# Inference Is All You Need

-

## What did Ilya see?

1. What does "Scaling Test Time Compute" mean?

2. Why does it matter for AI Engineers?

ellamind

# "Ask the same question a couple of times"

# "Ask the same question a couple of times"



Figure 3: Scaling inference time compute via repeated sampling leads to consistent coverage gains across a variety of model sizes (70M-70B), families (Llama, Gemma and Pythia) and levels of post-training (Base and Instruct models).

ellamınd
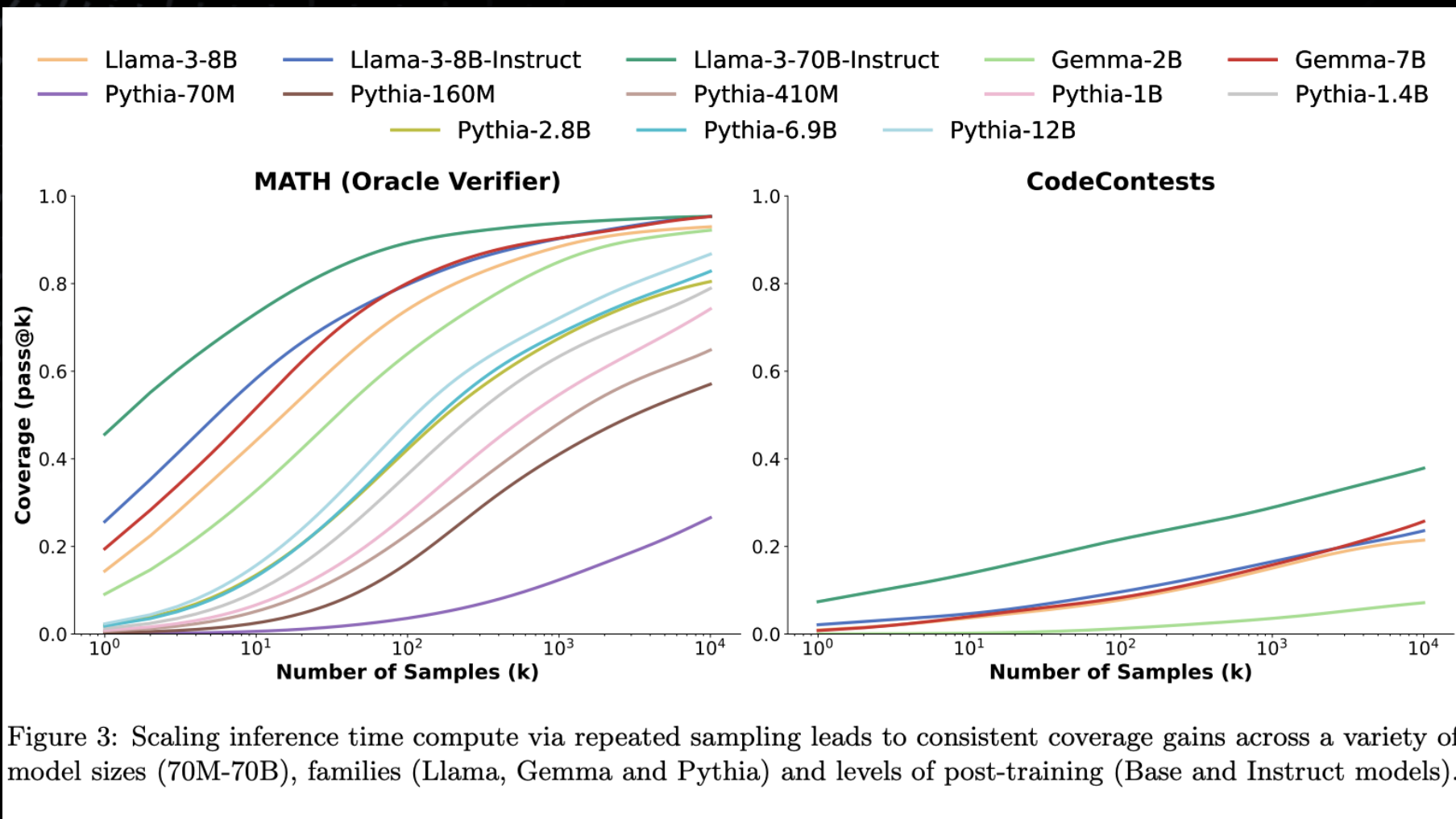
# Access to "Ground Truth"

## Golden Answer
## Unit Tests
## Formal Verfication
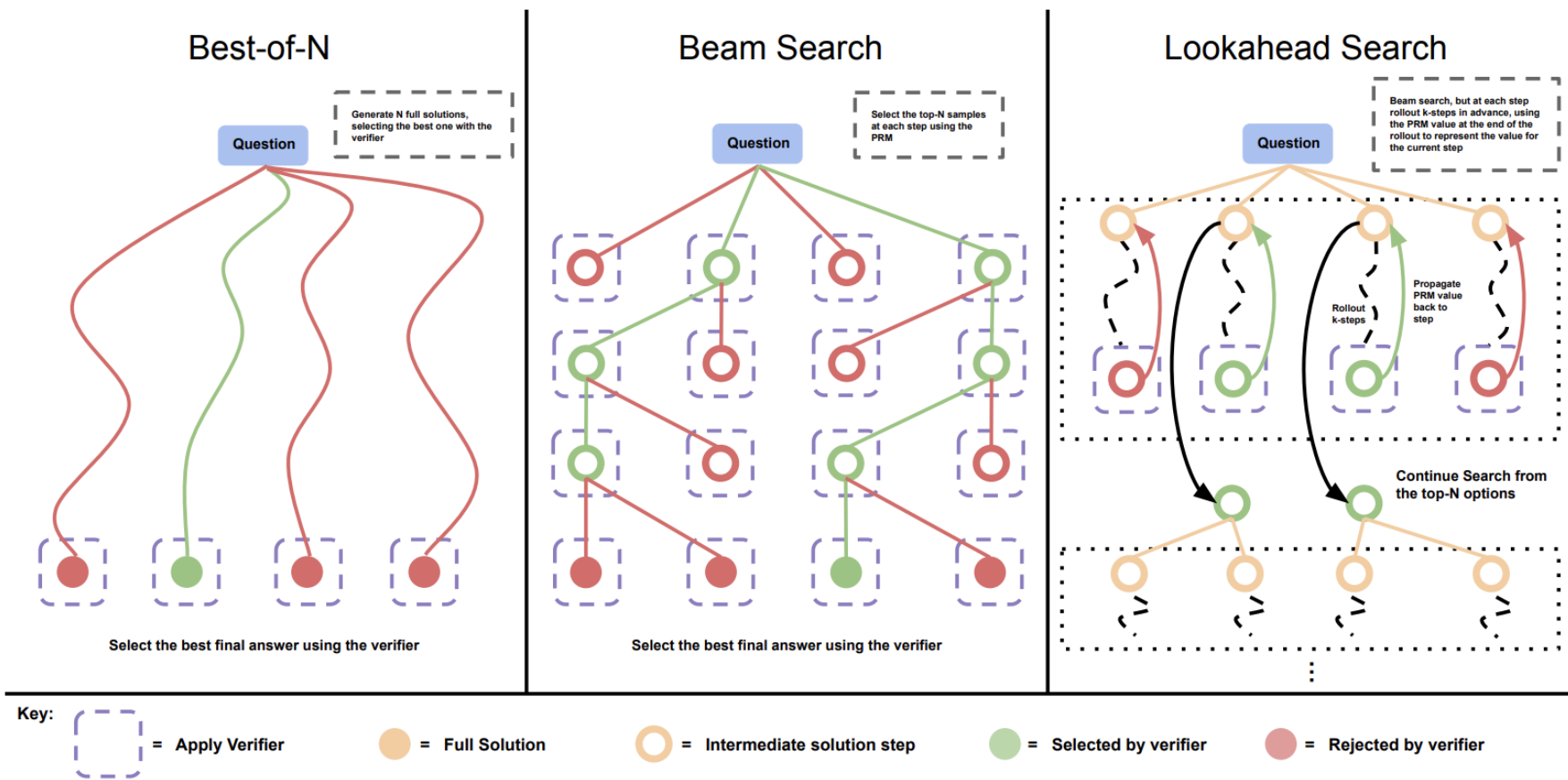## Any strict test

🎉 🎉 🎉

ellamind

# Reward Models!

ellamind

Figure 2 | *Comparing different PRM search methods.* **Left:** Best-of-N samples N full answers and then selects the best answer according to the PRM final score. **Center:** Beam search samples N candidates at each step, and selects the top M according to the PRM to continue the search from. **Right:** lookahead-search extends each step in beam-search to utilize a k-step lookahead while assessing which steps to retain and continue the search from. Thus lookahead-search needs more compute.

# AGI achieved internally? 🍓

## PRM800K

### Let's Verify Step by Step

Hunter Lightman[*]    Vineet Kosaraju[*]    Yura Burda[*]    Harri Edwards

Bowen Baker    Teddy Lee    Jan Leike    John Schulman    Ilya Sutskever

Karl Cobbe[*]

OpenAI

#### Abstract

ellamind

10

# Why does this matter for an AI engineer?

- **Sample on (idle) low VRAM hardware**
- **Inference is getting more optimized**
  - Groq, Cerebras
  - Quantization
  - Speculative Decoding
- **Collecting data for Instuction Finetuning can be harder than collecting preference data for RM**

ellamınd

**rasdani/inference-is-all-you-need**

ellamind

# ellamind

jan@ellamind.com

https://ellamind.com

ellamind GmbH
Konsul-Smidt-Straße 8p
28217 Bremen