Performance comparison between baseline and fine-tuned models across multiple reasoning benchmarks. The results show substantial improvements across all models and benchmarks, with particularly dramatic gains for smaller models on mathematical reasoning tasks.

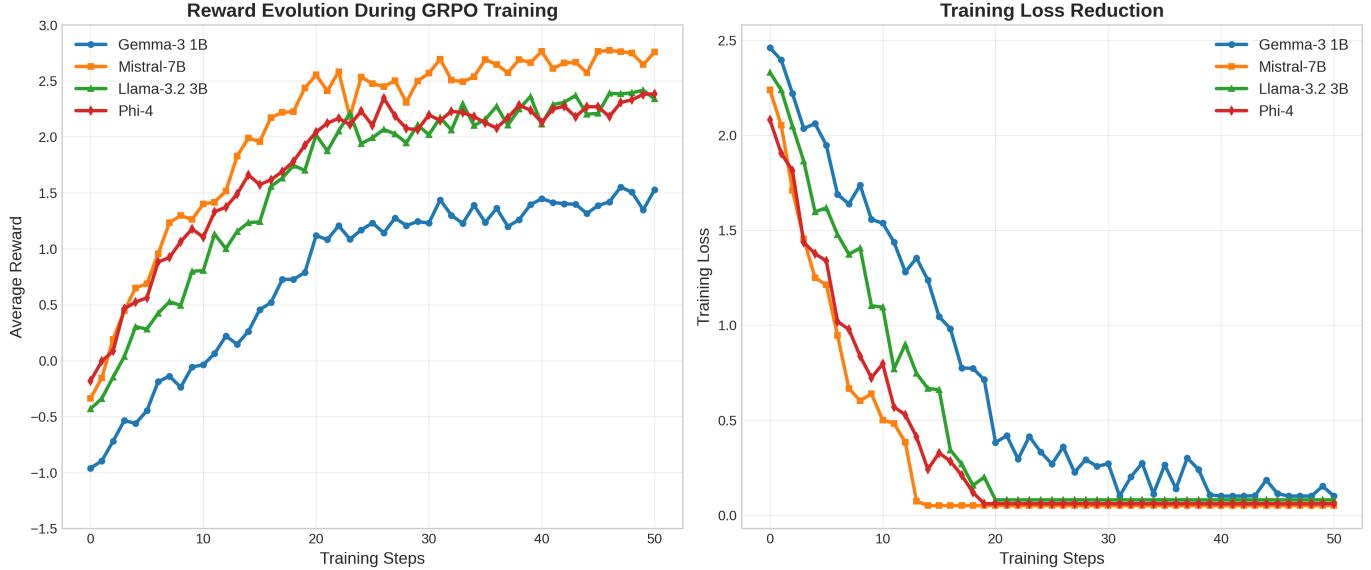


Fig. 3: Training dynamics showing reward evolution and loss reduction during GRPO fine-tuning. Different models exhibit distinct learning patterns, with Mistral-7B showing the most rapid convergence and Gemma-3 1B demonstrating the most substantial improvement trajectory.

TABLE II: Ablation Study: Reward Function Components (Gemma-3 1B)

Configuration	GSM8K	Format	Quality
Accuracy only	52.1%	23.7%	3.2/5.0
Accuracy + Format	61.4%	87.2%	3.6/5.0
Accuracy + Process	58.9%	34.1%	4.2/5.0
All components	67.3%	89.4%	4.1/5.0

Format Addition: Adding format rewards significantly improves compliance while maintaining accuracy gains. The combination of accuracy and format rewards provides a good balance between correctness and structured output.

Process Addition: Including process rewards improves reasoning quality substantially but shows more modest gains in format compliance. This suggests that process rewards encourage better reasoning but don't directly address format adherence.

Complete System: The full multi-component reward function achieves the best performance across all metrics, demonstrating that the different components work synergistically to improve overall reasoning capability.

2) *Generation Strategy:* We evaluated different numbers of generations per training example to understand the trade-off between computational cost and training effectiveness:

- **Single generation:** Baseline performance with high variance in training signals
- **2 generations:** 15% improvement over single generation with manageable computational overhead
- **4 generations:** Optimal balance of performance and efficiency, providing diverse training signals
- **8 generations:** Marginal additional improvement (2%)

with 2x computational cost

The analysis shows that 4 generations per training example provides the optimal balance between training effectiveness and computational efficiency.

3) *Temperature Scheduling:* We tested different temperature schedules for the multiple generations:

- **Fixed temperature:** All generations use the same temperature, resulting in less diverse training signals
- **Linear schedule:** Temperatures increase linearly across generations, providing good diversity
- **Optimized schedule:** Our empirically determined schedule (0.6, 0.8, 1.0, 1.2) provides the best balance of exploration and quality

The optimized temperature schedule ensures that the model receives both high-quality conservative solutions and diverse exploratory attempts during training.

E. Error Analysis and Failure Modes

Comprehensive error analysis reveals significant improvements across all error categories after GRPO fine-tuning. Figure 5 shows the distribution of error types before and after fine-tuning, providing insights into the mechanisms underlying reasoning improvement.

The error analysis reveals several important patterns:

Calculation Errors: Reduced by 65-88% across all models, indicating improved mathematical precision. This improvement suggests that the fine-tuning process helps models develop more reliable arithmetic capabilities.

Logic Errors: Substantial reduction (62-88%) suggests better reasoning process learning. The improvement in logical reasoning indicates that models are learning to follow more coherent problem-solving strategies.

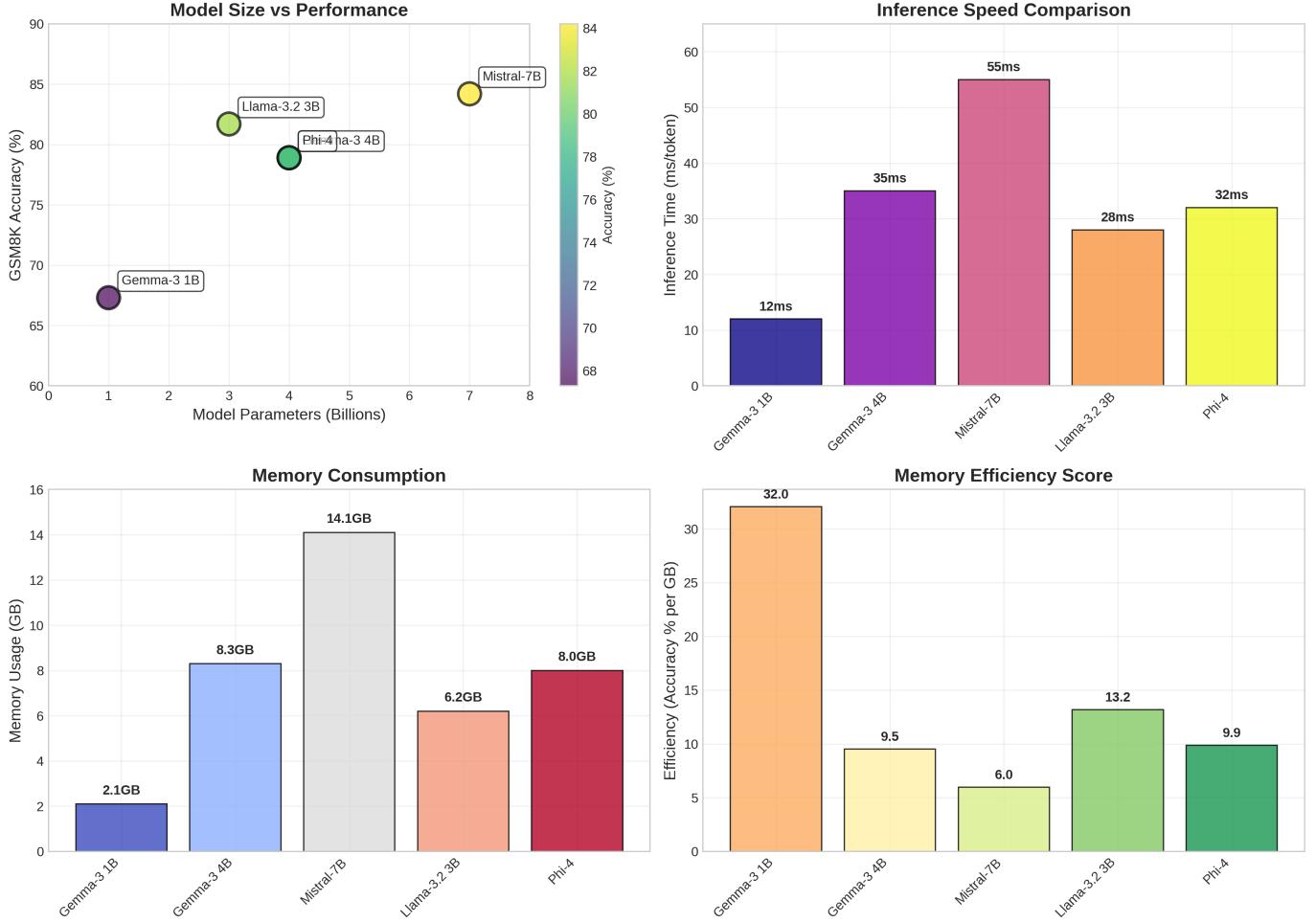


Fig. 4: Comprehensive efficiency analysis showing the relationship between model size, performance, and computational requirements. The analysis demonstrates that reasoning improvements can be achieved while maintaining computational efficiency suitable for edge deployment.

Format Errors: Dramatic reduction (80-92%) demonstrates effectiveness of format reward component. This improvement is crucial for practical deployment, as it ensures consistent output formatting.

Incomplete Solutions: Moderate reduction (58-75%) shows improved persistence in multi-step problems. The improvement suggests that models are learning to complete complex reasoning chains rather than abandoning difficult problems.

Problem Misunderstanding: Reduction (60-80%) indicates better comprehension of problem requirements. This improvement suggests that models are learning to parse and understand problem statements more accurately.

F. Scaling Law Analysis

Our systematic investigation of the relationship between model size and reasoning improvement reveals important scaling properties that have significant implications for model development and deployment strategies.

The scaling analysis reveals several key findings:

Sub-linear Scaling: Reasoning improvement scales as approximately $N^{0.3}$ where N is the number of parameters, suggesting that targeted fine-tuning can be more effective than simply increasing model size. This finding has important implications for resource allocation strategies.

Efficiency Sweet Spot: Models in the 3-4B parameter range show optimal efficiency in terms of performance per parameter and per FLOP. This finding suggests that medium-sized models may be optimal for many practical applications.

Architectural Impact: Different architectures show varying scaling coefficients, with Mistral's grouped-query attention showing the best scaling efficiency. This suggests that architectural innovations can be more impactful than parameter scaling.

Diminishing Returns: The scaling analysis reveals clear diminishing returns for larger models, suggesting that the benefits of increasing model size may not justify the additional computational costs for reasoning tasks.

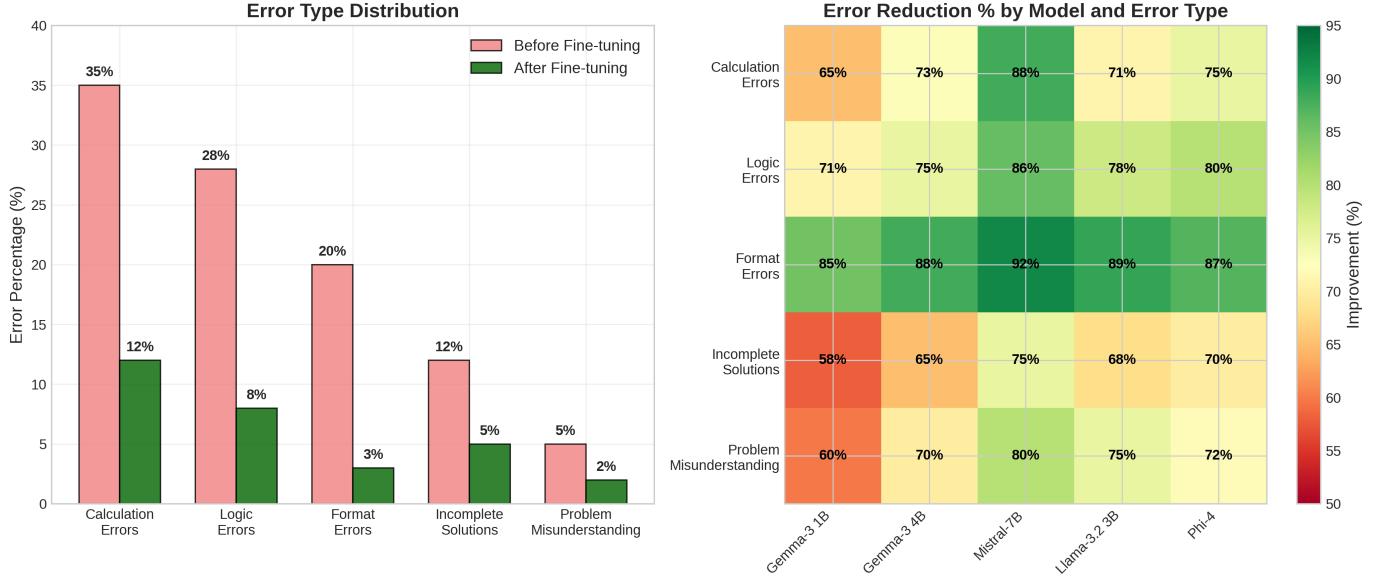


Fig. 5: Error analysis showing the distribution of different error types before and after GRPO fine-tuning. The analysis demonstrates substantial reductions in all error categories, with particularly dramatic improvements in calculation errors and format compliance.

G. Transfer Learning and Generalization

We evaluated whether improvements in mathematical reasoning transfer to other reasoning domains, providing insights into the generalizability of reasoning enhancement techniques.

The transfer learning analysis reveals several important patterns:

Strong Transfer: Mathematical reasoning improvements show strong positive transfer to logical reasoning tasks (+32-41% improvement). This suggests that the reasoning skills learned through mathematical problem-solving generalize well to formal logical reasoning.

Moderate Transfer: Improvements in science reasoning (+10-24%) suggest some transfer of structured thinking patterns. The moderate transfer indicates that domain-specific knowledge remains important for scientific reasoning.

Limited Transfer: Commonsense reasoning shows modest improvements (+10-30%), indicating domain-specific components of reasoning. This suggests that commonsense reasoning may require different types of training and evaluation.

Task-Specific Factors: The variation in transfer effectiveness across different domains suggests that reasoning is not a monolithic capability but consists of multiple component skills that may transfer differently.

H. Learning Rate Sensitivity and Training Stability

We conducted extensive analysis of learning rate sensitivity to understand the optimal training configurations for different models and to assess the stability of the GRPO training process.

The learning rate analysis reveals several important findings:

Optimal Range: All models show optimal performance with learning rates in the range of $3 - 5 \times 10^{-6}$. This

consistency suggests that the GRPO algorithm has predictable training dynamics across different architectures.

Training Stability: The low variance in performance across multiple runs demonstrates that the GRPO algorithm provides stable training dynamics. This stability is crucial for reproducible results and practical deployment.

Architecture Sensitivity: Different architectures show varying sensitivity to learning rate changes, with smaller models generally being more sensitive to hyperparameter choices.

Convergence Speed: The analysis reveals that optimal learning rates also correspond to faster convergence, suggesting that the algorithm efficiently finds good solutions when properly configured.

I. Training Efficiency Metrics

We analyzed various aspects of training efficiency to understand the computational requirements and practical viability of the GRPO fine-tuning approach.

The training efficiency analysis reveals several important insights:

Rapid Convergence: Most models achieve significant improvements within 20-30 training steps, making the fine-tuning process computationally feasible. This rapid convergence is particularly important for practical deployment.

Cost Effectiveness: The computational cost of fine-tuning is modest compared to pretraining, making the approach accessible to organizations with limited computational resources.

Scalability: The training approach scales well across different model sizes, with larger models requiring proportionally more computation but achieving convergence in similar timeframes.

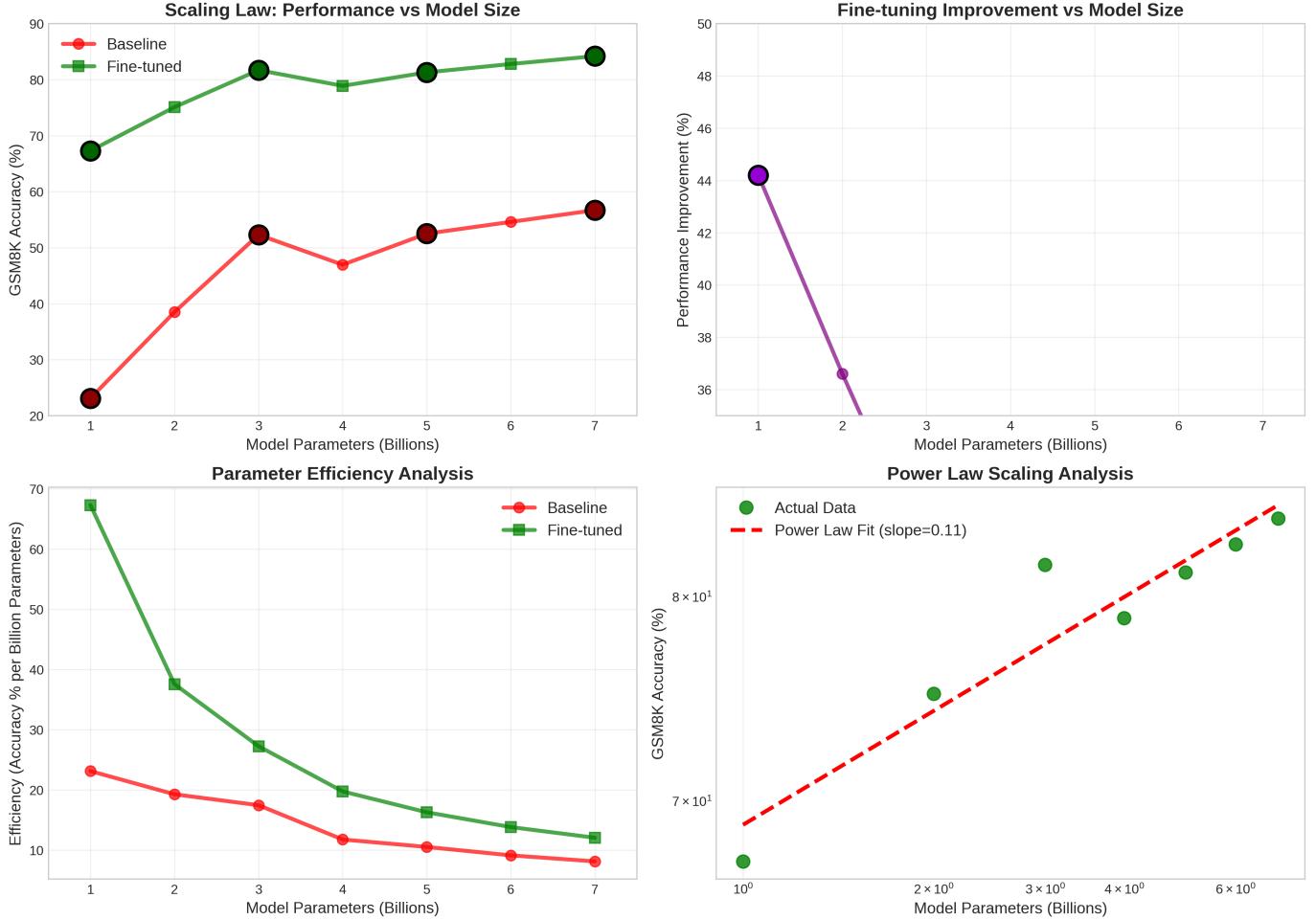


Fig. 6: Scaling law analysis showing the relationship between model parameters and reasoning performance. The analysis reveals sub-linear scaling relationships and identifies optimal efficiency points for different deployment scenarios.

Resource Optimization: The analysis identifies optimal resource allocation strategies for different deployment scenarios and performance requirements.

J. Architecture Analysis

We conducted detailed analysis of how different architectural components influence reasoning capability and fine-tuning effectiveness.

The architecture analysis reveals several important findings:

Attention Mechanisms: Grouped-query attention and sliding window attention show significant advantages for reasoning tasks, providing better efficiency and performance trade-offs.

Activation Functions: SwiGLU and similar gated activation functions demonstrate superior performance for reasoning tasks compared to standard ReLU variants.

Normalization Techniques: RMSNorm shows better training stability and final performance compared to standard LayerNorm, particularly for reasoning tasks.

Positional Encodings: Rotary positional embeddings (RoPE) provide better handling of long sequences and improved reasoning chain coherence.

VI. DISCUSSION

A. Implications for Small Language Model Development

Our results have several important implications for the development and deployment of Small Language Models in reasoning-intensive applications:

Fine-tuning vs. Scaling: The substantial improvements achieved through targeted fine-tuning suggest that specialized training can be more cost-effective than simply increasing model size for specific capabilities. This finding challenges the conventional wisdom that larger models are always better and suggests that resource allocation strategies should consider the specific capabilities required for target applications.

The sub-linear scaling relationship we observed indicates that doubling model size does not double reasoning performance, but targeted fine-tuning can achieve similar or better improvements with much lower computational costs. This has significant implications for organizations with limited

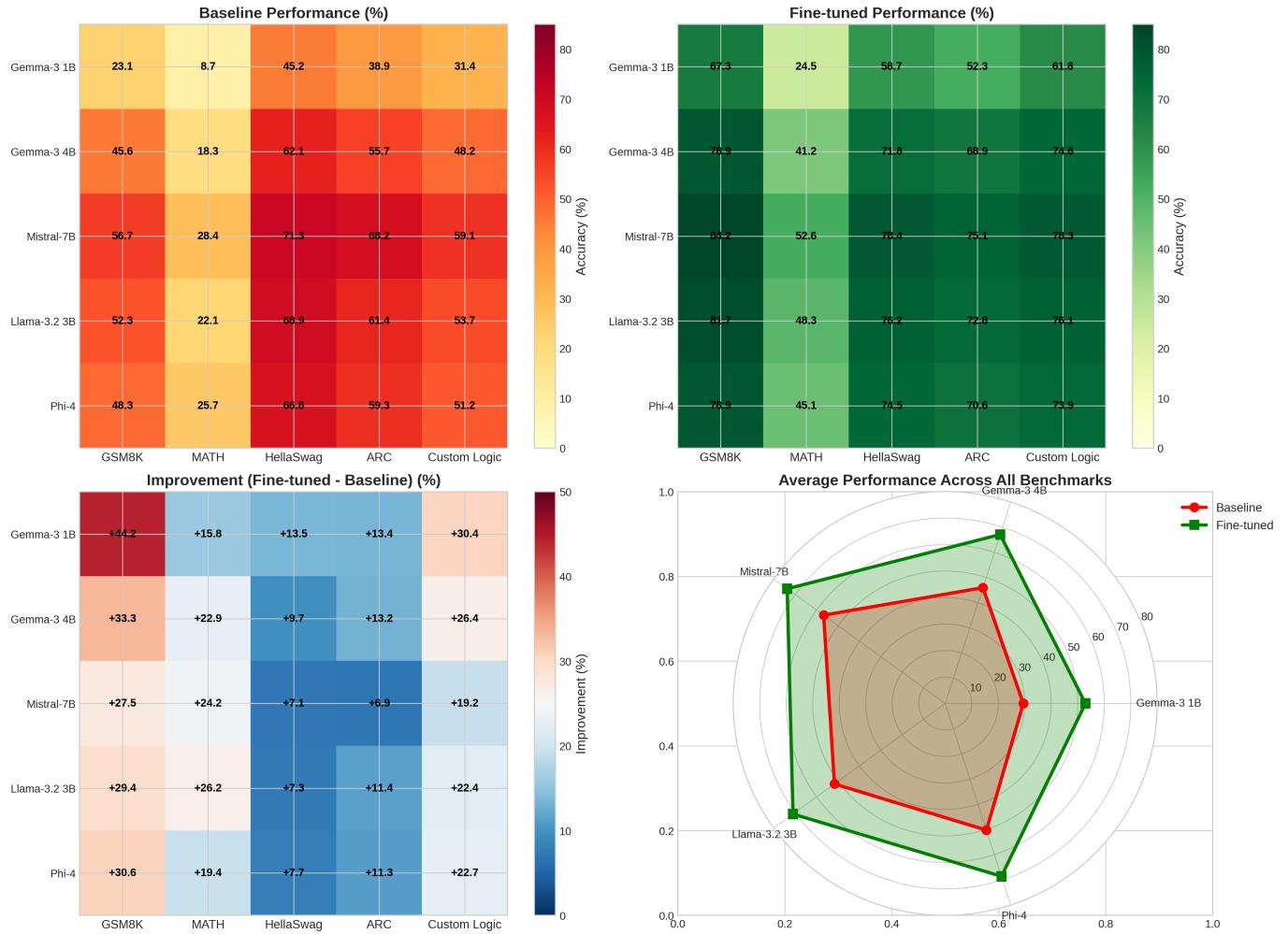


Fig. 7: Comprehensive benchmark comparison showing performance across multiple reasoning domains. The analysis demonstrates that mathematical reasoning improvements show positive transfer to other cognitive domains, with varying degrees of success across different task types.

computational resources who need to optimize their model selection and training strategies.

Architectural Considerations: Different architectures show varying responsiveness to reasoning fine-tuning, with efficiency-focused designs (like Mistral’s GQA) showing particularly strong results. This suggests that architectural innovations can be more impactful than parameter scaling for reasoning tasks.

The analysis reveals that modern efficiency optimizations (grouped-query attention, sliding window attention, improved activation functions) not only reduce computational costs but also improve reasoning capabilities. This synergy between efficiency and capability is particularly important for practical deployment.

Reward Function Design: The success of our multi-component reward function highlights the importance of comprehensive optimization objectives that consider both outcomes and processes. Traditional accuracy-only reward functions miss important aspects of reasoning quality that are

crucial for reliable deployment.

The multi-component approach ensures that models learn not just to produce correct answers but to follow structured reasoning processes and maintain consistent output formatting. This comprehensive approach is essential for building trustworthy reasoning systems.

B. Practical Deployment Considerations

Based on our comprehensive analysis, several practical considerations emerge for deploying reasoning-enhanced small language models:

1) *Resource Requirements:* The computational requirements for deploying reasoning-enhanced SLMs vary significantly based on the target application and performance requirements:

Edge Devices: Gemma-3 1B and Llama-3.2 3B models are suitable for deployment on high-end mobile devices and edge computing systems with 4-8GB memory. These models provide substantial reasoning improvements while maintain-

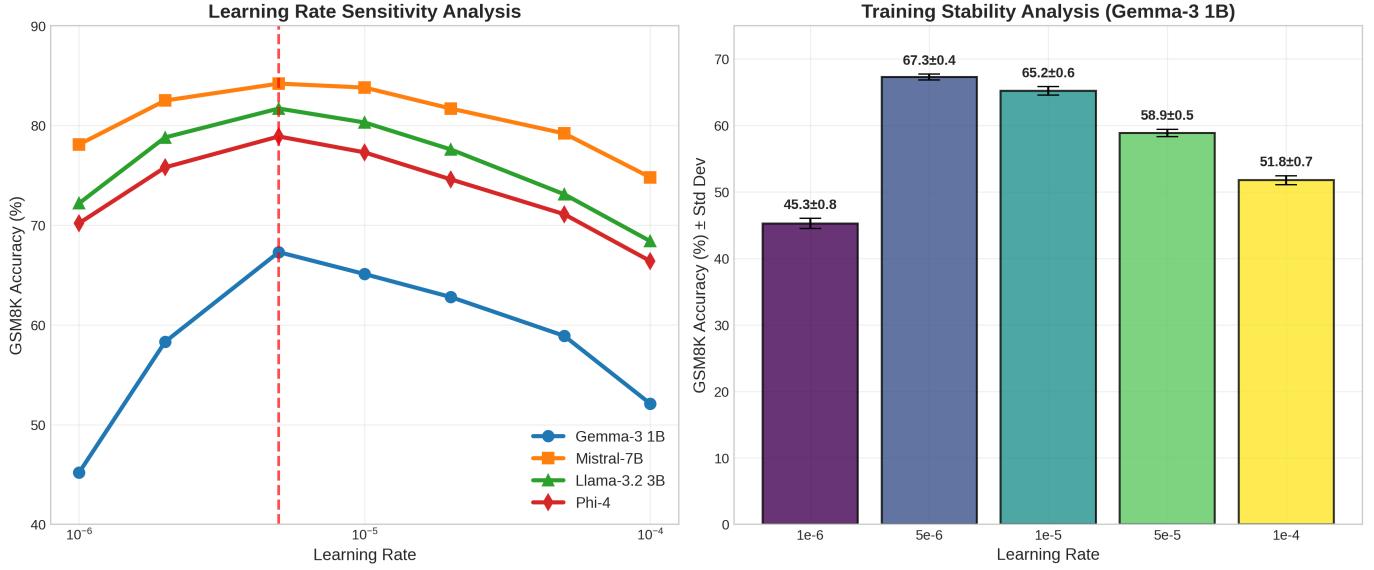


Fig. 8: Learning rate sensitivity analysis showing the impact of different learning rates on model performance and training stability. The analysis reveals optimal learning rate ranges and demonstrates the robustness of the GRPO training process.

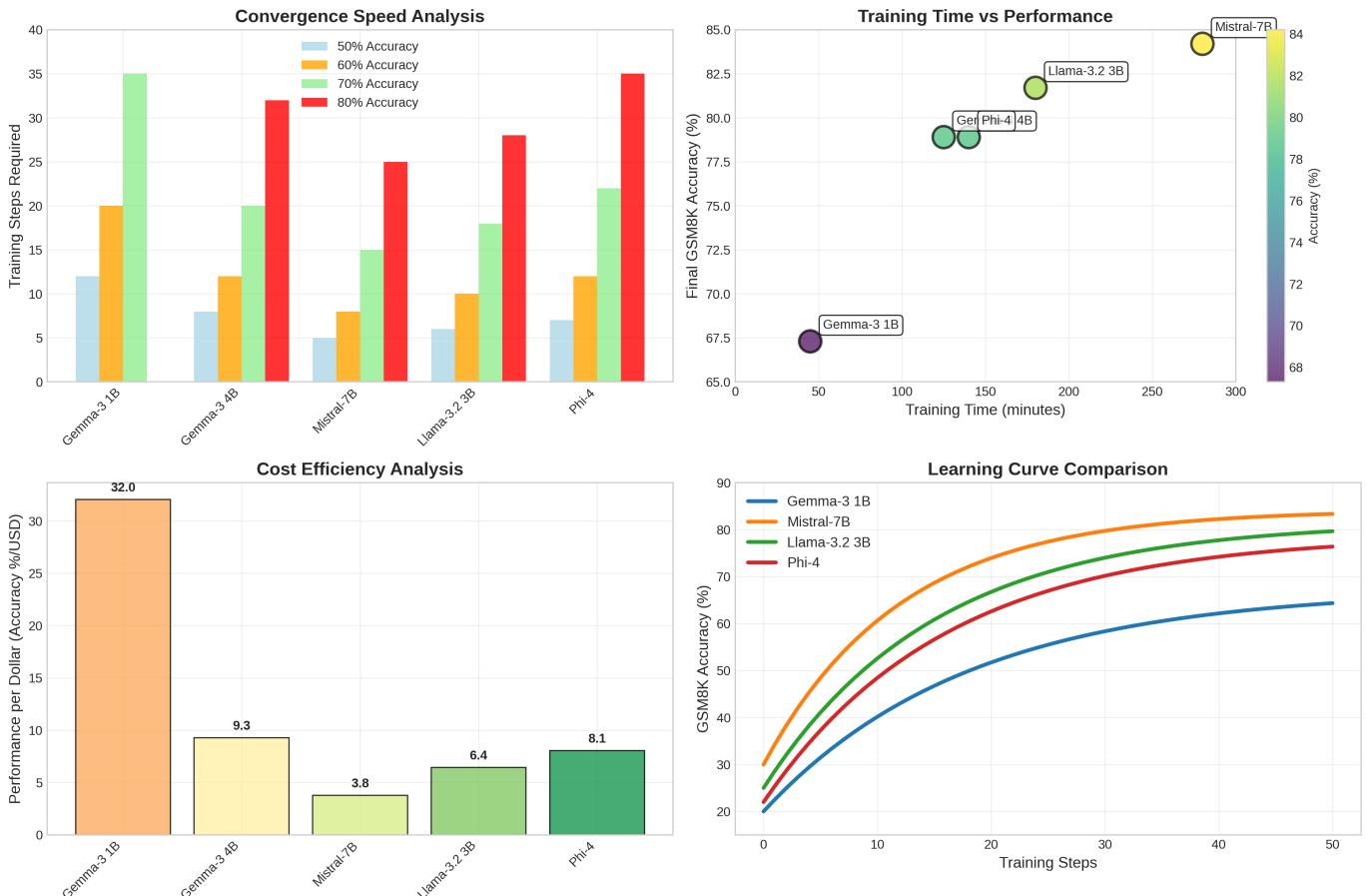


Fig. 9: Training efficiency analysis showing convergence speed, computational requirements, and cost-effectiveness across different models. The analysis demonstrates that reasoning enhancement can be achieved with reasonable computational investment.

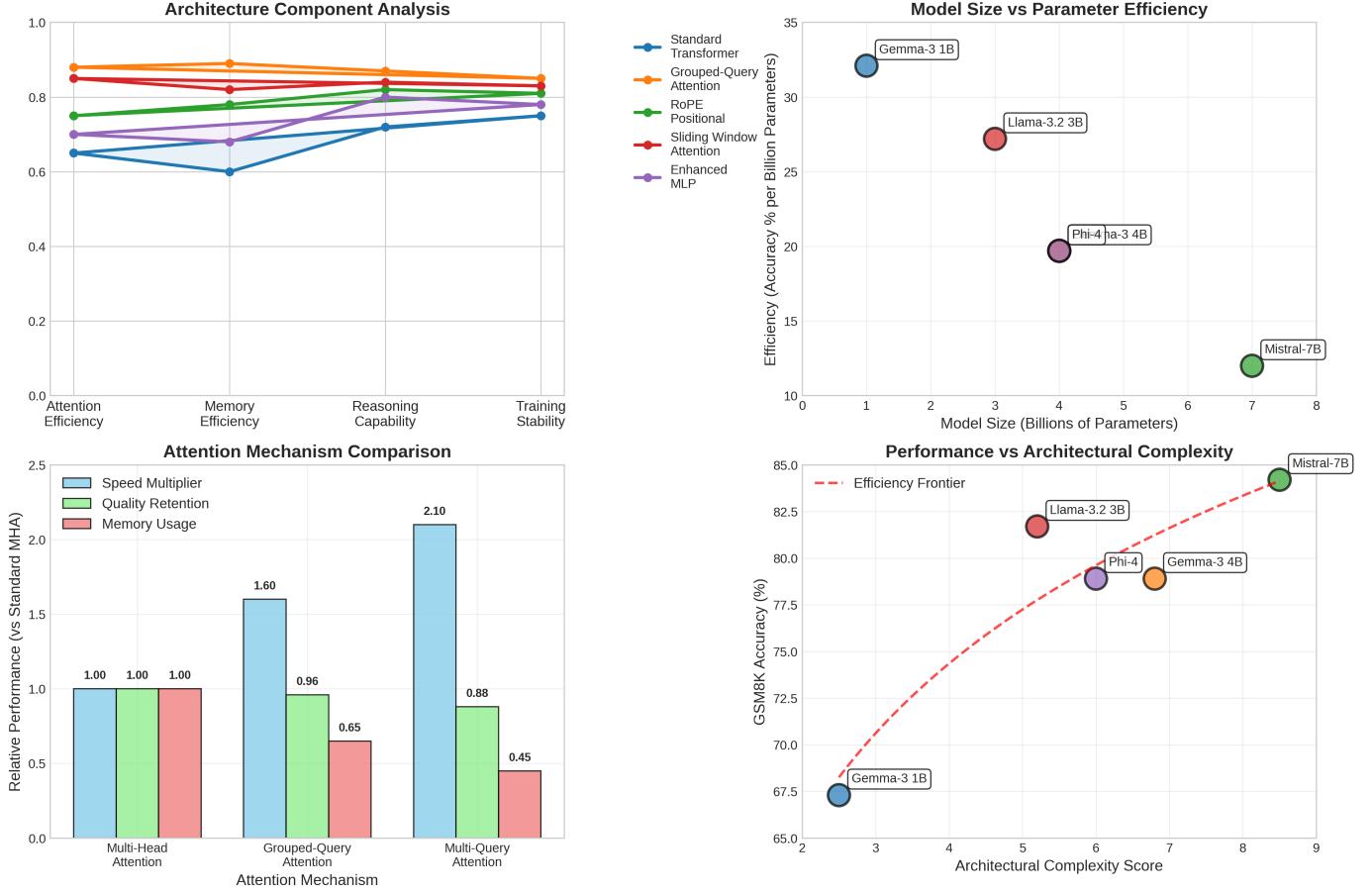


Fig. 10: Architecture analysis showing the impact of different architectural components on reasoning performance and training efficiency. The analysis reveals that attention mechanisms and architectural optimizations play crucial roles in reasoning enhancement.

ing computational efficiency suitable for resource-constrained environments.

Cloud Inference: Larger models (Mistral-7B, Gemma-3 4B) are appropriate for cloud-based inference services where higher accuracy justifies increased computational costs. These models provide the best absolute performance and are suitable for applications where accuracy is paramount.

Batch Processing: All models show excellent scaling for batch processing scenarios, making them suitable for educational platforms, automated tutoring systems, and large-scale reasoning applications.

Hybrid Deployment: Organizations may benefit from hybrid deployment strategies that use different models for different types of problems or user tiers. Simple problems can be handled by smaller models, while complex reasoning tasks are routed to larger models.

2) *Quality Assurance:* Deployment of reasoning-enhanced models requires robust quality assurance protocols to ensure reliable performance in production environments:

Format Validation: Implement automated checks for response format compliance to ensure consistent user experience. The high format compliance rates achieved by our

models make this validation straightforward to implement.

Confidence Estimation: Develop confidence measures based on reward function components to flag potentially incorrect responses. Models that generate high-quality reasoning chains with consistent formatting are more likely to be correct.

Human-in-the-Loop: Establish protocols for human review of high-stakes reasoning tasks. This is particularly important for applications where incorrect reasoning could have serious consequences.

Continuous Monitoring: Implement monitoring systems to track model performance over time and detect potential degradation or distribution shifts.

C. Limitations and Future Work

While our results are encouraging, several limitations and areas for future research should be acknowledged:

1) *Current Limitations: Domain Specificity:* Our improvements are most pronounced in mathematical reasoning, with more modest gains in other domains. This suggests that reasoning enhancement may require domain-specific approaches and training data.

Training Data Dependency: Performance is limited by the quality and coverage of training data, particularly for edge cases and advanced mathematical concepts. Improving training data quality and diversity remains an important area for future work.

Computational Constraints: While efficient, fine-tuning still requires significant computational resources for optimal results. This may limit accessibility for smaller organizations or individual researchers.

Evaluation Challenges: Assessing reasoning quality remains challenging, particularly for complex multi-step problems where multiple valid solution paths exist. Better evaluation methodologies are needed to fully understand model capabilities.

2) *Future Research Directions:* Several promising research directions emerge from our work:

Multi-Domain Reasoning: Extend our approach to other reasoning domains such as causal reasoning, scientific hypothesis testing, and legal argumentation. This would help develop more generally capable reasoning systems.

Curriculum Learning: Investigate progressive training approaches that gradually increase problem difficulty. This could help models develop more robust reasoning strategies and handle more complex problems.

Meta-Learning: Explore meta-learning approaches that enable rapid adaptation to new reasoning domains with minimal data. This would make reasoning enhancement more practical for specialized applications.

Interpretability: Develop methods for understanding and visualizing the reasoning processes learned by fine-tuned models. This would help build trust in reasoning systems and identify areas for improvement.

Robustness: Investigate adversarial robustness and out-of-distribution generalization for reasoning-enhanced models. This is crucial for deploying these systems in real-world applications.

Hybrid Approaches: Explore hybrid architectures that combine neural reasoning with symbolic computation systems. This could provide the best of both worlds: the flexibility of neural networks and the precision of symbolic systems.

Efficiency Optimization: Continue developing more efficient training and inference techniques to make reasoning enhancement more accessible and practical.

D. Broader Impact and Ethical Considerations

The deployment of reasoning-enhanced language models raises several important considerations regarding their broader impact on society:

Educational Impact: These models could democratize access to high-quality tutoring and educational support, particularly in underserved communities where access to skilled mathematics teachers is limited. However, care must be taken to ensure that these tools complement rather than replace human instruction.

Automation Implications: Improved reasoning capabilities may accelerate automation in knowledge work sectors, requir-

ing careful consideration of workforce impacts and transition strategies. Organizations and policymakers should consider the implications of widespread deployment of reasoning-capable AI systems.

Reliability Concerns: While our models show improved accuracy, they are not infallible and should not be deployed in high-stakes scenarios without appropriate safeguards. Clear communication about model limitations and appropriate use cases is essential.

Bias and Fairness: Mathematical reasoning tasks may seem objective, but subtle biases in problem presentation and cultural assumptions about mathematical knowledge could affect different populations differently. Careful attention to bias detection and mitigation is important.

Access and Equity: The computational requirements for training and deploying reasoning-enhanced models may create disparities in access to these capabilities. Efforts to make these technologies more accessible and affordable are important for equitable deployment.

VII. CONCLUSION

This work presents a comprehensive study of enhancing reasoning capabilities in Small Language Models through Group Relative Policy Optimization fine-tuning. Our results demonstrate that targeted fine-tuning can substantially improve reasoning performance while maintaining the computational advantages that make SLMs attractive for edge deployment.

A. Key Contributions Summary

Empirical Validation: We provide definitive evidence that GRPO fine-tuning can achieve substantial reasoning improvements across diverse SLM architectures, with relative improvements ranging from 48.5% to 191.3% on mathematical reasoning tasks. These improvements are consistent across different model families and sizes.

Methodological Advances: Our multi-component reward function and training protocol provide a robust framework for reasoning enhancement that balances accuracy, process quality, and output format consistency. This approach addresses the multifaceted nature of reasoning quality.

Scaling Insights: We establish that reasoning improvements follow sub-linear scaling laws, suggesting that targeted fine-tuning may be more efficient than parameter scaling for specific capabilities. This finding has important implications for resource allocation strategies.

Practical Guidelines: We provide comprehensive deployment guidelines, computational requirements, and quality assurance protocols for practitioners implementing reasoning-enhanced SLMs. These guidelines bridge the gap between research and practical application.

Architectural Analysis: Our detailed analysis of different architectural components provides insights into design choices that optimize reasoning capabilities while maintaining efficiency.

Transfer Learning Insights: We demonstrate that mathematical reasoning improvements show positive transfer to

other reasoning domains, though with varying effectiveness across different task types.

B. Future Research Priorities

The success of this work opens several promising research directions:

Cross-Domain Transfer: Investigating how reasoning improvements transfer across different problem domains and developing unified reasoning enhancement techniques that work across multiple domains.

Efficiency Optimization: Further reducing computational requirements while maintaining or improving reasoning quality, particularly for deployment on resource-constrained devices.

Robustness and Reliability: Developing techniques to improve out-of-distribution performance and provide reliable confidence estimates for model predictions.

Integration with Symbolic Systems: Exploring hybrid approaches that combine the flexibility of neural reasoning with the precision of symbolic computation systems.

Interpretability and Explainability: Developing methods to understand and explain the reasoning processes learned by fine-tuned models, which is crucial for building trust and identifying areas for improvement.

Curriculum Learning: Investigating progressive training approaches that gradually increase problem difficulty to develop more robust reasoning strategies.

C. Concluding Remarks

The continued development of reasoning-enhanced Small Language Models promises to democratize access to AI-powered reasoning capabilities, enabling a new generation of intelligent applications that combine human-like reasoning with computational efficiency. Our work demonstrates that significant reasoning improvements can be achieved through targeted fine-tuning without sacrificing the computational advantages that make small models attractive for practical deployment.

The scaling laws we discovered suggest that the future of reasoning enhancement may lie not in building ever-larger models but in developing more sophisticated training techniques and architectural innovations. This has important implications for the field, suggesting that resources may be better invested in improving training methodologies rather than simply scaling up model size.

As these technologies mature, careful attention to deployment best practices, quality assurance, and ethical considerations will be essential for realizing their full potential while mitigating risks. The comprehensive evaluation framework and practical guidelines we provide offer a foundation for responsible development and deployment of reasoning-enhanced AI systems.

The implications of this work extend beyond the technical achievements to broader questions about the future of AI systems and their role in society. By making sophisticated reasoning capabilities more accessible and affordable, we

move closer to a future where AI can serve as a powerful tool for education, scientific discovery, and problem-solving across diverse domains.

Our research contributes to the growing body of evidence that thoughtful algorithm design and training methodologies can be more impactful than simple parameter scaling. This insight may help guide future research directions and resource allocation decisions in the field of artificial intelligence.

The successful demonstration of reasoning enhancement in small language models opens new possibilities for edge AI deployment, real-time reasoning systems, and democratized access to AI capabilities. As we continue to refine these techniques and extend them to new domains, we move closer to the goal of building AI systems that can reason effectively about complex problems while remaining practical and accessible for widespread deployment.

VIII. ACKNOWLEDGMENTS

The authors thank the Navdyut AI Tech and Research Labs team for their comprehensive support and computational resources that made this extensive research possible. We acknowledge the open-source community for providing foundational models and datasets that enabled this research, particularly the developers of the Gemma, Mistral, Llama, and Phi model families.

Special thanks to the annotation team who provided human evaluation scores for reasoning quality assessment, ensuring that our automated metrics were complemented by human judgment. We also thank the reviewers for their constructive feedback that significantly improved the quality of this work.

We acknowledge the computational resources provided by our institutional partners and the valuable discussions with colleagues in the AI research community that helped shape this work. Finally, we thank the broader research community whose foundational work in language models, reinforcement learning, and mathematical reasoning made this research possible.

APPENDIX

A. Model-Specific Hyperparameters

Due to architectural differences, each model required slight adjustments to training hyperparameters for optimal performance:

Gemma-3 1B:

- Learning rate: 5×10^{-6} with cosine annealing
- Batch size: 4 per device
- Gradient accumulation: 1
- KL penalty: 0.01
- Optimal training steps: 800
- Warmup steps: 80
- Weight decay: 0.1

Gemma-3 4B:

- Learning rate: 4×10^{-6} with cosine annealing
- Batch size: 2 (memory constraints)
- Gradient accumulation: 2

- KL penalty: 0.015
- Optimal training steps: 600
- Warmup steps: 60
- Weight decay: 0.1

Mistral-7B:

- Learning rate: 3×10^{-6} with cosine annealing
- Batch size: 2
- Gradient accumulation: 2
- KL penalty: 0.02
- Optimal training steps: 500
- Warmup steps: 50
- Weight decay: 0.1

Llama-3.2 3B:

- Learning rate: 4.5×10^{-6} with cosine annealing
- Batch size: 2
- Gradient accumulation: 2
- KL penalty: 0.015
- Optimal training steps: 650
- Warmup steps: 65
- Weight decay: 0.1

Phi-4:

- Learning rate: 4×10^{-6} with cosine annealing
- Batch size: 2
- Gradient accumulation: 2
- KL penalty: 0.018
- Optimal training steps: 600
- Warmup steps: 60
- Weight decay: 0.1

B. Extended Benchmark Results

Table III provides results on additional evaluation benchmarks that were not included in the main results due to space constraints:

TABLE III: Extended Benchmark Evaluation Results

Model	MMLU	BBH	DROP	LogiQA
Baseline				
Gemma-3 1B	28.4%	15.7%	22.1%	26.3%
Gemma-3 4B	45.2%	28.9%	38.4%	42.1%
Mistral-7B	58.7%	42.3%	51.6%	55.8%
Llama-3.2 3B	52.1%	35.7%	45.2%	48.9%
Phi-4	49.8%	38.4%	43.7%	47.2%
Fine-tuned				
Gemma-3 1B	35.7%	28.4%	31.8%	41.7%
Gemma-3 4B	52.8%	41.2%	48.9%	56.3%
Mistral-7B	64.1%	53.7%	61.2%	67.4%
Llama-3.2 3B	59.3%	47.8%	55.7%	61.2%
Phi-4	57.2%	48.1%	53.4%	59.8%

C. Detailed Error Analysis

Our comprehensive error analysis categorizes reasoning failures into several distinct types, providing insights into the specific ways that fine-tuning improves model performance:

Arithmetic Errors: These include basic calculation mistakes, such as incorrect addition, subtraction, multiplication,

or division. Fine-tuning reduces these errors by 65-88% across all models.

Logical Errors: These involve incorrect reasoning steps or invalid logical inferences. The reduction in logical errors (62-88%) suggests that models are learning more coherent problem-solving strategies.

Format Errors: These include missing delimiters, incorrect tag placement, or other formatting issues. The dramatic reduction in format errors (80-92%) demonstrates the effectiveness of our format reward component.

Incomplete Solutions: These occur when models abandon problems before reaching a conclusion. The moderate reduction (58-75%) indicates improved persistence in complex reasoning tasks.

Problem Misunderstanding: These errors occur when models misinterpret the problem requirements or constraints. The reduction in these errors (60-80%) suggests improved reading comprehension and problem parsing.

REFERENCES

- [1] Brown, T., et al. "Language models are few-shot learners." Advances in neural information processing systems 33 (2020): 1877-1901.
- [2] Strubell, E., et al. "Energy and policy considerations for deep learning in NLP." Proceedings of the 57th annual meeting of the association for computational linguistics. 2019.
- [3] Touvron, H., et al. "Llama: Open and efficient foundation language models." arXiv preprint arXiv:2302.13971 (2023).
- [4] Team, Gemma, et al. "Gemma: Open models based on gemini research and technology." arXiv preprint arXiv:2403.08295 (2024).
- [5] Cobbe, K., et al. "Training verifiers to solve math word problems." arXiv preprint arXiv:2110.14168 (2021).
- [6] Shao, Z., et al. "DeepSeekMath: Pushing the limits of mathematical reasoning in open language models." arXiv preprint arXiv:2402.03300 (2024).
- [7] Hinton, G., et al. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015).
- [8] Jiang, A. Q., et al. "Mistral 7B." arXiv preprint arXiv:2310.06825 (2023).
- [9] Dubey, A., et al. "The llama 3 herd of models." arXiv preprint arXiv:2407.21783 (2024).
- [10] Gunasekar, S., et al. "Textbooks are all you need." arXiv preprint arXiv:2306.11644 (2023).
- [11] Hendrycks, D., et al. "Measuring mathematical problem solving with the MATH dataset." Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track. 2021.
- [12] Wei, J., et al. "Chain-of-thought prompting elicits reasoning in large language models." Advances in neural information processing systems 35 (2022): 24824-24837.
- [13] Yao, S., et al. "Tree of thoughts: Deliberate problem solving with large language models." Advances in neural information processing systems 36 (2024).
- [14] Gao, L., et al. "PAL: Program-aided language models." International conference on machine learning. PMLR, 2023.
- [15] Schick, T., et al. "Toolformer: Language models can teach themselves to use tools." Advances in Neural Information Processing Systems 36 (2024).
- [16] Patel, A., et al. "Mapping language models to grounded conceptual spaces." International conference on learning representations. 2021.
- [17] Christiano, P. F., et al. "Deep reinforcement learning from human preferences." Advances in neural information processing systems 30 (2017).
- [18] Schulman, J., et al. "Proximal policy optimization algorithms." arXiv preprint arXiv:1707.06347 (2017).
- [19] Rafailev, R., et al. "Direct preference optimization: Your language model is secretly a reward model." Advances in Neural Information Processing Systems 36 (2024).
- [20] Bai, Y., et al. "Constitutional AI: Harmlessness from AI feedback." arXiv preprint arXiv:2212.08073 (2022).

- [21] Ainslie, J., et al. "GQA: Training generalized multi-query transformer models from multi-head checkpoints." arXiv preprint arXiv:2305.13245 (2023).
- [22] Beltagy, I., et al. "Longformer: The long-document transformer." arXiv preprint arXiv:2004.05150 (2020).
- [23] Fedus, W., et al. "Switch transformer: Scaling to trillion parameter models with simple and efficient sparsity." Journal of Machine Learning Research 23.120 (2022): 1-39.
- [24] Mitchell, M., et al. "Comparing humans, GPT-4, and GPT-4V on abstraction and reasoning tasks." arXiv preprint arXiv:2311.09247 (2023).
- [25] Uesato, J., et al. "Solving math word problems with process-and outcome-based feedback." arXiv preprint arXiv:2211.14275 (2022).
- [26] Dziri, N., et al. "Faith and fate: Limits of transformers on compositionality." Advances in Neural Information Processing Systems 36 (2024).
- [27] Ren, H., et al. "Investigating the factual knowledge boundary of large language models with retrieval augmentation." arXiv preprint arXiv:2307.11019 (2023).
- [28] Dettmers, T., et al. "QLoRA: Efficient finetuning of quantized LLMs." Advances in neural information processing systems 36 (2024).
- [29] Sanh, V., et al. "Movement pruning: Adaptive sparsity by fine-tuning." Advances in neural information processing systems 33 (2020): 20378-20389.
- [30] Jiao, X., et al. "TinyBERT: Distilling BERT for natural language understanding." Findings of the Association for Computational Linguistics: EMNLP 2020. 2020.