# Scaling LLM Test-Time Compute Optimally can be More Effective than Scaling Model Parameters
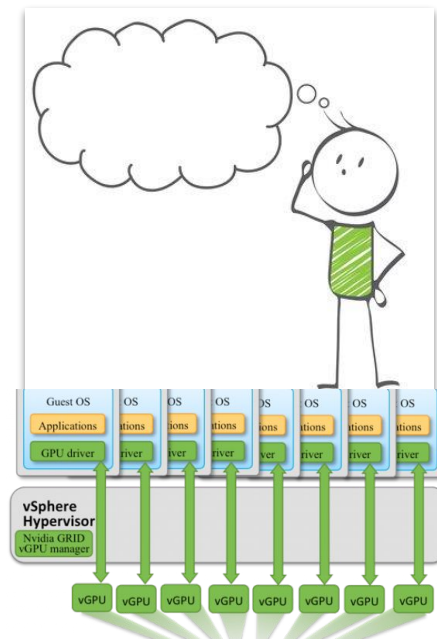
Elena Kote, Marvin Gao

September 23, 2024

# Question Definition

For humans, it's natural to think longer on difficult problems.

How to introduce this capability into LLMs?

Benefits:

- 💪 : Improve the performance on the testing phases

- ⏰ : Decrease the cost in training and pre-training phases;

- 💰 : Is possible to use extra, even spare, computation resource in testing steps to increase overall performance?

# Let's Think it Deeper

- 1. **Not only about** better performance on reasoning and inference;

- 2. The model will not need to be very large because of the improvement in testing steps, and we can **deploy smaller size** model in devices but reach the same, even better performance;

- 3. It's a path towards a **general self-improvement** LLMs without / weakly-with human supervision.

Please mathematically proof that in this world had chicken first and then had eggs

🤖 Thought for 8 seconds ⌃

**Mapping the problem**

The user seeks a mathematical proof that the chicken precedes the egg. This question traditionally falls under biology and evolution, but a mathematical logical argument might provide a solution.

**Revisiting assumptions**

I'm assembling a logical map for a mathematical proof to argue that chickens, not eggs, appeared first. Evolution suggests eggs preceded chickens, but I'll define parameters to challenge this idea.

**Identifying genetic origins**

I'm noting the genetic basis of a chicken, G_C, and how mutation during reproduction from proto-chickens led to the first true chicken, resulting from a non-chicken egg with a chicken embryo.

**Taking a closer look**

First, I'm proposing that a chicken egg is an egg laid by a chicken. This suggests the chicken came first, then laid the egg. I'm working through this logically and formally to clarify the sequence.

This explanation clarifies that the first chicken emerged due to genetic evolution and was born from an egg, which was laid by a non-chicken, disproving the common chicken-and-egg paradox.

**Logical Proof:**

1. **Ancestral Lineage:** Let's denote the immediate ancestors of the first true chicken as *proto-chickens*. These are organisms extremely similar to chickens but not classified as *Gallus gallus domesticus* due to slight genetic differences.

2. **Mutation Event:** Two proto-chickens mate and produce an offspring. Due to a specific genetic mutation during reproduction, this offspring has the exact genetic makeup to be classified as a chicken (C).

   - Mathematically, if $G_P$ represents the genetic code of proto-chickens and $G_C$ represents that of chickens, then:
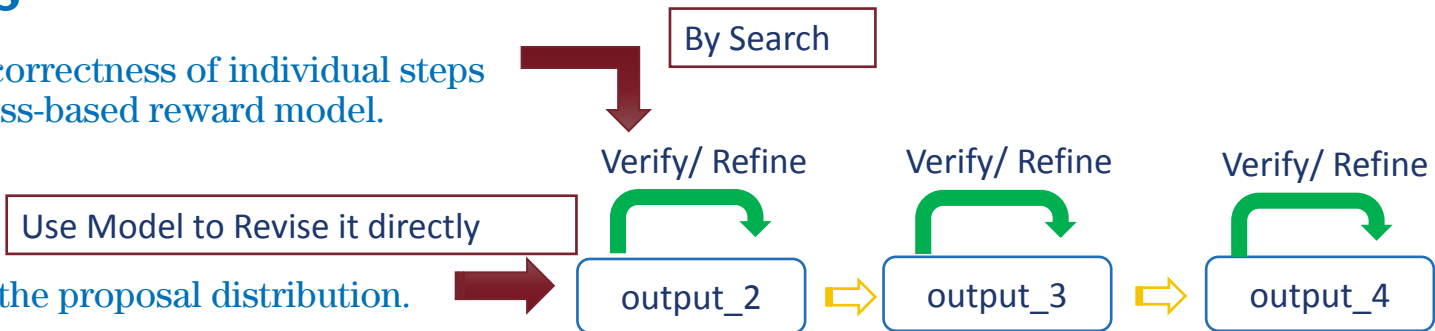
$$G_P + \Delta G = G_C$$

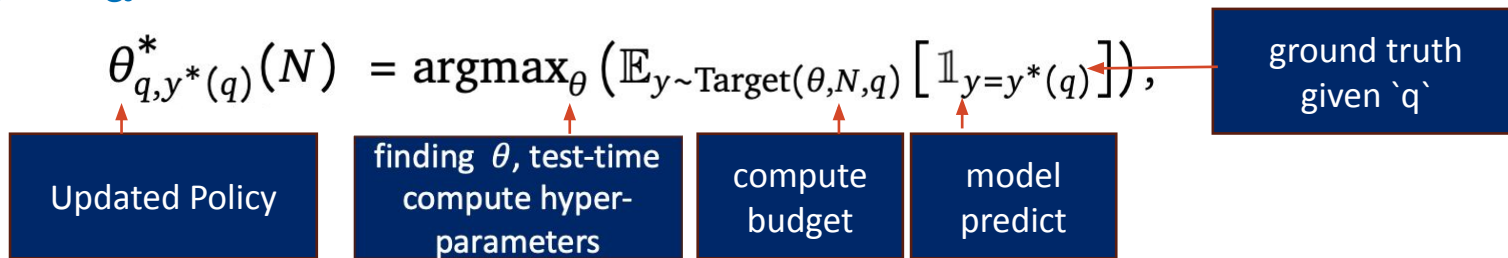   where $\Delta G$ is the genetic mutation occurring during reproduction.

3. **Egg Origin:** The egg from which this first chicken (C) hatches was produced by proto-chickens. According to our definition, since the egg was laid by proto-chickens (not chickens), it is **not** a "chicken egg" (E) but a "proto-chicken egg."

4. **First Chicken Egg:** The first true chicken egg (E) is the one laid by the first chicken (C). This egg contains the genetic material of a chicken and is laid by a chicken, satisfying our definition of a chicken egg.

5. **Conclusion:** Therefore, the sequence of events is:

   - Proto-chickens lay a proto-chicken egg ➔ First chicken (C) hatches.

   - First chicken (C) grows and lays the first chicken egg (E).

6. **Mathematical Representation:**

   - Let $t = 0$ be the time when the first chicken (C) exists.

   - At $t = -1$, only proto-chickens and proto-chicken eggs exist.

   - At $t = 1$, the first chicken egg (E) is produced by the first chicken (C).

# Methods

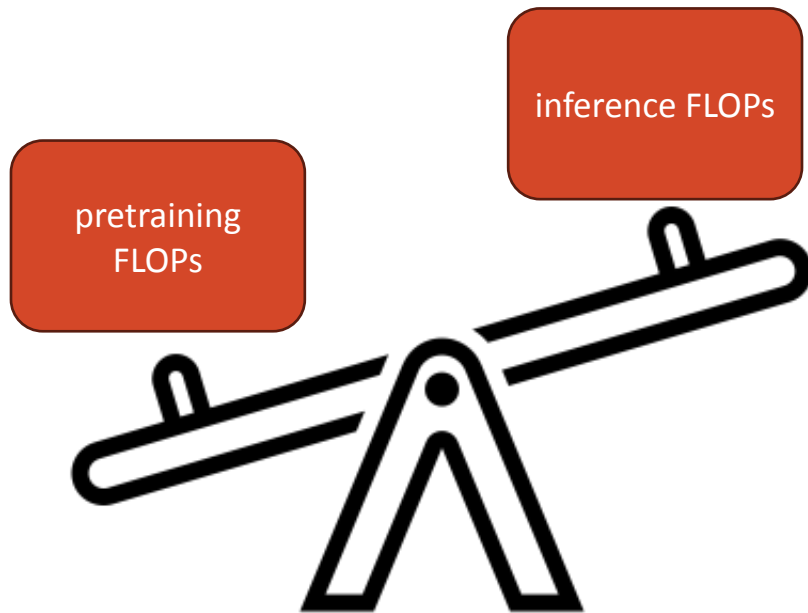- 1. Verify the correctness of individual steps using a process-based reward model.

By Search

Use Model to Revise it directly

- 2. Improving the proposal distribution.

Verify/ Refine    Verify/ Refine    Verify/ Refine

output_2    output_3    output_4

- What to know: **test-time compute-optimal scaling strategy**

$$\theta^*_{q,y^*(q)}(N) = \text{argmax}_\theta \left( \mathbb{E}_{y\sim\text{Target}(\theta,N,q)} \left[ \mathbb{1}_{y=y^*(q)} \right] \right),$$

ground truth given `q`

Updated Policy

finding $\theta$, test-time compute hyper-parameters
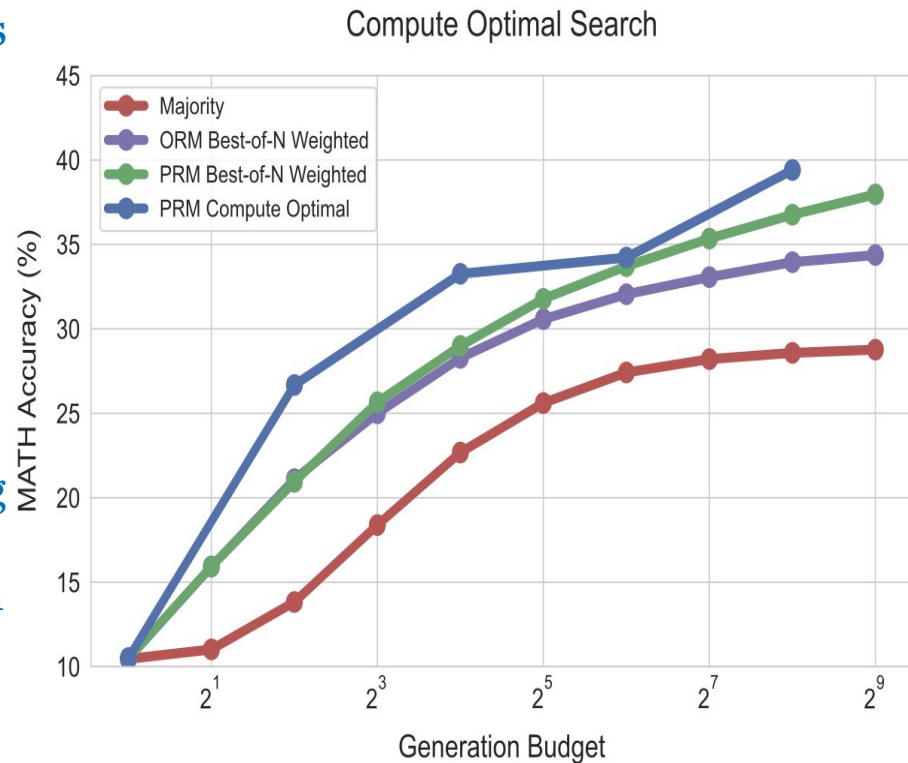
compute budget

model predict

# Efficiency Comparison: FLOPs

- comparison between a smaller model with additional test-time compute and pretraining a 14x larger model.

- scaling up test-time computation can be more preferable to scaling up pretraining, but not for all problems;

- For more difficult prompts, it will be less efficient when applying test time scaling up.

inference FLOPs

pretraining FLOPs

# Improve Accuracy by a Revisor

ORM: O

    Ev

PRM: Pr

    Ev

Which is

**ORM**
Offline Reward
Offline Reward Model

**Process Reward Model**



ORM Verses PRM

Figure 14 | We compare PRM and ORM models finetuned from PaLM 2-S* in a best-of-N evaluation. We use the PaLM 2-S* base LM to sample outputs, using a few-shot prompt. We see that the PRM greatly outperforms the ORM at a larg number of samples.

a single resuld-based answer

JOHNS HOPKINS
UNIVERSITY

# Experiment Settings and Methods

Difficulty of q: assign q to one of five difficulty levels. Use this as input to get optimal parameters settings. Getting from a LLM

$$\theta^*_{q,y^*(q)}(N) = \text{argmax}_\theta \left( \mathbb{E}_{y \sim \text{Target}(\theta,N,q)} \left[ \mathbb{1}_{y=y^*(q)} \right] \right),$$

ground truth given `q`

**Updated Policy**

finding $\theta$, test-time compute hyper-parameters

There are several different hyper parameters: e.g

**Search Algorithms**

OR

**Steps**

**Steps**

- **Dataset:** MATH benchmark, which consists of high-school competition level math problems.

- **Models:** PaLM 2-S*

- # Method ½: Scaling Test-Time Compute Via Verifiers

**Step1: Training Process Reward Model: supervise** PRMs without human labels, using estimates of per-step correctness obtained from running Monte Carlo rollouts from each step in the solution, corresponding to value estimated.

**Step2:** Score a set of solutions using PRM.



Compute Optimal Search

Legend:
- Majority
- ORM Best-of-N Weighted
- PRM Best-of-N Weighted
- PRM Compute Optimal

MATH Accuracy (%) vs Generation Budget ($2^1$ to $2^9$)

# PRM Search Methods



Rollout

Process Reward Model

Select By Verifier and Proceed

# Test-time Scaling for Search with Verifiers

- Sweep various search settings and hyperparameters

- Results
  - Smaller budget: beam search better
  - Greater budget: best-of-n
  - Lookahead underperforms for all budgets



Comparing PRM Search Methods

# Comparing search methods by question difficulty

- Four bars are increasing test-time compute
- Easy (1-2) : best of N better
  - Beam over-optimizes, exploits PRM signal
- Medium (3-4): beam better
- Intuition
  - On easy questions verifier will mostly be correct, by further optimizing we amplify spurious features thus degrading
  - On difficult questions model less likely to sample correct answer, search helps find it



Comparing Beam Search and Best-of-N by Difficulty Level

# Most powerful search methods underperform with high budget on easy questions

If 4 daps = 7 yaps, and 5 yaps = 3 baps, how many daps equal 42 baps?

```
So 4 daps is equal to 7 yaps, so that means that 1 dap is equal to 7/4 yaps or 1.75 yaps. [score = 0.21]
And 5 yaps is equal to 3 baps, so that means that 1 yap is equal to 3/5 baps or 0.6 baps. [score = 0.32]
And we can plug these values in to find the answer. [score = 0.31]
So 42 baps is equal to 42/0.6 yaps or 70 yaps. [score = 0.60]
And 70 yaps is equal to 70/1.75 daps or 40 daps. [score = 0.87]
So the answer is $\boxed{40}$. [score = 0.92]
And we're done! [score = 0.91]
And that is our final answer. [score = 0.90]
 [score = 0.91]
So 40 daps is equal to 42 baps. [score = 0.90]
And we're done! [score = 0.91]
And that is our final answer. [score = 0.90]
 [score = 0.64]
 [score = 0.78]
#### 40 [score = 0.83]
```

- Search causes model to generate low-information repetitive steps at end of solution

- Over-optimizing search can result in overly-short 1 or 2 step solutions

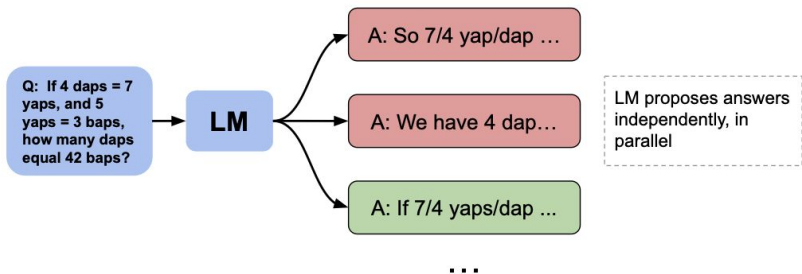# Comparing against baselines



Compute Optimal Search

- In both high and low generation budgets compute-optimal scaling can outperform Best-of-N by using **4x less compute**

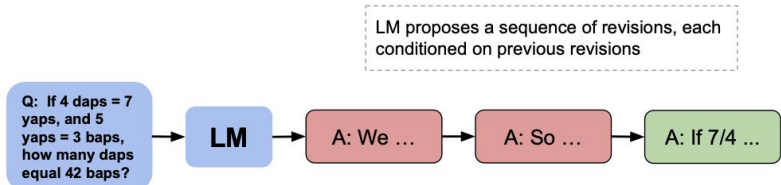- Curves for oracle and PRM-predicted difficulty bins generally overlap

# Modifying the Proposal Distribution

- Parallel Sampling
- Sequential Revision
  - Correct answer correlated to incorrect answers
  - Rollout generation
    - Sample in parallel at high temperature
    - Pair up post-hoc
    - Use character edit distance as correlation metric

# Combining the two strategies

- Each has its pros and cons so why not find a **balance**

- Parallel
  - A global search process
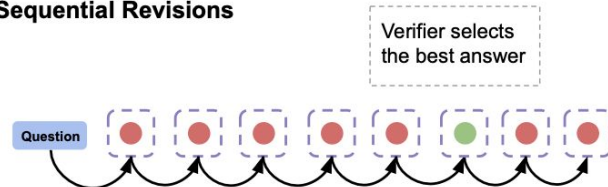  - Give different approaches of solving/answering something

- Sequential
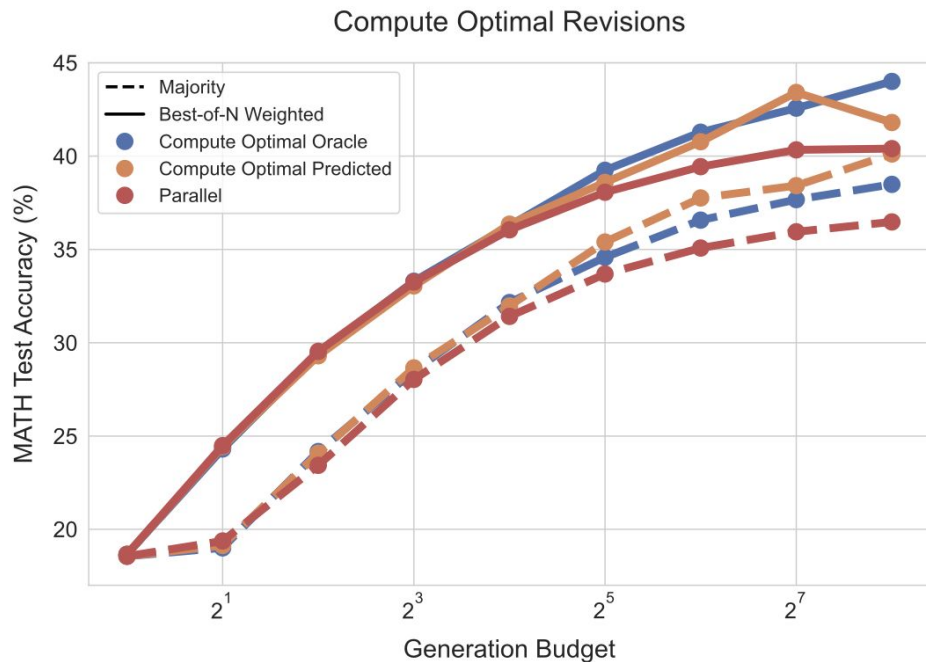  - Local refinement process

# Easier questions better with full sequential revision Harder questions better with combo of sequential - parallel

# Compute-Optimal Revisions vs baseline



Compute Optimal Revisions

- Can still outperform best-of-n by using up to **4x less compute**

# Exchanging Pre-Training and Test-Time Compute

Pretrain with X FLOP-s

Run Y FLOP-s on inference

Increase budget by factor of M

N model parameters

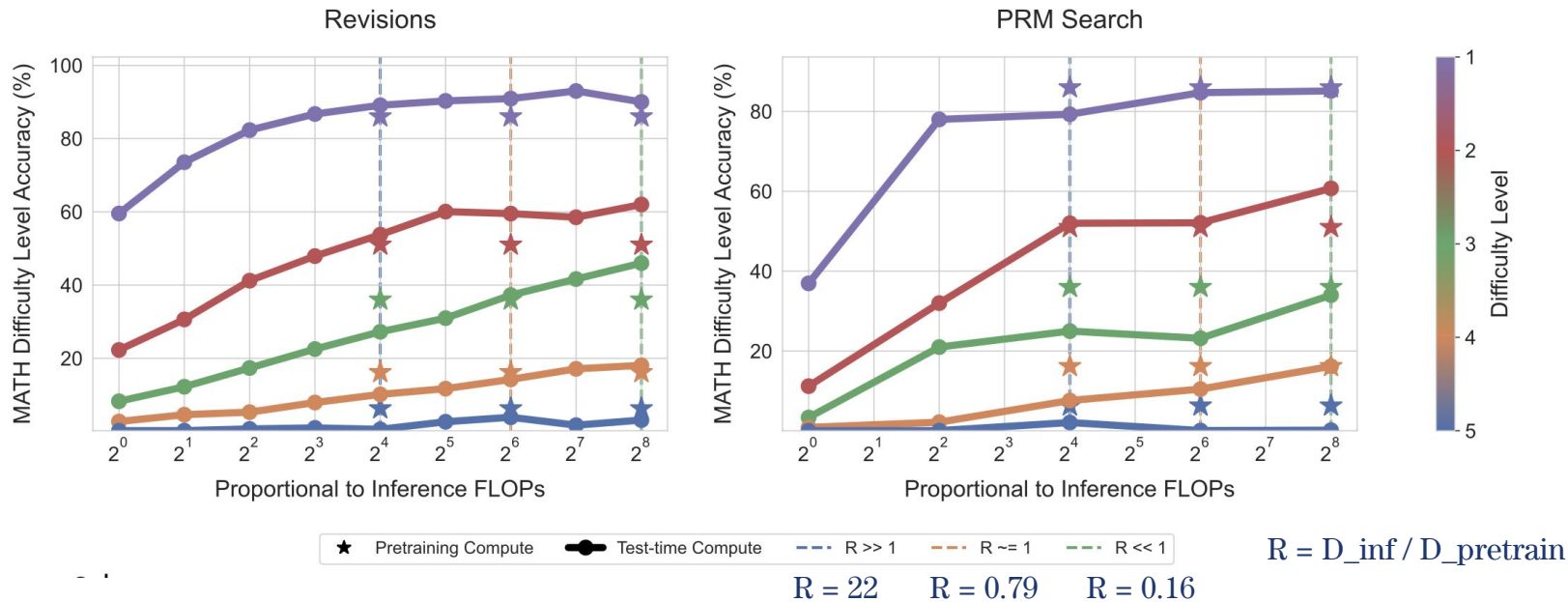D number of tokens used or generated

$X = 6ND_{pretrain}$
$Y = 2ND_{inference}$

Smaller model's inference compute can match larger pretrained model by a factor of

$$M + 3 \ (D_{pretain} / D_{inference}) (M-1)$$

$$R = D_{inference} / D_{pretrain}$$

**Comparing Test-time and Pretraining Compute**

$$R = D\_inf / D\_pretrain$$

R = 22     R = 0.79     R = 0.16

- Very difficult questions or high D_inf => more budget towards **pretraining**
- Easier questions or lower D_inf => more budget towards **test-time compute**

# Takeaways

Difficult of Prompts and Compute Budgets have strong effect on test time optimal efficacy

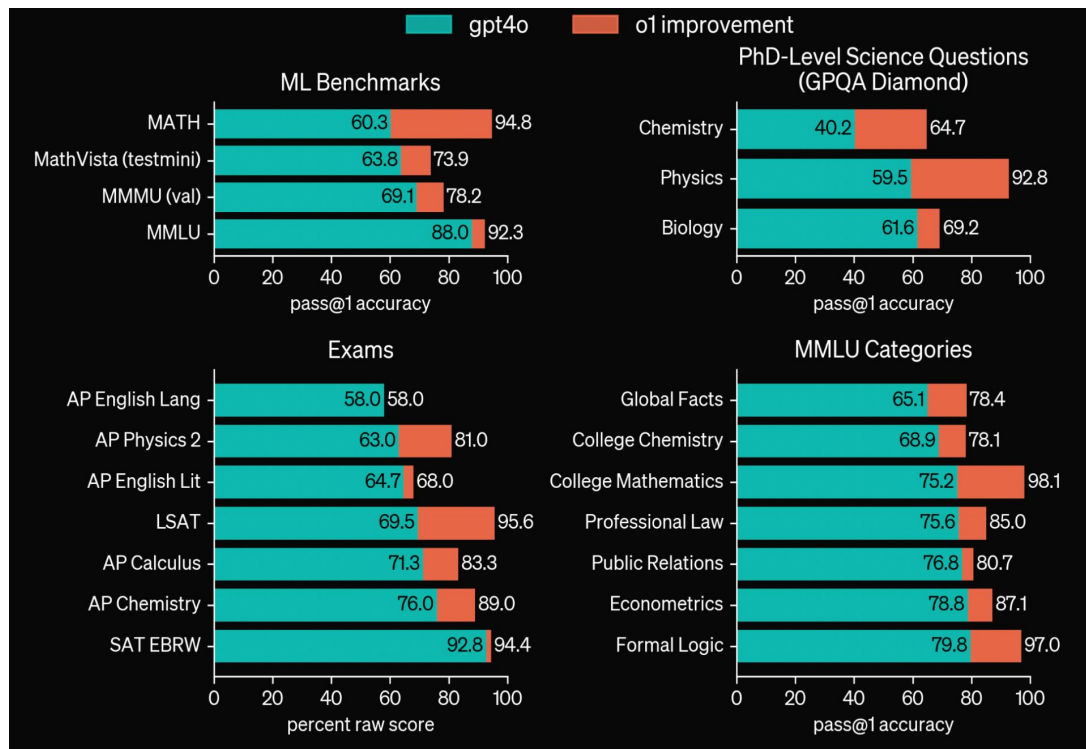Simple methods like search and revisions can scale very well on certain types of prompts

Utilize any resource you can use: By properly using test time computation resources, we can improve the end-2-end performance. (By a factor of 4)

There are a bunch of hyperparameters and variables that need to be accounted for, to get the optimal model

# OpenAI o1-preview

- Produces a long internal chain of thought before responding
- Through RL recognizes and corrects its mistakes
- Tries different approach when current isn't working
- Chain of thought reasoning robustly teaches human values and principles

# Thanks!

Elena Kote, Marvin Gao