# Cognitively Aligned Reasoning in AI via Classical Epistemic Structures using GRPO finetuning

## Abstract

We introduce a novel framework that enhances the reasoning capabilities of large language models by integrating classical Indian epistemology—Pramāṇa-śāstra—into both solution generation and evaluation. By leveraging six canonical sources of valid knowledge—Pratyakṣa (Perception), Anumāna (Inference), Upamāna (Comparison), Arthāpatti (Postulation), Anupalabdhi (Non-Apprehension), and Śabda (Testimony)—our approach annotates reasoning steps with explicit proof tags and enforces corresponding logical structures such as deductive, inductive, abductive, causal, and counterfactual reasoning. A custom verifier scores outputs by assessing alignment between proof types and reasoning modes, guided by a rigorously defined mapping informed by classical logic and contemporary epistemology.

To optimize model adherence to valid logical structures, we employ Group Relative Policy Optimization (GRPO), a reinforcement learning method that refines reasoning behaviors through groupwise reward comparisons, providing stable and efficient fine-tuning. Implemented with modern transformer architectures (e.g., Meta LLaMA) enhanced by LoRA fine-tuning and accelerated inference (vLLM), our system integrates epistemic proof types into the training loop, promoting logical consistency and interpretability.

This framework not only improves solution accuracy but also ensures transparency and auditable reasoning steps, facilitating pedagogical value and epistemic robustness. Our approach advances AI reasoning by grounding model decisions in provable epistemic sources, enabling models to select and apply appropriate reasoning schemes with greater cognitive alignment and reliability.

**Table of Contents**

## Background

Large language models (LLMs) have made significant strides in instruction following and reasoning, yet challenges remain in aligning their outputs with human intent and improving complex reasoning capabilities. Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022) successfully aligned models like InstructGPT (1.3B parameters) to outperform much larger models such as GPT-3 (175B parameters), improving truthfulness and reducing toxicity.

Mathematical reasoning remains a demanding domain for LLMs. Group Relative Policy Optimization (GRPO), a variant of PPO, was introduced by Shao et al. (2024) to optimize reasoning quality through groupwise reward comparisons. Their DeepSeekMath 7B model achieved 51.7% accuracy on the competition-level MATH benchmark and 64.2% on GSM8K, rivaling larger models like GPT-4 and outperforming Minerva 540B without external tools.

Recent theoretical work by Zhang et al. (2025) highlights limitations in Chain-of-Thought (CoT) prompting strategies. They show that fixed, universal CoT prompts (e.g., "think step by step") create high complexity in navigating the prompt space, hindering reasoning effectiveness. Their analysis demonstrates that task-specific prompting, guided by human supervision, outperforms unsupervised approaches and is essential to improve prompt-answer space navigation in reasoning tasks.

Together, these advances illustrate the importance of combining human-aligned reinforcement learning, optimized training methods like GRPO, and theoretically informed prompt design to enhance the reasoning and interpretability of LLMs.

## Problem Statement

This work aims to address the following research questions:

1. How can the principles of classical Indian epistemology, specifically Pramāṇa-śāstra, be systematically integrated into large language models to enhance their reasoning capabilities?

2. To what extent does annotating reasoning steps with explicit proof tags derived from canonical sources of valid knowledge improve the logical coherence and interpretability of model-generated solutions?

3. Can a custom verifier that evaluates the alignment between proof types and corresponding reasoning modes effectively enforce valid logical structures in language model outputs?

4. How can Group Relative Policy Optimization (GRPO) be utilized to optimize adherence to valid logical reasoning frameworks in language models via reinforcement learning?

5. Does incorporating epistemic proof types within the training process improve both the accuracy of model solutions and the transparency and auditability of their reasoning steps?

6. In what ways can grounding reasoning in provable epistemic sources enable models to reliably select and apply appropriate reasoning schemes, thereby achieving enhanced cognitive alignment and epistemic robustness?

## Research Questions

- How can formal proof types be integrated into the reasoning process of LLMs?

- Does the inclusion of epistemic proof structures improve model interpretability and accuracy?

- What are the effects of different reasoning types on model performance in math and common sense tasks?

## Aim and Objectives

- Aim: To develop a GRPO framework that incorporates epistemic proofs into LLM reasoning.

- Objectives:
    - Integrate proofs followed by reasoning types aligned with proofs for model training.
    - Implement LoRA fine-tuning and fast inference for efficient training.
    - Evaluate on benchmark datasets like GSM8K.
    - Develop reward functions enforcing logical ordering and proof-reasoning consistency.


## Significance of the Study

This work bridges classical epistemology and modern AI by integrating epistemic proof types into language model reasoning, enhancing interpretability and transparency. It introduces a proof-guided reinforcement learning method using GRPO to refine reasoning, improving both logical rigor and accuracy on complex reasoning tasks.

## Scope of the Study

This study focuses on formalizing six classical proof types alongside five reasoning types within language model reasoning. Implementation is carried out on the Meta LLaMA 3.1 8B model, (Low-Rank Adaptation) fine-tuning. The framework is applied specifically to mathematical problem solving and reasoning datasets. Multimodal perception and non-textual data modalities are excluded from this investigation.

## Research Methodology

1. Research Objectives
   - To design and implement a computational framework that incorporates classical Indian epistemic sources (Pramāṇas) and formal reasoning types into large language model (LLM) based problem-solving.

   - To fine-tune a state-of-the-art LLM (Meta's LLaMA 3) to generate stepwise reasoning explicitly annotated with proof types and reasoning forms.

   - To develop evaluation metrics and reward functions that assess the logical consistency, format adherence, and accuracy of the model's generated reasoning and final solutions.

2. Data Collection and Preparation
   - Dataset Selection: Use the GSM8K dataset, a benchmark of grade school math problems, to evaluate reasoning capabilities.

   - Data Annotation: Adapt dataset inputs to include system prompts that instruct the model to generate solutions following a prescribed format. This format requires:

     - Identification of the epistemic proof method (Perception, Inference, etc.)

     - Selection of the corresponding reasoning type (Deductive, Inductive, etc.)

     - Clear demarcation of the reasoning process and final answer via XML-like tags.

- Prompt : Create a detailed system prompt embedding classical proof concepts and reasoning classifications for guiding model output generation.

## 3. Model Setup and Fine-Tuning

- Model Choice: Use Meta LLaMA 3 (8B parameter) model leveraging:

    - 4-bit quantization for efficient memory utilization.

    - LoRA (Low-Rank Adaptation) fine-tuning with rank 32 for parameter-efficient adaptation.

    - PEFT (Parameter-Efficient Fine-Tuning) for faster convergence.

- Training Procedure:

    - Load pretrained LLaMA 3 model with tokenizer.

    - Fine-tune the model on the GSM8K dataset prompts emphasizing generation of explicit epistemic and reasoning annotations.

    - Use gradient checkpointing to manage memory during long sequence training.

## 4. Evaluation and Reward Functions

- Format Adherence:
    - Use regex-based reward functions (match_format_exactly, match_format_approximately) to check presence and correctness of *<start_working_out>...</end_working_out>* and *<SOLUTION>...</SOLUTION>* tags.

- Answer Accuracy:
    - Extract final numerical answers via regex and score exact matches and approximate numerical closeness.

    - Use domain-specific heuristics to reward near-correct answers.

- Logical Consistency:
    - Implement a check_logical_ordering function to verify correct pairing of proof types with allowable reasoning types (e.g., Perception → Inductive Reasoning).

- - - Penalize logically inconsistent or missing annotation structures.

- Composite Reward System:
  - Aggregate the above metrics during training to guide optimization, improving both answer correctness and epistemic/ reasoning coherence.

## Limitations and Future Work

- Limited to math problems; extension to other domains (e.g., legal, scientific reasoning) needed.

- Human expert evaluation of reasoning quality may be incorporated.

- Integration of multi-modal perception data as epistemic input could enhance realism.

## Requirements Resources

- Software: Python, Pydantic, HuggingFace datasets and transformers, unsloth library, vLLM.

- Hardware: GPU-enabled system with sufficient VRAM (e.g., 16GB+).

- Datasets: GSM8K for math reasoning tasks.

- Models: Meta LLaMA 3.1 8B with LoRA fine-tuning.

# Results

**Quantitative:**

- Track improvement in answer accuracy on GSM8K test splits.
- Measure increases in format adherence and logical consistency scores.

- Analyze correlations between proof-reasoning adherence and final answer correctness.

**Qualitative:**

- Inspect sample model outputs for coherent epistemic annotation.

- Validate logical soundness of reasoning chains per classical Indian epistemology.

| | GSM8K Question Summary | Answer | Proof | Reasoning | Correct Logical Pairing? | Score |
|---|---|---|---|---|---|---|
| 1 | House vs. Trailer loan monthly payment difference | 1500 | INFERENCE | DEDUCTIVE_REASONING | Yes | 0.5 |
| 2 | Janet: yearly piano vs. clarinet lesson cost difference | 1040 | INFERENCE | DEDUCTIVE_REASONING | Yes | 0.5 |
| 3 | Sabrina: total leaves needed for poultice | 29 | INFERENCE | DEDUCTIVE_REASONING | Yes | 0.5 |

| 4 | Average July 4th temp over 5 years in DC | 64 | PERCEPTION | INDUCTIVE_REASONING | Yes | 0.5 |
|---|---|---|---|---|---|---|
| 5 | Pages read by 3 people in 240 minutes | 328 | PERCEPTION | INDUCTIVE_REASONING | Yes | 0.5 |
| 6 | Martin rings big bell how many times | 36 | INFERENCE | DEDUCTIVE_REASONING | Yes | 0.5 |
| 7 | Bert: average words per daily crossword | 75 | INFERENCE | DEDUCTIVE_REASONING | Yes | 0.5 |

## Conclusion

Our proposed framework demonstrates the practical efficacy of integrating proofs and logic into large language model reasoning. By mapping proofs to corresponding reasoning types and enforcing logical coherence through structured annotations, the model is guided toward more interpretable and cognitively aligned outputs.

The empirical evaluation shows that the system reliably applies valid proof-reasoning combinations, as evidenced by a perfect consistency score across diverse examples. This confirms that our verifier and reasoning protocol not only align with epistemic theory but are also robust in applied settings.

Further, by incorporating Group Relative Policy Optimization (GRPO) during fine-tuning, the model internalizes these logical constraints as part of its decision-making process, leading to improved reasoning fidelity without sacrificing performance. The resulting system ensures both **epistemic transparency** and **logical rigor**, making it suitable for high-stakes domains such as education, legal reasoning, and scientific explanation.

## References

- **Training language models to follow instructions with human feedback** Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida
  https://doi.org/10.48550/arXiv.2203.02155

- **Why Does Your CoT Prompt (Not) Work? Theoretical Analysis of Prompt Space Complexity, its Interaction with Answer Space During CoT Reasoning with LLMs:** A Recurrent Perspective Xiang Zhang, Juntai Cao, Jiaqi Wei, Chenyu You, Dujian Ding
  https://doi.org/10.48550/arXiv.2503.10084

- **Classical epistemology texts (Nyāya logic).**

  https://doi.org/10.1007/978-1-4020-4425-0_9280

- **GSM8K dataset** - https://huggingface.co/datasets/openai/gsm8k