

# Tutorial: Test-Time Scaling for Mathematical Reasoning

---

Sean Welleck

May 28, 2025

Carnegie Mellon University

Algorithms for generating outputs with a language model

## Algorithms for generating outputs with a language model

Why? Use *test-time compute* to improve performance

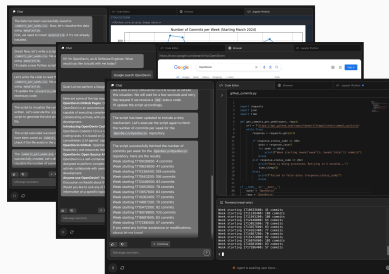
# Language models

RESEARCH

AI achieves silver-medal standard solving  
International Mathematical Olympiad  
problems

25 JULY 2024  
AlphaProof and AlphaGeometry teams

Solving olympiad problems



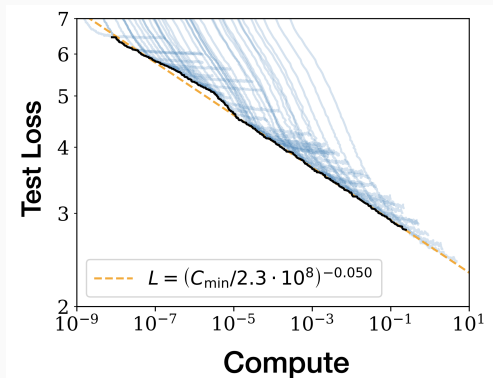
Writing code

Tasks framed as generating sequences: many other applications



# Approach 1: scale pretraining compute

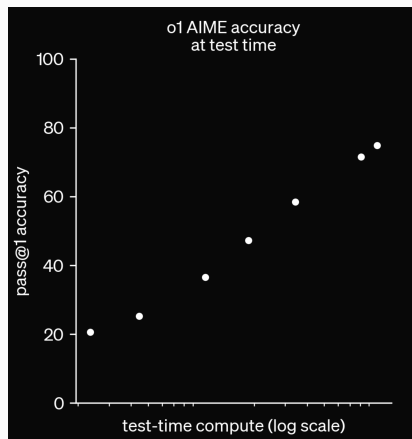
[2020-] Scale training-time compute: larger model, larger dataset



*Scaling Laws for Neural Language Models* [Kaplan et al., 2020]

## Approach 2: scale *test-time* compute

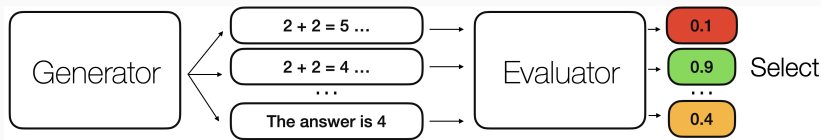
[Now] **Test-time scaling**: increase compute at generation time



Test-time compute vs. accuracy ([OpenAI, 2024])

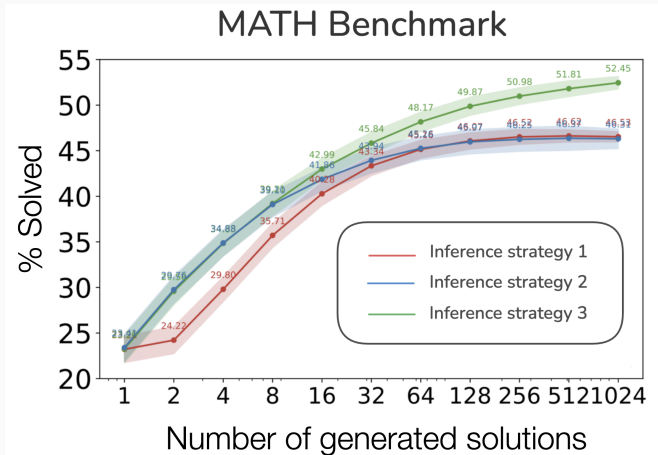
# Test-time scaling strategies

## 1. Generate multiple times



# Test-time scaling strategies

## 1. Generate multiple times



# Test-time scaling strategies

1. Generate multiple times
2. Generate longer outputs

input -> answer

Model Output

A: The answer is 27. ❌

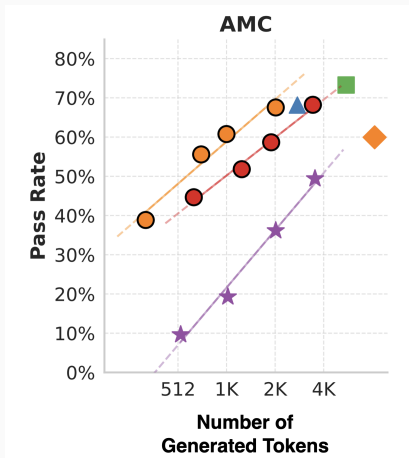
input -> **thought**, answer

Model Output

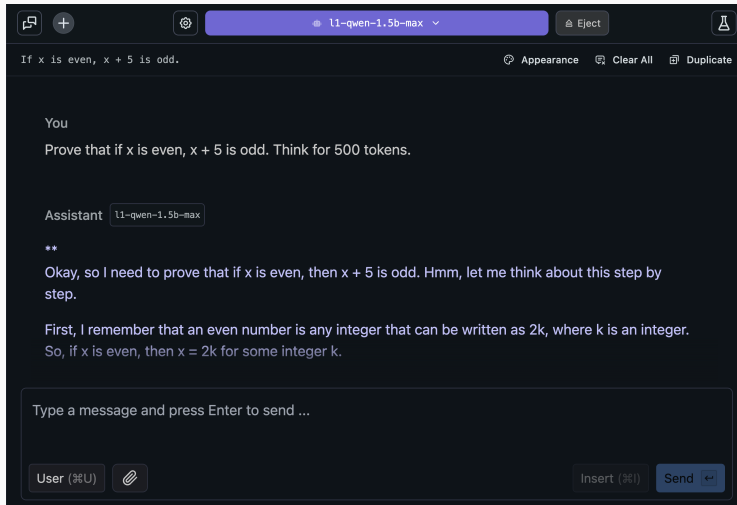
A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✅

# Test-time scaling strategies

1. Generate multiple times
2. Generate longer outputs



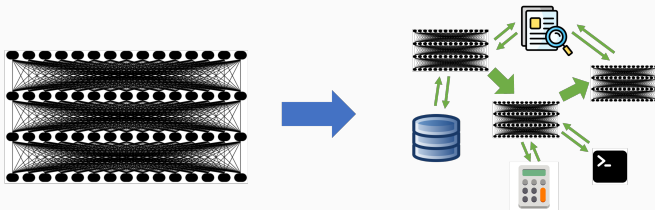
# Test-time scaling strategies



Demo: L1 reasoning model [Aggarwal and Welleck, 2025] on a laptop

# Test-time scaling strategies

1. Generate multiple times
2. Generate longer outputs
3. Incorporate other models/tools



[Zaharia et al., 2024]

Verifiers, code interpreters, search engines, ...



# Today's tutorial

1. **Part 1:** Generate multiple times
  - *Meta-generation*: chain, parallel, refinement, tree search
2. **Part 2:** Generate longer outputs
  - *Long chain-of-thought*

## 1. Part 1: Generate multiple times

- *Meta-generation*: chain, parallel, refinement, tree search

## 2. Part 2: Generate longer outputs

- *Long chain-of-thought*

## Recap: generation and decoding algorithms

**Generator:** Generates a sequence with a language model.

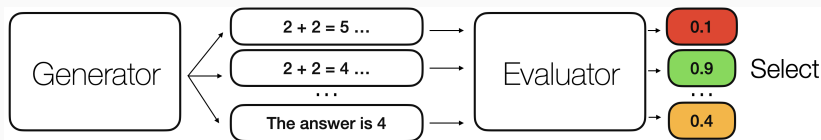


- Example: calling an LLM API
- Decoding algorithms
  - Temperature sampling
  - ...

$$y \sim g(p_{\theta}, x; \phi)$$

# Unified view: Meta-Generation

**Meta-generator:** Test-time strategies for calling a generator multiple times and incorporating external information<sup>1</sup>



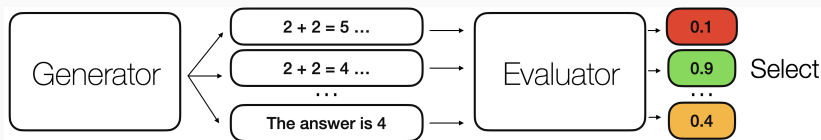
- Example: call API multiple times, select the best sequence with a separate model

---

<sup>1</sup>From *Decoding to Meta-Generation: Inference-time Algorithms for LLMs* [Welleck et al., 2024]

# Unified view: Meta-Generation

**Meta-generator:** Test-time strategies for calling a generator multiple times and incorporating external information<sup>1</sup>



- Example: call API multiple times, select the best sequence with a separate model

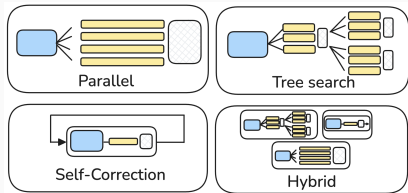
$$y \sim G(x, g; \Phi)$$

---

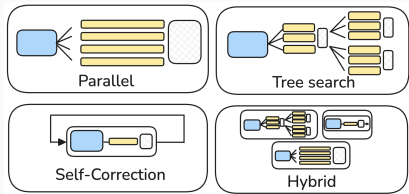
<sup>1</sup>From *Decoding to Meta-Generation: Inference-time Algorithms for LLMs* [Welleck et al., 2024]

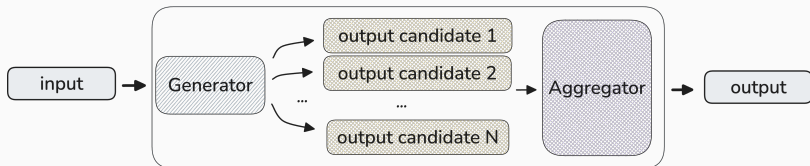
- Strategies

- Parallel
- Tree search
- Refinement/self-correction
- Chain (*Part II*)



- Strategies
  - Parallel
  - Tree search
  - Refinement/self-correction
  - Chain (Part II)





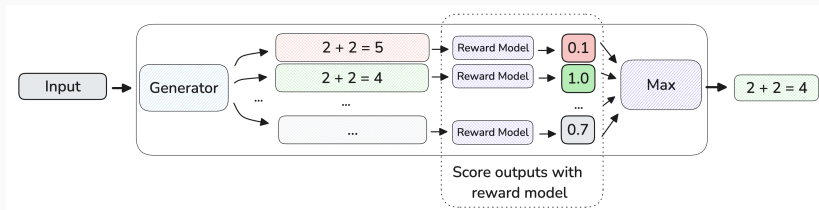
- Generate candidates:

$$\{y^{(1)}, \dots, y^{(N)}\} \sim G(\cdot|x)$$

- Aggregate:

$$y = h(y^{(1)}, \dots, y^{(N)})$$





$$\arg \max_{\{y^{(1)}, \dots, y^{(N)}\}} \underbrace{v(y)}_{\text{reward model}}$$

<sup>2</sup>[Stiennon et al., 2020, Nakano et al., 2022]

**X:** Input:  
Let  $f(r) = \sum_{j=2}^{2008} \frac{1}{j^r} = \frac{1}{2^r} + \frac{1}{3^r} + \dots + \frac{1}{2008^r}$ . Find  $\sum_{k=2}^{\infty} f(k)$ .

-----  
**LEMMA 34B solution:**

We have

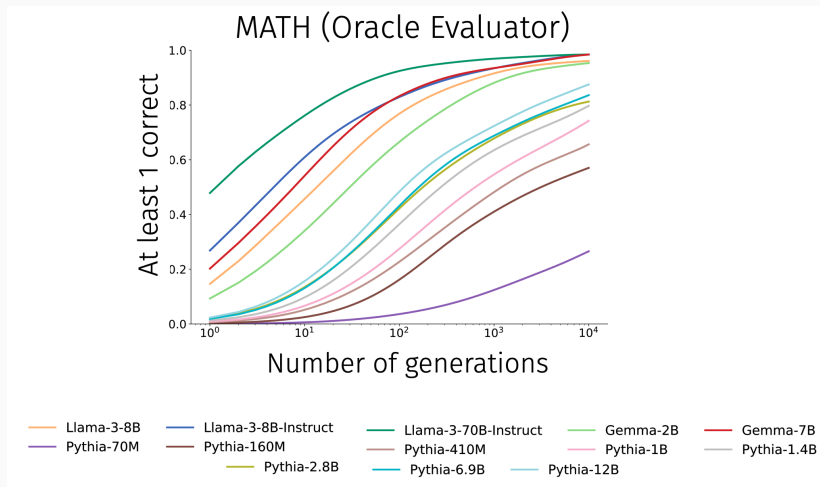
**Z:**

$$\begin{aligned}\sum_{k=2}^{\infty} f(k) &= \sum_{k=2}^{\infty} \sum_{j=2}^{2008} \frac{1}{j^k} = \sum_{j=2}^{2008} \sum_{k=2}^{\infty} \frac{1}{j^k} = \sum_{j=2}^{2008} \frac{1}{j^2} \sum_{k=0}^{\infty} \frac{1}{j^k} = \sum_{j=2}^{2008} \frac{1}{j^2} \frac{1}{1 - \frac{1}{j}} \\&= \sum_{j=2}^{2008} \frac{1}{j(j-1)} = \sum_{j=2}^{2008} \left( \frac{1}{j-1} - \frac{1}{j} \right) \\&= \left( \frac{1}{1} - \frac{1}{2} \right) + \left( \frac{1}{2} - \frac{1}{3} \right) + \dots + \left( \frac{1}{2007} - \frac{1}{2008} \right) \\&= 1 - \frac{1}{2008} \\&= \boxed{\frac{2007}{2008}}.\end{aligned}$$

**Y:** Final Answer: The final answer is  $\frac{2007}{2008}$ .

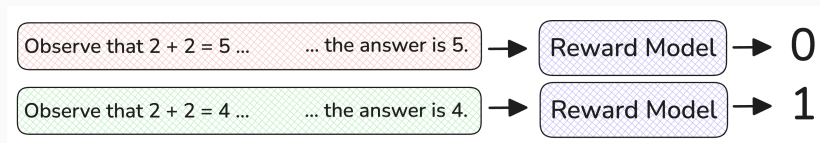
Example: solve a math problem

What if we had a perfect reward model  $v^*(y)$ ?



[Brown et al., 2024]

“Outcome reward model”  $v(y) \rightarrow [0, 1] \approx R(y)$ :

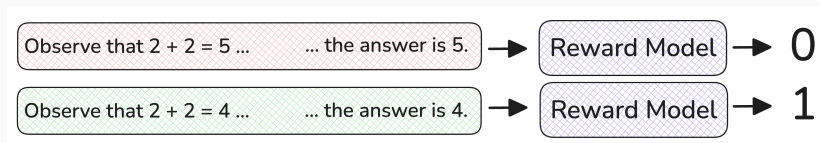


Train reward model with correct and incorrect examples.<sup>3</sup>

---

<sup>3</sup>E.g., [Cobbe et al., 2021]

“Outcome reward model”  $v(y) \rightarrow [0, 1] \approx R(y)$ :



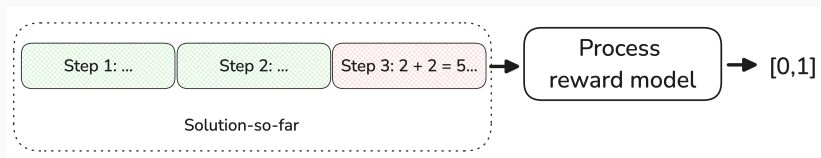
Train reward model with correct and incorrect examples.<sup>3</sup>

Terminology: Reward model  $\approx$  evaluator  $\approx$  critic  $\approx$  verifier  $\approx$  value  $\approx$  scoring model

---

<sup>3</sup>E.g., [Cobbe et al., 2021]

“Process reward model (PRM)”<sup>4</sup>

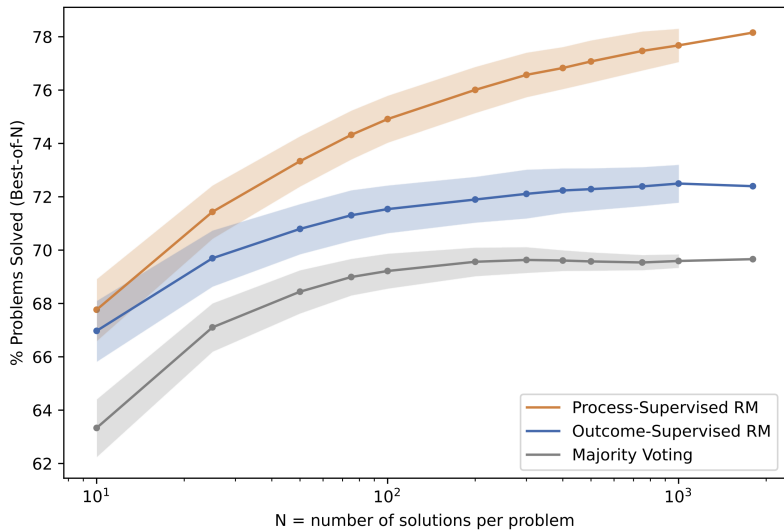


$$v(x, s_1, s_2, \dots, s_t) \rightarrow [0, 1]$$

---

<sup>4</sup>[Uesato et al., 2022, Lightman et al., 2024, Wang et al., 2024a]

# Best-of- $N$ | Reward Model



Lightman et al 2023

Note: Preference reward model

Hello, you are awesome

>

Hello, you are #&@#\*#@#

Train reward model with preference data.<sup>5</sup>

---

<sup>5</sup>E.g., [Stiennon et al., 2020]



Why Best-of-N?

- Approximates maximum (true) reward:

$$\begin{aligned}\text{Best-of-}N &= \arg \max_{y \in \{y^{(1)}, \dots, y^{(N)}\}} v(y) \\ &\approx \arg \max_y v(y)\end{aligned}\tag{1}$$

$$\approx \arg \max_y R(y)\tag{2}$$

Why Best-of-N?

- Approximates maximum (true) reward:

$$\begin{aligned}\text{Best-of-}N &= \arg \max_{y \in \{y^{(1)}, \dots, y^{(N)}\}} v(y) \\ &\approx \arg \max_y v(y)\end{aligned}\tag{1}$$

$$\approx \arg \max_y R(y)\tag{2}$$

(1) gets better as number of generations  $N$  increases!

Why Best-of-N?

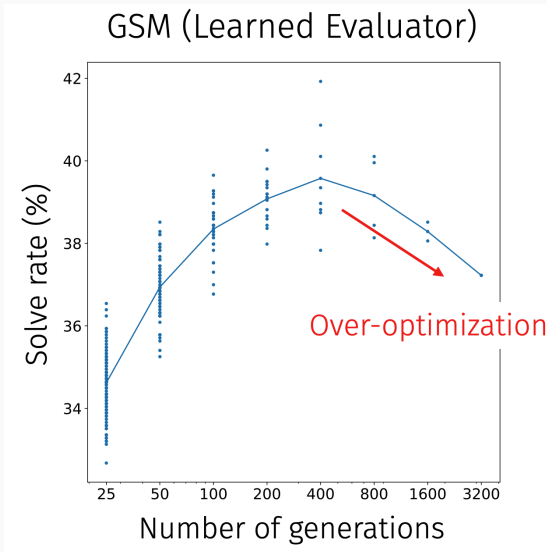
- Approximates maximum (true) reward:

$$\begin{aligned}\text{Best-of-}N &= \arg \max_{y \in \{y^{(1)}, \dots, y^{(N)}\}} v(y) \\ &\approx \arg \max_y v(y)\end{aligned}\tag{1}$$

$$\approx \arg \max_y R(y)\tag{2}$$

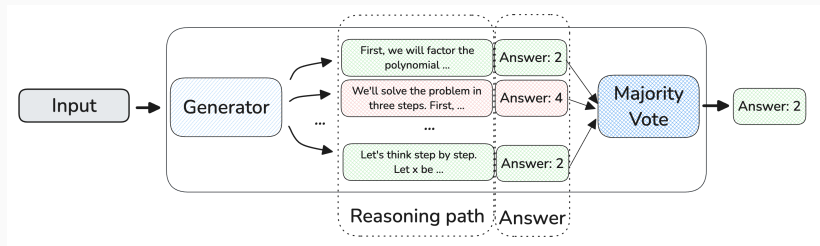
(1) gets better as number of generations  $N$  increases!

(2) can suffer from imperfect reward model, aka “over-optimization”



<sup>6</sup>Plot adapted from *Training Verifiers to Solve Math Word Problems* [Cobbe et al., 2021]

Voting aggregation:<sup>7</sup>

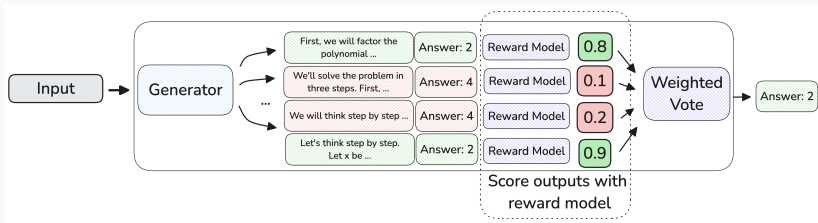


$$\arg \max_a \sum_{i=1}^N 1\{y^{(i)} = a\},$$

<sup>7</sup>Also called *self-consistency* [Wang et al., 2023]

# Weighted voting<sup>8</sup>

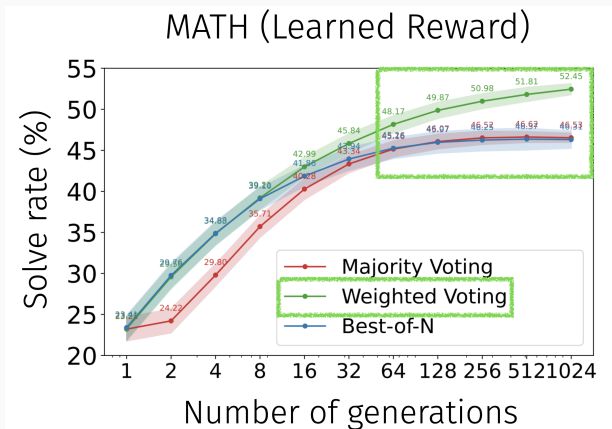
## Weighted Voting:



$$\arg \max_a \sum_{i=1}^N \underbrace{v(y^{(i)})}_{\text{reward model}} \cdot \mathbf{1}\{y^{(i)} = a\},$$

<sup>8</sup>[Li et al., 2023]

Can outperform Best-of-N, e.g.:<sup>9</sup>



<sup>9</sup>[Sun et al., 2024] *Easy-to-Hard Generalization: Scalable Alignment Beyond Human Supervision.*

# Why (weighted) voting?

As the number of candidates  $N \rightarrow \infty$ , voting accuracy converges to...<sup>10</sup>

$$\frac{1}{M} \sum_{i=1}^M \mathbb{I} \left[ a_i^* = \arg \max_a \underbrace{\sum_z v(x, z, a) g(z, a|x)}_{\text{"Marginalize out paths } z"} \right]$$

Notation:

- $(x, z, a)$ : (input, solution, answer)
- $M$ : number of test examples

---

<sup>10</sup>Theorem 2, [Wu et al., 2024] *Inference Scaling Laws*. Y. Wu, Z. Sun, S. Li, S. Welleck, Y. Yang.



# Why (weighted) voting?

As the number of candidates  $N \rightarrow \infty$ , voting accuracy converges to...<sup>10</sup>

$$\frac{1}{M} \sum_{i=1}^M \mathbb{I} \left[ a_i^* = \arg \max_a \underbrace{\sum_z v(x, z, a) g(z, a|x)}_{\text{"Marginalize out paths } z"} \right]$$

**Takeaway 1:** Will accuracy keep improving with more samples?

- **No**, it eventually converges to the accuracy shown above

---

<sup>10</sup>Theorem 2, [Wu et al., 2024] *Inference Scaling Laws*. Y. Wu, Z. Sun, S. Li, S. Welleck, Y. Yang.

# Why (weighted) voting?

As the number of candidates  $N \rightarrow \infty$ , voting accuracy converges to...<sup>10</sup>

$$\frac{1}{M} \sum_{i=1}^M \mathbb{I} \left[ a_i^* = \arg \max_a \underbrace{\sum_z v(x, z, a) g(z, a|x)}_{\text{"Marginalize out paths } z"} \right]$$

**Takeaway 2:** When is weighted voting better than voting?

- When  $v \cdot g$  assigns more total mass to correct answers than  $g$

---

<sup>10</sup>Theorem 2, [Wu et al., 2024] *Inference Scaling Laws*. Y. Wu, Z. Sun, S. Li, S. Welleck, Y. Yang.

# Why (weighted) voting?

As the number of candidates  $N \rightarrow \infty$ , voting accuracy converges to...<sup>10</sup>

$$\frac{1}{M} \sum_{i=1}^M \mathbb{I} \left[ a_i^* = \arg \max_a \underbrace{\sum_z v(x, z, a) g(z, a|x)}_{\text{"Marginalize out paths } z"} \right]$$

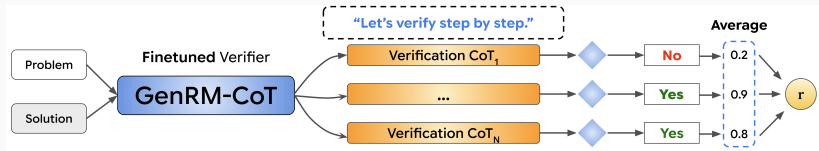
**Takeaway 3:** How do we improve performance further?

- Improve the reward model  $v$
- Improve the generator  $g$  (better model and/or better algorithm)

---

<sup>10</sup>Theorem 2, [Wu et al., 2024] *Inference Scaling Laws*. Y. Wu, Z. Sun, S. Li, S. Welleck, Y. Yang.

Improve the reward model:



Parallel generation *in the reward model too*<sup>11</sup>

Active area of research! (More later in this tutorial!!!)

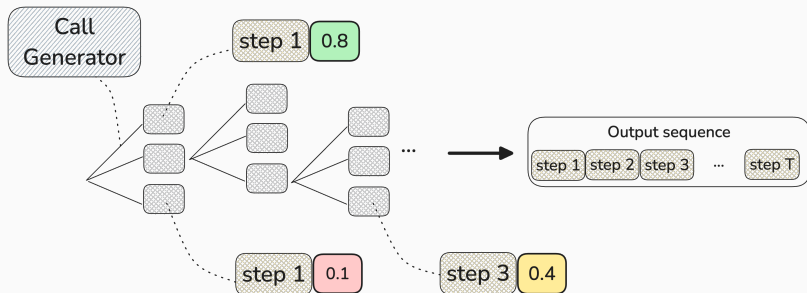
<sup>11</sup>[Zhang et al., 2024]

## Parallel

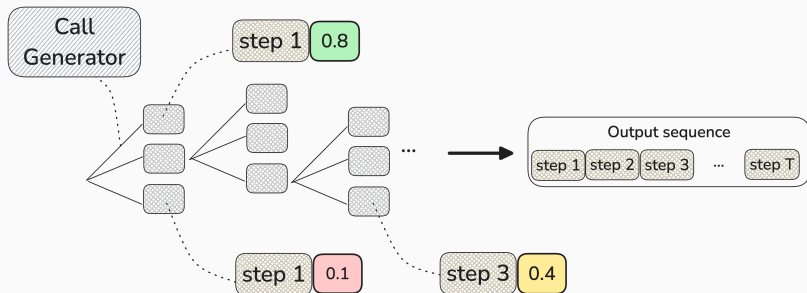
- Explores output space by generating full sequences
- Large performance gains in practice
- Bounded by the quality of the evaluator and generator
- Next: Can we better leverage *intermediate* evaluation?

- Strategies
  - Parallel
  - Tree search
  - Refinement

# Tree search | basic idea



# Tree search | basic idea



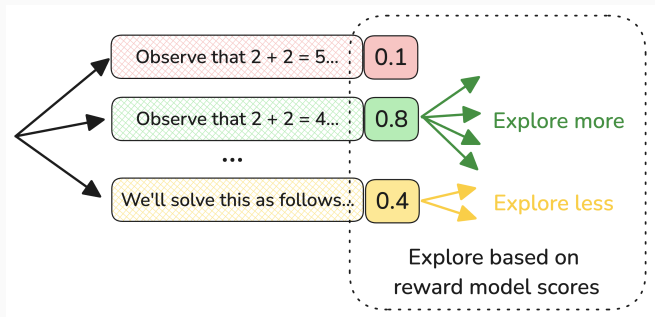
Design choices:

- States  $s$
- Transitions  $s \rightarrow s'$
- Scores  $v(s)$
- Strategy (breadth-first, depth-first, ...)



# Tree search | example (REBASE)

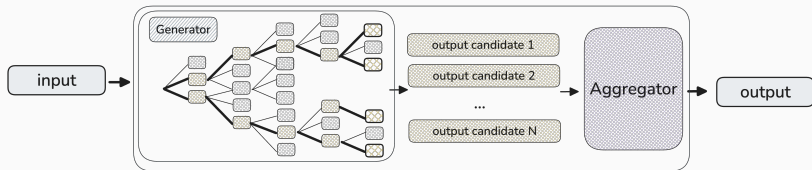
## 2. Reward Balanced Search (Rebase)<sup>12</sup>



$$\text{explore}_i = \text{Round} \left( \text{Budget} \frac{\exp(v(s_i)/\tau)}{\sum_j \exp(v(s_j)/\tau)} \right), \quad (3)$$

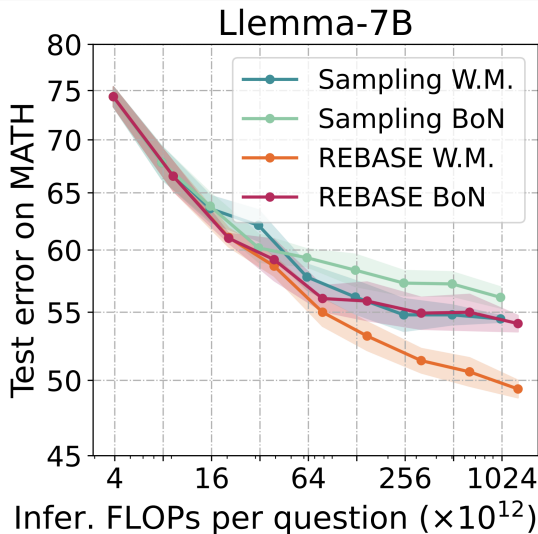
<sup>12</sup>[Wu et al., 2024] *Inference Scaling Laws: An Empirical Analysis of Compute-Optimal Inference.*

# Tree search | example



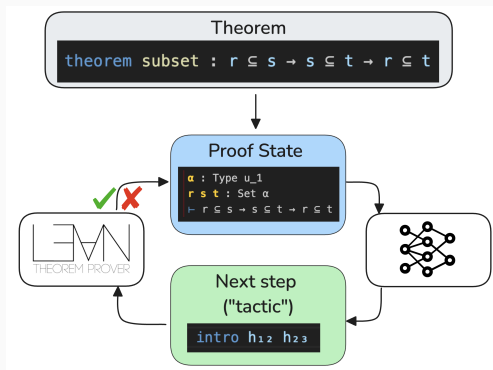
Run tree search to get candidates for aggregation (e.g., voting).

- **Key idea:** Leverages scores on *intermediate* states
  - Backtracking
  - Exploration



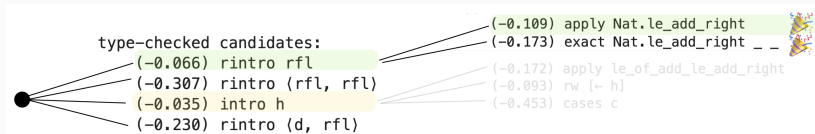
<sup>13</sup>[Wu et al., 2024] *Inference Scaling Laws: An Empirical Analysis of Compute-Optimal Inference.*

# Tree search | examples



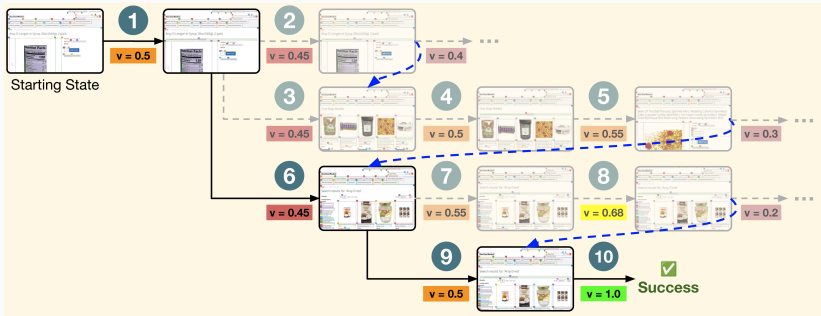
Formal theorem proving [Polu and Sutskever, 2020]

# Tree search | examples



Best-first search in formal theorem proving

# Tree search | examples



Best-first search in web agents [Koh et al., 2024]

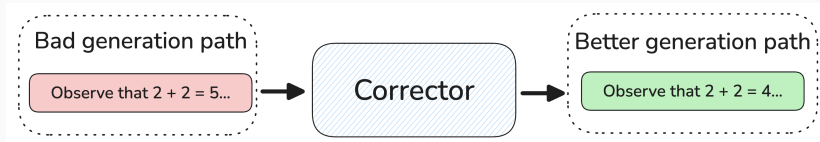
## Tree-search

- Can backtrack and explore using intermediate scores
- Requires a suitable environment and value function
  - Decomposition into states
  - Good reward signal

- Strategies
  - Parallel
  - Tree search
  - Refinement



# Refinement / self-correction

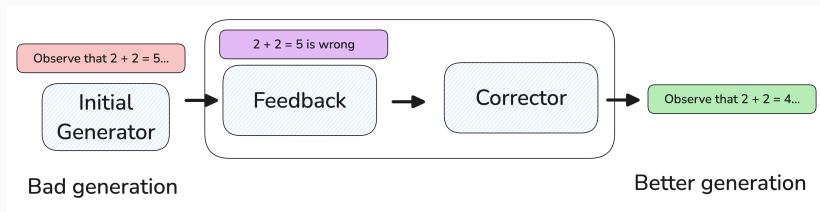


Improve a generation

Repeat:

$$\cdot y^{(i+1)} \sim g(x, y^{(i)})$$

# Refinement / self-correction

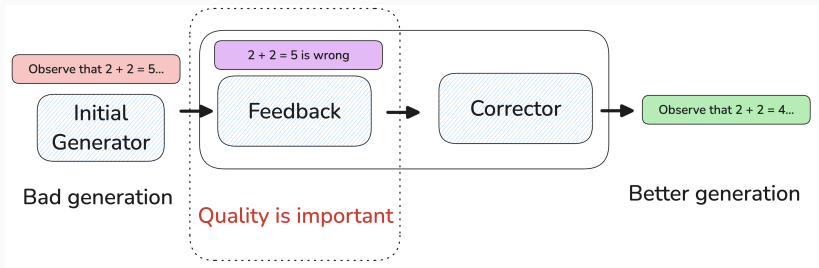


Improve a generation using feedback

Repeat:

$$\cdot y^{(i+1)} \sim g(x, y^{(i)}, F(y^{(i)}))$$

# Refinement / self-correction

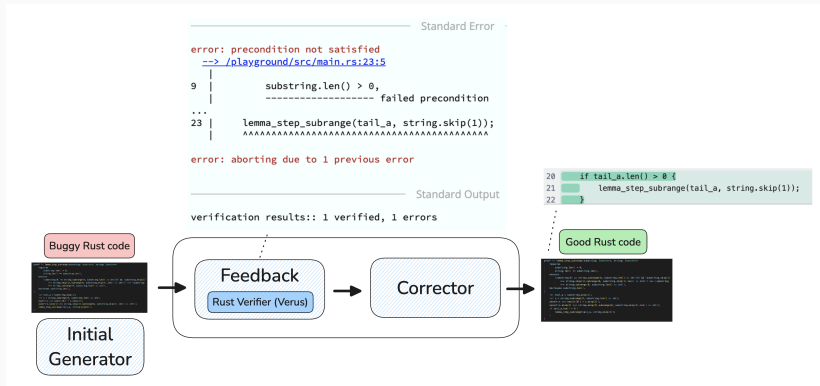


Improve a generation using feedback

In practice, the **quality and source of feedback** is crucial:

- **Extrinsic:** external information at inference time
- **Intrinsic:** no external information at inference time

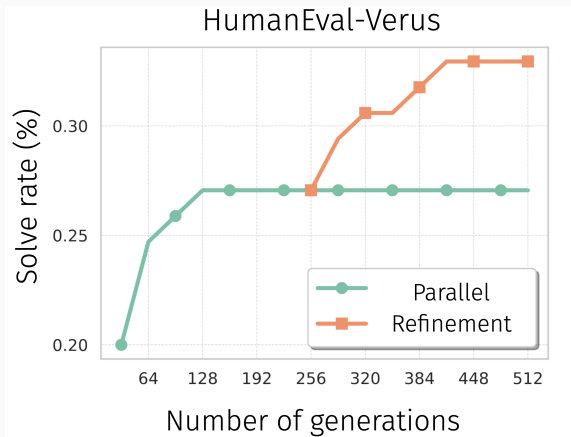
### 1. Extrinsic: external feedback



Feedback: external program verifier<sup>14</sup>

<sup>14</sup> [Aggarwal et al., 2024], *AlphaVerus*. P. Aggarwal, B. Parno, S. Welleck.

## 1. Extrinsic: external feedback



*AlphaVerus*. P. Aggarwal, B. Parno, S. Welleck.

## 1. Extrinsic: external feedback

Several **success cases**:

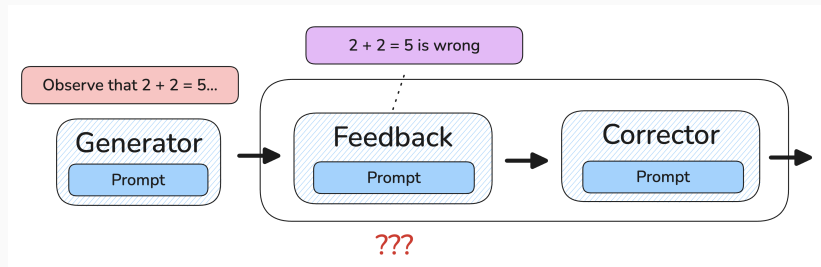
- Verifiers [Aggarwal et al., 2024]
- Code interpreters [Chen et al., 2024]
- Retrievers [Asai et al., 2024]
- Tools + agent environment<sup>14</sup>
- ...

Intuition: adds new information, can detect and localize errors

---

<sup>14</sup><https://x.com/gneubig/status/1866172948991615177>

## 2. Intrinsic: Re-prompt the same model:



Re-prompt a single LLM, e.g. [Madaan et al., 2023]

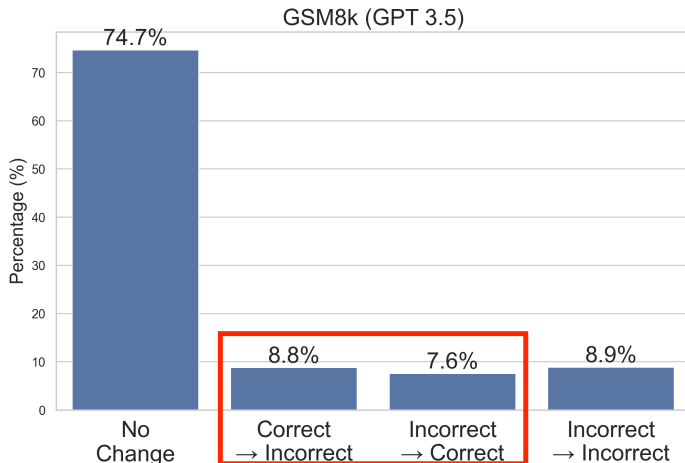


## Mixed results:

- Easy to evaluate tasks: **positive** [Wang et al., 2024b]
  - E.g., missing info [Asai et al., 2024]
- Mathematical reasoning: **mixed**<sup>15</sup>

---

<sup>15</sup>E.g., [Huang et al., 2024] *Large Language Models Cannot Self-Correct Reasoning Yet*



**Takeaway:** feedback is too noisy From [Huang et al., 2024]

# Meta-generators | refinement

Generate “TAYLORSWIFT”

- Generator:
  - $p(\text{character})$
- Feedback:
  - Incorrect characters
- Corrector:
  - Regenerate incorrect



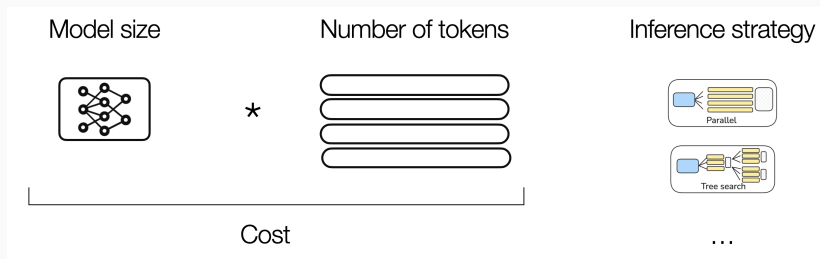
## Refinement / self-correction

- Extrinsic
  - **Positive results** for environments that detect or localize errors
- Intrinsic
  - **Mixed results**, depends on difficulty of verification

- Strategies
  - Parallel
  - Tree search
  - Refinement
- *Inference scaling laws*

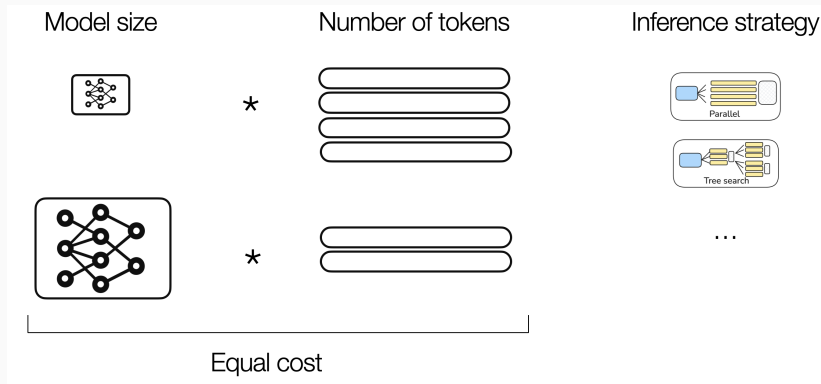
# Inference scaling laws [Wu et al., 2024]

Compute is a function of model size and number of generated tokens

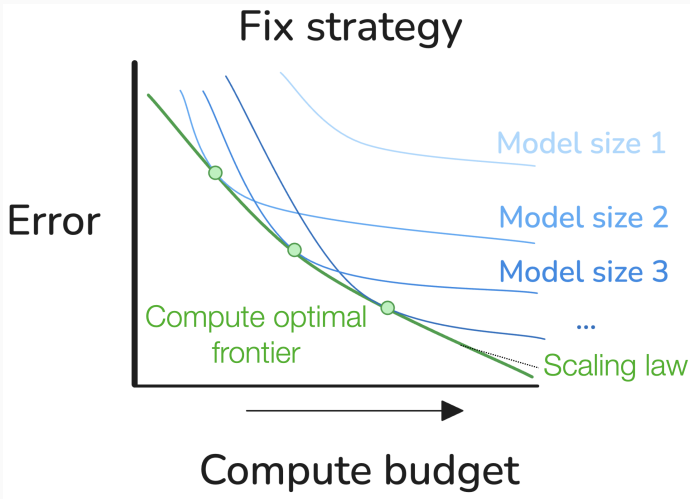


# Inference scaling laws [Wu et al., 2024]

We can choose to increase model size or number of tokens

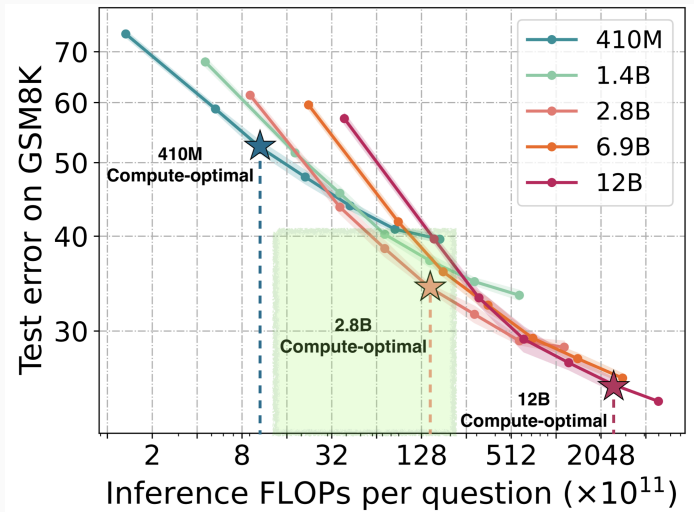


# Inference scaling laws [Wu et al., 2024]





# Inference scaling laws [Wu et al., 2024]



Using a smaller model and generating more is often best [Wu et al., 2024].

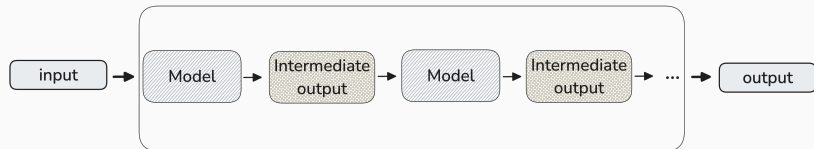
## Designing better strategies

- Example: design a better tree search [Wu et al., 2024]
- Example: select inference strategy based on problem difficulty [Snell et al., 2024]
- ...

- When allocated optimally, performance improves with compute
- Best model size and strategy varies with the budget
  - Sometimes smaller models are better!

- Strategies for generating multiple sequences
- Parallel, tree search, refinement
- Choose methods based on task performance and cost

1. Part 1: Generating multiple sequences
2. Part 2: Generating a single long sequence
  - *Long chain-of-thought*



Compose generators:

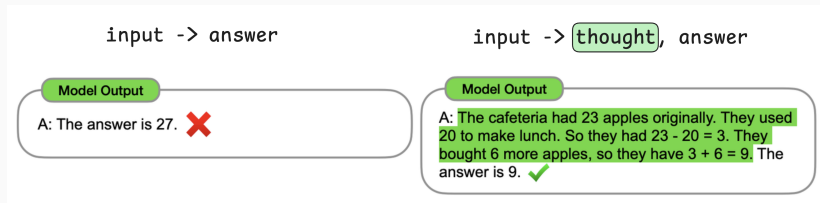
$$y_1 \sim g_1(x)$$

$$y_2 \sim g_2(x, y_1)$$

$$y_3 \sim g_3(x, y_2)$$

$\vdots$

Simple example: *Chain-of-thought* [Wei et al., 2022]:

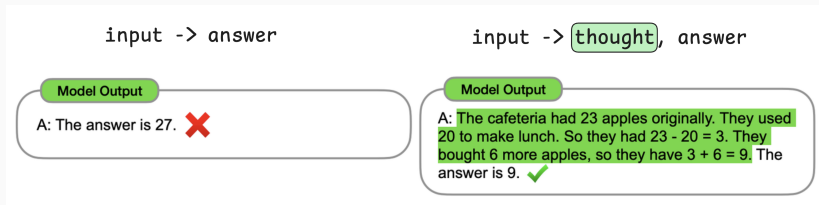


A simple decomposition:

- Generate a thought,  $z \sim g(\cdot|x)$
- Generate an answer,  $a \sim g(\cdot|x, z)$

# Meta-generators | chain

Simple example: *Chain-of-thought* [Wei et al., 2022]:



Increases expressivity<sup>15</sup>

- Variable output length, analogous to a writeable tape
- Idea: train a model that searches on its own in the thought!

---

<sup>15</sup>E.g., [Feng et al., 2023, Merrill and Sabharwal, 2024, Nowak et al., 2024]



# Training a model to search: Basic idea

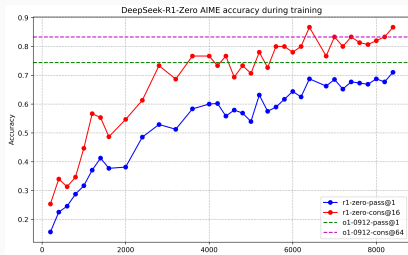
- Train a model to generate a “thought” prior to a final output

$$p_{\theta}(\underbrace{y}_{\text{“output”}}, \underbrace{z}_{\text{“thought”}} | x)$$

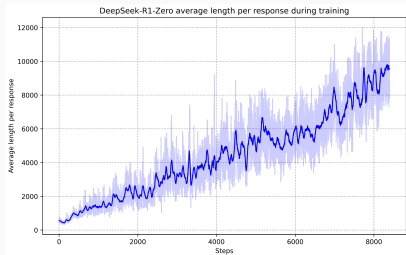
- At inference time, just sample a thought+output
- In principle, the model can learn to try alternatives, perform refinement, backtrack **within the thought**

- Approach 1: reinforcement learning
  - Policy: given a math problem  $x$ , generate a thought + answer
  - Reward: is the answer correct

# Training for long chain-of-thought [DeepSeek-AI et al., 2025]



Accuracy improves during training



Response length increases to > 10,000

# Training for long chain-of-thought [DeepSeek-AI et al., 2025]

---

Question: If  $a > 1$ , then the sum of the real solutions of  $\sqrt{a - \sqrt{a + x}} = x$  is equal to

---

Response: <think>

To solve the equation  $\sqrt{a - \sqrt{a + x}} = x$ , let's start by squaring both ...

$$\left(\sqrt{a - \sqrt{a + x}}\right)^2 = x^2 \implies a - \sqrt{a + x} = x^2.$$

Rearrange to isolate the inner square root term:

$$(a - x^2)^2 = a + x \implies a^2 - 2ax^2 + (x^2)^2 = a + x \implies x^4 - 2ax^2 - x + (a^2 - a) = 0$$

...

**Wait, wait. Wait. That's an aha moment I can flag here.**

Let's reevaluate this step-by-step to identify if the correct sum can be ...

We started with the equation:

$$\sqrt{a - \sqrt{a + x}} = x$$

First, let's square both sides:

$$a - \sqrt{a + x} = x^2 \implies \sqrt{a + x} = a - x^2$$

Next, I could square both sides again, treating the equation: ...

...

---

Example response

## 1. Uncertainty

- *Wait... / Hold on...*
- *Wait-actually, does this formula apply here?*

## 2. Branching, backtracking, retrying

- *Alternatively, generating functions could model this problem...*
- *Revisiting...*
- *Wait, I'm overthinking. Let's try again...*

## 3. Verification

- *Let's check if we made an error. We should verify...*
- *This is a contradiction, so we must have made a mistake.*
- *Let's test this with...*

## 4. Key Points

- *Key takeaway... / It's worth noting...*

## 5. Clarification

- *In other words... / To clarify...*

## 6. Synthesis

- *Ultimately... / Putting it all together...*

## 1. Classification Criteria Identification

N criteria

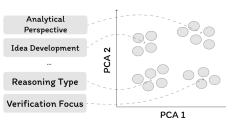
Analytical Perspective: Top-down vs Bottom-up

Idea Development: Multiple vs Single

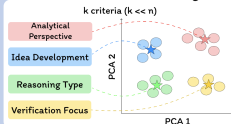
Reasoning Type: Inductive vs Deductive

Verification Focus: Hypothesis-driven vs Data-driven

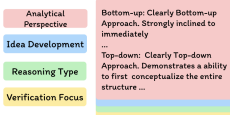
## 2. Classification Criteria Embedding



## 3. Criteria Compression via Hierarchical Clustering



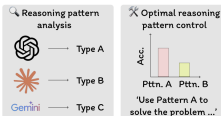
## 4. Rubric Generation



## 5. Pattern Analysis Report Generation

The reasoning pattern shows a clear top-down approach ...  
 The response demonstrates multiple thinking ways ...  
 The given response demonstrates strong deductive reasoning ...  
 Verification is driven by data rather than hypotheses ...

## Use cases



## The CoT ENCYCLOPEDIA: Analyzing, Predicting, and Controlling how a Reasoning Model will Think

Seungyun Lee<sup>1,3\*</sup> Seungone Kim<sup>2\*</sup> Minju Seo<sup>1</sup> Yongrae Jo<sup>3</sup>

Dongyoung Go<sup>4,5</sup> Hyeonbin Hwang<sup>1</sup> Jinho Park<sup>1</sup>

Xiang Yue<sup>2</sup> Sean Wellick<sup>2</sup> Graham Neubig<sup>2</sup> Moontae Lee<sup>1</sup> Minjoon Seo<sup>1</sup>

KAIST AI<sup>1</sup> Carnegie Mellon University<sup>2</sup> LG AI Research<sup>3</sup>

NAVER Search US<sup>4</sup> Cornell University<sup>5</sup>

{seungyun, minjoon}@kaist.ac.kr seungone@cmu.edu

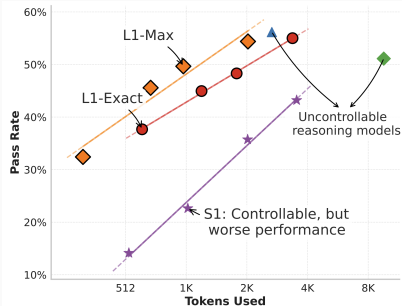
# Controlling the length: L1 [Aggarwal and Welleck, 2025]

- Train model with reinforcement learning to adhere to length constraints
- E.g. “use up to 1000 tokens” provided in the prompt
- Reward: correctness and length constraint penalty

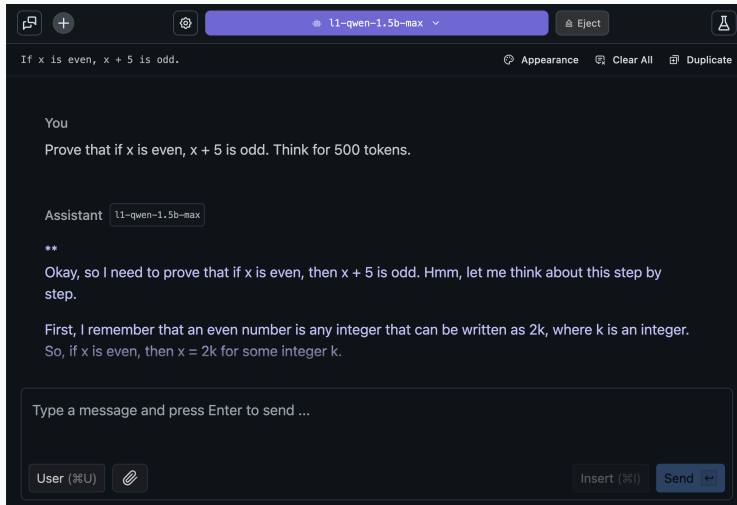
## L1: Controlling How Long A Reasoning Model Thinks With Reinforcement Learning

Pranjal Aggarwal  
Carnegie Mellon University

Sean Welleck  
Carnegie Mellon University



# Controlling the length: L1 [Aggarwal and Welleck, 2025]

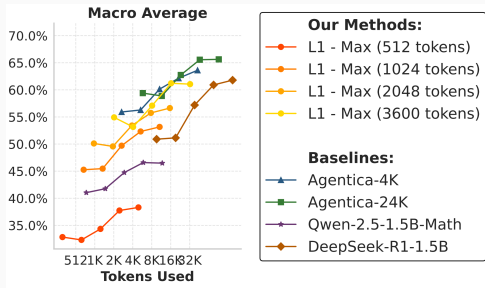


Demo: Using L1 on a laptop in LM studio



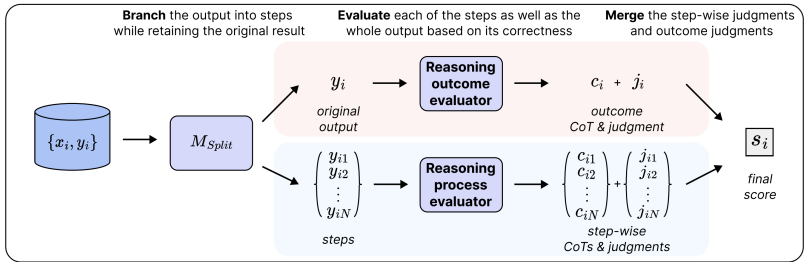
# Long chain-of-thought | sequential vs. parallel

- **Sequential:** long chain-of-thought
- **Parallel:** majority voting (multiple long COTs)



[Aggarwal and Welleck, 2025]

# Reasoning models as evaluators in Best-of-N [Kim et al., 2025]

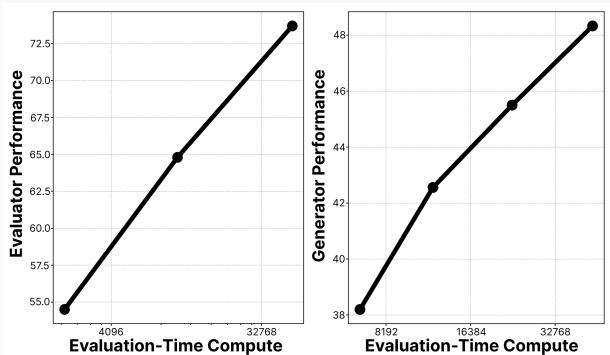


## Scaling Evaluation-time Compute with Reasoning Models as Process Evaluators

Seungone Kim<sup>1\*</sup> Ian Wu<sup>2\*</sup> Jinu Lee<sup>3\*</sup> Xiang Yue<sup>1</sup>  
Seongyun Lee<sup>4</sup> Mingyeong Moon<sup>2</sup> Kiril Gashteovski<sup>5,6</sup> Carolin Lawrence<sup>5</sup>  
Julia Hockenmaier<sup>3</sup> Graham Neubig<sup>1</sup> Sean Welleck<sup>1</sup>

<sup>1</sup>CMU <sup>2</sup>Independent Researcher <sup>3</sup>UIUC <sup>4</sup>KAIST AI  
<sup>5</sup>NEC Laboratories Europe <sup>6</sup>Ss.Cyril and Methodius University of Skopje

# Reasoning models as evaluators in Best-of-N [Kim et al., 2025]



## Scaling Evaluation-time Compute with Reasoning Models as Process Evaluators

Seungone Kim<sup>1\*</sup> Ian Wu<sup>2\*</sup> Jinu Lee<sup>3\*</sup> Xiang Yue<sup>1</sup>  
Seongyun Lee<sup>4</sup> Mingyeong Moon<sup>2</sup> Kiril Gashteovski<sup>5,6</sup> Carolin Lawrence<sup>5</sup>  
Julia Hockenmaier<sup>3</sup> Graham Neubig<sup>1</sup> Sean Welleck<sup>1</sup>

<sup>1</sup>CMU <sup>2</sup>Independent Researcher <sup>3</sup>UIUC <sup>4</sup>KAIST AI  
<sup>5</sup>NEC Laboratories Europe <sup>6</sup>Ss.Cyrl and Methodius University of Skopje

- Train a model to generate a long sequence, then use a simple inference algorithm
- Internally can perform backtracking, self-correction, etc.
- Emerging area of research!

1. Test-time inference strategies take a trained model and improve performance by:
  - Generating tokens according to a strategy
  - Incorporate external information
    - Reward models
    - Environment feedback

Very active and evolving research area!

1. Test-time inference strategies take a trained model and improve performance by:
  - Generating tokens according to a strategy
  - Incorporate external information
    - Reward models
    - Environment feedback
2. Two complementary meta-generation strategies
  - Call a generator to generate a thought prior to an answer
    - Long chain-of-thought
  - Call a generator multiple times in a structured way
    - parallel, tree search, refinement

Very active and evolving research area!

Published in Transactions on Machine Learning Research (11/2024)

## From Decoding to Meta-Generation: Inference-time Algorithms for Large Language Models

Sean Welleck  
Carnegie Mellon University

wellecks@cmu.edu

Amanda Bertsch\*  
Carnegie Mellon University

abertsch@cs.cmu.edu

Matthew Finlayson\*  
University of Southern California

mfinlayso@usc.edu

Hailey Schoelkopf\*  
EleutherAI

hailey@eleuther.ai

Alex Xie  
Carnegie Mellon University

alex@cs.cmu.edu

Graham Neubig  
Carnegie Mellon University

gneubig@cmu.edu

Ilya Kulikov  
Meta FAIR

ikuliko@meta.com

Zaid Harchaoui  
University of Washington

zaid@uw.edu

\*Co-second authors

Reviewed on OpenReview: <https://openreview.net/forum?id=esk9tc1tM5>

## NeurIPS 2024 Tutorial: Beyond Decoding: Meta-Generation Algorithms for Large Language Models



Sean Welleck<sup>1</sup>



Amanda Bertsch<sup>1</sup>



Matthew Finlayson<sup>2</sup>



Alex Xie<sup>1</sup>



Graham Neubig<sup>1</sup>



Konstantin Golobokov<sup>5</sup>



Hailey Schoelkopf<sup>3</sup>



Ilya Kulikov<sup>4</sup>



Zaid Harchaoui<sup>5</sup>

<sup>1</sup>Carnegie Mellon University

<sup>2</sup>University of Southern California

<sup>3</sup>Work done while at EleutherAI

<sup>4</sup>Meta AI

<sup>5</sup>University of Washington

NeurIPS 2024 LLM Inference Tutorial:  
<https://cmu-l3.github.io/neurips2024-inference-tutorial/>

TMLR Survey [Welleck et al., 2024]



Home About Call for Papers Speakers Schedule Organizers

## THE FIRST WORKSHOP ON TEST-TIME SCALING AND REASONING MODELS (SCALR @ COLM 2025)

October 10, 2025, Montreal, Canada

### About The Workshop

The ScalR workshop focuses on the emerging challenges and opportunities in scaling and reasoning capabilities of large language models at test time. As models continue to grow in size and capability, understanding how to effectively scale their performance and enhance their reasoning abilities during inference has become increasingly important.

### WHERE & WHEN

📍 Palais des Congrès  
Montreal, Canada

📅 October 10, 2025

1st Workshop on Test-Time Scaling and Reasoning Models at COLM 2025  
<https://scalr-workshop.github.io/>



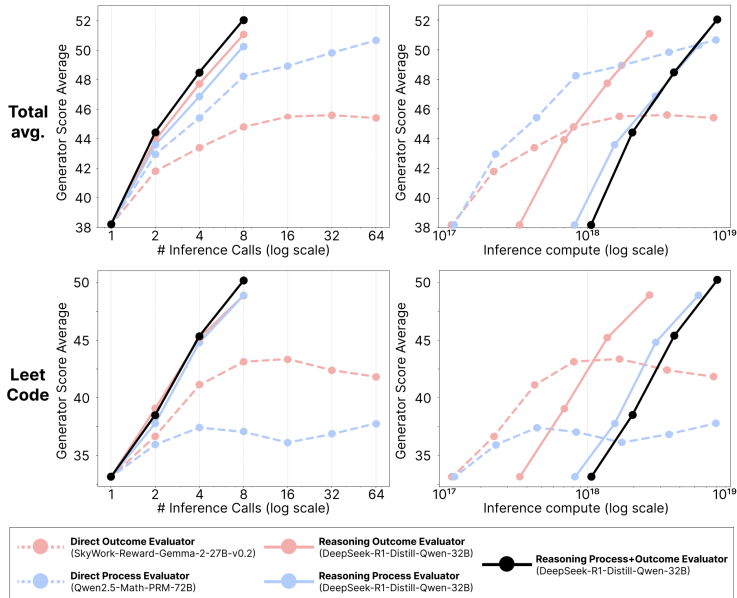
# Thank you!

Sean Welleck  
Carnegie Mellon University  
Learning, Language, and Logic (L3) Lab  
[www.wellecks.com](http://www.wellecks.com)

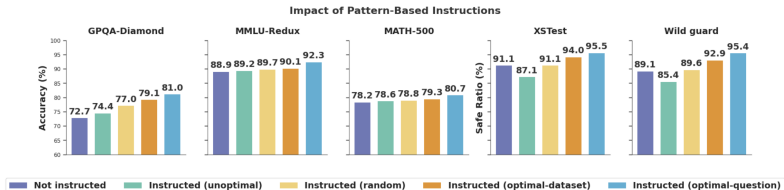
# Appendix

---

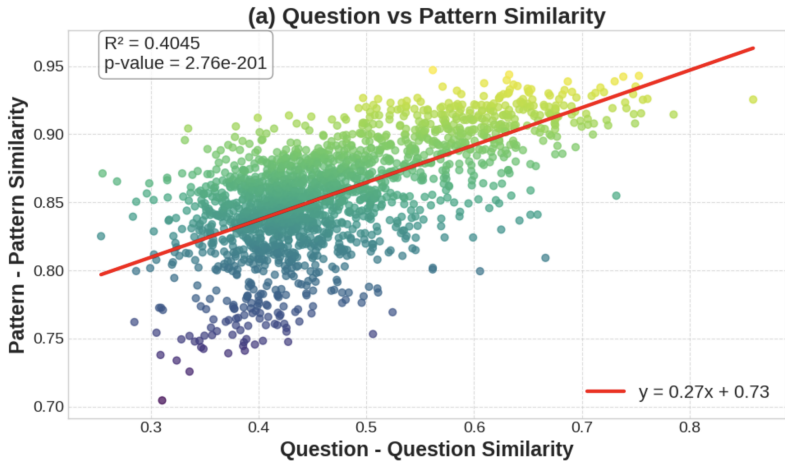
# Additional results



# Additional results



# Additional results





Aggarwal, P., Parno, B., and Welleck, S. (2024).

**Alphaverus: Bootstrapping formally verified code generation through self-improving translation and treefinement.**

<https://arxiv.org/abs/2412.06176>.



Aggarwal, P. and Welleck, S. (2025).

**L1: Controlling how long a reasoning model thinks with reinforcement learning.**

<https://arxiv.org/abs/2503.04697>.



Asai, A., He\*, J., Shao\*, R., Shi, W., Singh, A., Chang, J. C., Lo, K., Soldaini, L., Feldman, Tian, S., Mike, D., Wadden, D., Latzke, M., Minyang, Ji, P., Liu, S., Tong, H., Wu, B., Xiong, Y., Zettlemoyer, L., Weld, D., Neubig, G., Downey, D., Yih, W.-t., Koh, P. W., and Hajishirzi, H. (2024).

**OpenScholar: Synthesizing scientific literature with retrieval-augmented language models.**

*Arxiv.*



Brown, B., Juravsky, J., Ehrlich, R., Clark, R., Le, Q. V., Ré, C., and Mirhoseini, A. (2024).

**Large language monkeys: Scaling inference compute with repeated sampling.**

<https://arxiv.org/abs/2407.21787>.



Chen, X., Lin, M., Schärli, N., and Zhou, D. (2024).

**Teaching large language models to self-debug.**

*In The Twelfth International Conference on Learning Representations.*



Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. (2021).

**Training verifiers to solve math word problems.**

<https://arxiv.org/abs/2110.14168>.





DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., Zhang, X., Yu, X., Wu, Y., Wu, Z. F., Gou, Z., Shao, Z., Li, Z., Gao, Z., Liu, A., Xue, B., Wang, B., Wu, B., Feng, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Ding, H., Xin, H., Gao, H., Qu, H., Li, H., Guo, J., Li, J., Wang, J., Chen, J., Yuan, J., Qiu, J., Li, J., Cai, J. L., Ni, J., Liang, J., Chen, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L., Zhang, M., Zhang, M., Tang, M., Li, M., Wang, M., Li, M., Tian, N., Huang, P., Zhang, P., Wang, Q., Chen, Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen, R. J., Jin, R. L., Chen, R., Lu, S., Zhou, S., Chen, S., Ye, S., Wang, S., Yu, S., Zhou, S., Pan, S., Li, S. S., Zhou, S., Wu, S., Ye, S., Yun, T., Pei, T., Sun, T., Wang, T., Zeng, W., Zhao, W.,

Liu, W., Liang, W., Gao, W., Yu, W., Zhang, W., Xiao, W. L., An, W., Liu, X., Wang, X., Chen, X., Nie, X., Cheng, X., Liu, X., Xie, X., Liu, X., Yang, X., Li, X., Su, X., Lin, X., Li, X. Q., Jin, X., Shen, X., Chen, X., Sun, X., Wang, X., Song, X., Zhou, X., Wang, X., Shan, X., Li, Y. K., Wang, Y. Q., Wei, Y. X., Zhang, Y., Xu, Y., Li, Y., Zhao, Y., Sun, Y., Wang, Y., Yu, Y., Zhang, Y., Shi, Y., Xiong, Y., He, Y., Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Ou, Y., Wang, Y., Gong, Y., Zou, Y., He, Y., Xiong, Y., Luo, Y., You, Y., Liu, Y., Zhou, Y., Zhu, Y. X., Xu, Y., Huang, Y., Li, Y., Zheng, Y., Zhu, Y., Ma, Y., Tang, Y., Zha, Y., Yan, Y., Ren, Z. Z., Ren, Z., Sha, Z., Fu, Z., Xu, Z., Xie, Z., Zhang, Z., Hao, Z., Ma, Z., Yan, Z., Wu, Z., Gu, Z., Zhu, Z., Liu, Z., Li, Z., Xie, Z., Song, Z., Pan, Z., Huang, Z., Xu, Z., Zhang, Z., and Zhang, Z. (2025).

**Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.**



Feng, G., Zhang, B., Gu, Y., Ye, H., He, D., and Wang, L. (2023).  
**Towards revealing the mystery behind chain of thought: A theoretical perspective.**

*In Thirty-seventh Conference on Neural Information Processing Systems.*



Huang, J., Chen, X., Mishra, S., Zheng, H. S., Yu, A. W., Song, X., and Zhou, D. (2024).

**Large language models cannot self-correct reasoning yet.**




*In The Twelfth International Conference on Learning Representations.*



Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. (2020).

**Scaling laws for neural language models.**

<https://arxiv.org/abs/2001.08361>.

-  Kim, S., Wu, I., Lee, J., Yue, X., Lee, S., Moon, M., Gashteovski, K., Lawrence, C., Hockenmaier, J., Neubig, G., and Welleck, S. (2025). **Scaling evaluation-time compute with reasoning models as process evaluators.**  
<https://arxiv.org/abs/2503.19877>.
-  Koh, J. Y., McAleer, S., Fried, D., and Salakhutdinov, R. (2024). **Tree search for language model agents.**  
*arXiv preprint arXiv:2407.01476*.
-  Lee, S., Kim, S., Seo, M., Jo, Y., Go, D., Hwang, H., Park, J., Yue, X., Welleck, S., Neubig, G., Lee, M., and Seo, M. (2025). **The cot encyclopedia: Analyzing, predicting, and controlling how a reasoning model will think.**  
<https://arxiv.org/abs/2505.10185>.



Li, Y., Lin, Z., Zhang, S., Fu, Q., Chen, B., Lou, J.-G., and Chen, W. (2023).

**Making language models better reasoners with step-aware verifier.**

In Rogers, A., Boyd-Graber, J., and Okazaki, N., editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5315–5333, Toronto, Canada. Association for Computational Linguistics.



Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. (2024).

**Let’s verify step by step.**

In *The Twelfth International Conference on Learning Representations*.



Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N., Prabhume, S., Yang, Y., Gupta, S., Majumder, B. P., Hermann, K., Welleck, S., Yazdanbakhsh, A., and Clark, P. (2023).

**Self-refine: Iterative refinement with self-feedback.**

In *Thirty-seventh Conference on Neural Information Processing Systems*.



Merrill, W. and Sabharwal, A. (2024).

**The expressive power of transformers with chain of thought.**

In *The Twelfth International Conference on Learning Representations*.



Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., Jiang, X., Cobbe, K., Eloundou, T., Krueger, G., Button, K., Knight, M., Chess, B., and Schulman, J. (2022).

**Webgpt: Browser-assisted question-answering with human feedback.**

<https://arxiv.org/abs/2112.09332>.



Nowak, F., Svete, A., Butoi, A., and Cotterell, R. (2024).

**On the representational capacity of neural language models with chain-of-thought reasoning.**

In Ku, L.-W., Martins, A., and Srikumar, V., editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12510–12548, Bangkok, Thailand. Association for Computational Linguistics.



OpenAI (2024).

**Learning to reason with llms.**

<https://openai.com/index/learning-to-reason-with-llms/>.



Polu, S. and Sutskever, I. (2020).

**Generative language modeling for automated theorem proving.**



Snell, C., Lee, J., Xu, K., and Kumar, A. (2024).

**Scaling llm test-time compute optimally can be more effective than scaling model parameters.**

<https://arxiv.org/abs/2408.03314>.





Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. (2020).

**Learning to summarize with human feedback.**

In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3008–3021. Curran Associates, Inc.



Sun, Z., Yu, L., Shen, Y., Liu, W., Yang, Y., Welleck, S., and Gan, C. (2024).

**Easy-to-hard generalization: Scalable alignment beyond human supervision.**

In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.



Uesato, J., Kushman, N., Kumar, R., Song, F., Siegel, N., Wang, L., Creswell, A., Irving, G., and Higgins, I. (2022).

**Solving math word problems with process- and outcome-based feedback.**



Wang, P., Li, L., Shao, Z., Xu, R., Dai, D., Li, Y., Chen, D., Wu, Y., and Sui, Z. (2024a).

**Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations.**

In Ku, L.-W., Martins, A., and Srikumar, V., editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9426–9439, Bangkok, Thailand. Association for Computational Linguistics.



Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi, E. H., Narang, S., Chowdhery, A., and Zhou, D. (2023).

**Self-consistency improves chain of thought reasoning in language models.**

*In The Eleventh International Conference on Learning Representations.*



Wang, Y., Wu, Y., Wei, Z., Jegelka, S., and Wang, Y. (2024b).

**A theoretical understanding of self-correction through in-context alignment.**

<https://arxiv.org/abs/2405.18634>.



Wei, J., Wang, X., Schuurmans, D., Bosma, M., brian ichter, Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. (2022).

**Chain of thought prompting elicits reasoning in large language models.**

In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K., editors, *Advances in Neural Information Processing Systems*.



Welleck, S., Bertsch, A., Finlayson, M., Schoelkopf, H., Xie, A., Neubig, G., Kulikov, I., and Harchaoui, Z. (2024).

**From decoding to meta-generation: Inference-time algorithms for large language models.**

<https://arxiv.org/abs/2406.16838>.



Wu, Y., Sun, Z., Li, S., Welleck, S., and Yang, Y. (2024).  
**Inference scaling laws: An empirical analysis of  
compute-optimal inference for problem-solving with language  
models.**

<https://arxiv.org/abs/2408.00724>.



Zaharia, M., Khattab, O., Chen, L., Davis, J. Q., Miller, H., Potts, C.,  
Zou, J., Carbin, M., Frankle, J., Rao, N., and Ghodsi, A. (2024).  
**The shift from models to compound ai systems.**

[https://bair.berkeley.edu/blog/2024/02/18/  
compound-ai-systems/](https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/).



Zhang, L., Hosseini, A., Bansal, H., Kazemi, M., Kumar, A., and  
Agarwal, R. (2024).

**Generative verifiers: Reward modeling as next-token prediction.**