# CSCI 5541: Natural Language Processing

**Lecture 20: Reasoning**
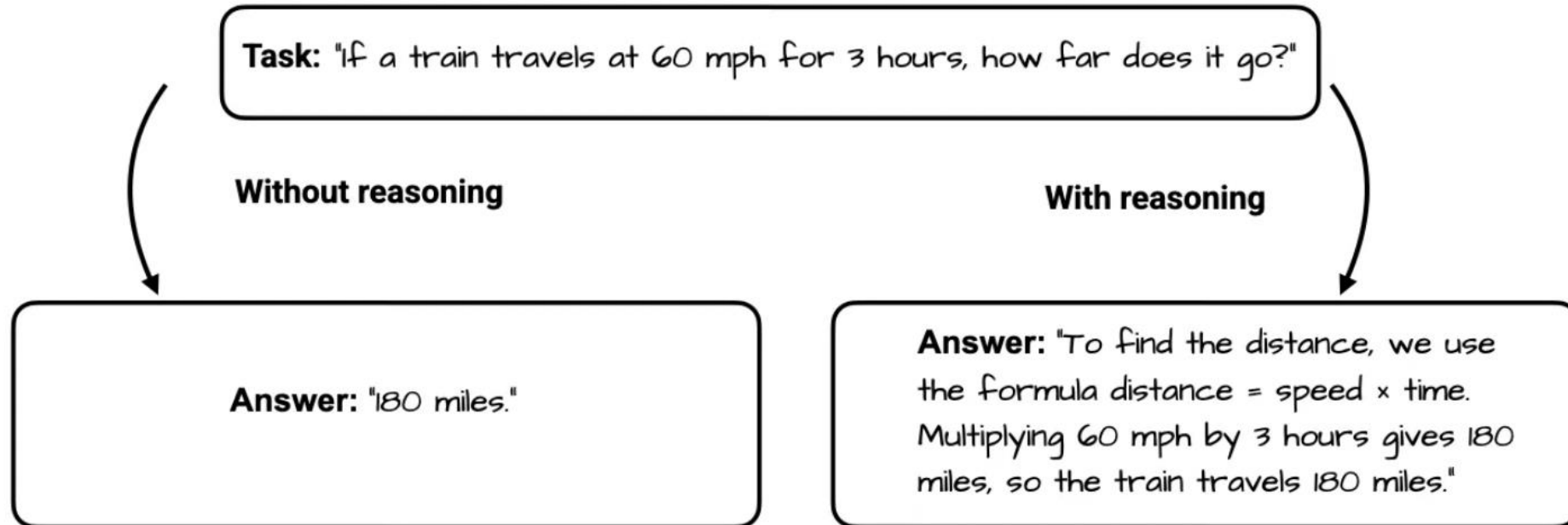
Slides borrowed from Sean Welleck (CMU), and Sebastian Raschka (UW-Madison)

UNIVERSITY OF MINNESOTA
Driven to Discover®

# Announcements (4/22)

❑ Presentations (see canvas/slack after class for more details)
- o Location (Shepherd Drone Lab)
- o Poster Format and Printing Information
- o Presentation Dates
  - ✔ Group B → April 29
  - ✔ Group A → May 1
- o All students *must* show up for both days to give feedback to other students

❑ HW5 Due Today (4/22)

❑ HW6 Out (Due May 5)

# What is Reasoning



Task: "If a train travels at 60 mph for 3 hours, how far does it go?"

**Without reasoning**

**With reasoning**

Answer: "180 miles."

Answer: "To find the distance, we use the formula distance = speed × time. Multiplying 60 mph by 3 hours gives 180 miles, so the train travels 180 miles."
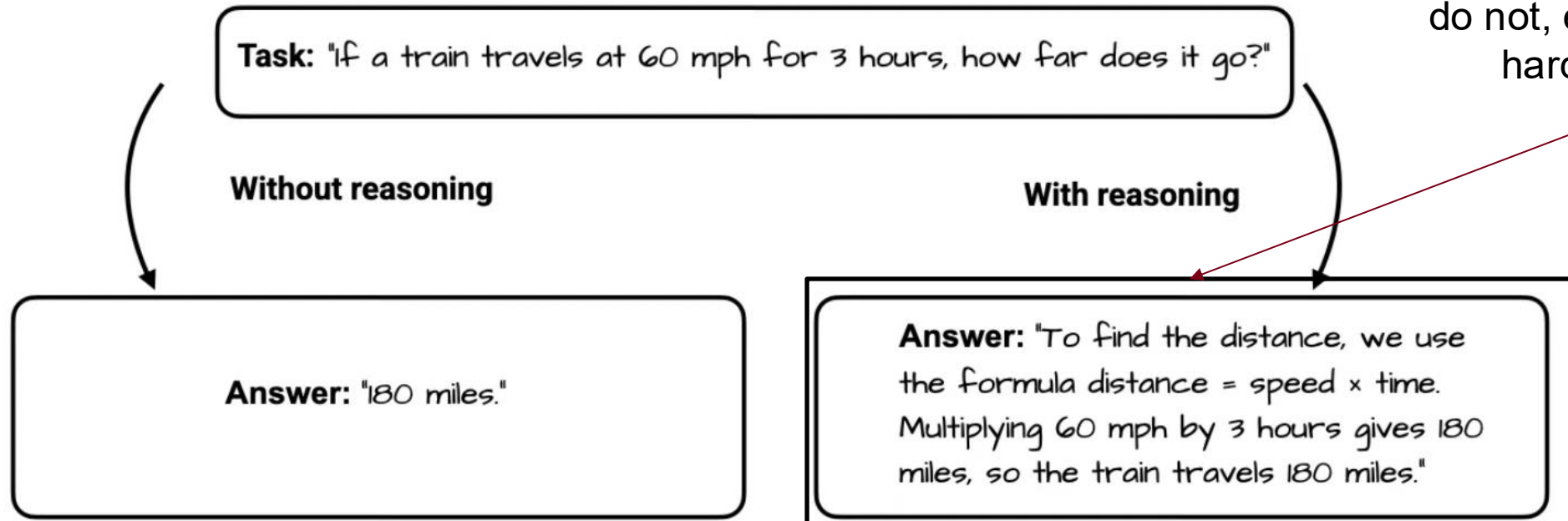
Sebastian Raschka

The State of LLM Reasoning Model Inference (Raschka, 2025)

# What is Reasoning

Models that use intermediate steps to answer questions oftentimes perform better than those which do not, especially on harder tasks
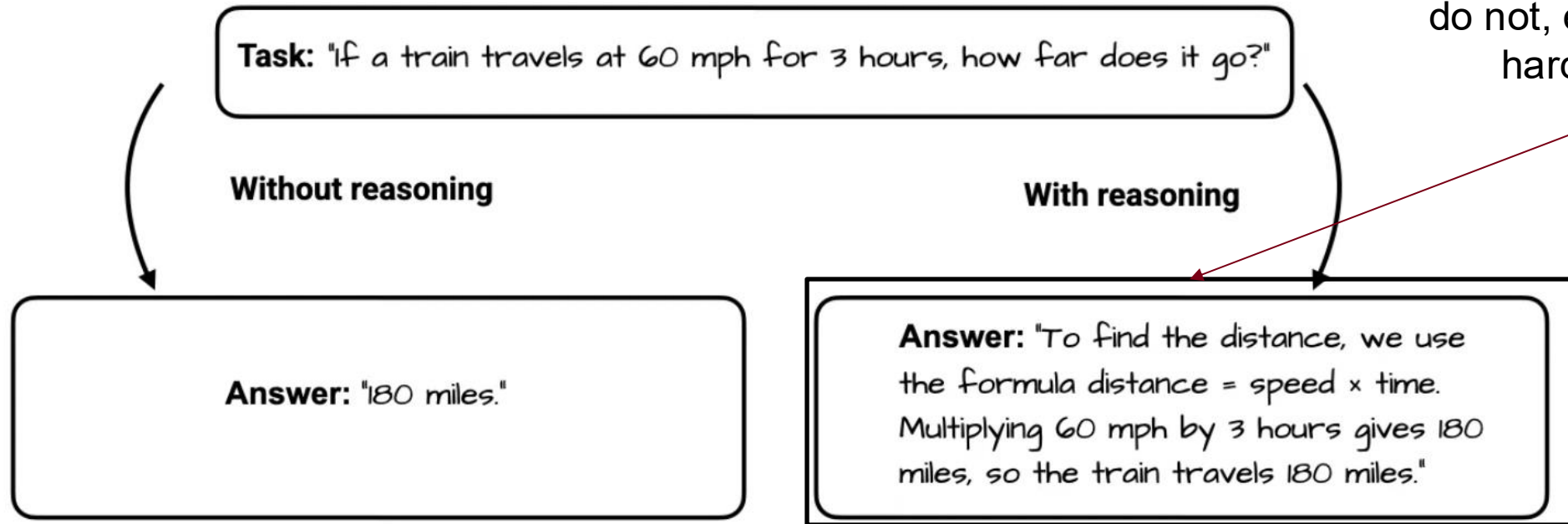
**Task:** "If a train travels at 60 mph for 3 hours, how far does it go?"

**Without reasoning**

**With reasoning**

**Answer:** "180 miles."

**Answer:** "To find the distance, we use the formula distance = speed × time. Multiplying 60 mph by 3 hours gives 180 miles, so the train travels 180 miles."

Sebastian Raschka

The State of LLM Reasoning Model Inference (Raschka, 2025)

# What is Reasoning

Models that use intermediate steps to answer questions oftentimes perform better than those which do not, especially on harder tasks
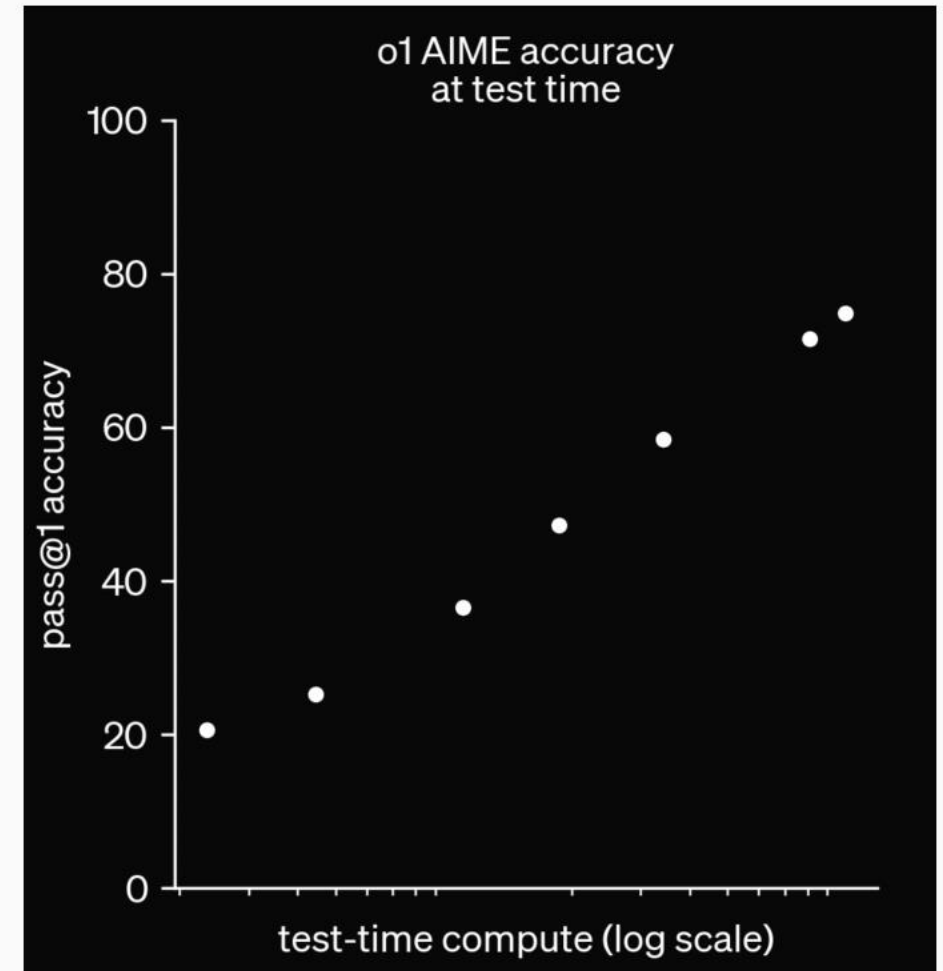


**Task:** "If a train travels at 60 mph for 3 hours, how far does it go?"

**Without reasoning**

**With reasoning**

**Answer:** "180 miles."

**Answer:** "To find the distance, we use the formula distance = speed × time. Multiplying 60 mph by 3 hours gives 180 miles, so the train travels 180 miles."

Sebastian Raschka

How can we get models to reason more in order to get to the right solution?

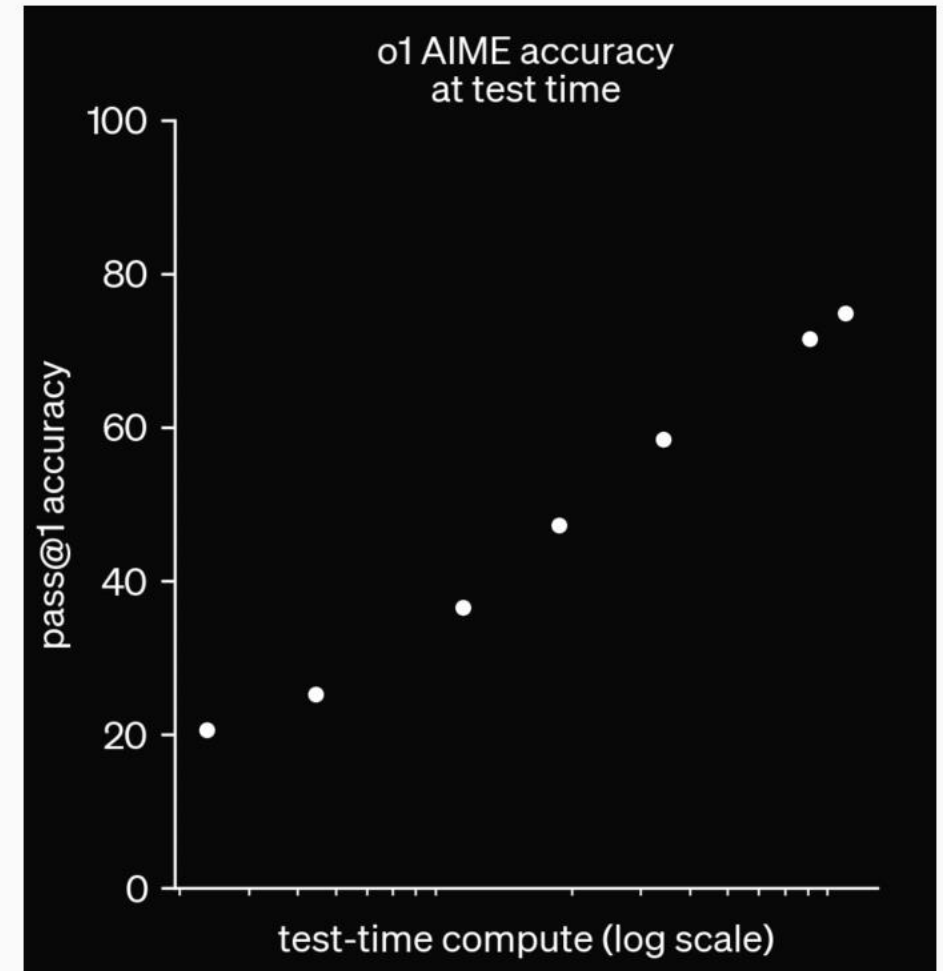The State of LLM Reasoning Model Inference (Raschka, 2025)

# What is Reasoning

❑ In a more general sense, reasoning increases the number of generated tokens to improve model performance

❑ We examine several approaches that achieve significant performance improvement by increasing this 'reasoning' (i.e. increasing the number of tokens generated)
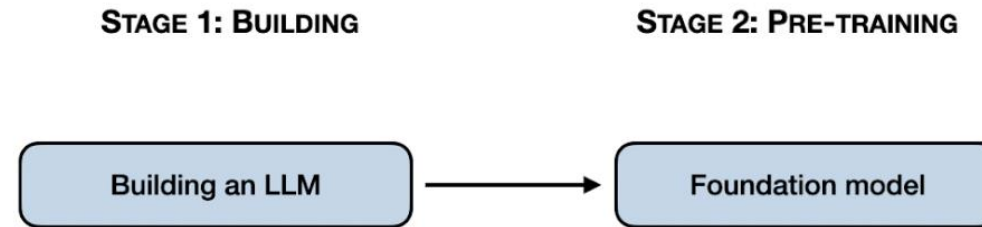


o1 AIME accuracy at test time

Compute ∝ Model size × Generated tokens

# What is Reasoning

❑ In a more general sense, reasoning increases the number of generated tokens to improve model performance

❑ We examine several approaches that achieve significant performance improvement by increasing this 'reasoning' (i.e. increasing the number of tokens generated)



o1 AIME accuracy at test time

pass@1 accuracy vs test-time compute (log scale)

$$\text{Compute} \propto \text{Model size} \times \text{Generated tokens}$$

# Different Kinds of Reasoning

**STAGE 1: BUILDING**

Building an LLM

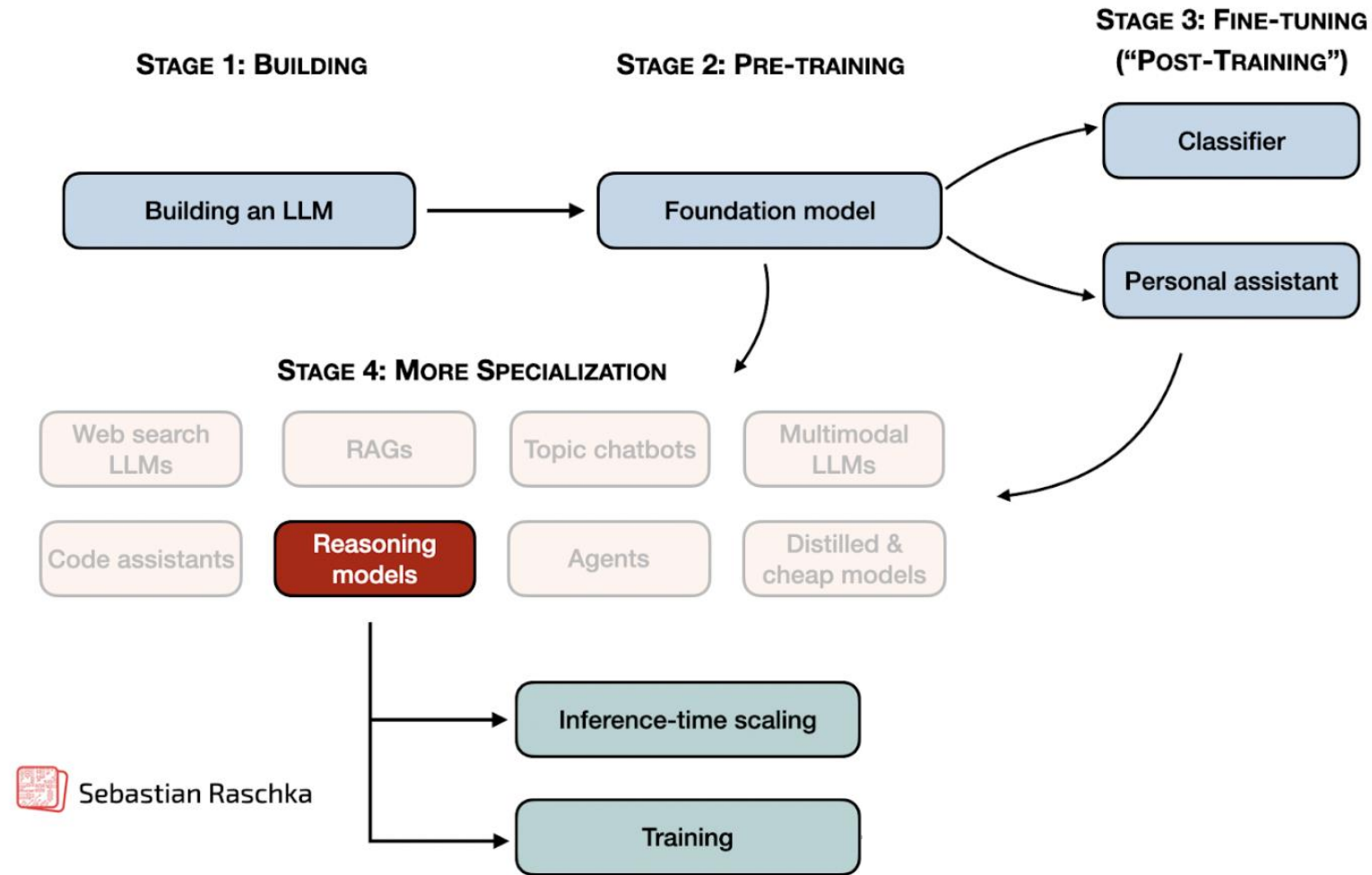# Different Kinds of Reasoning

**STAGE 1: BUILDING**　　　　**STAGE 2: PRE-TRAINING**

Building an LLM　→　Foundation model

# Different Kinds of Reasoning

STAGE 3: FINE-TUNING
("POST-TRAINING")

STAGE 1: BUILDING

STAGE 2: PRE-TRAINING

Building an LLM → Foundation model

Classifier

Personal assistant

# Different Kinds of Reasoning



**STAGE 1: BUILDING**

Building an LLM

**STAGE 2: PRE-TRAINING**

Foundation model

**STAGE 3: FINE-TUNING ("POST-TRAINING")**

Classifier

Personal assistant

**STAGE 4: MORE SPECIALIZATION**

| | | | |
|---|---|---|---|
| Web search LLMs | RAGs | Topic chatbots | Multimodal LLMs |
| Code assistants | Reasoning models | Agents | Distilled & cheap models |

# Different Kinds of Reasoning



**STAGE 1: BUILDING**

**STAGE 2: PRE-TRAINING**

**STAGE 3: FINE-TUNING ("POST-TRAINING")**

Building an LLM → Foundation model

Classifier

Personal assistant

**STAGE 4: MORE SPECIALIZATION**

Web search LLMs

RAGs

Topic chatbots

Multimodal LLMs

Code assistants

Reasoning models

Agents

Distilled & cheap models

Sebastian Raschka

Inference-time scaling

Training

# Different Kinds of Reasoning



STAGE 1: BUILDING

STAGE 2: PRE-TRAINING

STAGE 3: FINE-TUNING ("POST-TRAINING")

Building an LLM → Foundation model → Classifier / Personal assistant

STAGE 4: MORE SPECIALIZATION

Web search LLMs | RAGs | Topic chatbots | Multimodal LLMs

Code assistants | **Reasoning models** | Agents | Distilled & cheap models

Inference-time scaling

Training

Sebastian Raschka

# Different Kinds of Reasoning

❑ Test-time Scaling
  ○ Parallel
  ○ Tree-Search
  ○ Refinement

❑ Training Reasoners
  ○ Revisiting RLHF & PPO
  ○ From PPO to GRPO
  ○ RLHF to RLVR
  ○ Distillation from reasoning models
  ○ DeepSeek deepdive

❑ Latent Space Reasoning
  ○ Methods (Inner Thinking/Latent Reasoning/CoCoNut)

# Test-Time Scaling

# Test-Time Scaling

# Test-Time Scaling (Prompting)

**Regular prompting**

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: The answer (arabic numerals) is

*(Output) 8* **X**

**Chain-of-thought prompting**

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: **Let's think step by step.**

*(Output)* *There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.* ✓

# Test-Time Scaling (Prompting)



Good initial step for better results – but not as scalable. We have less control over increasing the number of generated tokens.

# Test-Time Scaling Overview

❏ Test-Time Scaling Approaches
  - Parallel
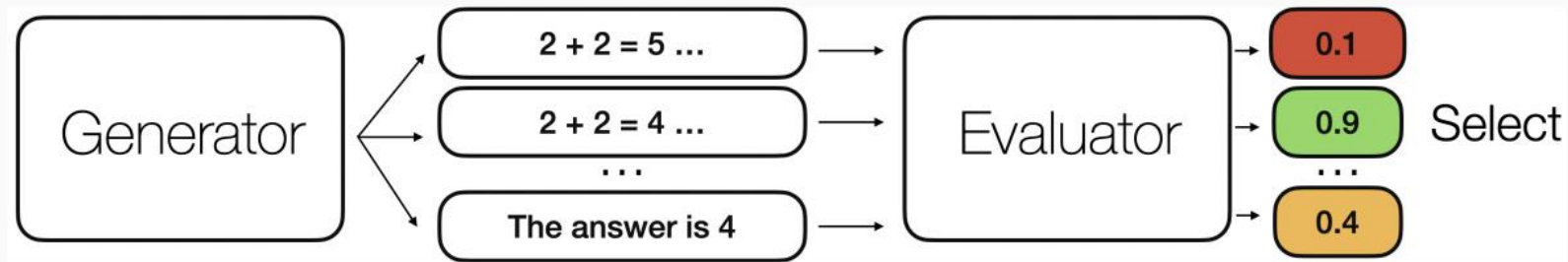  - Tree Search
  - Refinement

# Test-Time Scaling

# Test-Time Scaling



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

# Test-Time Scaling
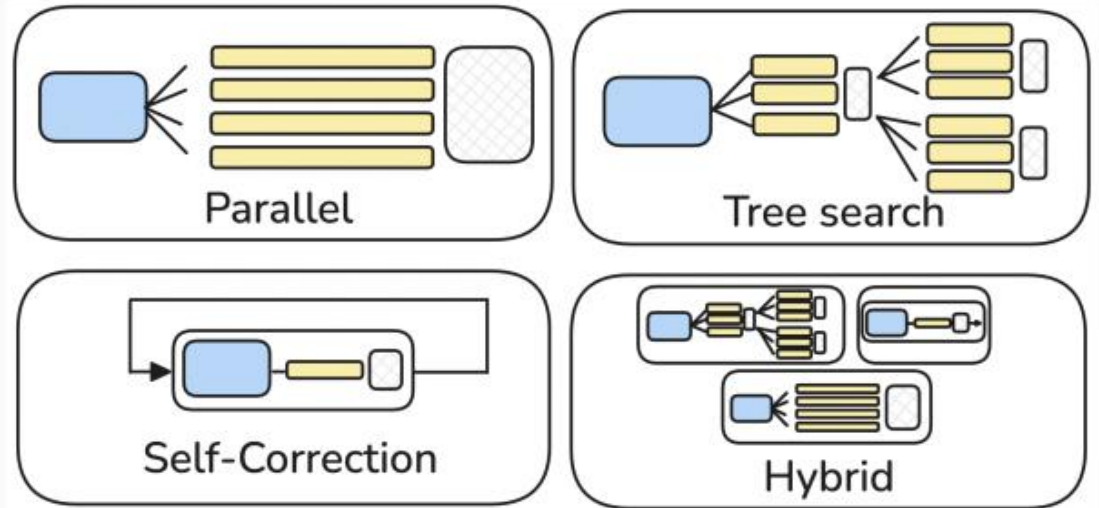
**Meta-generator:** Strategies for calling a generator multiple times



- Example: call API multiple times, select the best sequence with a separate model
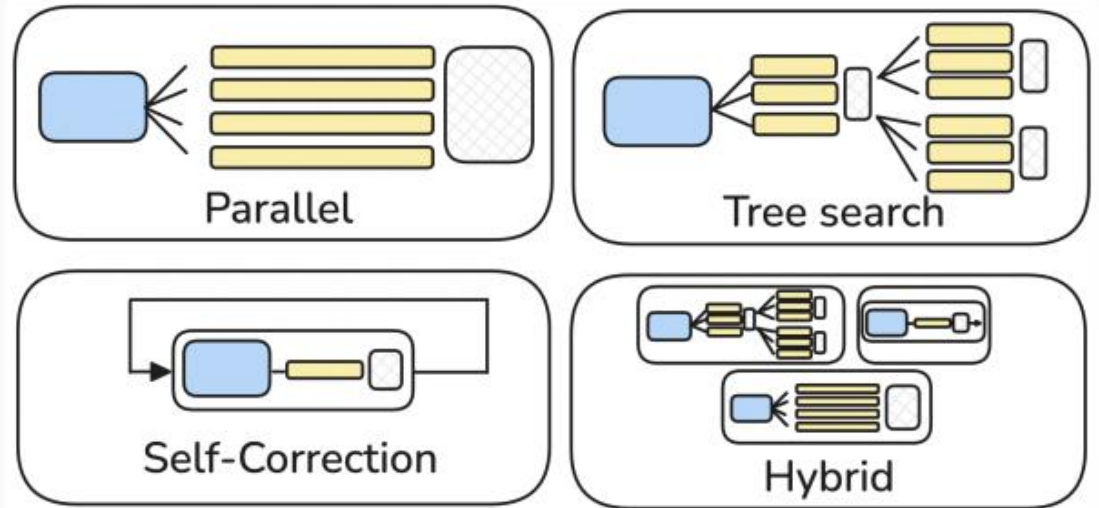
# Test-Time Scaling Approaches

- **Strategies**
  - Parallel
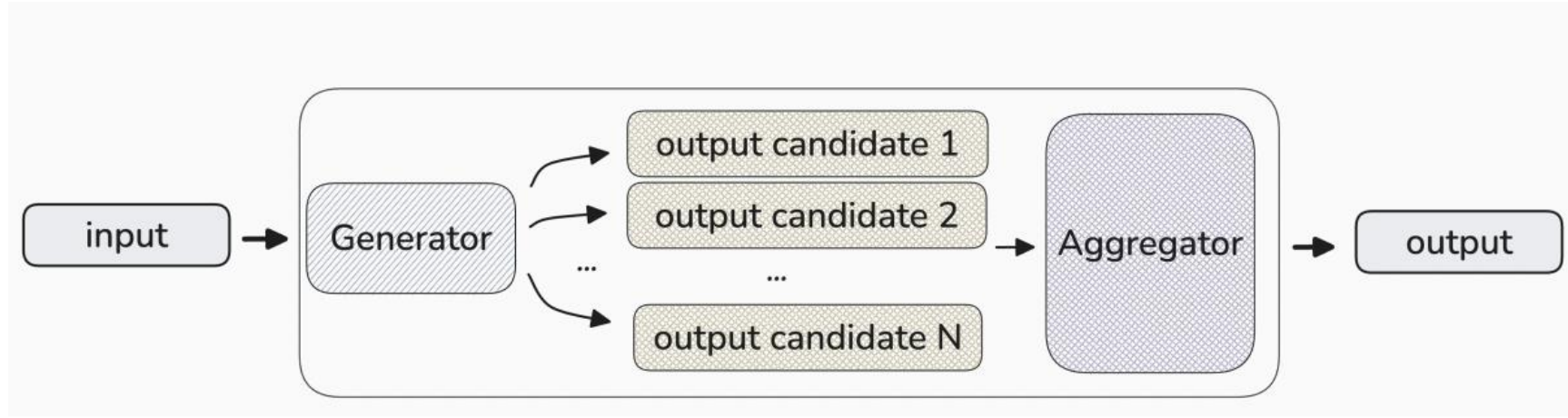  - Tree search
  - Refinement/self-correction

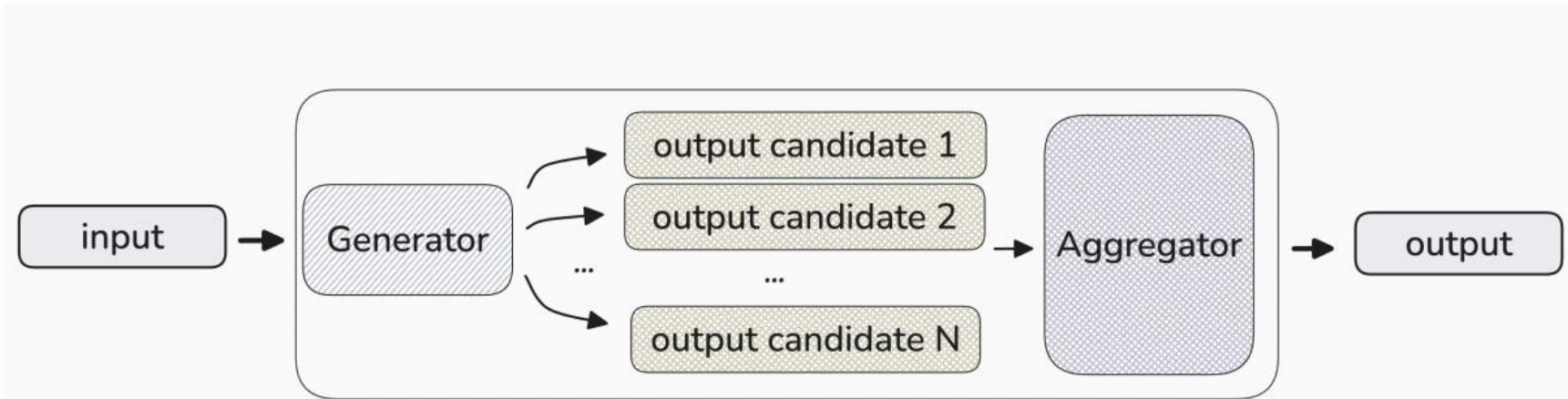# Test-Time Scaling Approaches

- **Strategies**
  - Parallel
  - Tree search
  - Refinement/self-correction



Parallel

Tree search

Self-Correction
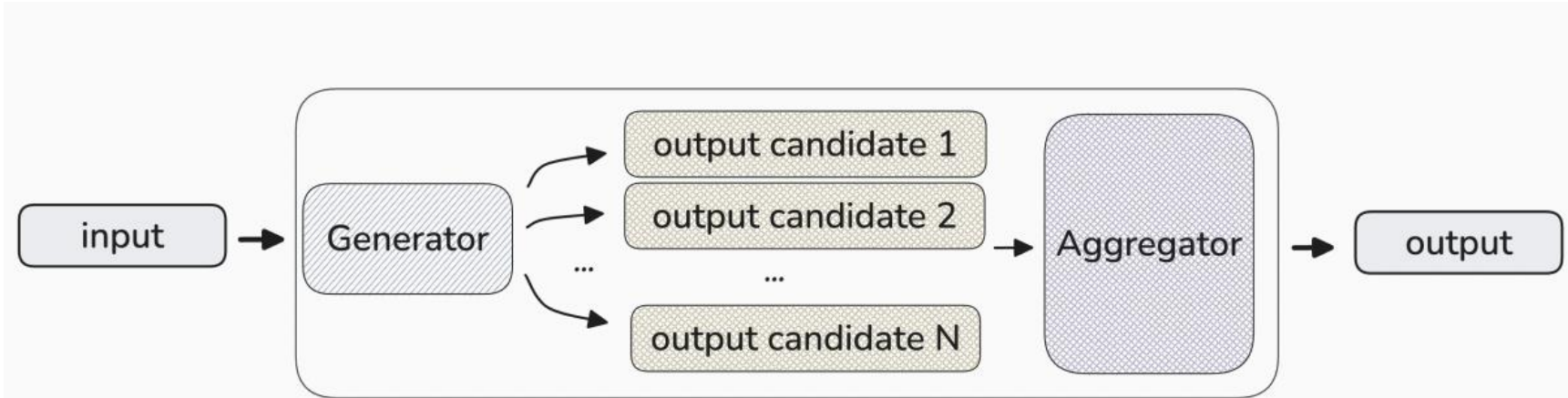
Hybrid

# Parallel Generation

# Parallel Generation



- Generate candidates:

$$\{y^{(1)}, \ldots, y^{(N)}\} \sim G(\cdot|x)$$
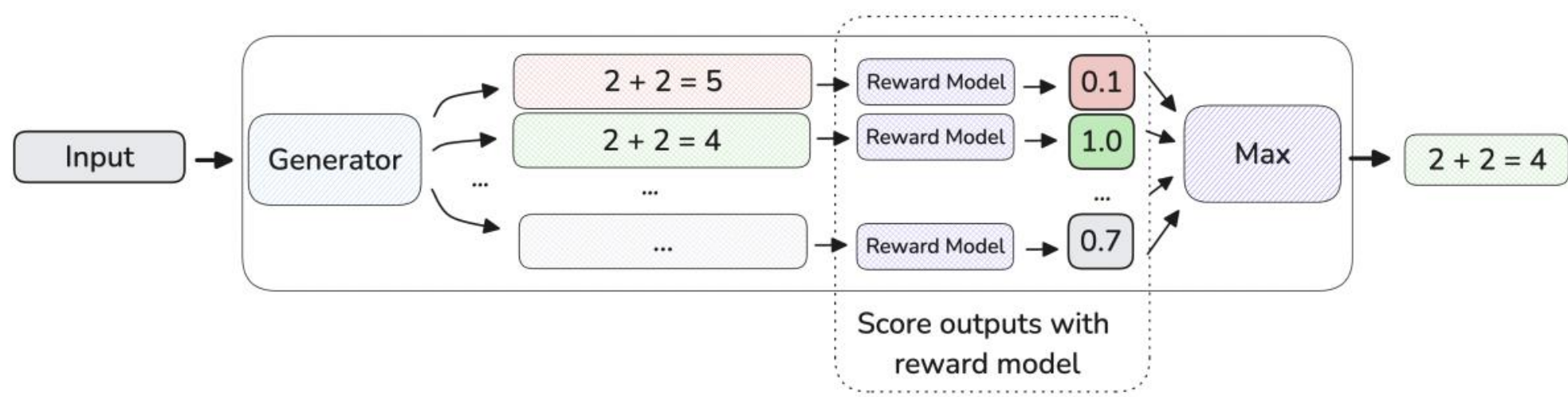
# Parallel Generation



- Generate candidates:

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

- Aggregate:

$$y = h(y^{(1)}, \dots, y^{(N)})$$
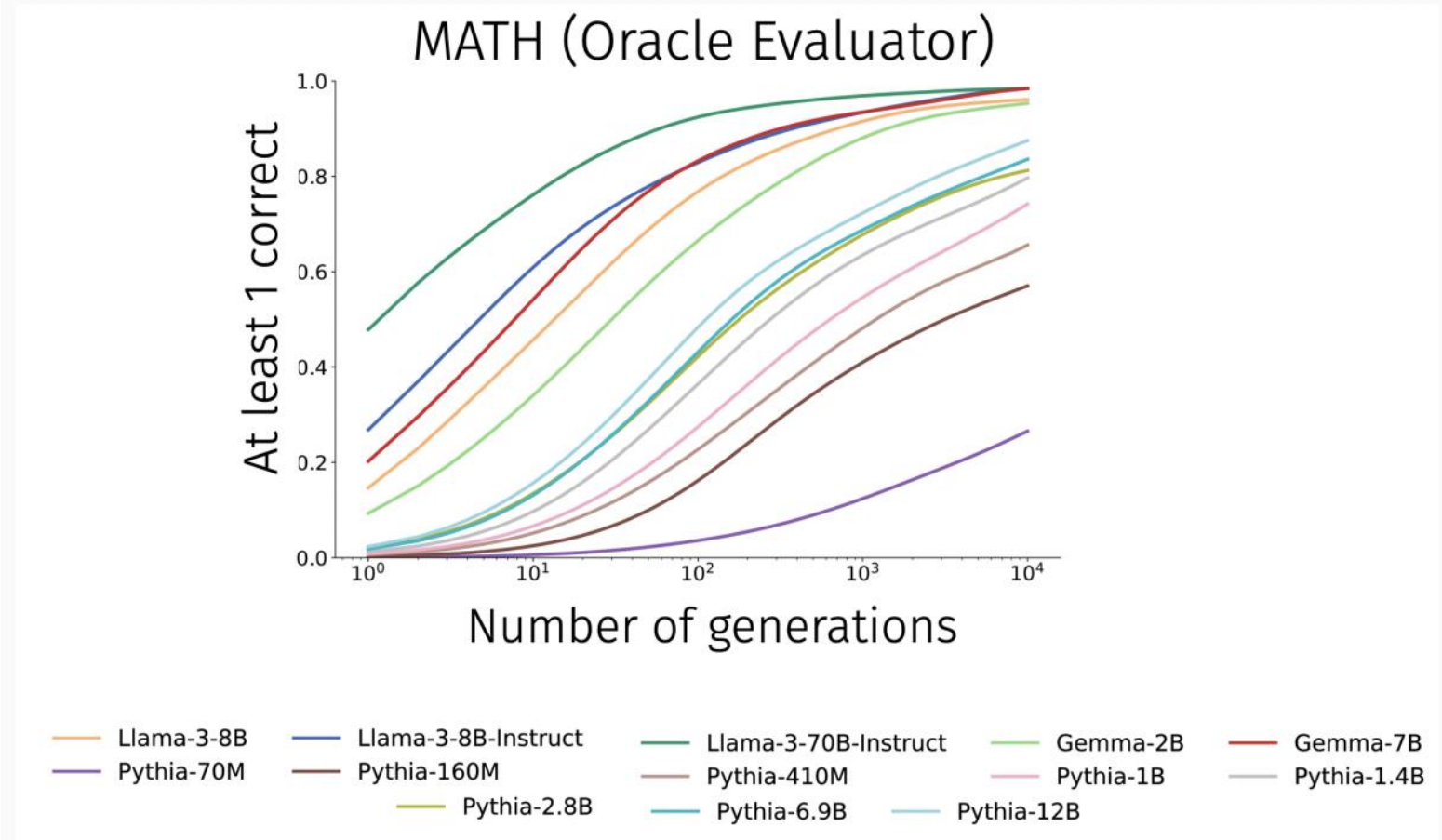
# Parallel Generation (Best-of-N)



$$\underset{\{y^{(1)},...,y^{(N)}\}}{\arg\max} \underbrace{v(y)}_{\text{reward model}}$$

Stiennon, et. al (2020), Nakano et. al (2022)
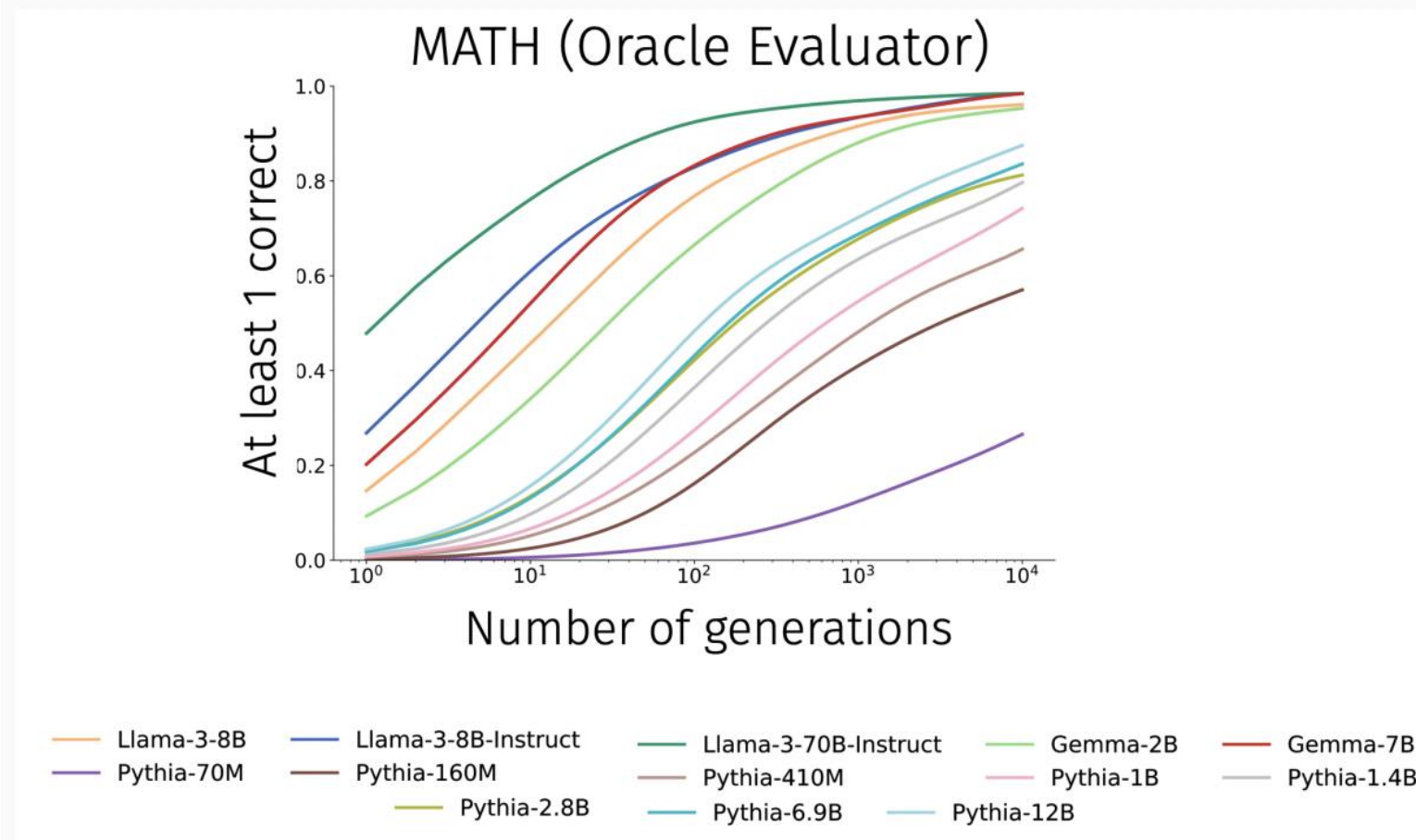
# Parallel Generation (Best-of-N)

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



MATH (Oracle Evaluator)

At least 1 correct vs Number of generations

Llama-3-8B
Llama-3-8B-Instruct
Llama-3-70B-Instruct
Gemma-2B
Gemma-7B
Pythia-70M
Pythia-160M
Pythia-410M
Pythia-1B
Pythia-1.4B
Pythia-2.8B
Pythia-6.9B
Pythia-12B

# Parallel Generation (Best-of-N)

Q: What is a 'perfect reward model'?

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



MATH (Oracle Evaluator)

Legend:
- Llama-3-8B
- Llama-3-8B-Instruct
- Llama-3-70B-Instruct
- Gemma-2B
- Gemma-7B
- Pythia-70M
- Pythia-160M
- Pythia-410M
- Pythia-1B
- Pythia-1.4B
- Pythia-2.8B
- Pythia-6.9B
- Pythia-12B

# Parallel Generation (Best-of-N)

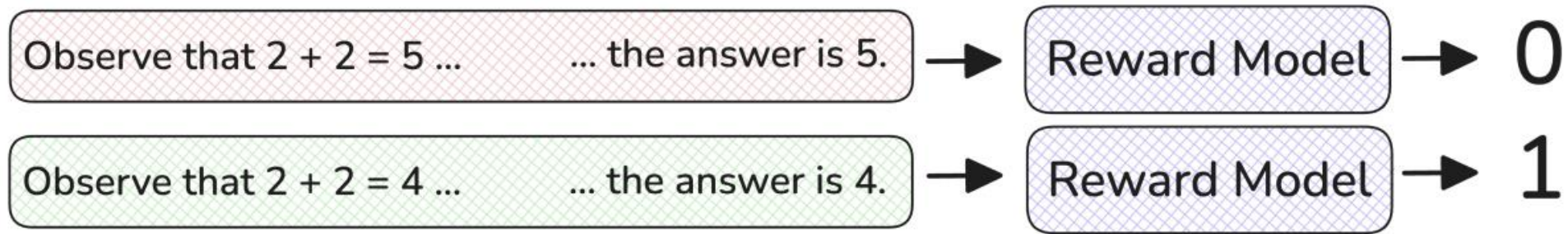Q: What is a 'perfect reward model'?

A: A 'perfect reward model' allows us to always select the correct answer if it is given.

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



MATH (Oracle Evaluator)

At least 1 correct vs Number of generations

Legend:
- Llama-3-8B
- Llama-3-8B-Instruct
- Llama-3-70B-Instruct
- Gemma-2B
- Gemma-7B
- Pythia-70M
- Pythia-160M
- Pythia-410M
- Pythia-1B
- Pythia-1.4B
- Pythia-2.8B
- Pythia-6.9B
- Pythia-12B

# Best of N: Finding a reward model in practice



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

Observe that 2 + 2 = 5 ... ... the answer is 5. → Reward Model → 0

Observe that 2 + 2 = 4 ... ... the answer is 4. → Reward Model → 1

Train reward model with correct and incorrect examples.[2]

Cobbe et. al (2021)

# Best of N: Finding a reward model in practice

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

Hello, you are awesome $>$ Hello, you are #&@#*@#

Train reward model with preference data.[2]

Stiennon et. al (2020)

# Why Best of N?

Why Best-of-*N*?

- Approximates maximum (true) reward:

$$\text{Best-of-}N = \underset{y\in\{y^{(1)},\ldots,y^{(N)}\}}{\arg\max}\ v(y)$$

$$\approx \underset{y}{\arg\max}\ v(y) \qquad\qquad (1)$$

$$\approx \underset{y}{\arg\max}\ R(y) \qquad\qquad (2)$$

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

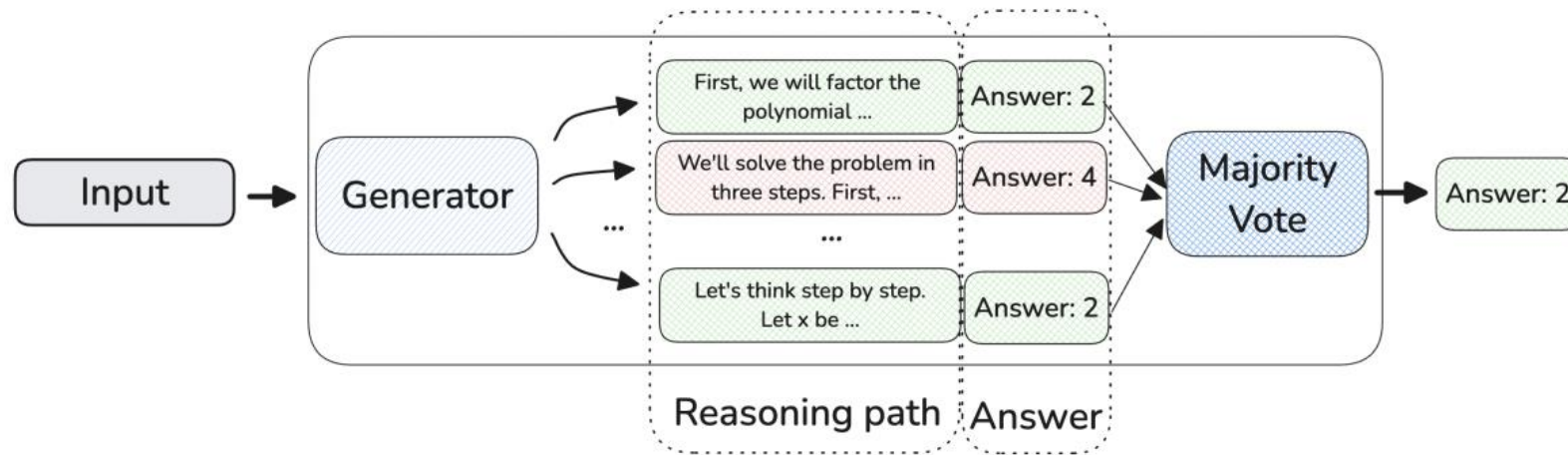(2) Suffers from imperfect reward model, aka "over-optimization"

# Best of N (example)



GSM (Learned Evaluator)

Over-optimization

Cobbe et. al (2020)

# Parallel Generation (Voting)



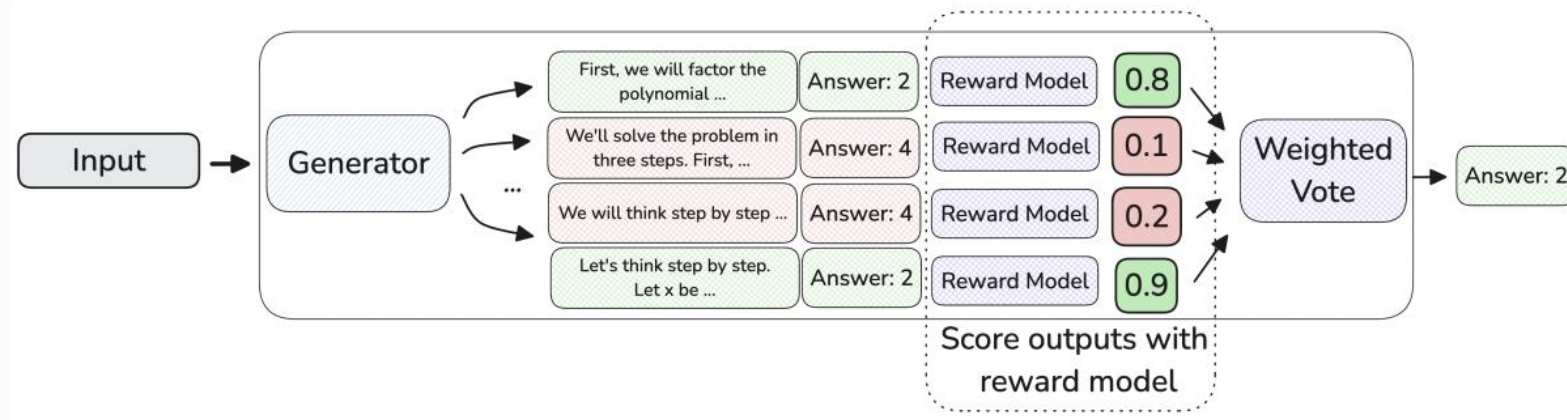Voting aggregation:[4]

Reasoning path | Answer

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

Wang et. al (2023)

# Parallel Generation (Weighted Voting)



Weighted Voting:
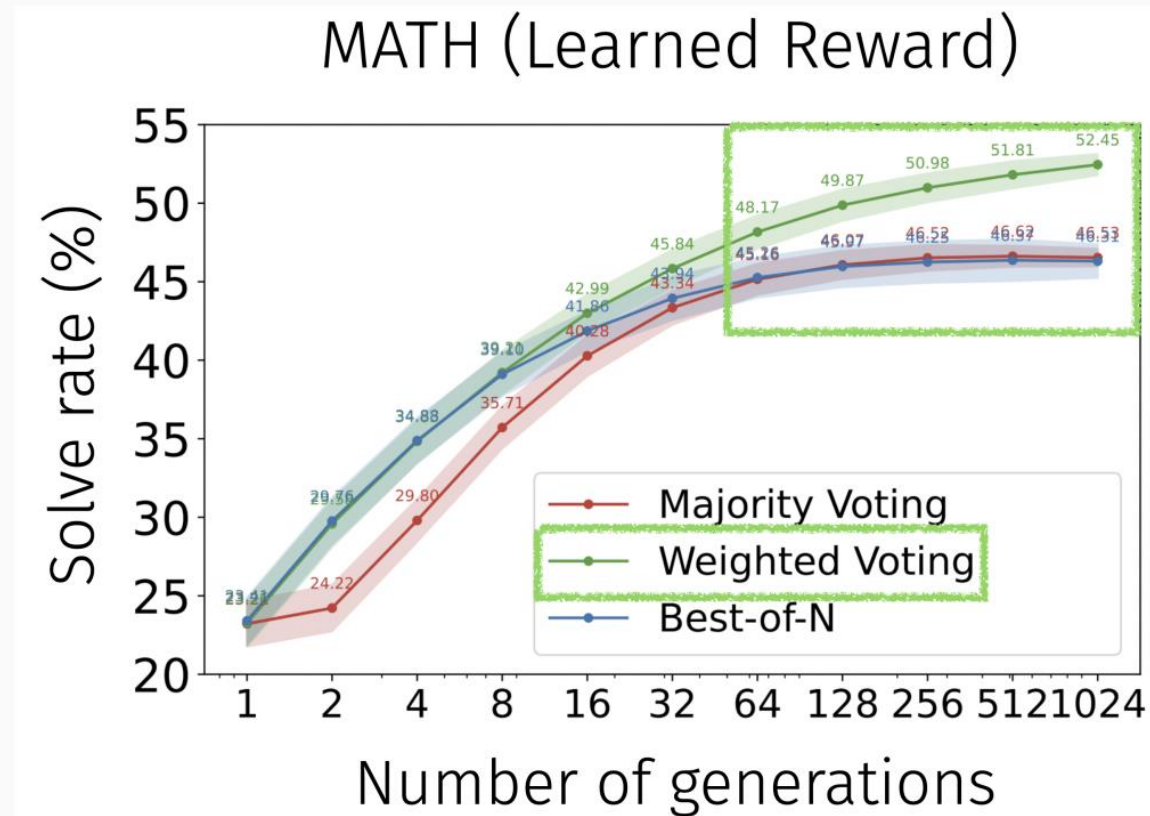
$$\arg\max_a \sum_{i=1}^{N} \underbrace{v(y^{(i)})}_{\text{reward model}} \cdot 1\{y^{(i)} = a\},$$
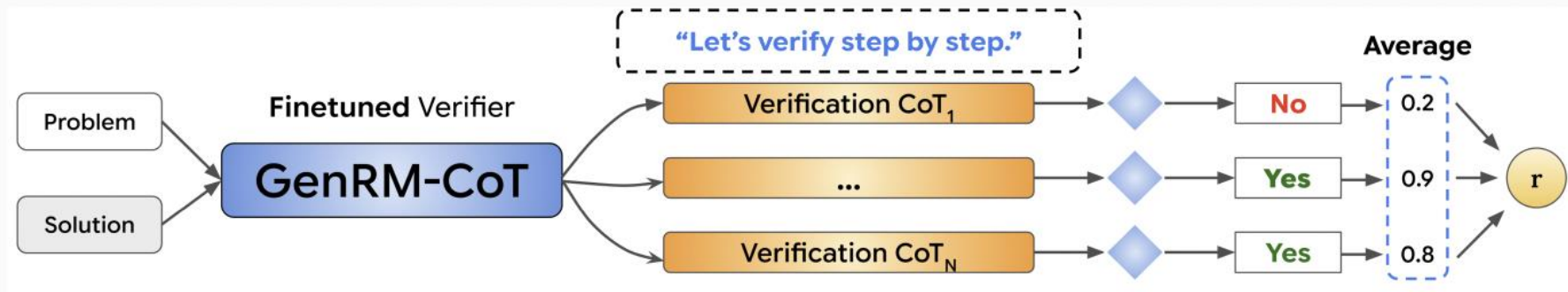
Li et. al (2023)

# Parallel Generation (Weighted Voting)



Can outperform Best-of-N, e.g.:[6]

MATH (Learned Reward)

[Sun et al., 2024] Easy-to-Hard Generalization: Scalable Alignment Beyond Human Supervision.

# Parallel Generation



Improve the reward model:

"Let's verify step by step."

Problem → Finetuned Verifier **GenRM-CoT** → Verification CoT$_1$ → No → Average 0.2

Solution → ... → Yes → 0.9 → r

Verification CoT$_N$ → Yes → 0.8

Parallel generation *in the reward model too*[8]

Active area of research!

Zhang et. al (2024)

# Parallel Generation

❏ Parallel

- o Explores output space by generating full sequences
- o Large performance gains in practice
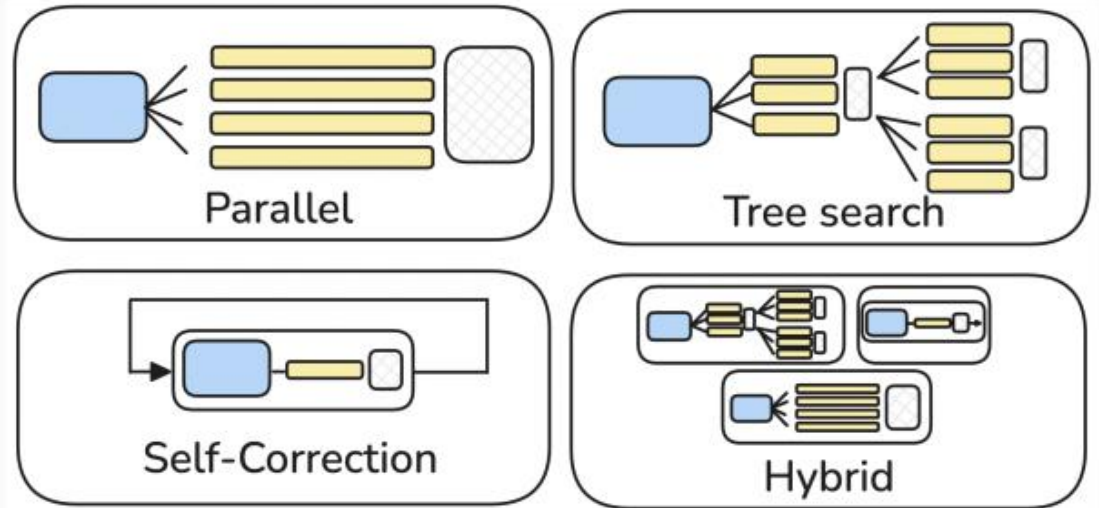- o Bounded by the quality of the evaluator and generator

❏ Insight: only uses the verifier at the end (full sequence outputs)

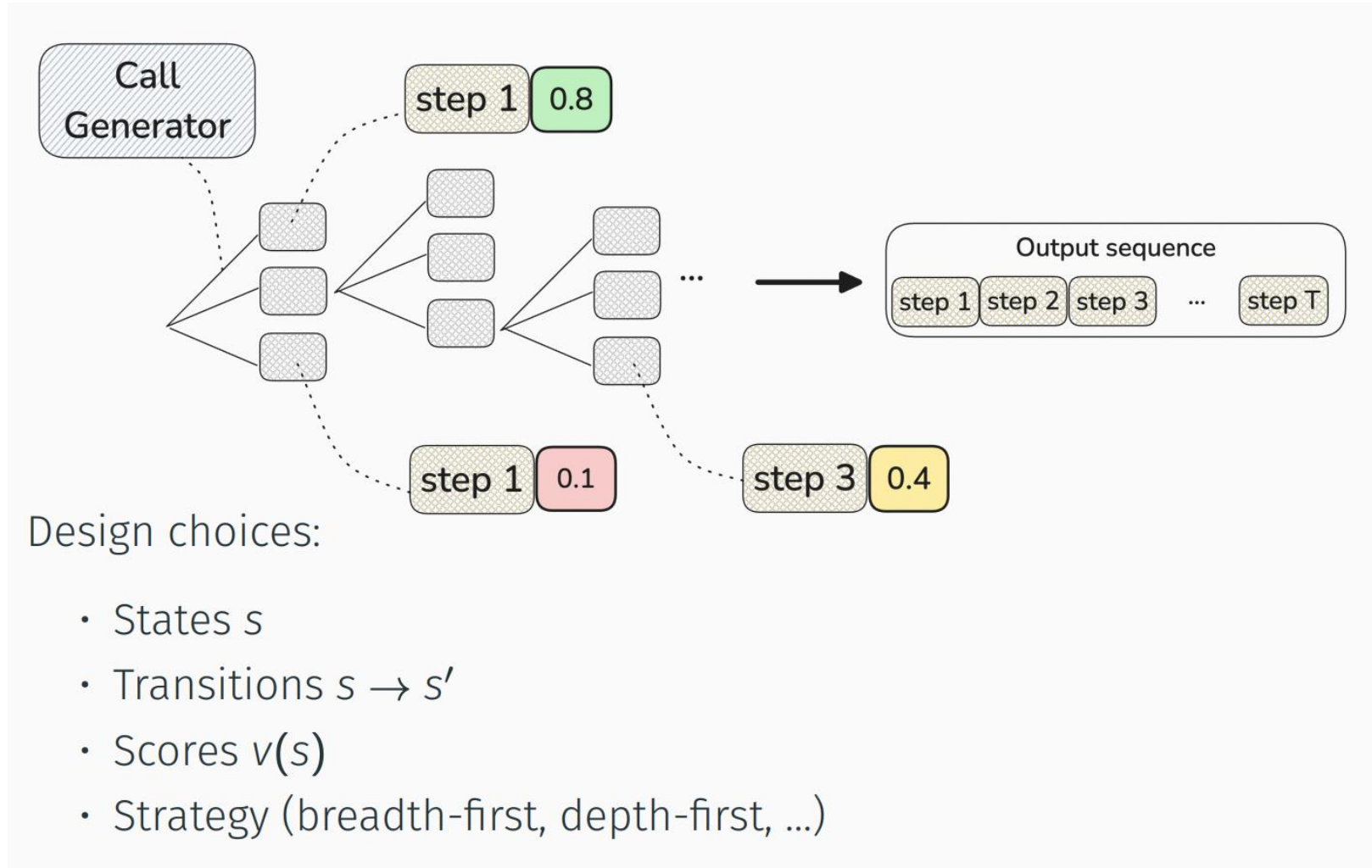- o Next: How can intermediate evaluation improve on this approach

# Test-Time Scaling Approaches

- Strategies
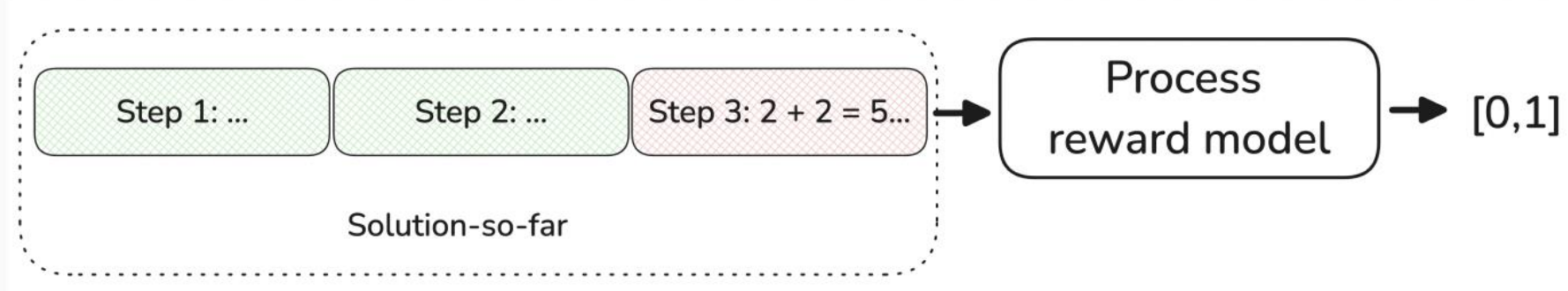  - Parallel
  - Tree search
  - Refinement/self-correction



Parallel

Tree search

Self-Correction

Hybrid

# Tree Search



Design choices:

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

# Tree Search



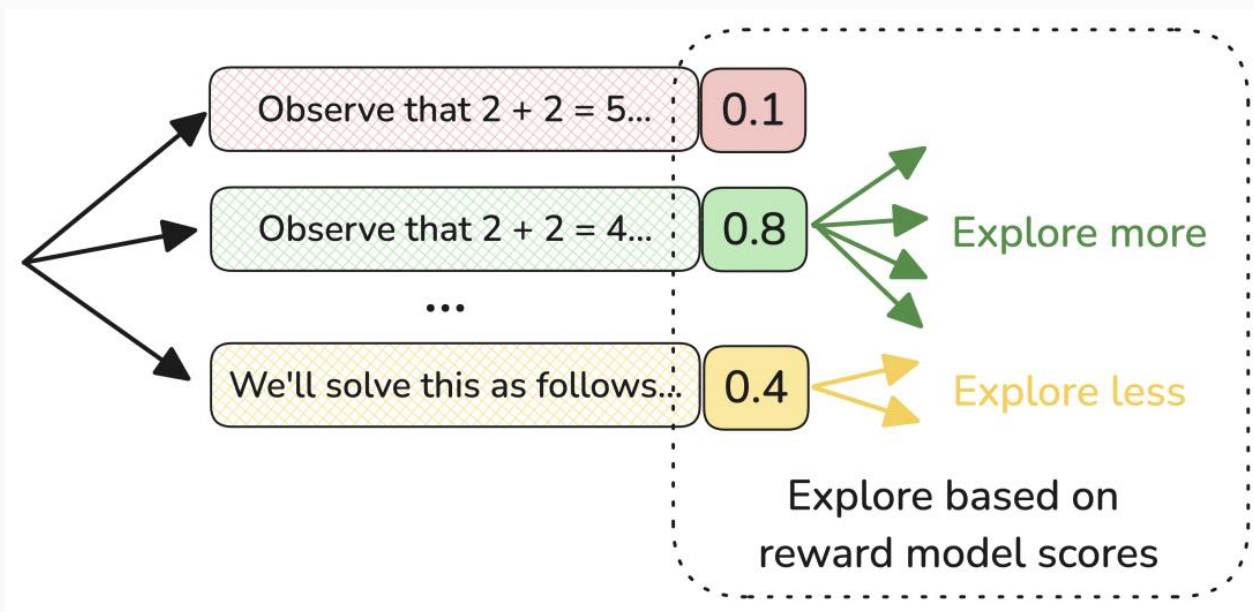1. Scores: "process reward model (PRM)"[9]

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

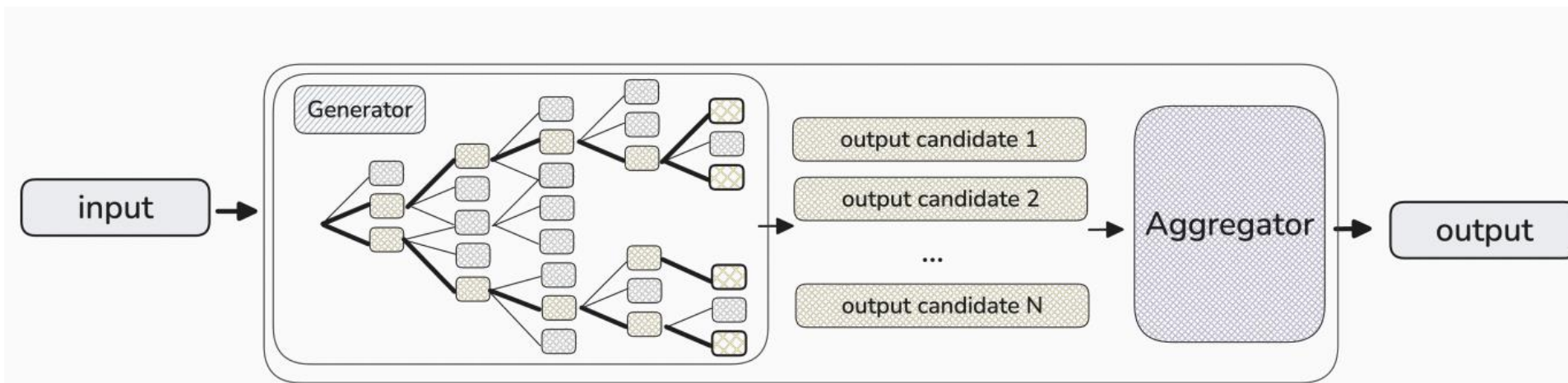[Uesato et al., 2022, Lightman et al., 2024, Wang et al., 2024a]

# Tree Search



2. Reward Balanced Search (Rebase)[10]

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

[Wu et al., 2024b] Inference Scaling Laws: An Empirical Analysis of Compute-Optimal Inference.
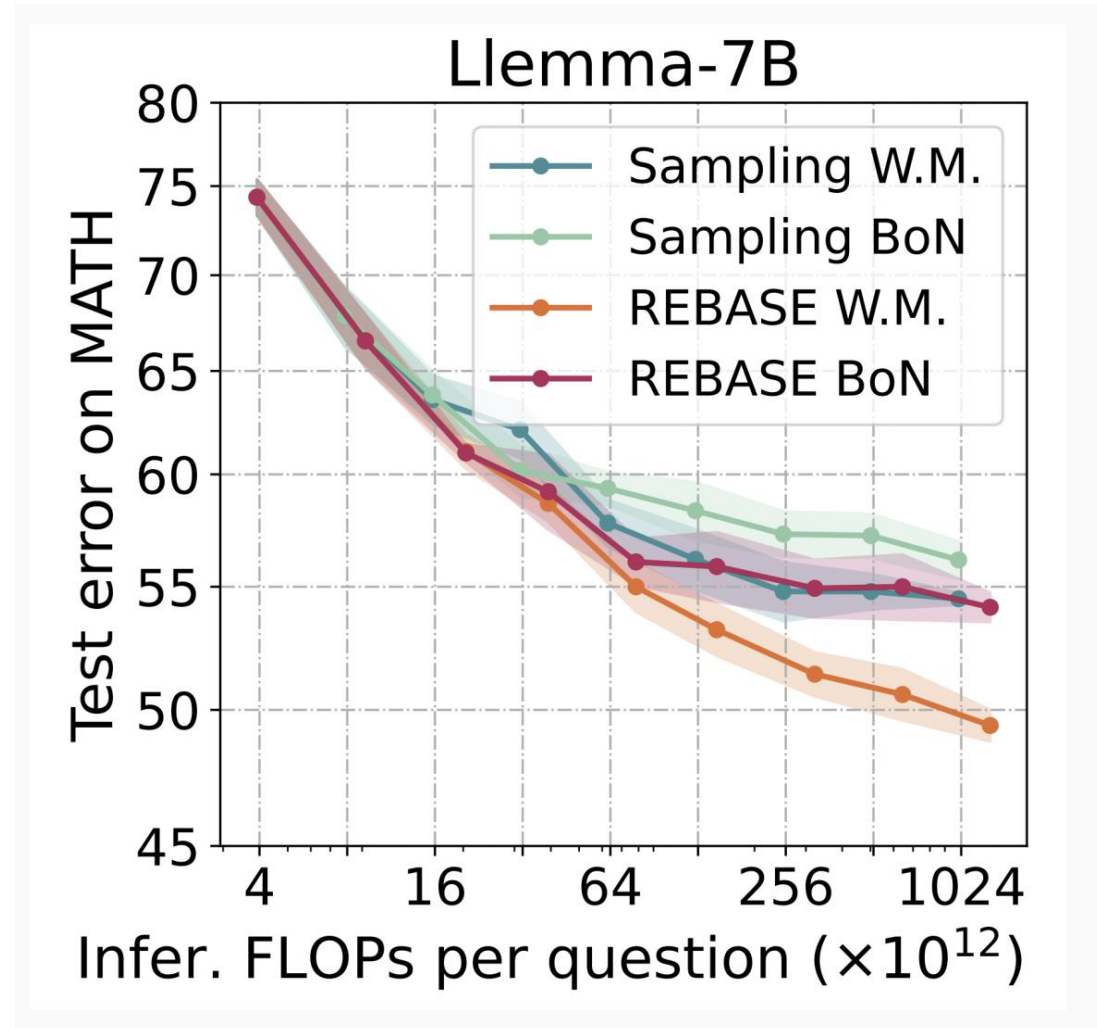
# Tree Search



Run tree search to get candidates for aggregation (e.g., voting).
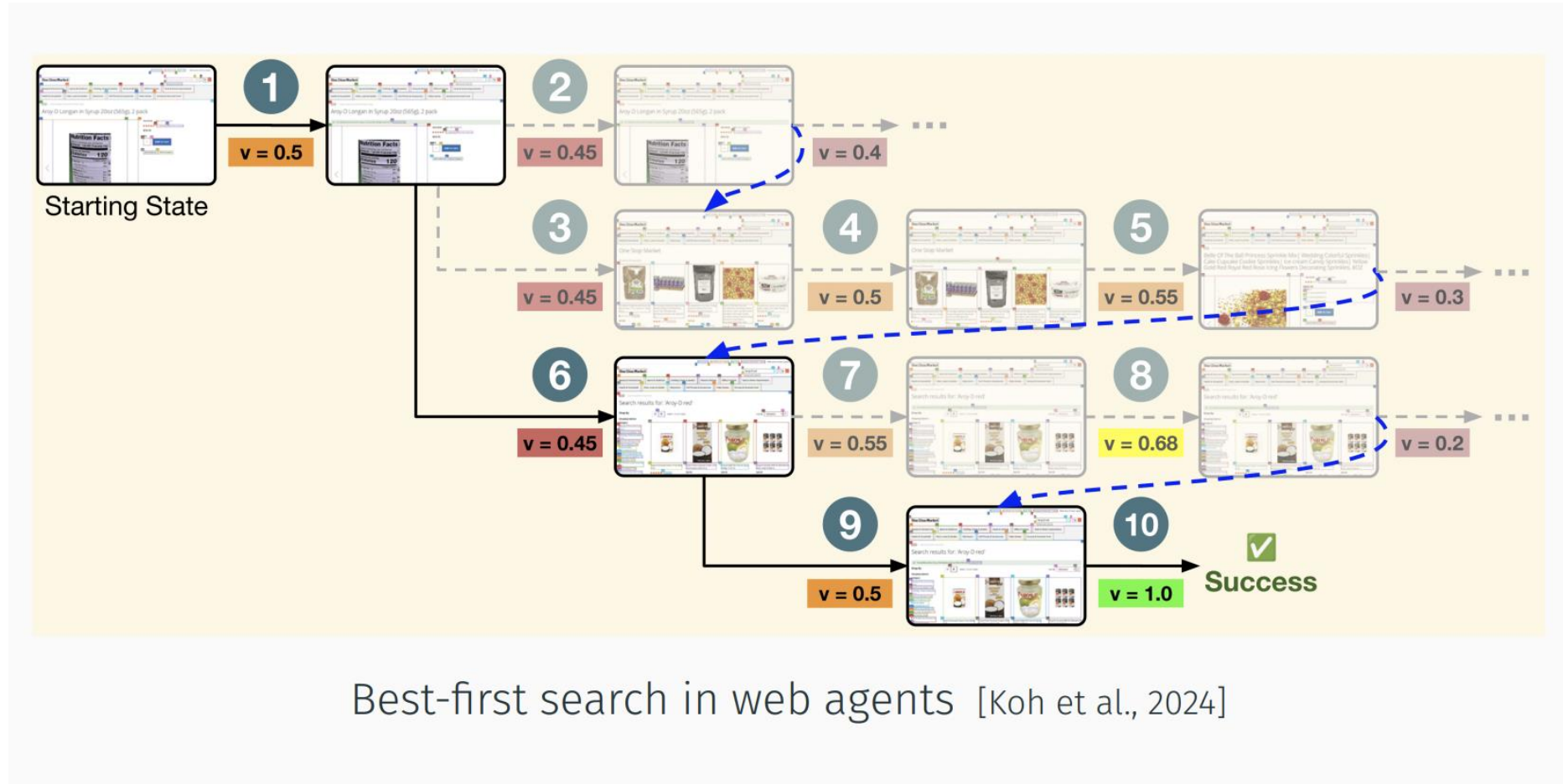
- **Key idea:** Leverages scores on *intermediate* states
    - Backtracking
    - Exploration

# Tree Search



[Wu et al., 2024b] Inference Scaling Laws: An Empirical Analysis of Compute-Optimal Inference.

# Tree Search (Example)



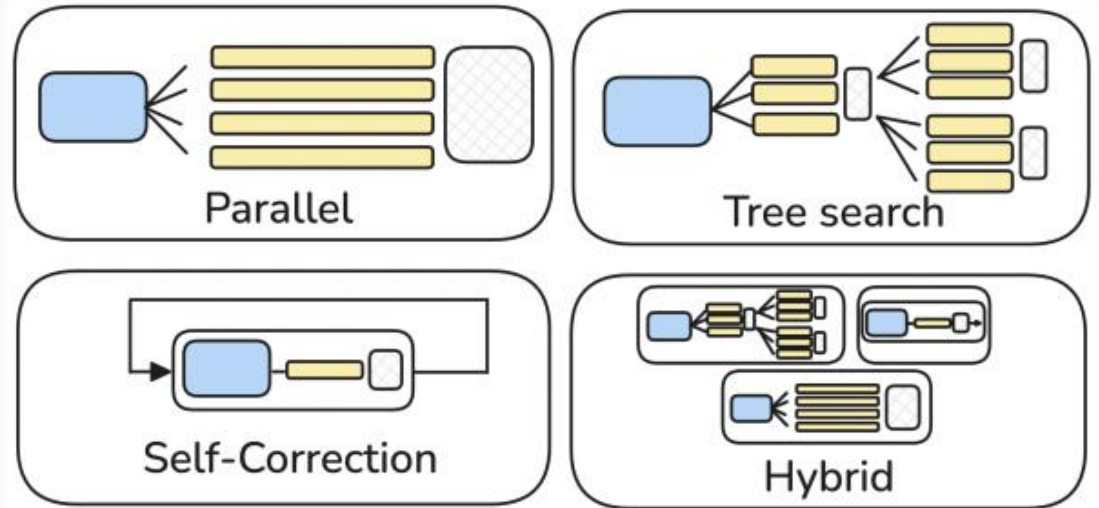Best-first search in web agents [Koh et al., 2024]

# Tree Search

❑ Can backtrack and explore using intermediate scores

❑ Requires a suitable environment and value function

- Decomposition into states
- Good reward signal

# Refinement/self-correction



- Strategies
  - Parallel
  - Tree search
  - Refinement/self-correction

Parallel

Tree search

Self-Correction

Hybrid

# Refinement/self-correction



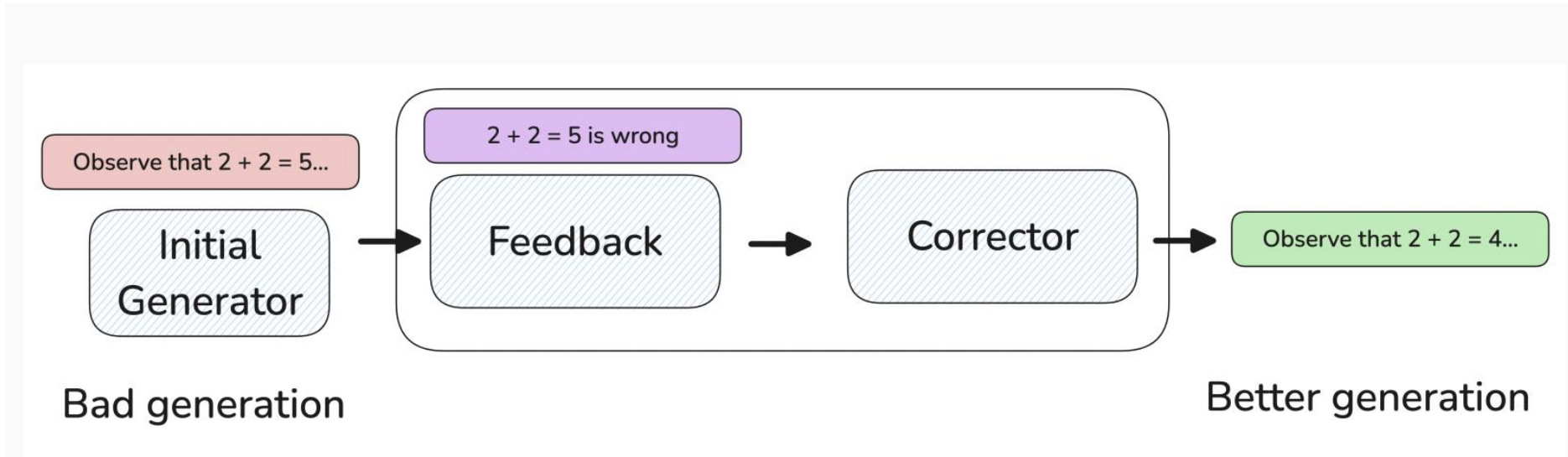Improve a generation

Repeat:

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

# Refinement/self-correction



Improve a generation using feedback

Repeat:

- $y^{(i+1)} \sim g(x, y^{(i)}, F(y^{(i)}))$
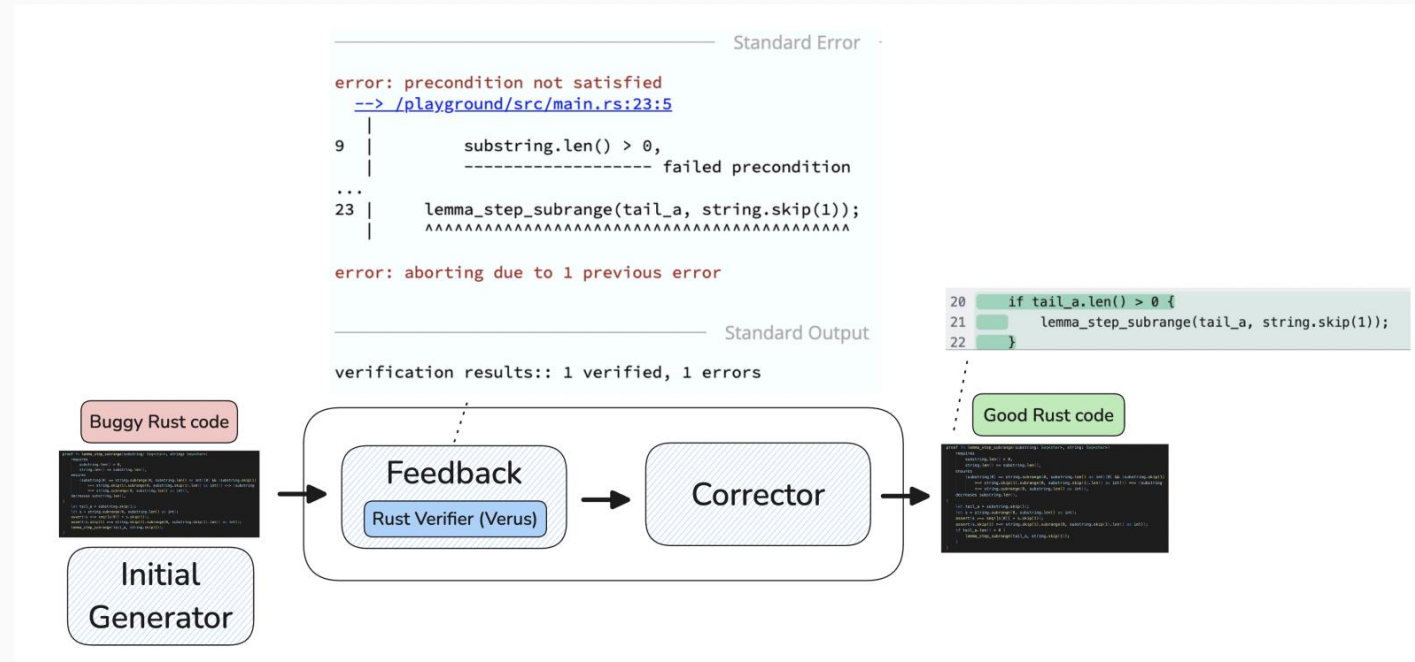
# Refinement/self-correction

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

❑ Extrinsic: external information at inference time

❑ Intrinsic: no external information at inference time

# Refinement (Extrinsic)



**1. Extrinsic**: external feedback

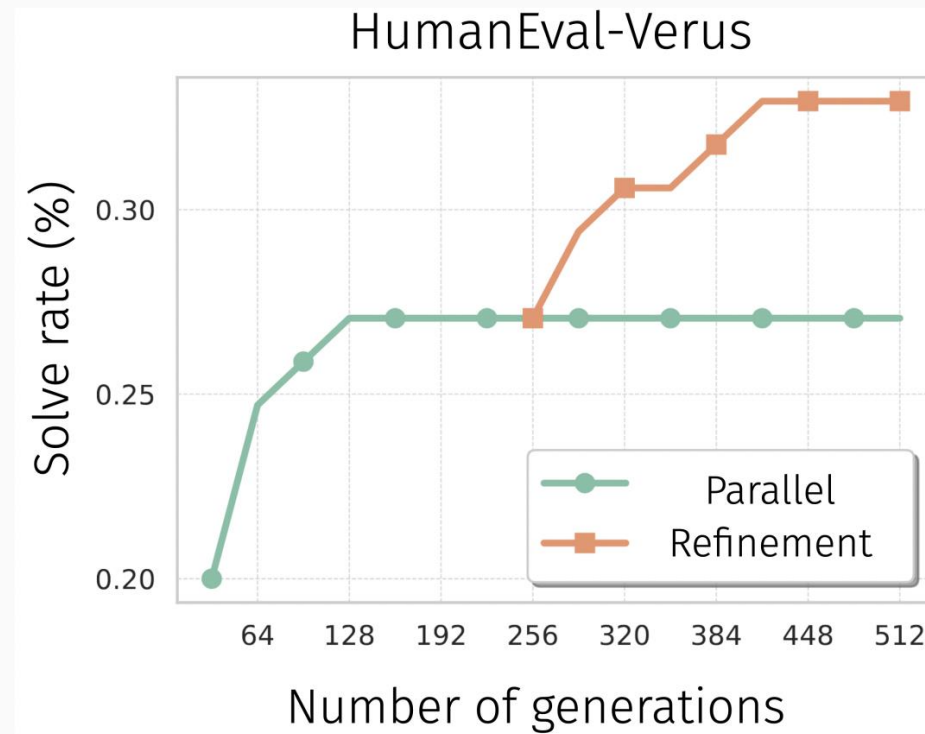Feedback: external program verifier[12]

---

[12] [Aggarwal et al., 2024], *AlphaVerus*. P. Aggarwal, B. Parno, S. Welleck.

# Refinement (Extrinsic)



**1. Extrinsic**: external feedback

HumanEval-Verus

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

# Refinement (Extrinsic)

❑ Extrinsic: External Feedback

❑ Succuess cases

  ○ Verifiers [Aggarwal et al., 2024]

  ○ Code interpreters [Chen et al., 2024]

  ○ Retrievers [Asai et al., 2024]

  ○ ...

❑ Intuition: adds new information that allows detection and localization of erros

# Refinement (Intrinsic)



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

2 + 2 = 5 is wrong

Observe that 2 + 2 = 5...

Generator
Prompt

Feedback
Prompt

Corrector
Prompt

???

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

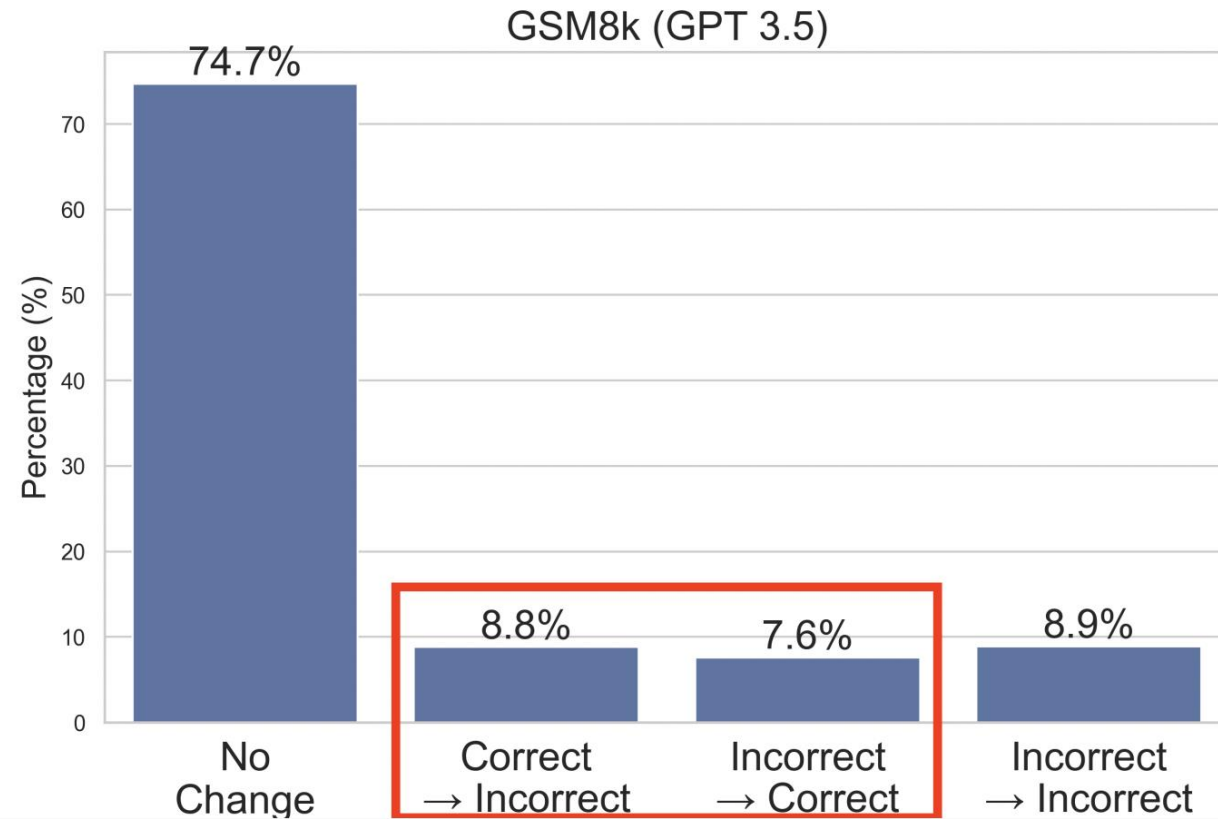# Refinement (Intrinsic)

Mixed results

❑ Easy to evaluate tasks: positive [Wang et al., 2024]

❑ Mathematical reasoning: mixed [Huang et al., 2024]

# Refinement (Intrinsic)



GSM8k (GPT 3.5)

74.7% — No Change
8.8% — Correct → Incorrect
7.6% — Incorrect → Correct
8.9% — Incorrect → Incorrect

Percentage (%)

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

# Refinement

Refinement / self-correction

❑ Extrinsic

  ○ Positive results for environments that detect or localize errors

❑ Intrinsic

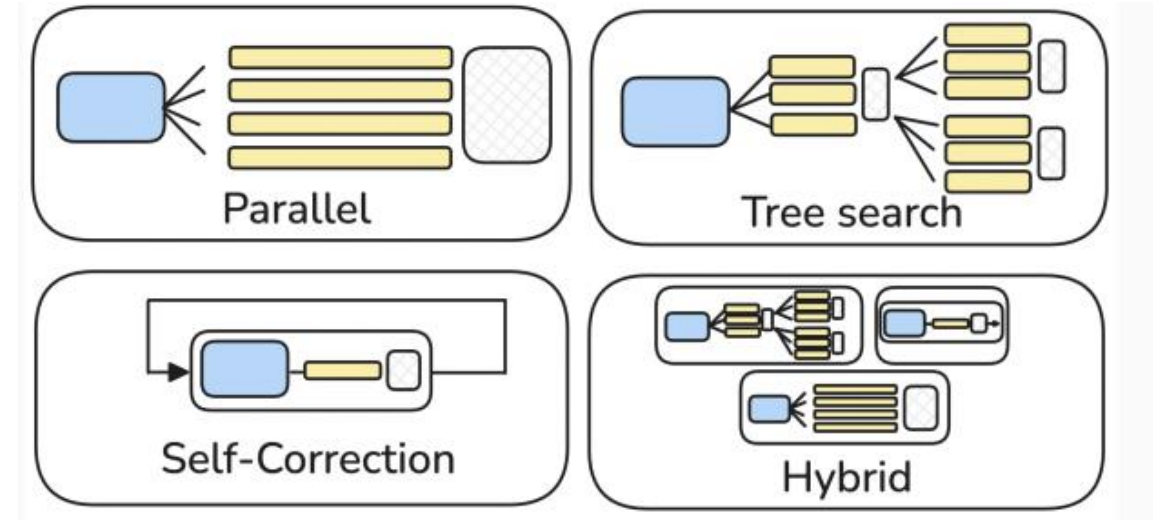  ○ Mixed results, depends on difficulty of verification

# Training Reasoners
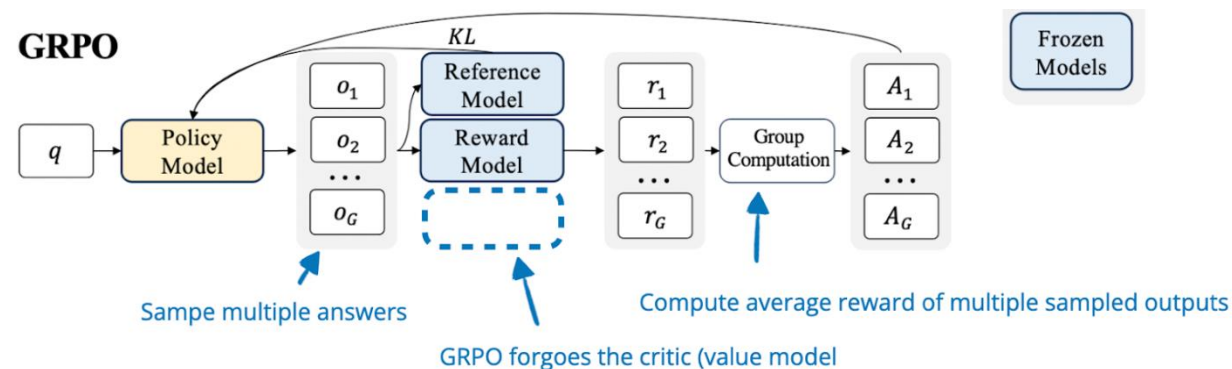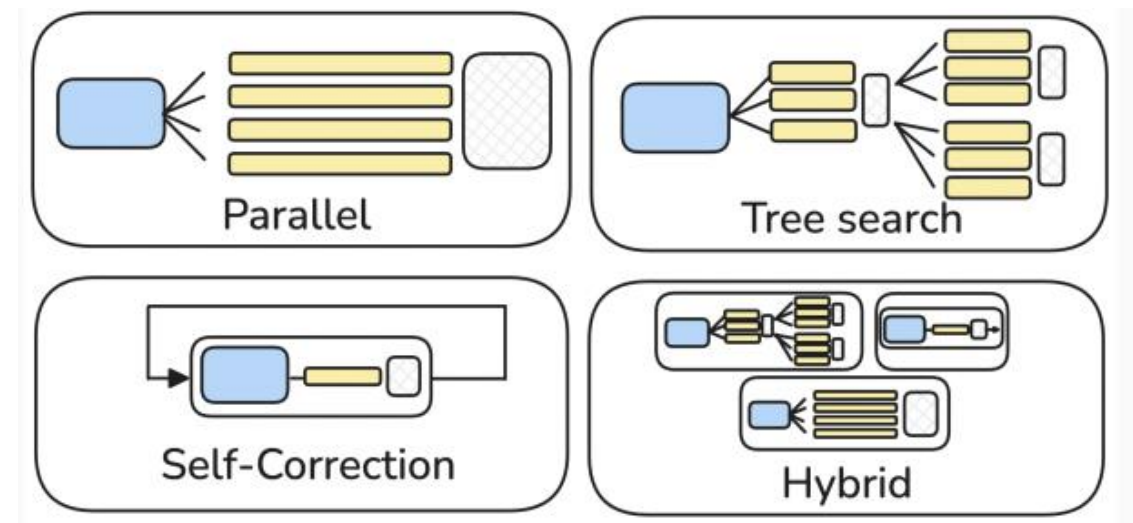
# Latent Space Reasoning

❏ Test-time strategies rely on already trained models in order to get improvements with more token generation



Parallel

Tree search

Self-Correction

Hybrid

# Latent Space Reasoning

❑ Test-time strategies rely on already trained models in order to get improvements with more token generation

❑ What if we just trained the models to reason directly?



Parallel

Tree search

Self-Correction

Hybrid



**GRPO**

KL

$q$ → Policy Model → $o_1$, $o_2$, ..., $o_G$ → Reference Model / Reward Model → $r_1$, $r_2$, ..., $r_G$ → Group Computation → $A_1$, $A_2$, ..., $A_G$

Frozen Models

Sampe multiple answers

GRPO forgoes the critic (value model

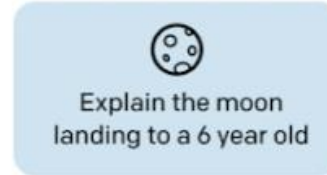Compute average reward of multiple sampled outputs

# Training Reasoners

❑ Revisiting RLHF & PPO

❑ From PPO to GRPO

❑ RLHF to RLVR

❑ Distillation from Reasoning Models

❑ DeepSeek deepdive

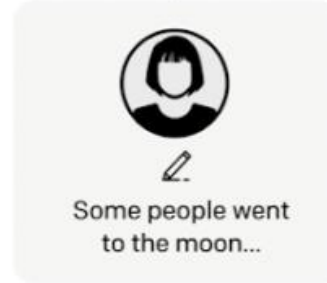# Revisiting RLHF & PPO

**RLHF** Step 1

Sample prompt

Explain the moon landing to a 6 year old

Human writes response

Some people went to the moon...

Supervised fine-tuning of pre-trained LLM

SFT

Time & labor intensive

# Revisiting RLHF & PPO

**RLHF Step 2**

LLM fine-tuned in step 1:

SFT

Sample prompt

Explain the moon landing to a 6 year old

| | |
|---|---|
| **A** Explain gravity... | **B** Explain war... |
| **C** Moon is natural satellite of... | **D** People went to the moon... |

Collect model responses

**Human ranks responses**

D > C > A = B

Time & labor intensive

RM

**Train reward model** (Another LLM)

D > C > A = B

# Revisiting RLHF & PPO



**RLHF Step 3**

Sample prompt — Write a story about frogs

Generate output — PPO — Once upon a time...

Fine-tuned LLM (step 1) — SFT

Reward LLM (step 2) — RM — D > C > A = B

RM — $r_k$

Proximal policy optimization algorithm

Calculate reward to update model

# From PPO to GRPO



Guo et al., (2025)

# RLHF to RLVR

RLHF with **PPO**

**Original policy**
(LLM from SFT stage)

**New policy**
(LLM being trained with PPO)

**Critic**
(Value model estimating expected reward)

**Reward model**
(Trained from human feedback; predicts reward scores)

# RLHF to RLVR

### RLHF with **PPO**

**Original policy**
(LLM from SFT stage)

**New policy**
(LLM being trained with PPO)

**Critic**
(Value model estimating expected reward)

**Reward model**
(Trained from human feedback; predicts reward scores)

### RLHF with **GRPO**

**Original policy**
(LLM from SFT stage)

**New policy**
(LLM being trained with PPO)

**Critic**
(Value model estimating expected reward)

**Reward model**
(Trained from human feedback; predicts reward scores)

# RLHF to RLVR

### RLHF with **PPO**

**Original policy**
(LLM from SFT stage)

**New policy**
(LLM being trained with PPO)

**Critic**
(Value model estimating expected reward)

**Reward model**
(Trained from human feedback; predicts reward scores)

### RLHF with **GRPO**

**Original policy**
(LLM from SFT stage)

**New policy**
(LLM being trained with PPO)

**Critic**
(Value model estimating expected reward)

**Reward model**
(Trained from human feedback; predicts reward scores)

### **RLVR** with GRPO

**Original policy**
(LLM from SFT stage)

**New policy**
(LLM being trained with PPO)

**Critic**
(Value model estimating expected reward)

**Reward model**
(Trained from human feedback; predicts reward scores)

# RLHF to RLVR (RL with Verifiable Rewards)



**Input:** Solve $5 \times (3 + 4)$

**LLM output:**
$$5 \times (3 + 4) =$$
$$5 \times 7 =$$
$$= 35$$

**Verifier output:**

$5 \times (3+4)$

$35$

**Correct**

# Distillation from Reasoning Models

1. Train a large, very capable reasoning language model



**DeepSeed-R1 Model Distillation**

DeepSeek-R1 model(*671B*)

Model Distillation

DeepSeek-R1-Distill-Llama(*8B*)     DeepSeek-R1-Distill-Qwen(*7B*)

Other Distilled Models ....
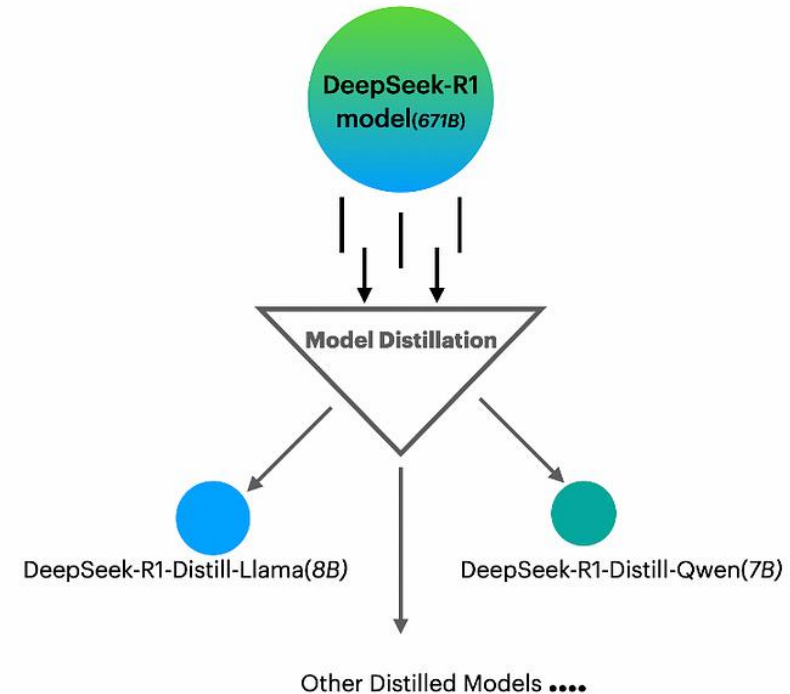
# Distillation from Reasoning Models

1. Train a large, very capable reasoning language model

2. Get a number of outputs from this reasoning model (i.e. curate a reasoning dataset



**DeepSeed-R1 Model Distillation**

DeepSeek-R1 model(*671B*)

Model Distillation

DeepSeek-R1-Distill-Llama(*8B*)          DeepSeek-R1-Distill-Qwen(*7B*)

Other Distilled Models ....

# Distillation from Reasoning Models

1. Train a large, very capable reasoning language model

2. Get a number of outputs from this reasoning model (i.e. curate a reasoning dataset)
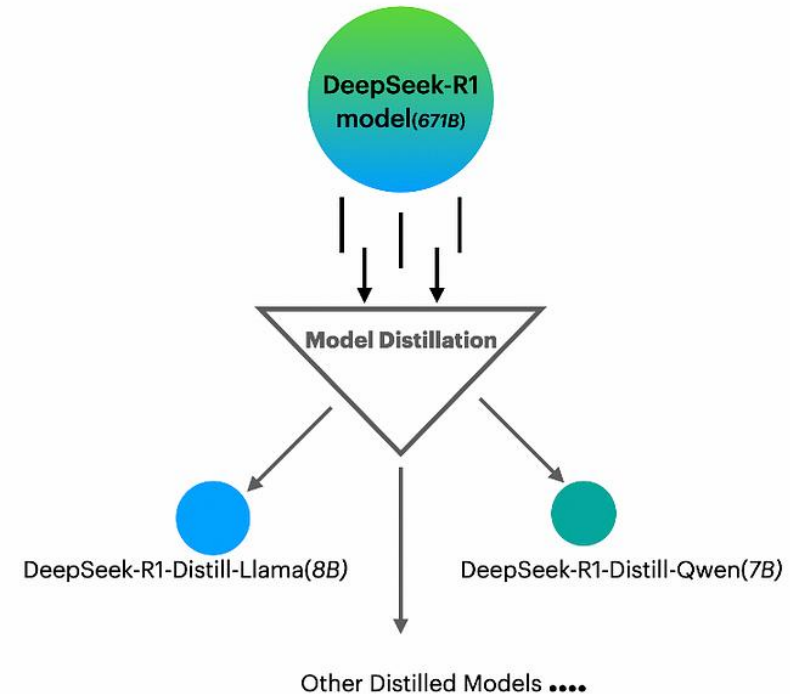
3. Train smaller language models with SFT on this reasoning dataset
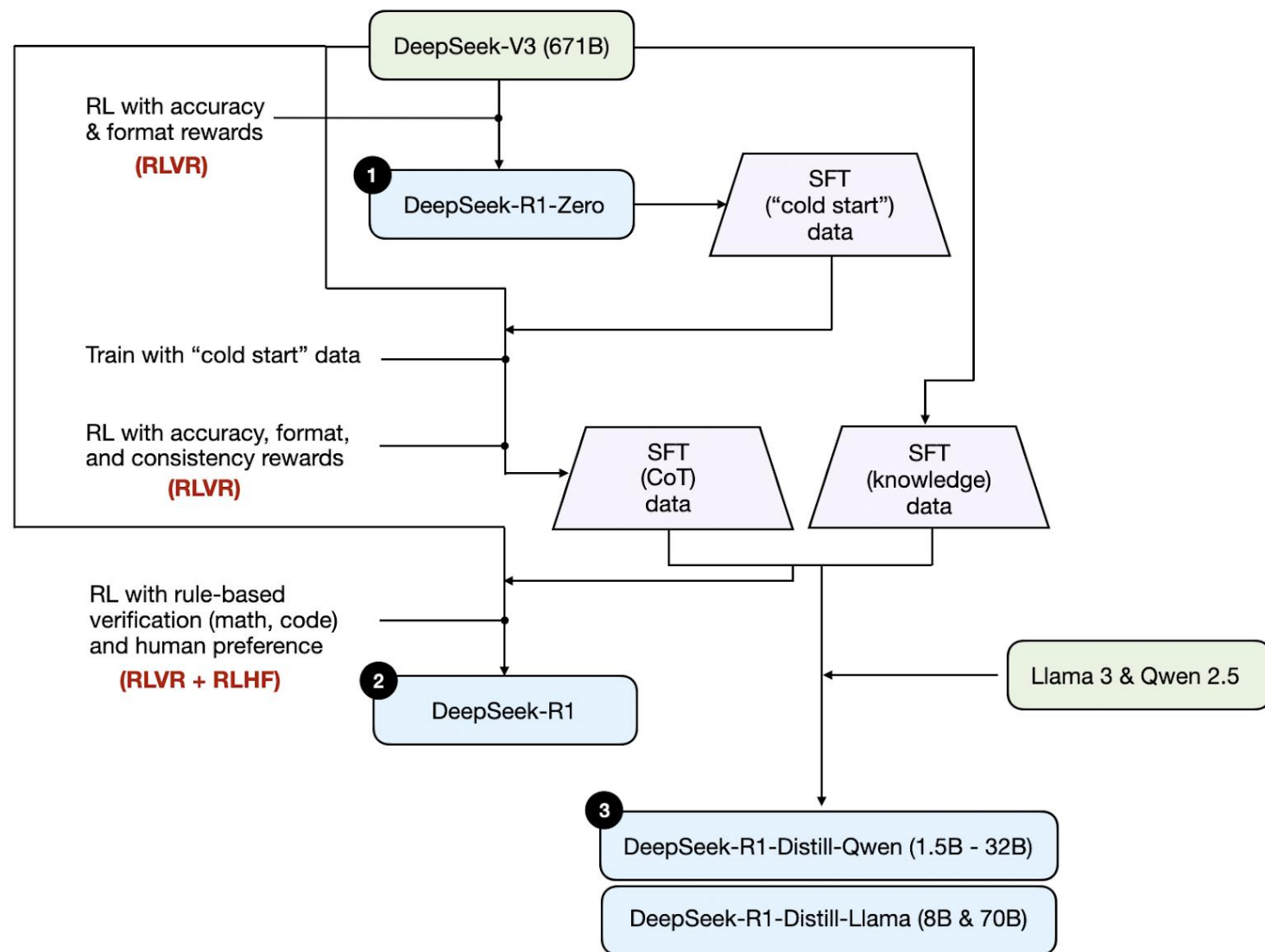


**DeepSeed-R1 Model Distillation**

DeepSeek-R1 model(*671B*)

Model Distillation

DeepSeek-R1-Distill-Llama(*8B*)          DeepSeek-R1-Distill-Qwen(*7B*)

Other Distilled Models ....

# Training (DeepSeek R1 deepdive)

# DeepSeek R1

DeepSeek R1 was one of the first efforts to open source models trained explicitly to reason



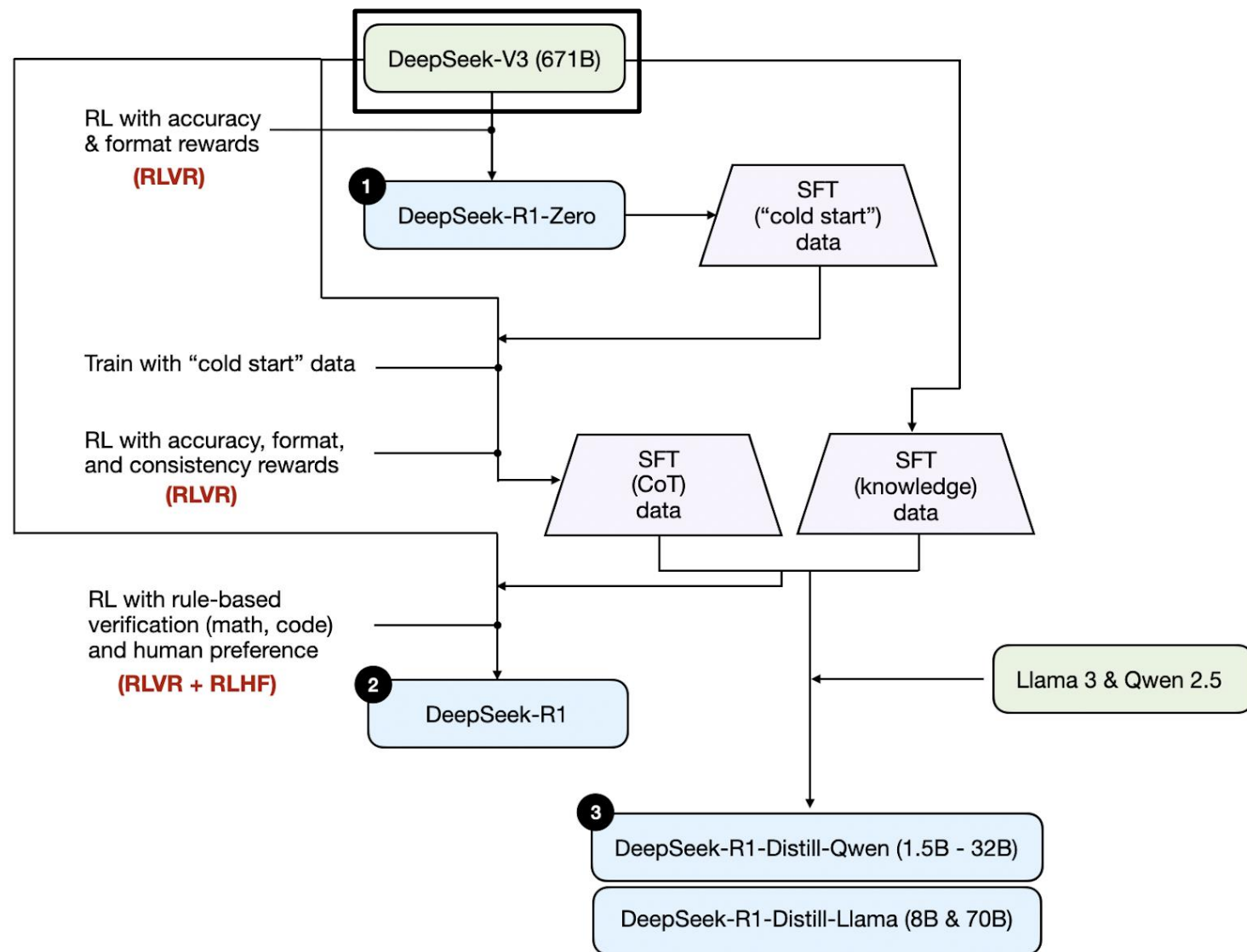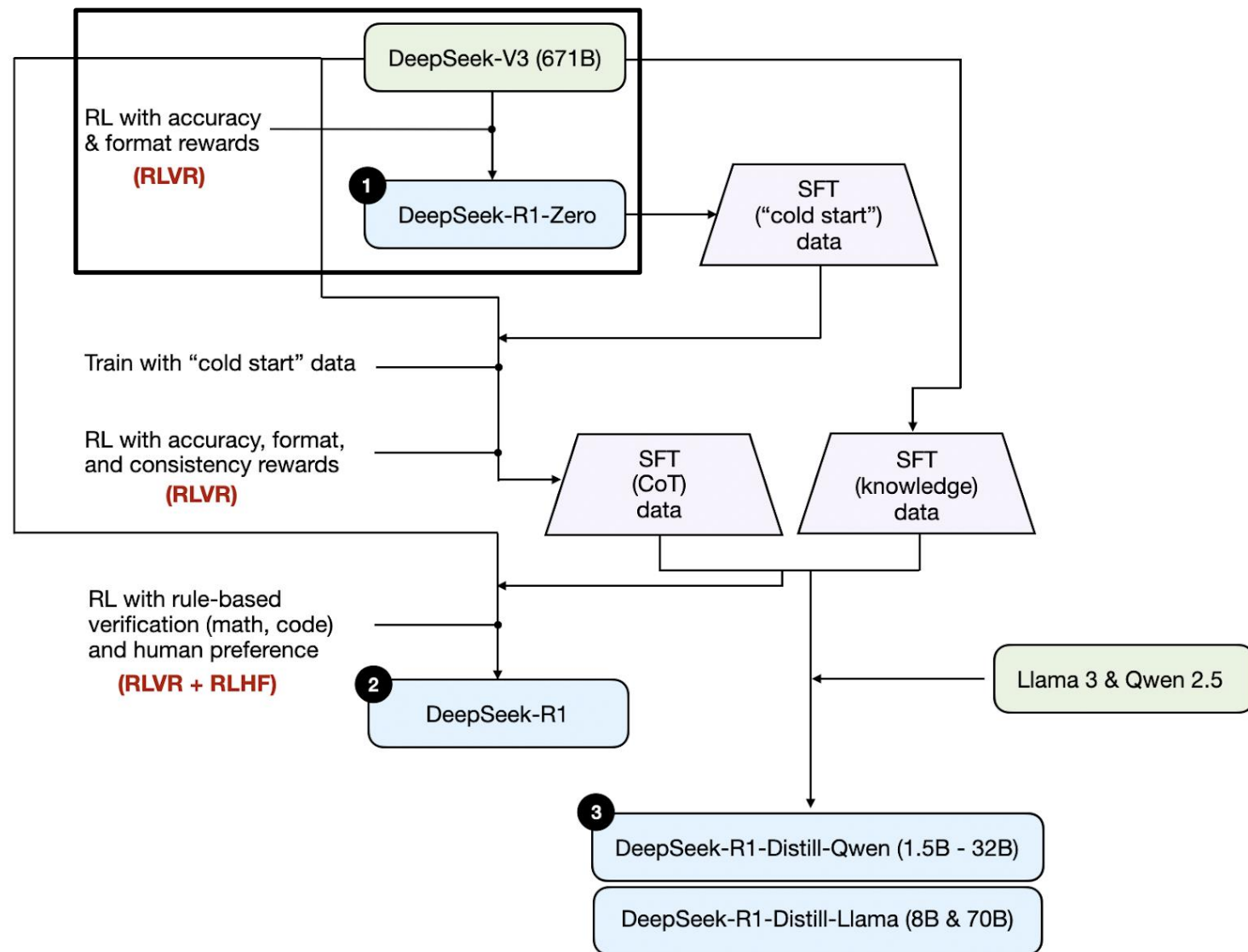Guo et al., (2025)

# Base Model

Base Model from pretraining



- DeepSeek-V3 (671B)
- RL with accuracy & format rewards **(RLVR)**
- ① DeepSeek-R1-Zero
- SFT ("cold start") data
- Train with "cold start" data
- RL with accuracy, format, and consistency rewards **(RLVR)**
- SFT (CoT) data
- SFT (knowledge) data
- RL with rule-based verification (math, code) and human preference **(RLVR + RLHF)**
- ② DeepSeek-R1
- Llama 3 & Qwen 2.5
- ③ DeepSeek-R1-Distill-Qwen (1.5B - 32B)
- DeepSeek-R1-Distill-Llama (8B & 70B)

Guo et al., (2025)

# DeepSeek-R1-Zero

- ❑ DeepSeek-R1-Zero is trained with RLVR
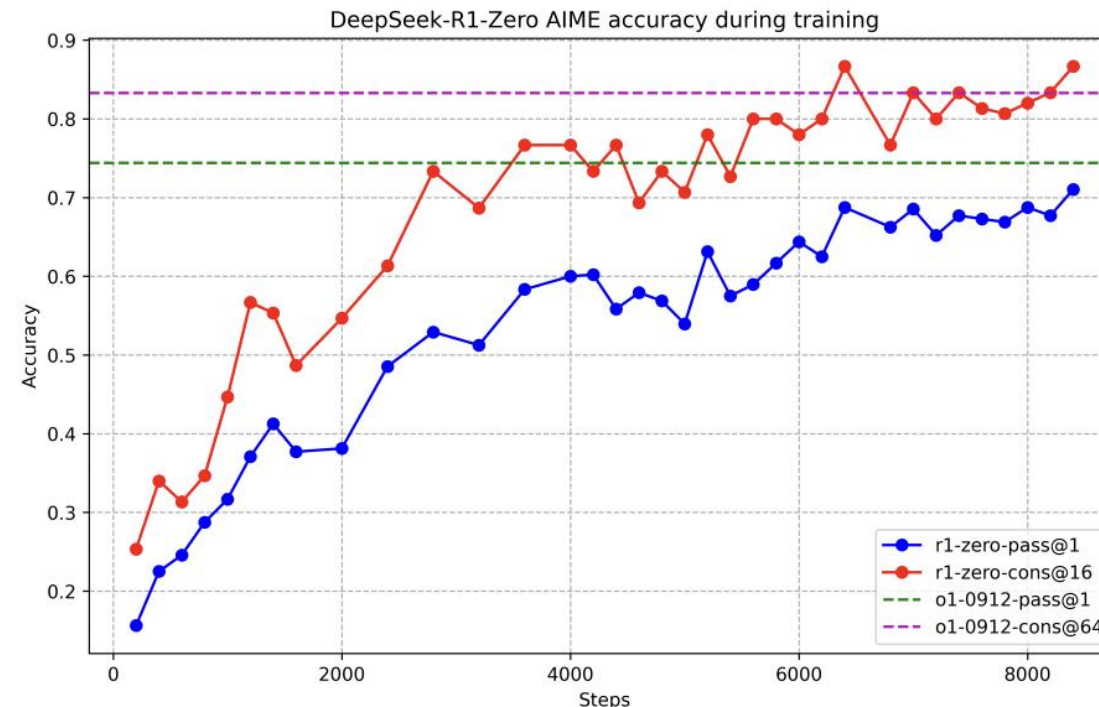- ❑ No SFT was applied, only RL



Guo et al., (2025)

# DeepSeek R1-Zero (Rewards)

- **Accuracy rewards**: The accuracy reward model evaluates whether the response is correct. For example, in the case of math problems with deterministic results, the model is required to provide the final answer in a specified format (e.g., within a box), enabling reliable rule-based verification of correctness. Similarly, for LeetCode problems, a compiler can be used to generate feedback based on predefined test cases.
- **Format rewards**: In addition to the accuracy reward model, we employ a format reward model that enforces the model to put its thinking process between '<think>' and '</think>' tags.

# DeepSeek R1-Zero (Performance)
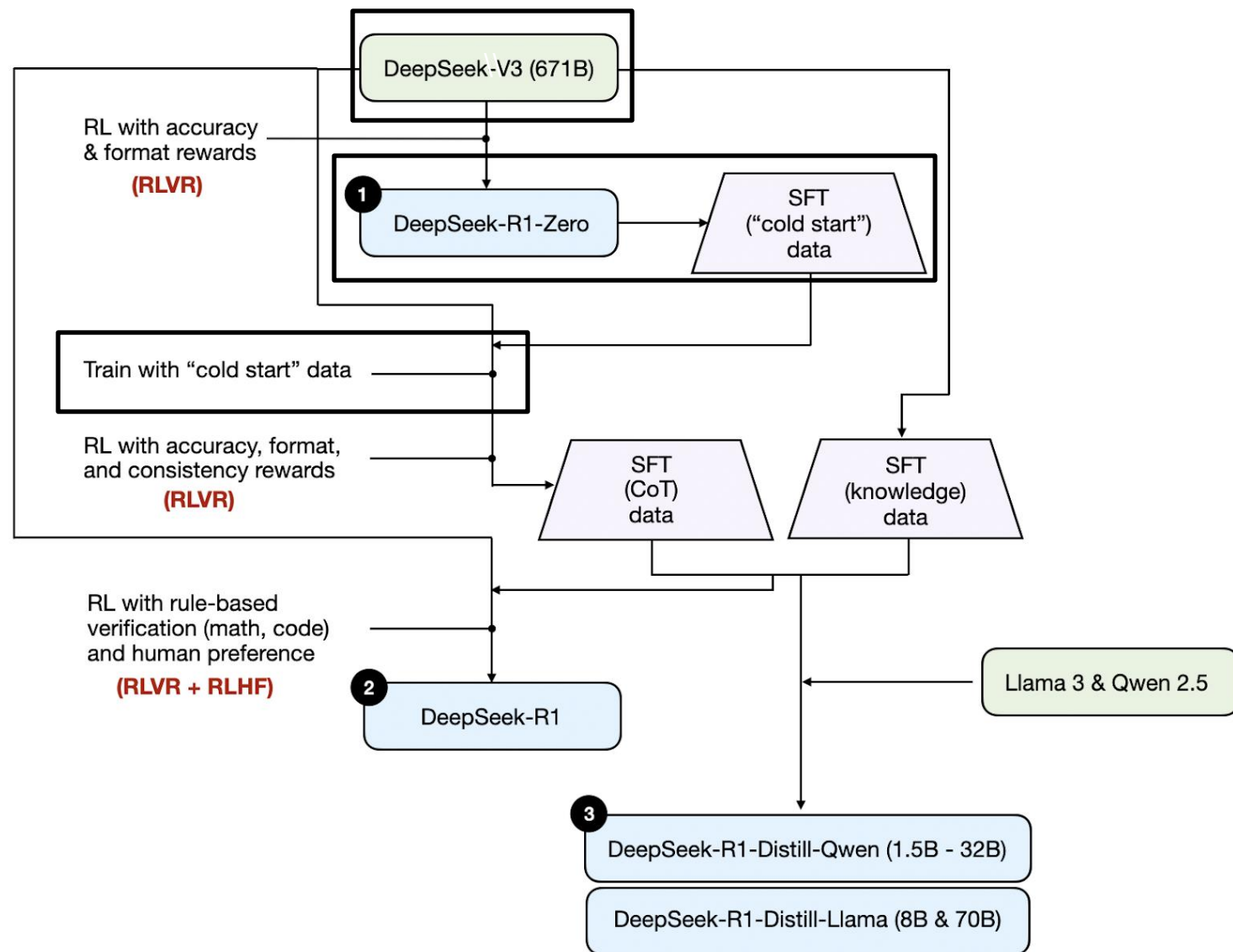


DeepSeek-R1-Zero AIME accuracy during training

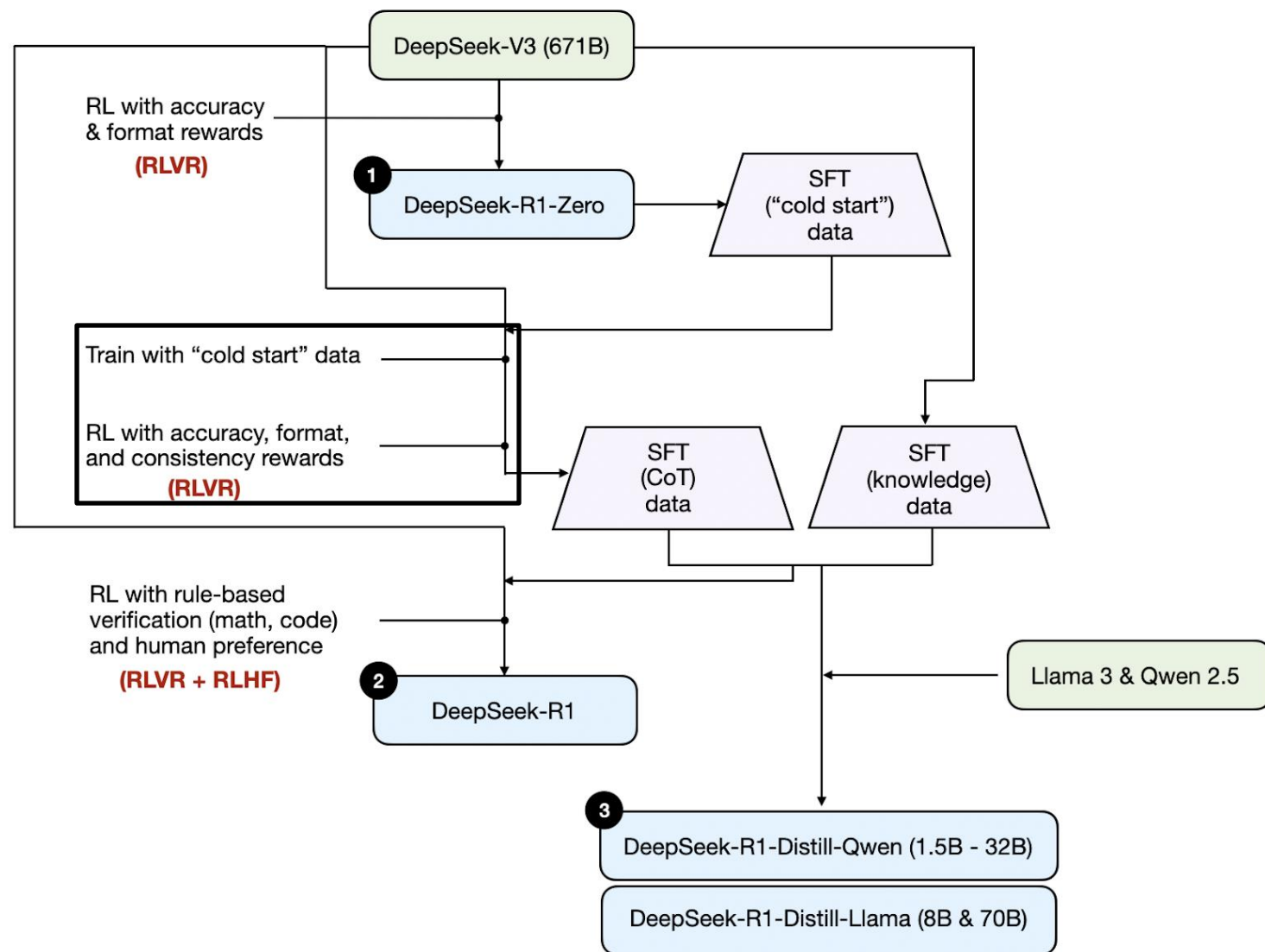| Model | AIME 2024 | | MATH-500 | GPQA Diamond | LiveCode Bench | CodeForces |
|---|---|---|---|---|---|---|
| | pass@1 | cons@64 | pass@1 | pass@1 | pass@1 | rating |
| OpenAI-o1-mini | 63.6 | 80.0 | 90.0 | 60.0 | 53.8 | 1820 |
| OpenAI-o1-0912 | 74.4 | 83.3 | 94.8 | 77.3 | 63.4 | 1843 |
| DeepSeek-R1-Zero | 71.0 | 86.7 | 95.9 | 73.3 | 50.0 | 1444 |

Guo et al., (2025)

# DeepSeek-R1

❏ Collect cold start data from R1-Zero using CoT, and altered prompts

❏ Use human annotated post-processing on this data

❏ Use SFT on base model with this post-processed dataset
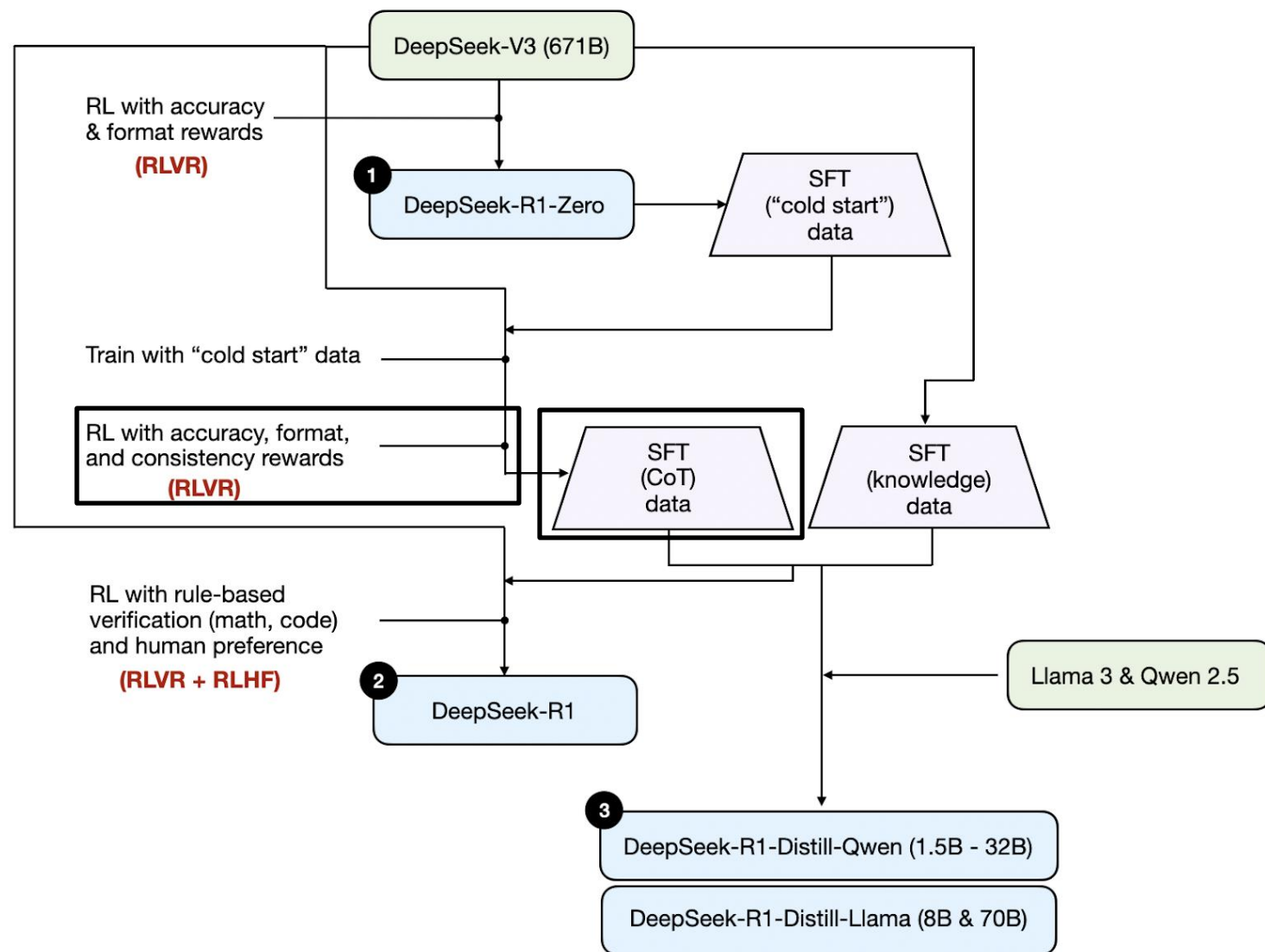
Guo et al., (2025)

# DeepSeek-R1

❏ Apply RLVR to model trained with SFT on "cold start" data

❏ Use same accuracy, format rewards as was used for R1-Zero

❏ Add consistency reward to prevent language-switching
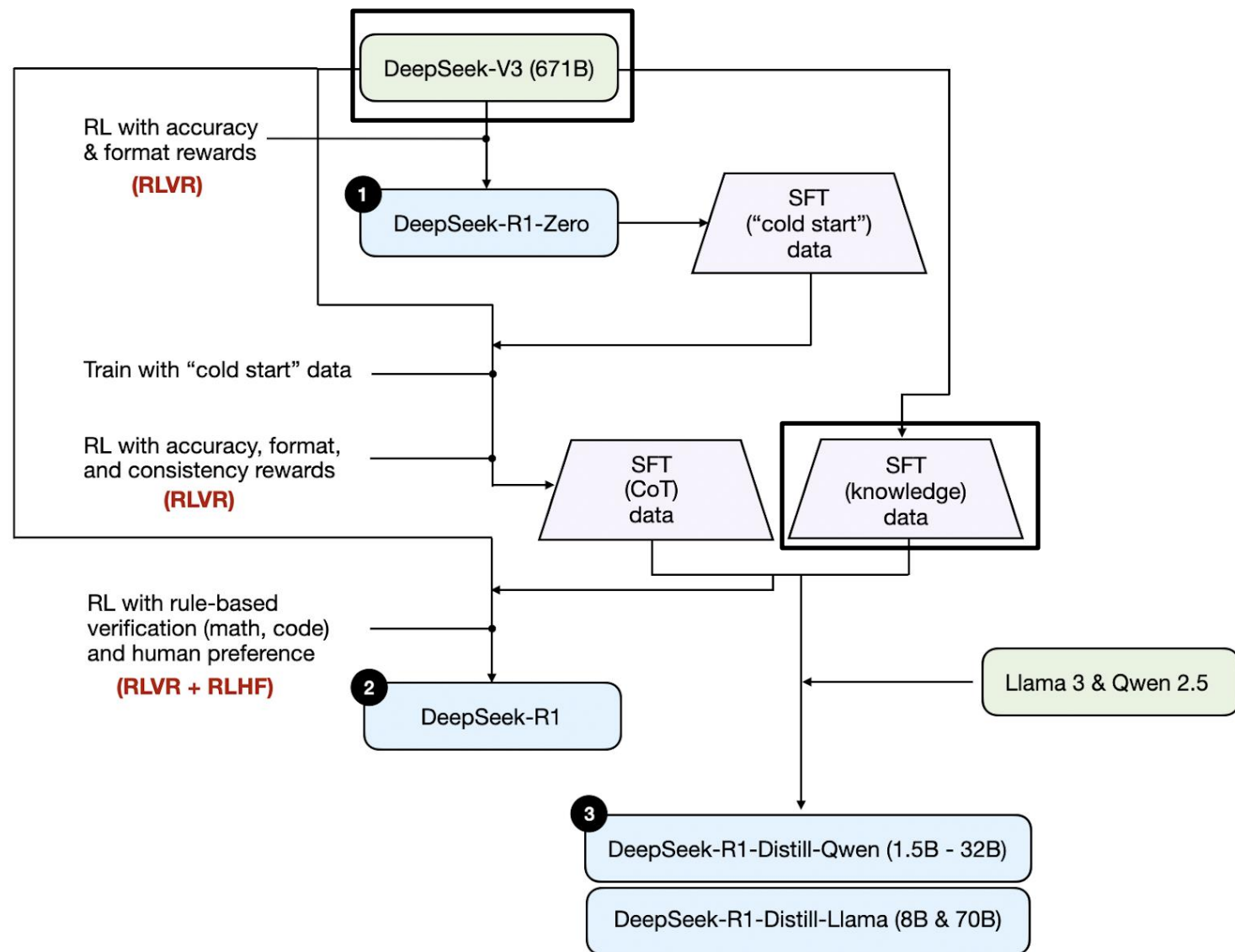
Guo et al., (2025)

# DeepSeek-R1

❑ Using model from previous slide, curate a reasoning dataset (CoT data)

❑ Sample from this model only when correct and apply simple reward based filters
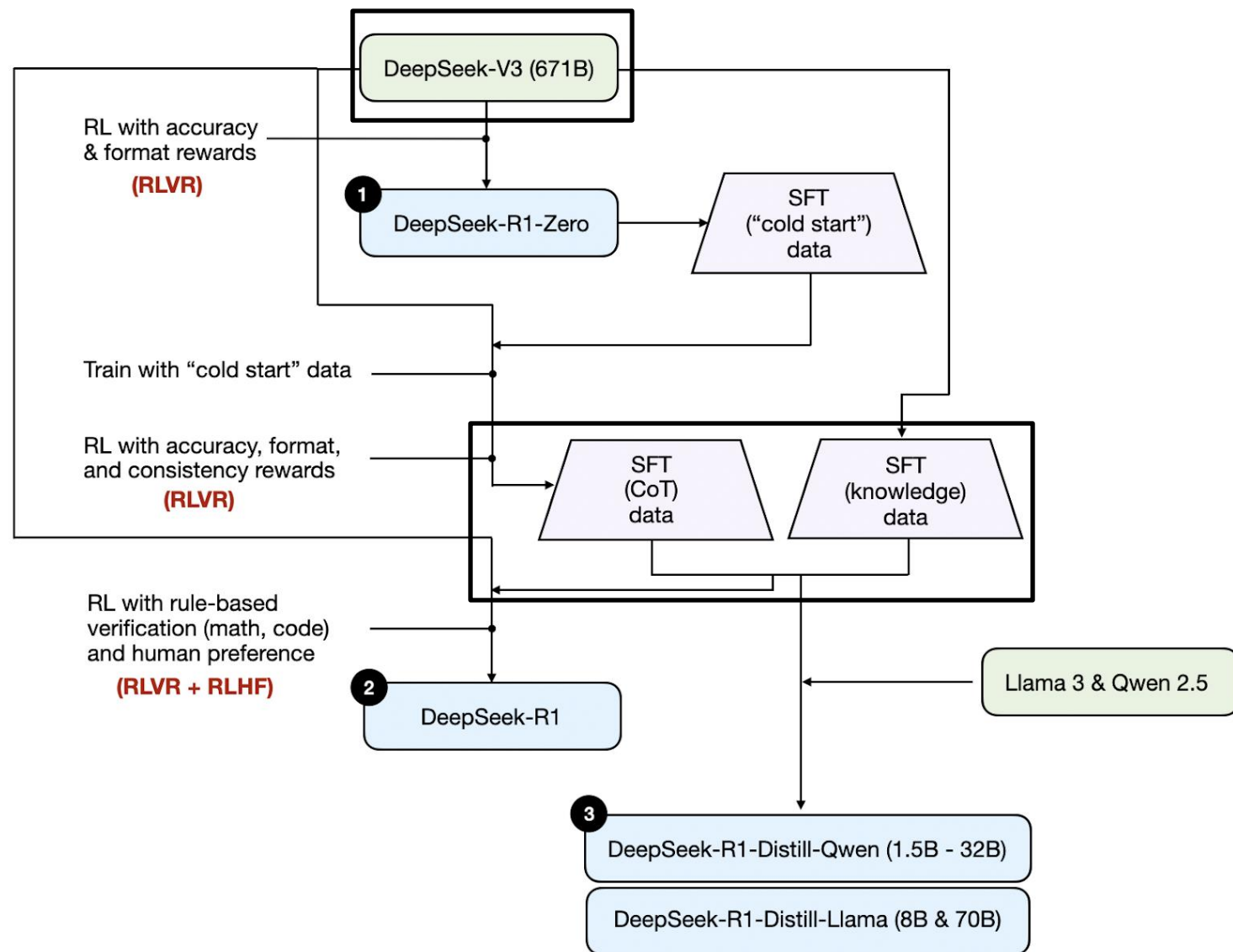
Guo et al., (2025)

# DeepSeek-R1

❑ Sample non-reasoning data from the base model

❑ Sample for tasks related to writing, factual QA, self-cognition, and translation, etc.
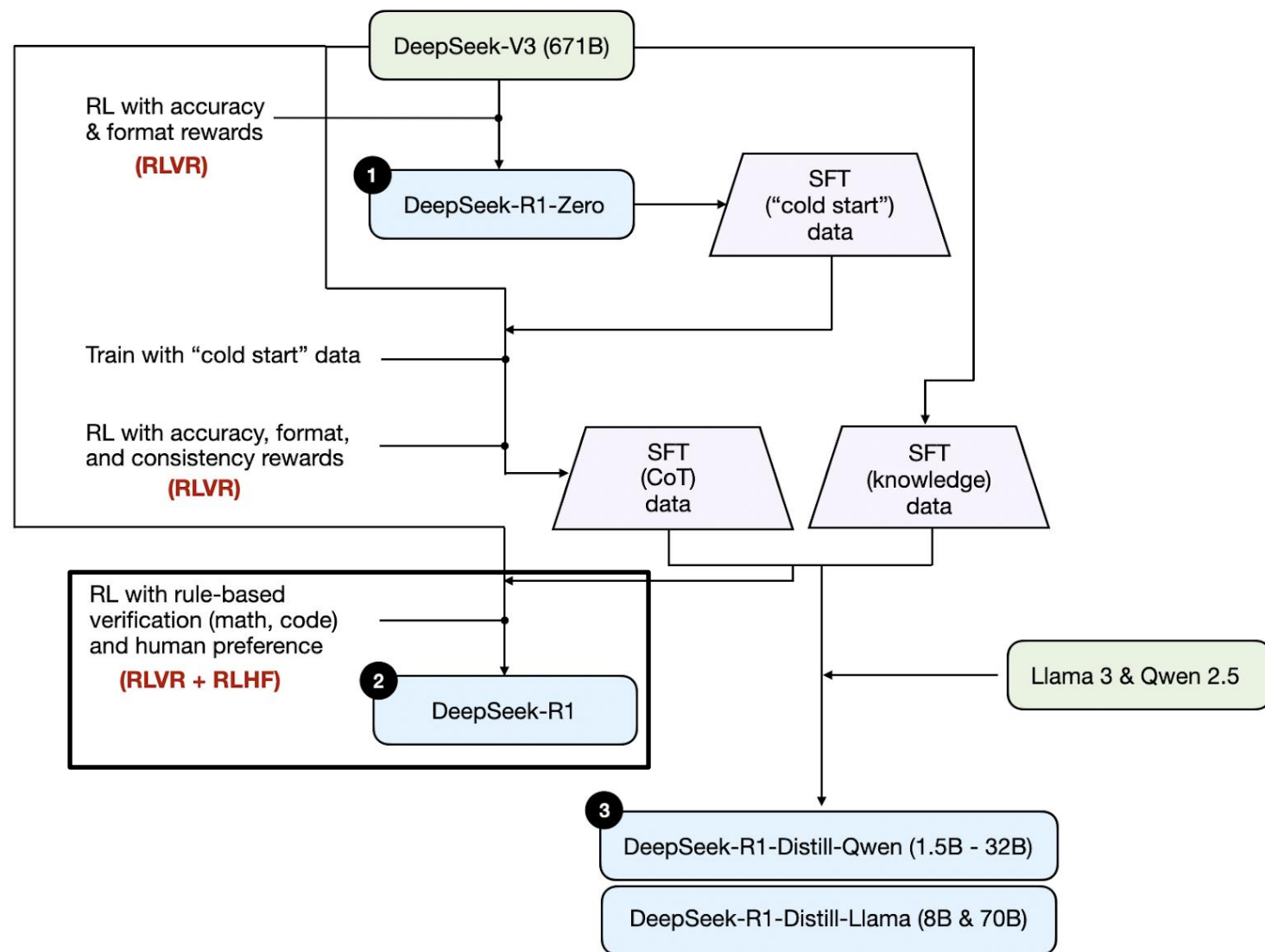


Guo et al., (2025)

# DeepSeek-R1

❑ Using the two curated datasets, perform SFT on the base model

❑ Combined, this is a dataset of around 800K generated, rather than human created, instances



Guo et al., (2025)

# DeepSeek-R1

❏ After applying SFT to the base model, combine RLVR and RLHF methods to produce the final model (DeepSeek-R1)

Guo et al., (2025)

# DeepSeek-R1

- ❑ Very strong model performance
- ❑ Note that the base model is a MoE model
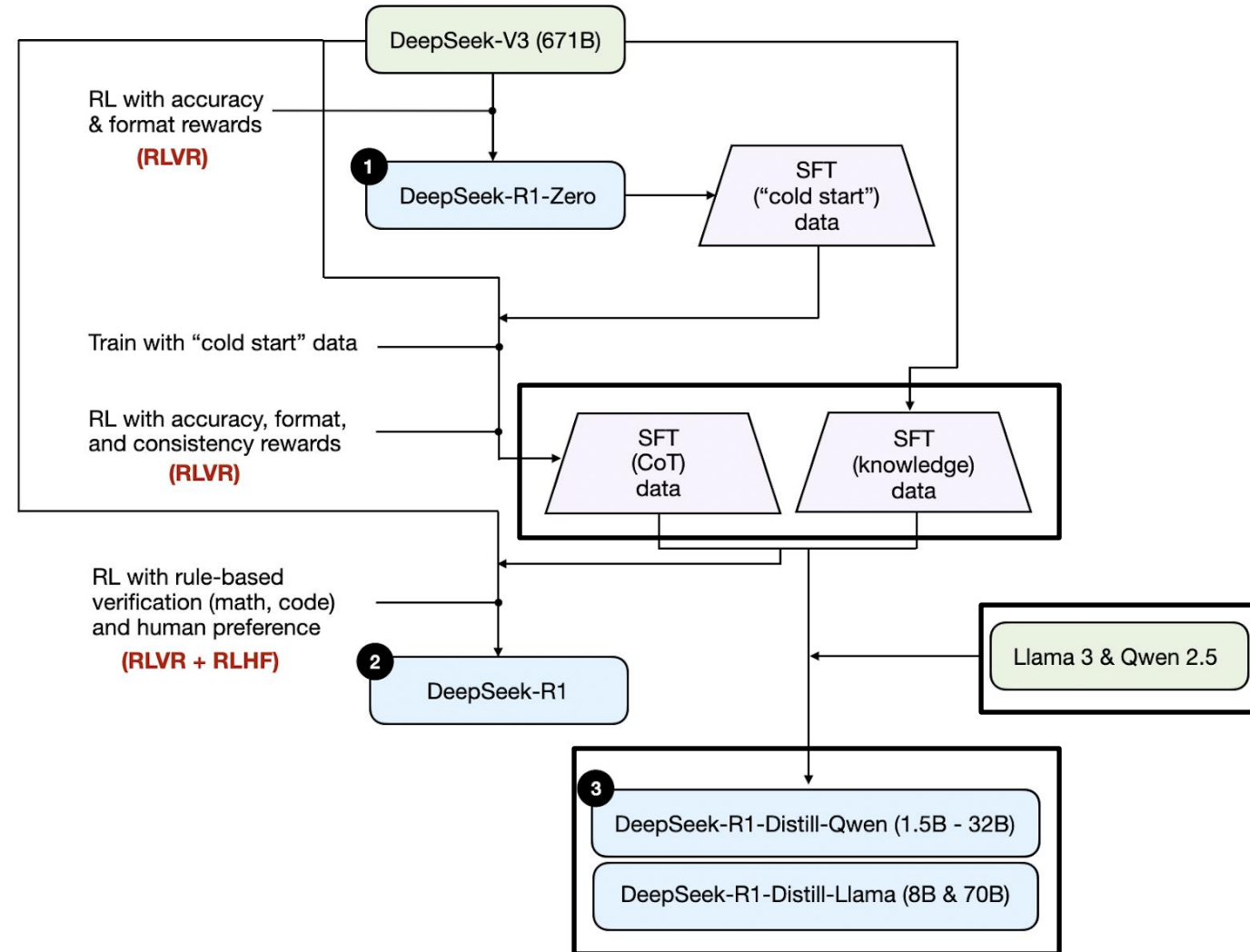- ❑ Strongest open source result up to that point

| Benchmark (Metric) | Claude-3.5-Sonnet-1022 | GPT-4o 0513 | DeepSeek V3 | OpenAI o1-mini | OpenAI o1-1217 | DeepSeek R1 |
|---|---|---|---|---|---|---|
| Architecture | - | - | MoE | - | - | MoE |
| # Activated Params | - | - | 37B | - | - | 37B |
| # Total Params | - | - | 671B | - | - | 671B |
| **English** MMLU (Pass@1) | 88.3 | 87.2 | 88.5 | 85.2 | **91.8** | 90.8 |
| MMLU-Redux (EM) | 88.9 | 88.0 | 89.1 | 86.7 | - | **92.9** |
| MMLU-Pro (EM) | 78.0 | 72.6 | 75.9 | 80.3 | - | **84.0** |
| DROP (3-shot F1) | 88.3 | 83.7 | 91.6 | 83.9 | 90.2 | **92.2** |
| IF-Eval (Prompt Strict) | **86.5** | 84.3 | 86.1 | 84.8 | - | 83.3 |
| GPQA Diamond (Pass@1) | 65.0 | 49.9 | 59.1 | 60.0 | **75.7** | 71.5 |
| SimpleQA (Correct) | 28.4 | 38.2 | 24.9 | 7.0 | **47.0** | 30.1 |
| FRAMES (Acc.) | 72.5 | 80.5 | 73.3 | 76.9 | - | **82.5** |
| AlpacaEval2.0 (LC-winrate) | 52.0 | 51.1 | 70.0 | 57.8 | - | **87.6** |
| ArenaHard (GPT-4-1106) | 85.2 | 80.4 | 85.5 | 92.0 | - | **92.3** |
| **Code** LiveCodeBench (Pass@1-COT) | 38.9 | 32.9 | 36.2 | 53.8 | 63.4 | **65.9** |
| Codeforces (Percentile) | 20.3 | 23.6 | 58.7 | 93.4 | **96.6** | 96.3 |
| Codeforces (Rating) | 717 | 759 | 1134 | 1820 | **2061** | 2029 |
| SWE Verified (Resolved) | **50.8** | 38.8 | 42.0 | 41.6 | 48.9 | 49.2 |
| Aider-Polyglot (Acc.) | 45.3 | 16.0 | 49.6 | 32.9 | **61.7** | 53.3 |
| **Math** AIME 2024 (Pass@1) | 16.0 | 9.3 | 39.2 | 63.6 | 79.2 | **79.8** |
| MATH-500 (Pass@1) | 78.3 | 74.6 | 90.2 | 90.0 | 96.4 | **97.3** |
| CNMO 2024 (Pass@1) | 13.1 | 10.8 | 43.2 | 67.6 | - | **78.8** |
| **Chinese** CLUEWSC (EM) | 85.4 | 87.9 | 90.9 | 89.9 | - | **92.8** |
| C-Eval (EM) | 76.7 | 76.0 | 86.5 | 68.9 | - | **91.8** |
| C-SimpleQA (Correct) | 55.4 | 58.7 | **68.0** | 40.3 | - | 63.7 |

Guo et al., (2025)

# DeepSeek-R1-{Qwen, Llama} (*B)

❑ Apply SFT to Llama3 and Qwen2.5 models using the ~800k SFT data from the pipeline

❑ Note: **Only** SFT is applied at this point, no RL is used



Guo et al., (2025)

# DeepSeek-R1-{Qwen, Llama} (*B)

❑ Even without any reinforcement learning, SFT with a reasoning dataset is sufficient to achieve very good performance with these other open Small Language Models (SLMs)

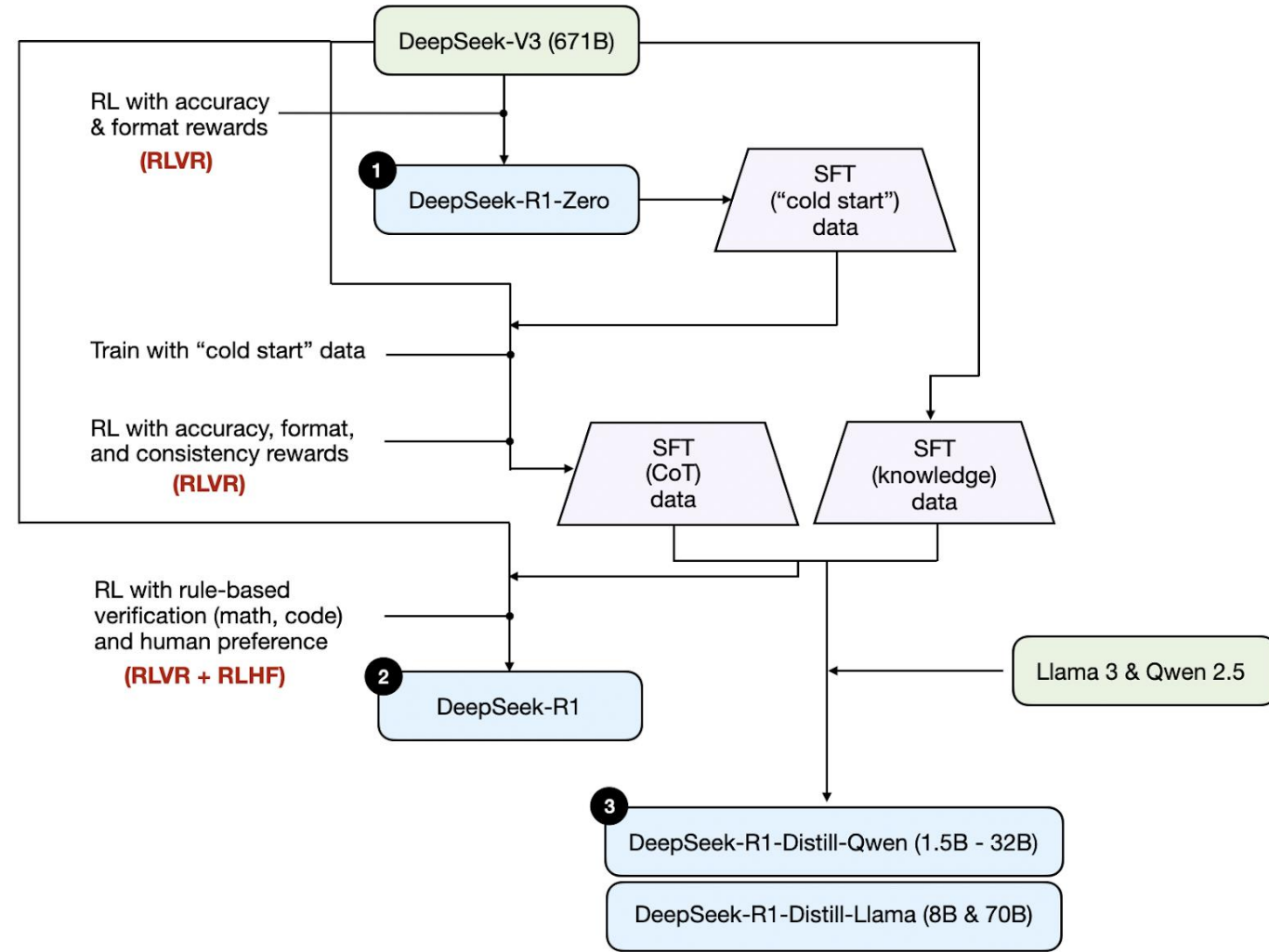❑ Opens possibility of improving SLMs by curating better reasoning datasets

| Model | AIME 2024 | | MATH-500 | GPQA Diamond | LiveCode Bench | CodeForces |
|---|---|---|---|---|---|---|
| | pass@1 | cons@64 | pass@1 | pass@1 | pass@1 | rating |
| GPT-4o-0513 | 9.3 | 13.4 | 74.6 | 49.9 | 32.9 | 759 |
| Claude-3.5-Sonnet-1022 | 16.0 | 26.7 | 78.3 | 65.0 | 38.9 | 717 |
| OpenAI-o1-mini | 63.6 | 80.0 | 90.0 | 60.0 | 53.8 | **1820** |
| QwQ-32B-Preview | 50.0 | 60.0 | 90.6 | 54.5 | 41.9 | 1316 |
| DeepSeek-R1-Distill-Qwen-1.5B | 28.9 | 52.7 | 83.9 | 33.8 | 16.9 | 954 |
| DeepSeek-R1-Distill-Qwen-7B | 55.5 | 83.3 | 92.8 | 49.1 | 37.6 | 1189 |
| DeepSeek-R1-Distill-Qwen-14B | 69.7 | 80.0 | 93.9 | 59.1 | 53.1 | 1481 |
| DeepSeek-R1-Distill-Qwen-32B | **72.6** | 83.3 | 94.3 | 62.1 | 57.2 | 1691 |
| DeepSeek-R1-Distill-Llama-8B | 50.4 | 80.0 | 89.1 | 49.0 | 39.6 | 1205 |
| DeepSeek-R1-Distill-Llama-70B | 70.0 | **86.7** | **94.5** | **65.2** | **57.5** | 1633 |

Guo et al., (2025)

# DeepSeek Main Contributions (3 key areas)

1. Possible to train reasoning model with RL alone on math/coding (DeepSeek-R1-Zero)

2. It is possible to curate a reasoning dataset that enables SFT for reasoning models (DeepSeek-R1-{Qwen, Llama} (*B))

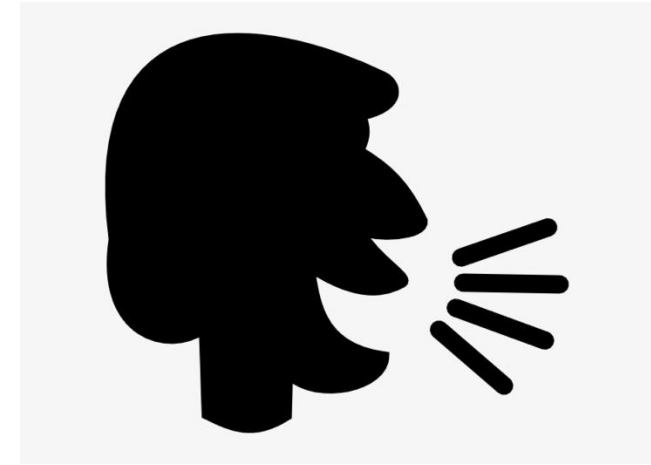3. Open Sources RL + SFT solution which is competitive with closed-source models (DeepSeek-R1)
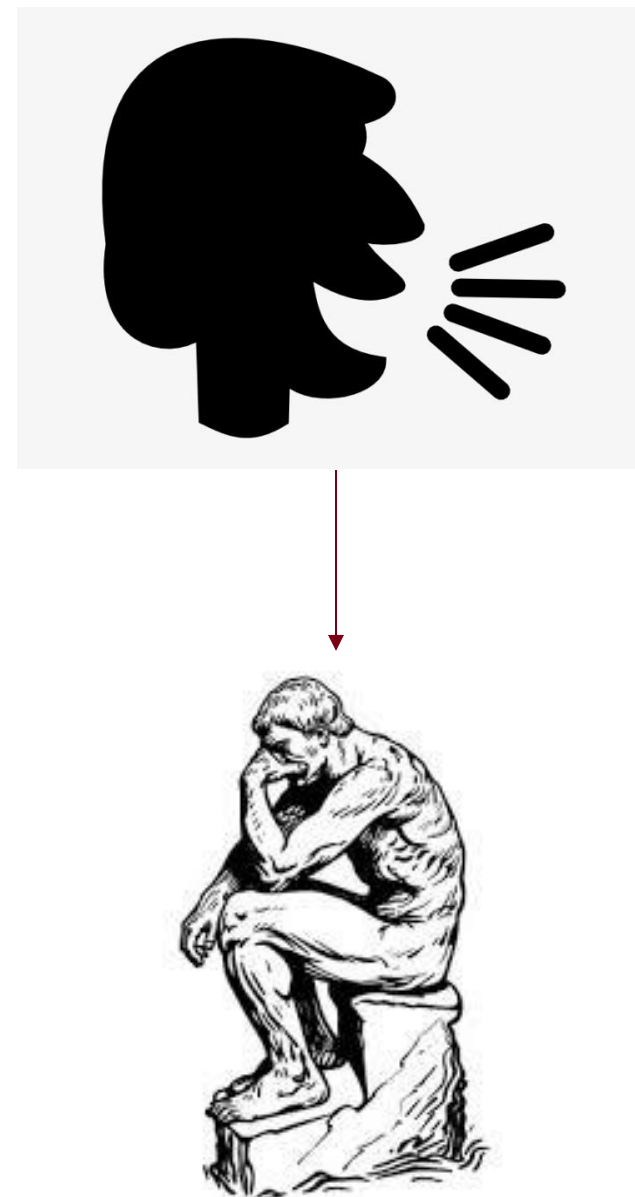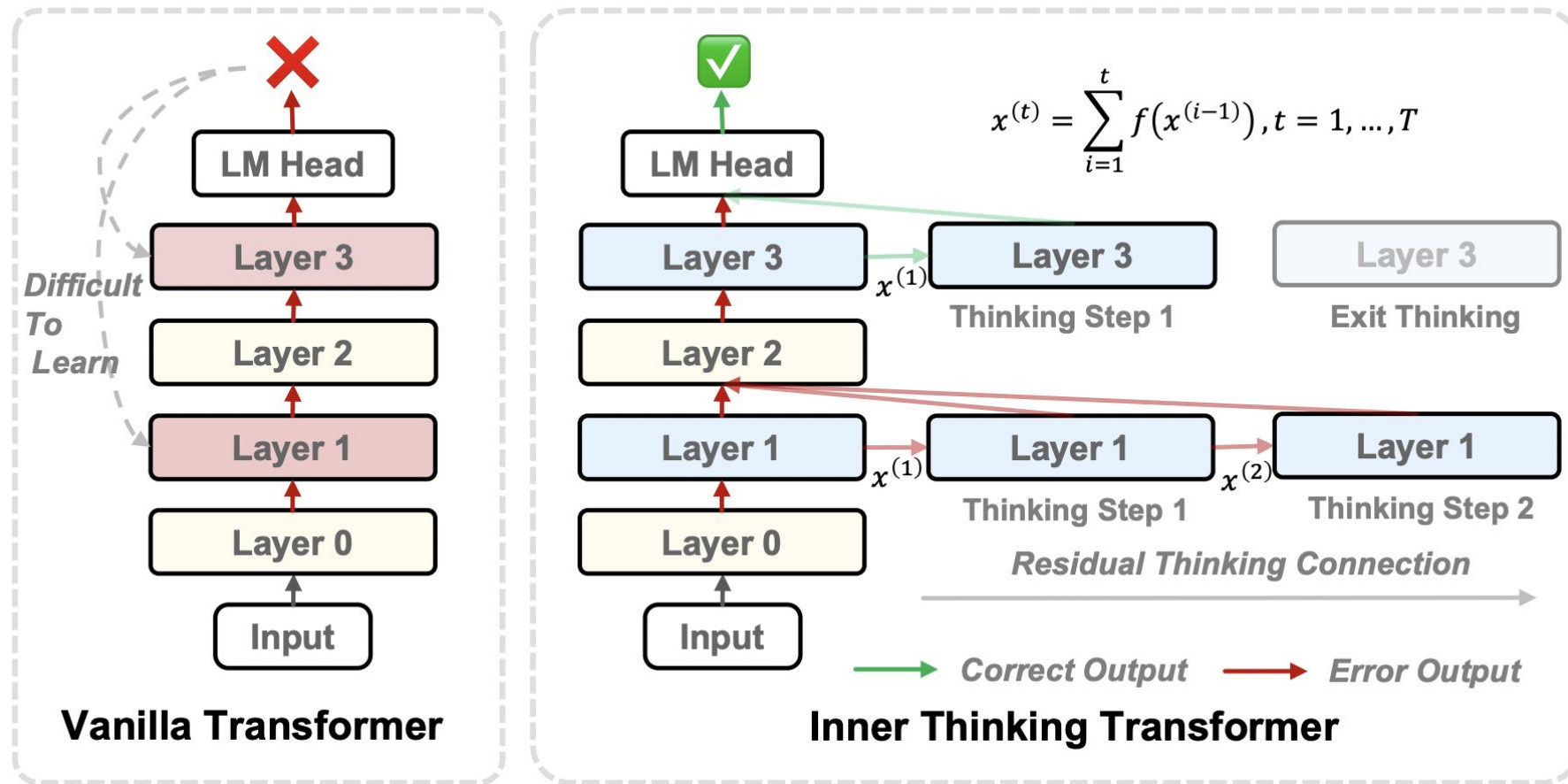


Guo et al., (2025)

# Latent Space Reasoning

# Latent Space Reasoning



❑ What if, instead of generating more tokens at inference, we just used more compute by altering the hidden states of the transformer itself

# Latent Space Reasoning

❑ What if, instead of generating more tokens at inference, we just used more compute by altering the hidden states of the transformer itself
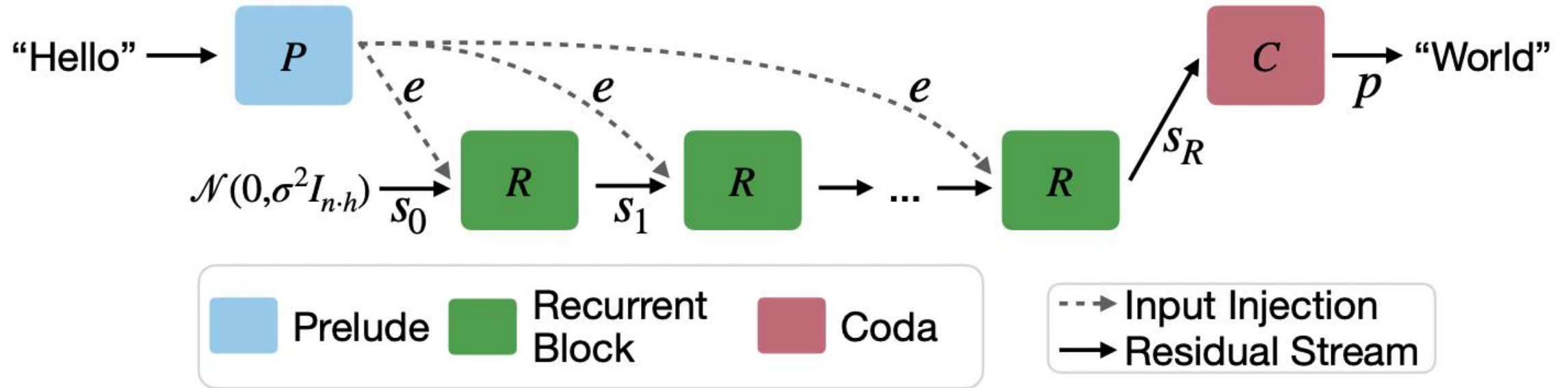
❑ In analogical terms, less talking, more thinking
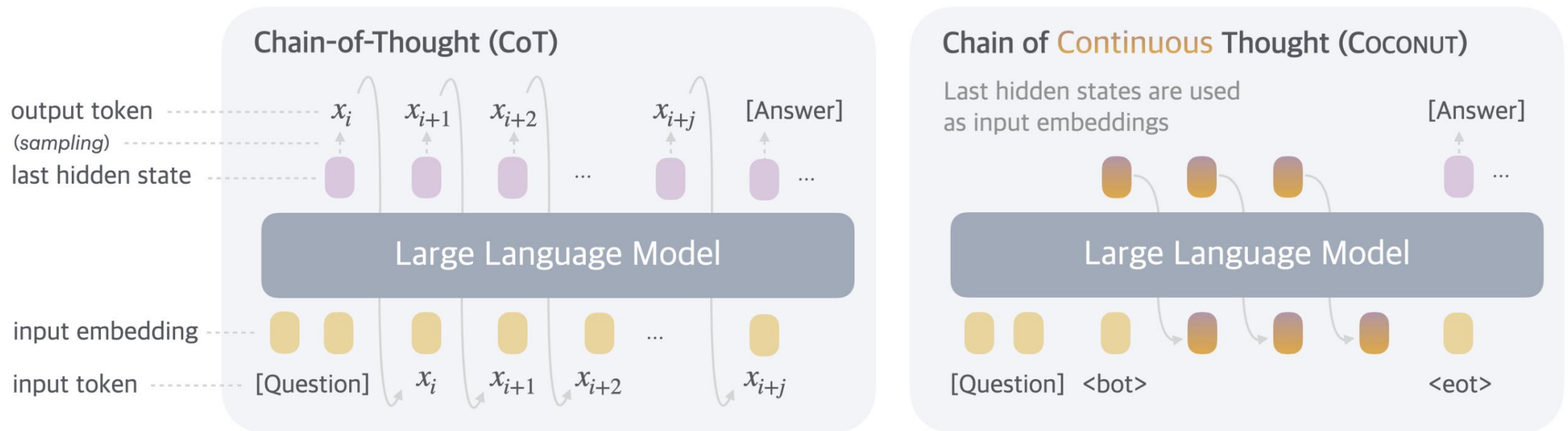
# Methods (Inner Thinking Transformer)



$$x^{(t)} = \sum_{i=1}^{t} f(x^{(i-1)}), t = 1, \ldots, T$$

**Vanilla Transformer**

**Inner Thinking Transformer**

Chen et al., (2025)

# Methods (Scaling up Test-Time Compute with Latent Reasoning)



Geiping et al., (2025)

# Training Large Language Models to Reason in a Continuous Latent Space



Hao et al., (2024)