

Naga Rohan Kumar Bayya



### Agenda



01. The Problem: LLMs vs. Formal Reasoning

02. Project Goal & Objectives

03. **Background: Key Concepts** 

04. **Preliminary Work** 

05. Main Experiments & Results

06. Infrastructure & Demo

# The Problem: LLMs vs. Formal Reasoning

## The Problem: Formal Reasoning

- Large Language Models (LLMs) are trained on linguistic patterns, not formal reasoning.
- They produce plausible-sounding but often incorrect answers for math and logic.
- Struggle with multi-step deductions and symbolic manipulation.
- Lack the ability to verify the correctness of their own steps.



## Project Goal & Objectives

## To bridge the gap between language models and formal mathematical reasoning.



#### **Primary Focus**

 Enhance LLM capabilities in formal theorem proving, a comparatively under-studied domain.



#### **Preliminary Objectives**

- Improve mathematical problem-solving through specialized fine-tuning.
- Integrate symbolic tools (Mathics, Lean) to ensure correctness.



#### **Experiments**

 Investigate advanced architectures (MoE) and training methods (OBT, GRPO, SFT, etc).

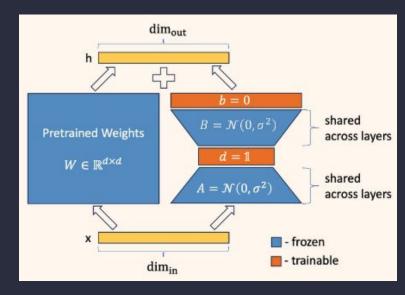
## Background: Key Concepts

## Parameter-Efficient Fine-Tuning (PEFT)

- Fine-tune models efficiently without retraining all parameters.
  - LoRA (Low-Rank Adaptation): Freezes original weights and trains small "adapter" matrices. Drastically reduces memory and compute.
  - DoRA (Decomposed Low-Rank Adaptation): An improvement upon LoRA to achieve better accuracy..

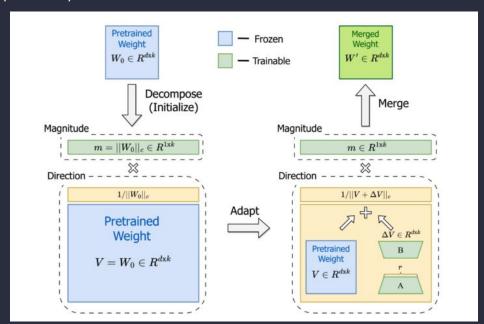
#### LoRA

- Freezes the original, large pre-trained weights (W)
- Injects two small, trainable "low-rank" matrices (A and B)
- $W_new = W + B.A$



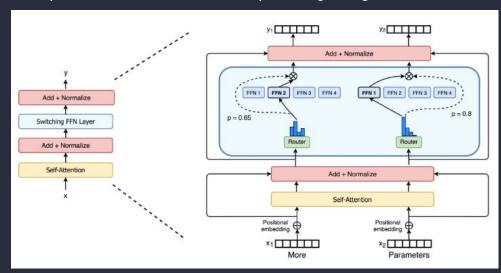
#### **DoRA**

- Achieves performance closer to full fine-tuning
- Splits original weight (W) in 2 parts: magnitude(m) and direction (D)
- LoRA is applied only to D



## Mixture-of-Experts (MoE) Architecture

- Analogy: A team of specialists, not a single generalist.
- A "gating network" routes each input token to a few specialized "expert" subnetworks.
- Allows for a massive number of total parameters with low computational cost.
- Theoretically more resistant to catastrophic forgetting.



## Post-Training Techniques

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

A prompt and several model outputs are sampled.

A labeler ranks

the outputs from best to worst.

This data is used

to train our reward model.



Explain the moon

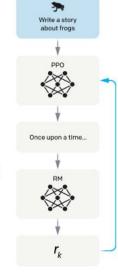
landing to a 6 year old

Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

The policy generates an output.

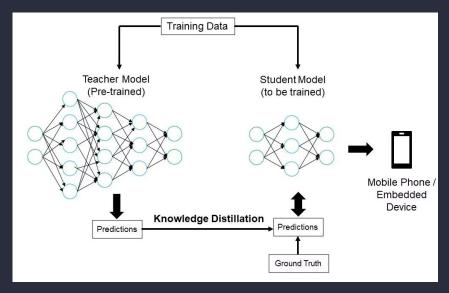


The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

#### **Model Distillation**

- Student-teacher technique for transferring knowledge to smaller model
- Working
  - Teacher Model: A large, highly capable (often proprietary) LLM.
  - Teacher generates a massive, high-quality dataset with CoT
  - Student is fine-tued on dataset
- Efficient, cost effective, and fast



### Chain-of-Thought (CoT) **Prompting**

- Break down complex problems into intermediate steps
- Improved reasoning and transparency

#### **Standard Prompting**

#### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

#### **Model Output**

A: The answer is 27.



#### **Chain-of-Thought Prompting**

#### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

#### **Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. <

## Tool-Augmented Reasoning

An LLM acts as an orchestrator, offloading tasks to reliable external tools.



LLM (Brain)

Parses the problem



Mathics (Calculator)

For precise symbolic math



Lean (Verifier)

For ensuring logical correctness of proofs.

#### **Mathics Tool**

```
@tool
def perform_math(expression: str) -> int:
    """ Use this tool when you need to calculate a complex mathematical expression.
    from mathics.session import MathicsSession
    session = MathicsSession(add_builtin=True, catch_interrupt=True)
    result = session.evaluate(expression).to_python()
    return result
```

> Entering new AgentExecutor chain...

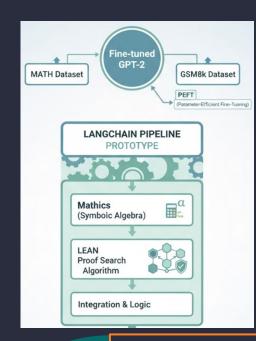
```
Invoking: `perform_math` with `{'expression': 'Integrate[Sin[x]^2, {x, 0, Pi/2}]'}`
System`Times[1/4, System`Pi] The integral of sin@(x) from 0 to pi/2 is equal to 1/4 * pi.
```

> Finished chain.

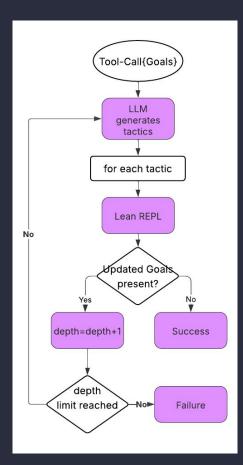
## Preliminary Work

## Small-Scale Validation Validate core hypotheses using a small model

- Fine-tuned GPT-2 on MATH and GSM8k datasets using PEFT
- Implemented LEAN Proof search algorithm
- Prototyped a LangChain pipeline to integrate Mathics and Lean proof search



### **Proof Search**



## Key Findings & Learnings Promising Results, Clear Limitations

- PEFT successfully improved mathematical formatting and basic reasoning.
- LLMs show learning capabilities when finetuned on different modalities
- Basic tool integration was not robust enough/feasible for conversational use.
- GPT-2's small size was a major bottleneck for complex problems.

## Transition to Main Experiments

- Move from GPT-2 to more capable models like Llama-3.2-3B and Phi-tiny-MoE
- Focus on theorem proving domain
- Implement more sophisticated training pipelines inspired by TheoremLlama and open-rl.

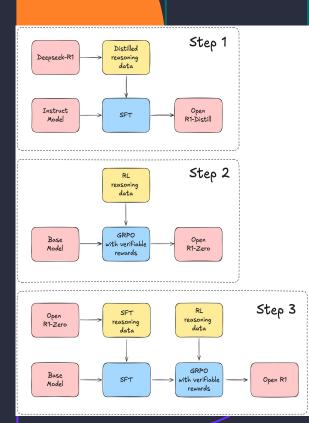
Quarte

## Experiment-1

Open-r1 + data supervision



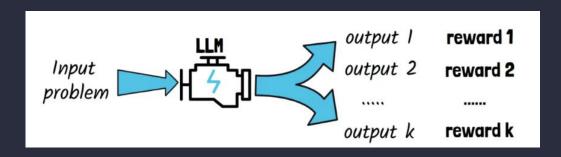
- Checkpoint-1 SFT (Supervised Fine-Tuning)
  with Distillation: Teach the model the basic task
  on text-guided data.
- Checkpoint-2 GRPO: Refine the model. It generates multiple proofs, a custom reward function picks the "best" one, and the model is trained only on that winner.



#### **GRPO**

- Group Relative Policy Optimization
- Samples a group of outputs <reasoning, answer>
- Assigns reward to outputs
  - Accuracy
  - Format

A conversation between User and Assistant. The user asks a question, and the Assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within <think> </think> and <answer> </answer> tags, respectively, i.e., <think> reasoning process here </think> <answer> answer here </answer>. User: prompt. Assistant:



#### **GRPO - Custom Reward Fn**

- Normalization
  - Extracts code block from model's completion and gold-standard proof
  - o Removes comments and normalizes whitespace
- Jaccard similarity between the tactics
- Penalties: sorry, by skip, etc
- Format bonus: <think> tags

## The Impact of Data Supervision Style

Does the style of training data matter for theorem proving?

- Regime 1 (Formal-Only): Fine-tuned on DeepSeek-Prover-V1 (Lean code only).
- Regime 2 (Text-Guided): Fine-tuned on DeepSeek-Prover-V2-generation (NL rationale + Lean code).

```
let total students := 900
let selected students := 10
let A students := 8
let B students := 2
let A correct := 7
let B correct := 8
let total correct := A correct + B correct
let total samples := Nat.choose
selected students total correct
let A_samples := Nat.choose A_correct
selected students
let B samples := Nat.choose B correct
selected students
let prob A := A samples / total samples
let prob B := B samples / total samples
simp [prob_A, prob_B, Nat.choose, factorial,
Nat.sub zero, Nat.mul zero, Nat.mul one]
```

```
First, we need to prove that (5^2 + 12^2 = 13^2). This is a
straightforward arithmetic problem:
1. Calculate (5^2 = 25).
2. Calculate \(12^2 = 144\).
3. Calculate (13^2 = 169).
4. Add \(25 + 144 = 169\), which is \(13^2\).
Thus, the equation holds.
### Step 1: Abstract Plan
1. Calculate (5^2 = 25).
2. Calculate \(12^2 = 144\).
3. Calculate (13^2 = 169).
4. Add \(25 + 144 = 169\) to verify that \(5^2 + 12^2 = 13^2\).
### Step 2: Lean 4 'have' Statements
```lean4
theorem
have pythag: 5 ^2 + 12 ^2 = 13 ^2 = 10
sorry
### Complete Lean 4 Proof
```lean4
theorem
have pythag : 5 ^2 + 12 ^2 = 13 ^2 = by
norm_num
<;> rfl
<;> simp [pow_two]
<;> norm_num
<;> rfl
```

## Results

		SFT-Non CoT		GRPO+SF
Theorem	Base Model	dataset	SFT-CoT dataset	Т
imp_chain	5	10	93	93
succ25_mem_of_zero_and _succ_closed	40	10	55	80
_succ_closed	40	10	<u> </u>	80
add_zero_eq_self	40	10	98	98
succ_add_one_eq_add_two	40	5	40	40
succ_pred_self	5	10	40	45
two_divides_double (Prompt 1)	10	5	5	94
two_divides_double				
(Prompt 2)	40	10	45	93

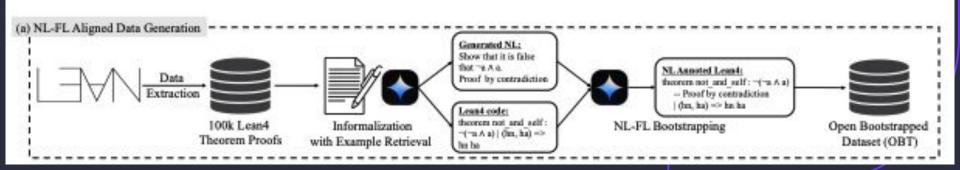
Max Points
60
25
15

## Experiment-2

TheoremLlama style Fine-tuning

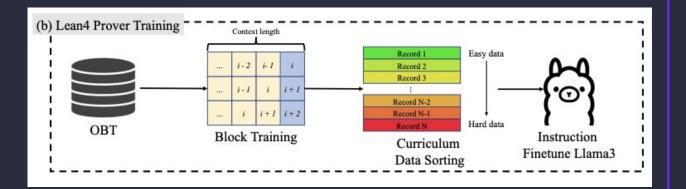
## Open Bootstrapped Theorems Dataset

- Lean 4 code with NL reasoning
- Extraction & Retrieval
- Informalization
- Bootstrapping



#### **Advanced Training with OBT Data**

- OBT: Open Bootstrapped Theorems
- Curriculum Sorting
  - Train on examples from easy to hard.
  - Prevents the model from being overwhelmed
- Few-Shot Block Training
  - Provide the model with 3 solved examples in the prompt before asking it to solve the current problem.
  - The model learns to generate more concise and human-like proofs.



## **Results**

Theorem Statement (Simplified)	Base Model Score	OBT-Finetuned Model Score
theorem imp_chain	35	50
theorem		
succ25_mem_of_ze		
ro	20	40
theorem		
add_zero_eq_self	5	55
theorem		
succ_add_one_eq_		
add	10	100
theorem		
succ_pred_self	5	20
theorem		
two_divides_double	5	25

Main Category	Max Points
Proof Score	60
Explanation Score	25
Format and Integrity Score	15

Quarter

## Experiment-3

MoE and Catastrophic Forgetting

#### Sequential vs Mixed

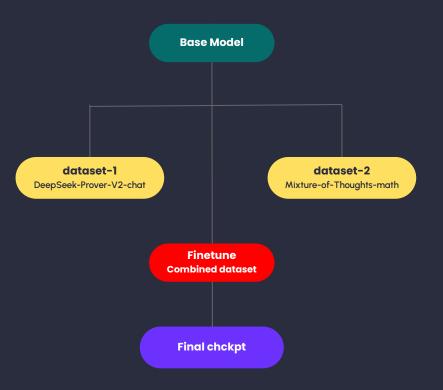
- Backbone: microsoft/phi-mini-moe-instruct
- Two fine tuning regimes
  - Sequential
  - Mixed
- Catastrophic forgetting



#### Sequential

### **Base Model Finetune** DeepSeek-Prover-V2-chat Intermediate chkpt **Finetune** Mixture-of-Thoughts-math Final chckpt

#### Mixed



## Results

Theorem	Base Model	Mixed FT	Sequential FT
imp_chain	35	84	5
succ25_mem_of_zero_and _succ_closed	20	60	5
add_zero_eq_self	5	97	0
succ_add_one_eq_add_two	0	96	30
succ_pred_self	5	92	5

Main Category	Max Points
Proof Score	60
Explanation Score	25
Format and Integrity Score	15

## Infrastructure: Evaluation Framework

## Infrastructure for Systematic Evaluation

#### Tech Stack



Ol. Frontend (React)
chat UI for side-by-side model comparison



#### **Backend (FastAPI)**

An orchestration layer to route requests to the O2correct model



#### Inference Layer (TGI on H200 GPUs)

O3. High-performance model serving using Text Generation Inference.



### The UI in Action

#### Conclusion

- Data-centric fine-tuning is a powerful strategy.
- Training pipelines like OBT and GRPO are key in niche domains like formal theorem proving.
- Training implications on catastrophic forgetting

# Any questions? Ask away!

## Thank you