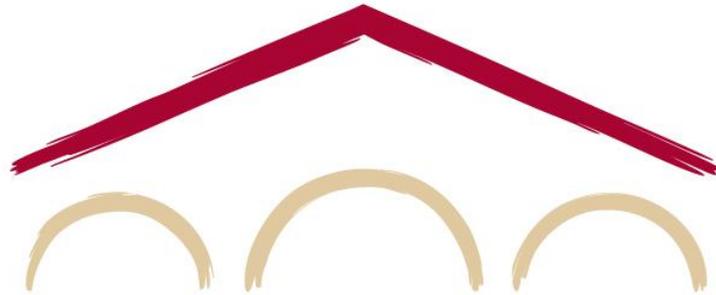


Natural Language Processing with Deep Learning

CS224N/Ling284



Yejin Choi

Lecture 12: Reasoning 1/2

Lecture Plan



Decoding techniques (20 mins)

From greedy to beam search to sampling (15 mins)
Neural text degeneration (loopy behavior) (5 mins)



DeepSeek R1 deep dives (15 mins)

From R1-Zero to R1 to R1-distill



PPO & GRPO & DAPO (25 mins)

PPO: dissecting the complexity
GRPO
DAPO

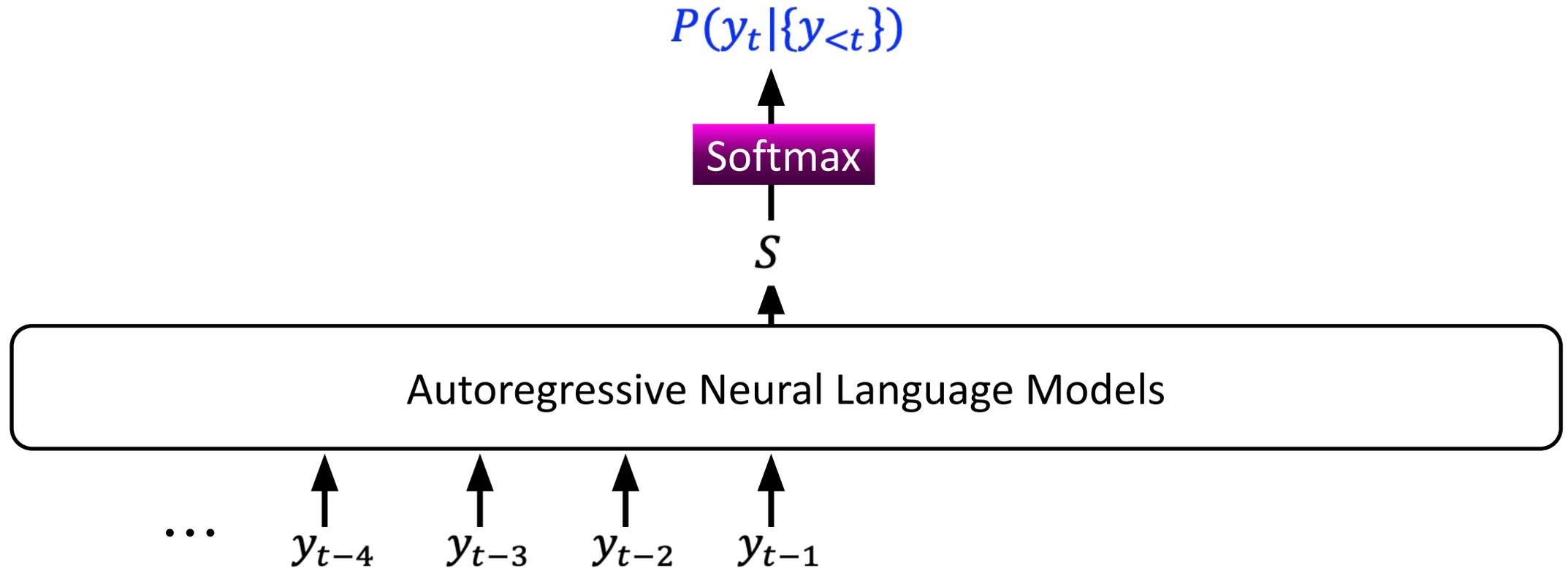


The nature of “reasoning” (20 min)

What works, why it works, and when it fails

How to decode a token during “inference”

- At each time step t , our model computes a vector of scores for each token in our vocabulary, $S \in \mathbb{R}^V$. Then, we compute a probability distribution P over $w \in V$ using these scores:



How to decode a token during “inference”

- At each time step t , our model computes a vector of scores for each token in our vocabulary, $\mathbf{s} \in \mathbb{R}^V$:

$$\mathbf{s} = f(\{y_{<t}\}; \theta) \longrightarrow f(\cdot; \theta) \text{ is your model}$$

- Then, we compute a probability distribution P over $w \in V$ using these scores:

$$P(y_t = w | \{y_{<t}\}) = \frac{\exp(\mathbf{s}_w)}{\sum_{w' \in V} \exp(\mathbf{s}_{w'})}$$

- Decoding algorithm defines a function to select a token from this distribution:

$$\hat{y}_t = g(P(y_t | \{y_{<t}\})) \longrightarrow g(\cdot) \text{ is your decoding algorithm}$$

Basic decoding algorithms

- Decoding algorithm defines a function to select a token from this distribution:

$$\hat{y}_t = g(P(y_t|\{y_{<t}\})) \quad \longrightarrow g(\cdot) \text{ is your decoding algorithm}$$

- Greedy decoding

$$\hat{y}_t = \mathit{argmax}_{w \in V} P(y_t = w | y_{<t})$$

- **Limitation:** Greedy decoding is myopic—it maximizes local probability at each step but may miss globally optimal sequences.
- Beam search
 - maintains the top-k most probable partial sequences (hypotheses) at each step, where k is the beam width.

Basic decoding algorithms

- Decoding algorithm defines a function to select a token from this distribution:

$$\hat{y}_t = g(P(y_t | \{y_{<t}\})) \quad \longrightarrow \quad g(\cdot) \text{ is your decoding algorithm}$$

- **Beam search**

- maintains the top-k most probable partial sequences (hypotheses) at each step, where k is the beam width.

- at each time step t , for each hypothesis h in the beam B_{t-1} :

- $\text{score}(h \circ y_t) = \text{score}(h) + \log P(y_t | h)$

- final sequence output $y^* = \operatorname{argmax}_{h' \in B_T} \log P(h')$

a partial sequence until previous time t-1

$$B_t = \operatorname{top-k}_{h \in B_{t-1}, y \in V} \{(h \circ y, \text{score}(h \circ y))\}$$

Basic decoding algorithms

- when $k=1$, beam search reduces to greedy decoding
- Beam search was the defacto standard decoding algorithm for the classical NLP models, but it is no longer used often in modern day language models. Why?

Aspect	Greedy	Beam Search
Search space	1 path	
Complexity		
Quality	Local optimum	Better approximation of global optimum
Diversity	None	Limited (within beam)

Circa 2019, GPT-2 was introduced with remarkable text quality

SYSTEM PROMPT
(HUMAN-WRITTEN)

In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

MODEL COMPLETION
(MACHINE-WRITTEN,
10 TRIES)

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

This notion of "sampling" was surprising at that time!

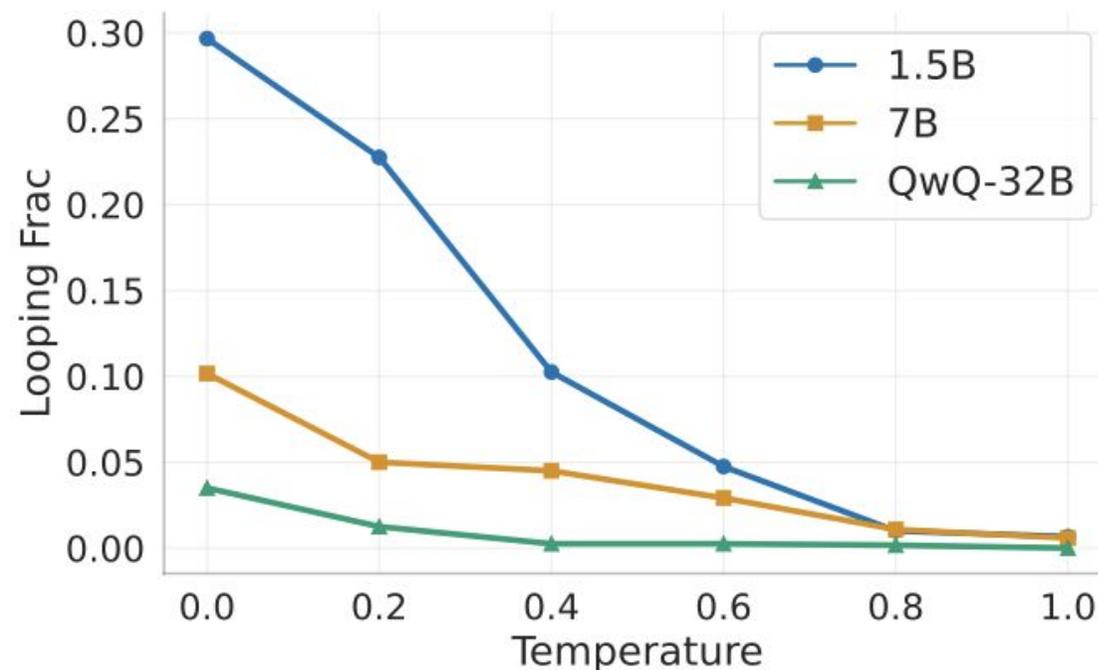
Circa 2025, is looping gone with reasoning models?

- Looping in reasoning LLMs isn't an edge case!
- On AIME, authors flag a response as looping if a 30-gram repeats >20 times.
- [When looping more likely?](#)
 - Across open reasoning LLMs (qwen, openthinker, phi, llama-r1):
 - Low temps => more looping
 - Smaller models => more looping
 - Harder problems => more looping

WAIT, WAIT, WAIT... WHY DO REASONING MODELS LOOP?

Charilaos Pipis^{*1}, Shivam Garg^{*2}, Vasilis Kontonis²,
Vaishnavi Shrivastava², Akshay Krishnamurthy², Dimitris Papailiopoulos^{2,3}

¹ MIT ² Microsoft Research ³ University of Wisconsin-Madison



Circa 2015, is looping gone with reasoning models?

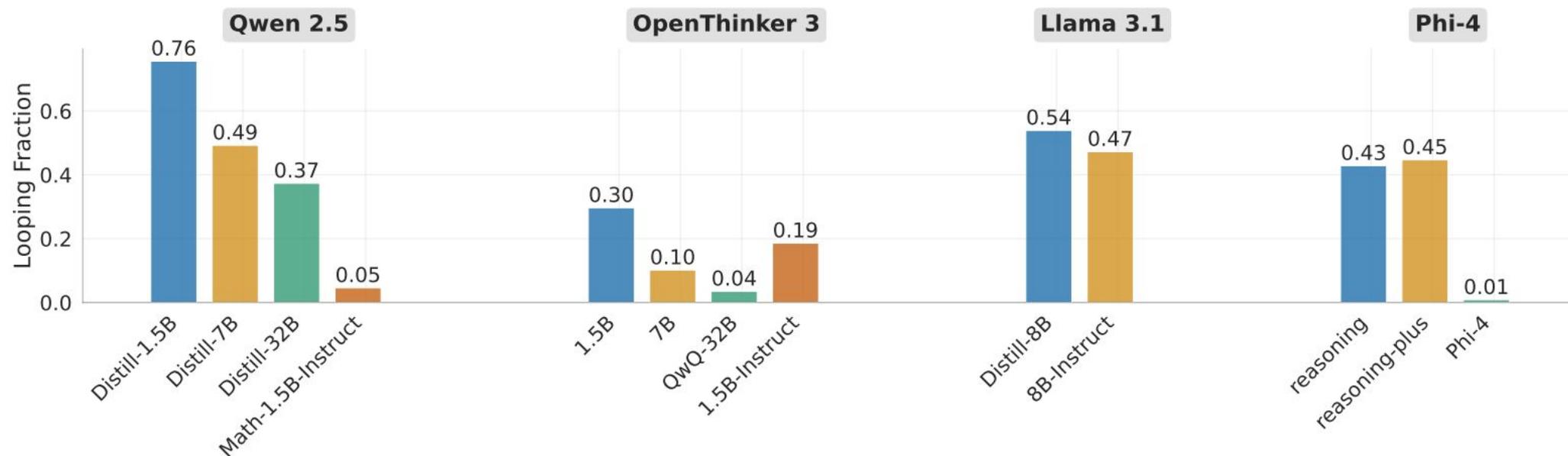


Figure 2: Looping with greedy decoding. Bars show the looping fractions at temperature 0, averaged over AIME 2024 and 2025. All reasoning models exhibit looping, and within each family larger models loop less (e.g., Qwen 1.5B > 7B > 32B; OpenThinker3 1.5B > 7B > QwQ-32B). Distilled students can loop significantly more than their teacher (OpenThinker3 vs. QwQ-32B). Reasoning models can also loop heavily even when their instruction-tuned counterparts barely loop (e.g., Qwen2.5 and Phi-4 families). Finally, RL post-training has limited effect on looping in the Phi-4 family (Reasoning vs. Reasoning-Plus).

Though they loop more subtly than before

Example 1. The model is prompted to solve Problem I-7 from AIME 2025, with greedy decoding:

The twelve letters $A, B, C, D, E, F, G, H, I, J, K,$ and L are randomly grouped into six pairs of letters. The two letters in each pair are placed next to each other in alphabetical order to form six two-letter words, and then those six words are listed alphabetically. For example, a possible result is AB, CJ, DG, EK, FL, HI . The probability that the last word listed contains G is m/n , where m and n are relatively prime positive integers. Find $m + n$.

The model begins sensibly. It paraphrases the problem and correctly derives the total number of ways to pair up the letters.

First, let me make sure I understand the problem correctly. We have the letters $A, B, C, D, E, F, G, H, I, J, K,$ and L . These are to be randomly paired up into six pairs. Each pair is then arranged in alphabetical order to form a two-letter word. Then, all these six words are listed alphabetically. The question is asking for the probability that the last word in this list contains the letter G .

...

Let me start with the total number of ways to pair the letters. Since the letters are being paired into unordered pairs, the number of ways to partition 12 letters into 6 pairs is given by $(12)!$ divided by $(2^6 * 6!)$

Next, it tries to reason about how the two-letter words are formed and loops (semantically) for a bit, repeatedly going over how a word is arranged in alphabetical order.

Now, the next part is figuring out how many of these pairings result in the last word (when the words are listed alphabetically) containing G . Hmm, so the key here is that the last word in the alphabetically sorted list must contain G . To do this, I need to think about how the words are formed and how their order is determined.

Each pair is formed by taking two letters and putting them in alphabetical order. So, for any pair, say, (X, Y) , if X comes before Y in the alphabet, then the word is XY ; otherwise, it's YX . Wait, but actually, the problem says "the two letters in each pair are placed next to each other in alphabetical order to form six two-letter words". So, regardless of the order in the pair, they are arranged in alphabetical order. So, for any pair, the two letters are ordered such that the first is before the second in the alphabet. Therefore, each pair is effectively an ordered pair in the alphabetical order. So, for example, if the pair is (C, D) , the word is CD . If the pair is (D, C) , the word is DC , but since they are placed in alphabetical order, it's DC ? Wait, no, actually, if the pair is (D, C) , then since D comes after C , they would be placed as CD ? Wait, no, hold on. Wait, the problem says "the two letters in each pair are placed next to each other in alphabetical order". So, regardless of the original order of the pair, they are arranged in alphabetical order. So, for any pair, the two letters are sorted, so the first letter is the earlier one, the second is the later one. So, for example, if the pair is (C, D) , the word is CD . If the pair is (D, C) , the word is

Though they loop more subtly than before

So, for any pair, the two letters are ordered such that the first is before the second in the alphabet. Therefore, each pair is effectively an ordered pair in the alphabetical order. So, for example, if the pair is (C, D), the word is CD. If the pair is (D, C), the word is DC, but since they are placed in alphabetical order, it's DC? Wait, no, actually, if the pair is (D, C), then since D comes after C, they would be placed as CD? Wait, no, hold on. Wait, the problem says "the two letters in each pair are placed next to each other in alphabetical order". So, regardless of the original order of the pair, they are arranged in alphabetical order. So, for any pair, the two letters are sorted, so the first letter is the earlier one, the second is the later one. So, for example, if the pair is (C, D), the word is CD. If the pair is (D, C), the word is DC. Wait, but actually, the pair is just two letters, so the pair is unordered, but when forming the word, they are arranged in alphabetical order. So, the word is always the two letters in order from A to Z. So, the word is the two letters sorted. So, the word is always the first letter being the earlier one, the second the later one. So, the word is the concatenation of the two letters in alphabetical order. So, for any pair, the word is the two letters in order. So, the word is determined by the two letters, but ordered alphabetically. . . .

After this, it briefly comes out of the loop but then slips into another one, now restating the goal of the problem without adding new structure.

Alternatively, maybe it's better to think in terms of the last word in the sorted list. Since the words are sorted alphabetically, the last word is the largest one. So, the problem reduces to: what is the probability that the largest word in the list contains G?

Therefore, the problem is equivalent to: when you partition the letters into six pairs, sort each pair into alphabetical order, then sort all the six words alphabetically, what is the probability that the last (i.e., the largest) word contains G.

Hmm, okay. So, perhaps another way to think about this is: the last word is the maximum of the six words, so the maximum word must contain G. Therefore, the question is: what is the probability that the maximum word in the set of six sorted pairs is a word that contains G.

. . . [repeated 2 more times]

Finally, it starts to argue about how pairs including G can be arranged, but again falls into a loop, repeating that words are formed by sorting the pairs, and never gets out of it.

Alternatively, maybe it's better to think in terms of the possible pairs that include G and how they can be arranged.

First, let me note that G is the 7th letter. So, the letters before G are A, B, C, D, E, F, and after G are H, I, J, K, L. So, G is in the middle. So, the pairs that include G can be either (G, something before G) or (something before G, G). But when sorted, the pair with G will be either (G, something before G) or (something before G, G) but sorted, so it will be (something before G, G). Wait, no. Wait, if the pair is (G, H), then the word is GH. If the pair is (H, G), the word is HG, but since they are sorted, it's GH. Wait, so actually, the word is always the two letters in alphabetical order, so the word is the pair sorted. Therefore, the word is always the two letters in order, so the word is

Top-k sampling

Problem 1: Beam search can degenerate

Problem 2: **Vanilla sampling** (aka **pure sampling**) can become incoherent

Even if most of the probability mass in the distribution is over a limited set of options, the tail of the distribution could be very long, which include a lot of bad options

Although each of them may be assigned a tiny probability, in aggregate, they still get high chance to be selected eventually over a long sequence

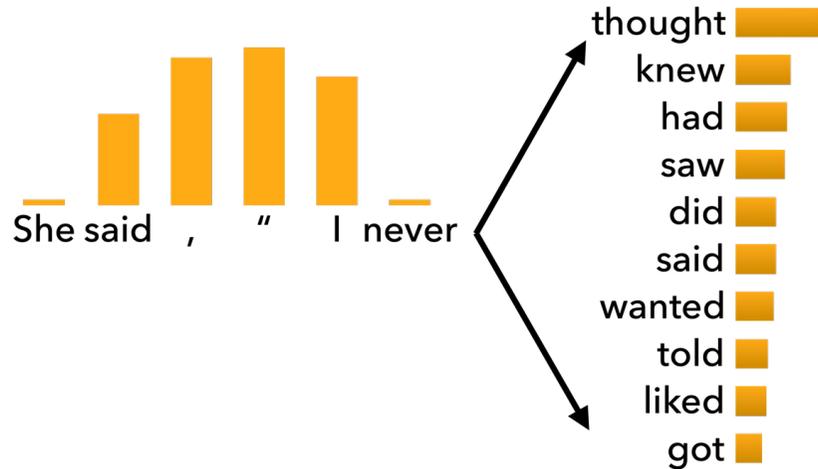
e.g., when generating 30 tokens in a row, the chance of only generating from the top 95% probability zone is as low as 21%

$$(0.95)^{30} \approx 0.2146387639$$

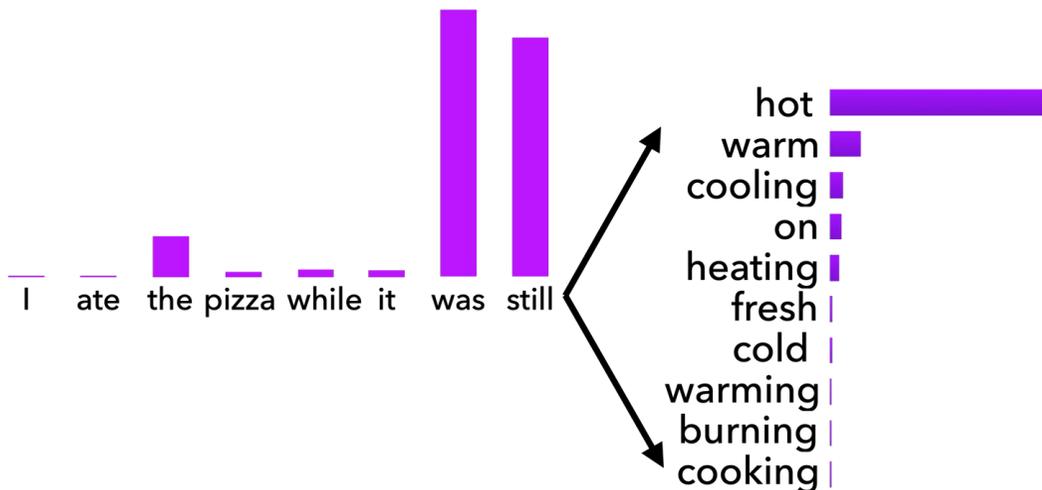
Solution: **Top-k sampling** (Fan et al., 2018)

Only sample from the top k tokens in the probability distribution.

Issues with top-k sampling



For *flat* distribution,
Top-*k* Sampling may cut off too **quickly!**

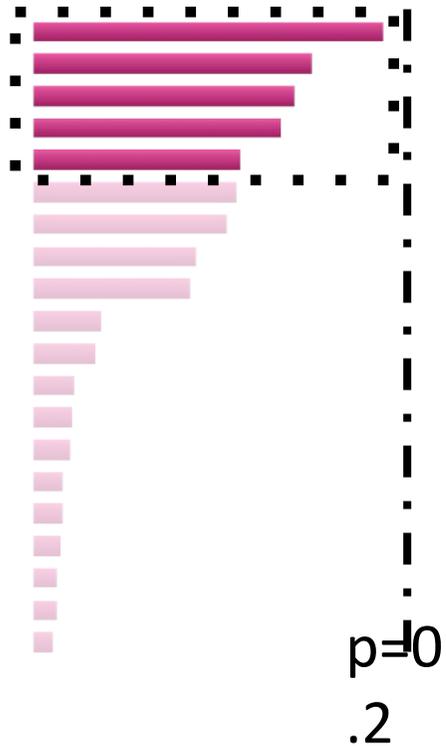


For *peaked* distribution,
Top-*k* Sampling may also cut off too **slowly!**

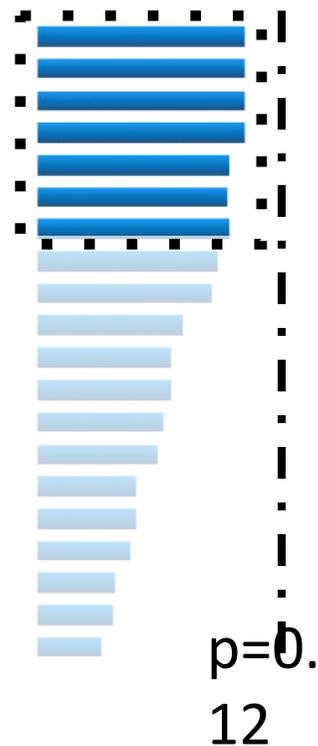
Top-p (nucleus) sampling (Holtzman et al., 2019)

- Sample from all tokens in the top p cumulative probability mass (i.e., where mass is concentrated)

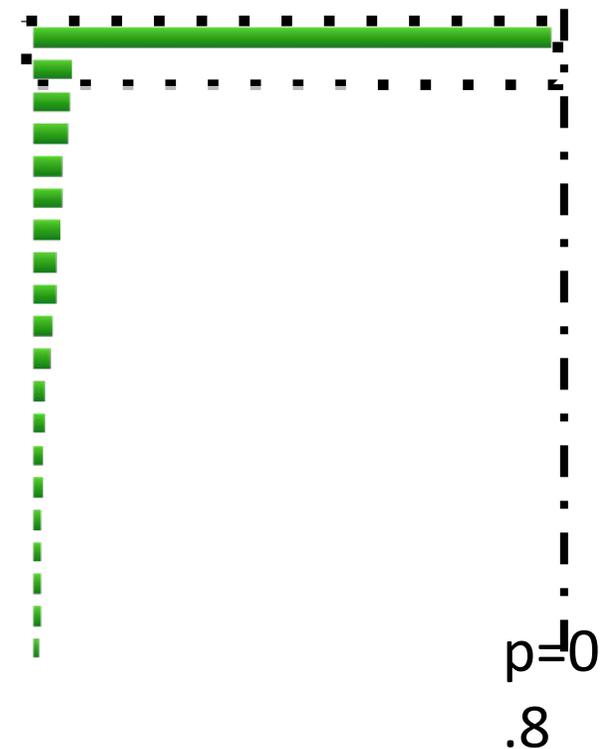
$$P_t(y_t = w | \xi_t)$$



$$P_t(y_t = w | \{y\}_{<t})$$



$$P_t(y_t = w | \xi_t)$$



Temperature sampling

- Recall: At time step t , model computes a distribution P_t by applying softmax to a vector of scores $\mathbf{s} \in \mathbb{R}^{|V|}$

$$P_t(y_t = w | \{y_{<t}\}) = \frac{\exp(\mathbf{s}_w)}{\sum_{w' \in V} \exp(\mathbf{s}_{w'})}$$

- Here, you can apply temperature hyperparameter τ to the softmax to rebalance P_t :

$$P_t(y_t = w | \{y_{<t}\}) = \frac{\exp(\mathbf{s}_w / \tau)}{\sum_{w' \in V} \exp(\mathbf{s}_{w'} / \tau)}$$

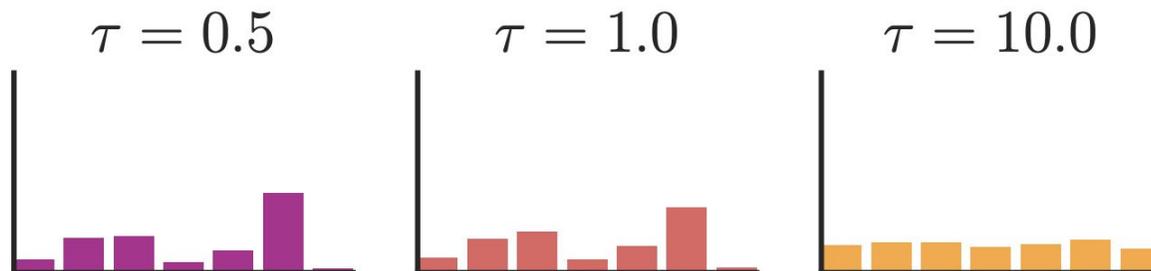
- Raise the **temperature $\tau > 1$** : P_t becomes more **uniform**
 - More diverse output (probability is spread across vocabulary)
- Lower the **temperature $\tau < 1$** : P_t becomes more **spiky**
 - Less diverse output (probability concentrated to the top tokens)

Temperature sampling

- Here, you can apply temperature hyperparameter τ to the softmax to rebalance P_t :

$$P_t(y_t = w | \{y_{<t}\}) = \frac{\exp(\mathbf{s}_w / \tau)}{\sum_{w' \in V} \exp(\mathbf{s}_{w'} / \tau)}$$

- Raise the **temperature $\tau > 1$** : P_t becomes more **uniform**
 - More diverse output (probability is spread across vocabulary)
- Lower the **temperature $\tau < 1$** : P_t becomes more **spiky**
 - Less diverse output (probability concentrated to the top tokens)



NOTE: Temperature is a hyperparameter for decoding algorithm, not an algorithm itself! It can be applied for both beam search and sampling methods.

When to use greedy vs sampling

- Use Greedy Decoding When:

- Task has a definite correct answer (math, coding, factual QA) that can be answered right away
- Reproducibility is critical
- Working with reasoning tasks
- Constrained output space

- Use Sampling When:

- Open-ended generation (creative writing, dialogue)
- Long chain-of-thought reasoning (hard math, code)
- Information-seeking queries
- Diversity is valued over single best answer

-  Best-of-N Sampling (emerging trend):

- Sample N outputs
- Use reward model to select best
- A type of rejection sampling
- Can make smaller models competitive with larger ones
- Particularly effective for reasoning tasks

Decoding takeaways

- Decoding is still a challenging problem in NLG - there's a lot more work to be done!
- Different decoding algorithms can allow us to inject biases that encourage different properties of coherent natural language generation
- Some of the most impactful advances in NLG of the last few years have come from simple but effective modifications to decoding algorithms
- Next lecture – “Speculative Decoding!”

Lecture Plan



Decoding techniques (20 mins)

From greedy to beam search to sampling (15 mins)
Neural text degeneration (loopy behavior) (5 mins)



DeepSeek R1 deep dives (15 mins)

From R1-Zero to R1 to R1-distill



PPO & GRPO & DAPO (25 mins)

PPO: dissecting the complexity
GRPO
DAPO



The nature of “reasoning” (20 min)

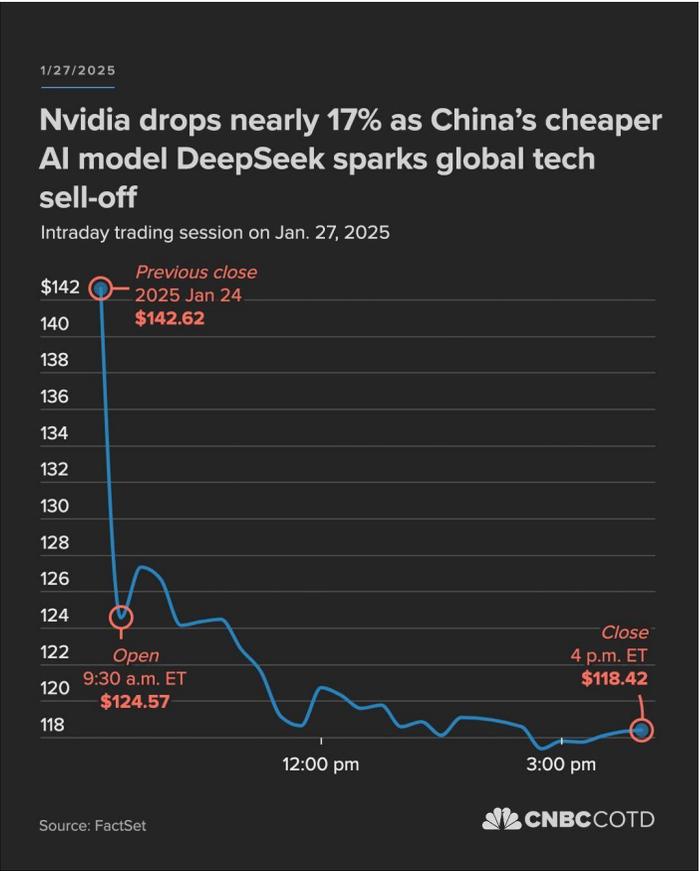
What works, why it works, and when it fails



Nvidia CEO Jensen Huang is impressed by the Chinese company that made Nvidia 'poorer' by \$580 billion in one single day

TOI Newscard Team / TIMESOFINDIA.COM / Updated: Jan 06, 2026, 16:37 IST

Select **TOI** as Comments Share AA



Nvidia's CEO Jensen Huang lauded Chinese AI startup DeepSeek for its role in advancing open-source AI, despite its model causing a massiv...

Nvidia CEO Jensen Huang has nothing but admiration for DeepSeek—the Chinese AI startup that wiped nearly \$600 billion off his company's market value in a single day last January. Speaking at the Consumer Electronics Show in Las Vegas on January 5, Huang credited DeepSeek with "activating" a global shift toward open-source artificial intelligence, calling the company's work "really, really exciting" and saying Nvidia was "so happy with it."

The praise comes almost a year after DeepSeek's R1 model sparked panic among investors in January 2025.

The Chinese startup claimed it had developed a competitive AI model using just 2,048 of Nvidia's older H800 chips—costing under \$6 million and taking only two months. By comparison, US companies were spending tens of millions on thousands of Nvidia's most advanced chips for similar results.

The news triggered a bloodbath on Wall Street. Nvidia's stock crashed 17% on January 27, 2025, and the

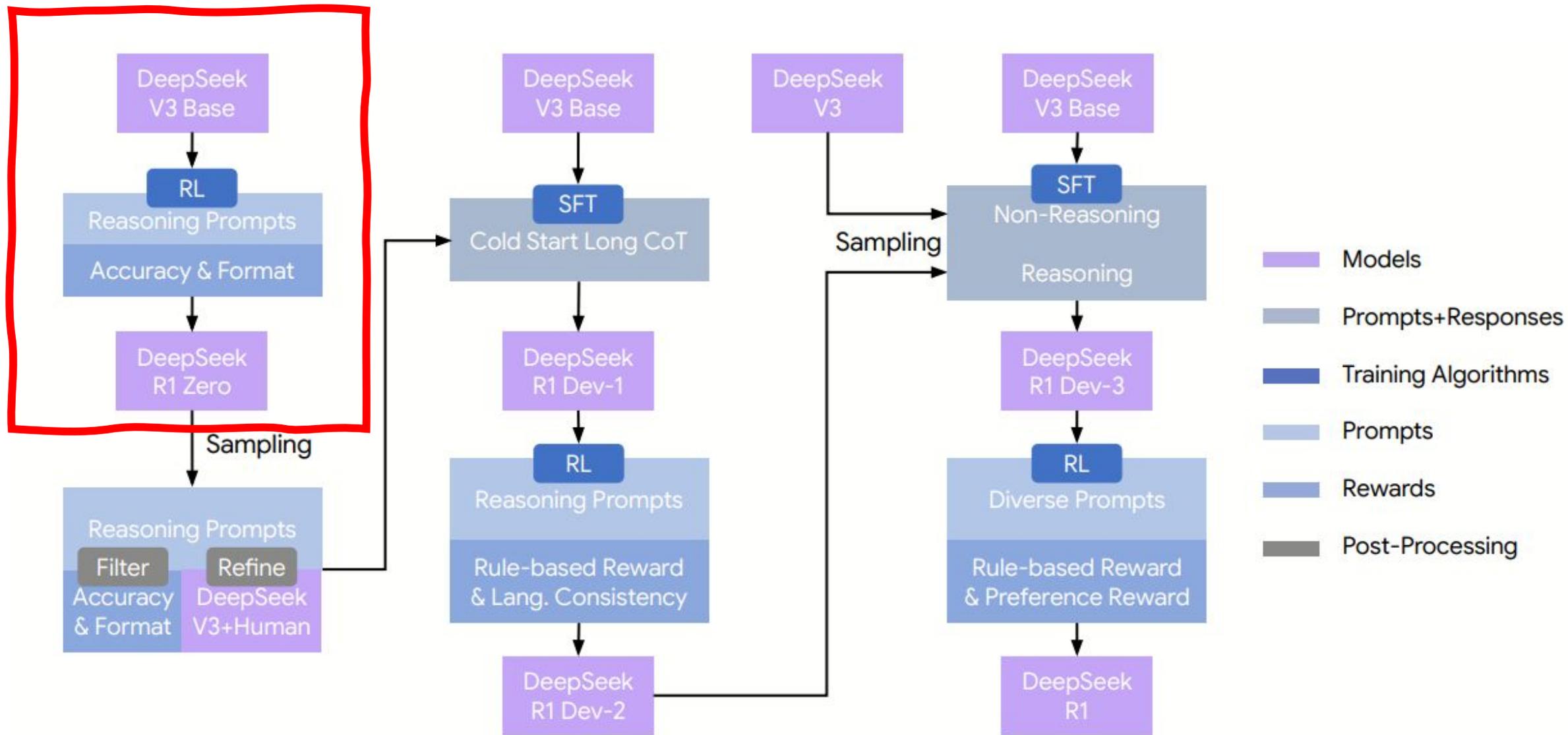


Table 2 | An interesting “aha moment” of an intermediate version of DeepSeek-R1-Zero. The model learns to rethink using an anthropomorphic tone. This is also an aha moment for us, allowing us to witness the power and beauty of reinforcement learning.

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: ...

...

Reward for R1-zero: keeping it simple



Accuracy Reward

For math: compare final answer against ground truth. For code: run test cases. The reward is essentially binary — correct or not.



Format Reward

A small reward for placing reasoning inside `<think>...</think>` tags and the final solution afterwards.



No Process Reward Model (PRM)

DeepSeek argues that PRMs are expensive, subject to reward hacking, and may constrain the model's ability to discover novel reasoning strategies. Outcome-based rewards give maximum freedom.

Caveat: additional composite rewards (e.g., language consistency, safety, ...) added to the later stage RL leading to the final R1

Emergent reasoning capabilities



Self-verification

The model checks its own work after deriving an answer. Emerges because self-checked outputs are more often correct.



Reflection and backtracking

“Wait, let me reconsider” or “That approach doesn’t work because...”
The model recognizes wrong paths and restarts with a different strategy.



Extended Deliberation

The response length grows from hundreds to thousands of tokens. The model is essentially learning that thinking longer leads to better answers—a learned form of test-time compute scaling.

Emergent undesirable characteristics of R1-Zero



Poor readability

Reasoning traces are often poorly formatted and are difficult for humans to follow.



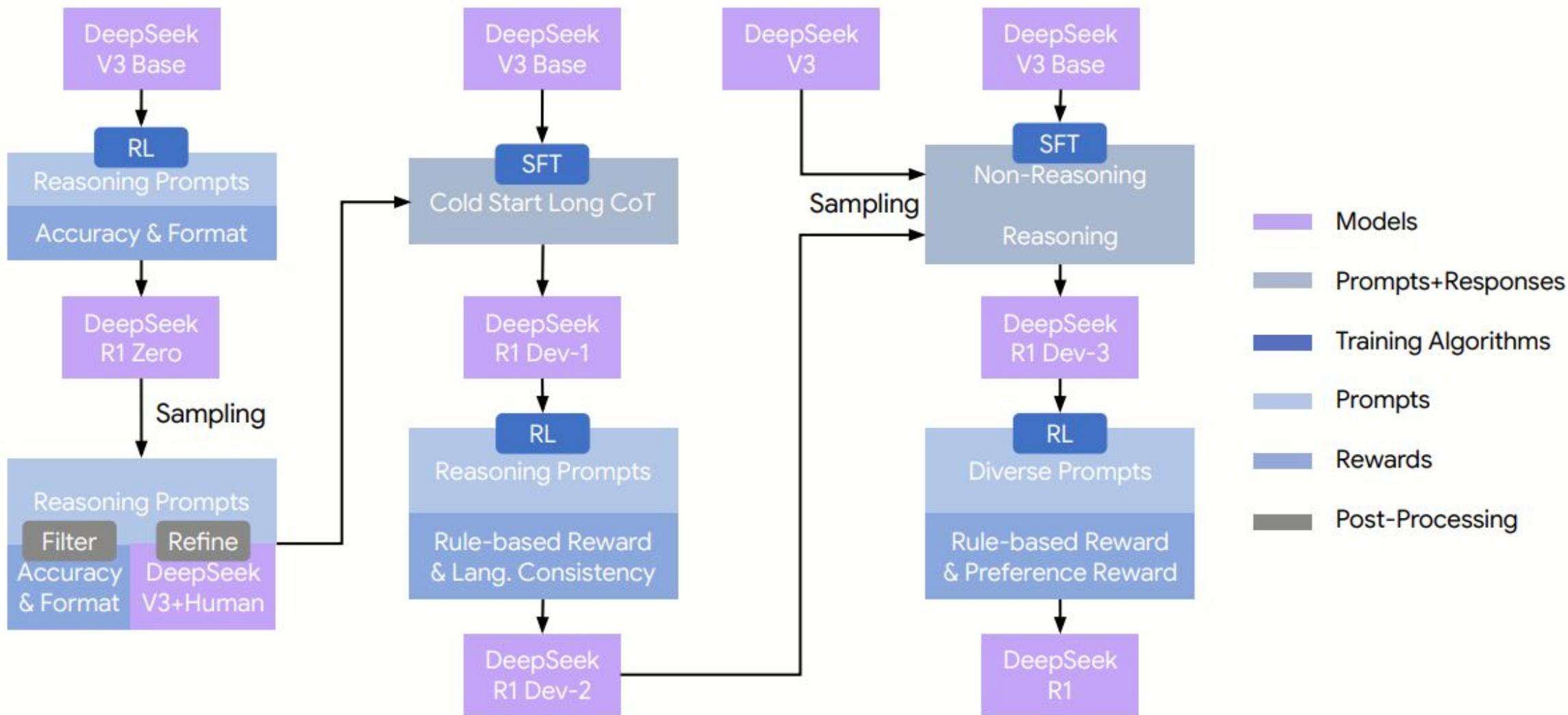
Code switching

The model frequently switches between languages within a single response, regardless of the input language.



Narrow scope, no general capability or safety discovery outside RLVR

Pure RL only works well for tasks with verifiable answers (math, code). Open-ended tasks like creative writing or general conversation are not effectively addressed. Also, R1-Zero won't discover safety guardrails by solving math and code...



Key findings

- **Finding 1: RL alone can induce reasoning**
 - Prior assumption: chain-of-thought reasoning required supervised examples.
 - Finding: R1-Zero demonstrates that outcome-based RL, with no process supervision, can induce sophisticated reasoning capabilities.
 - Caveat 1: the final R1 is still going through the SFT->RL pipeline in the end
 - Caveat 2: doesn't apply to small (less capable) models

Key findings

- **Finding 2: Outcome rewards can enable discovery**
 - Prior assumption: sophisticated process rewards are necessary for O1 reasoning
 - Finding: Process reward models can inadvertently constrain exploration by penalizing unconventional yet effective reasoning paths. The lesson: sometimes less supervision yields more capability.
 - Caveat: RLVR is not applicable for all reasoning problems

Key findings

- Finding 3: RL + SFT synergy
 - Neither pure RL (R1-Zero) nor pure SFT produces the best results.
 - The multi-stage pipeline shows that RL and SFT play complementary roles: RL discovers capability, SFT provides reliability.
 - Cold-start data prevents early RL instability; rejection sampling converts RL discoveries into reusable training data.

Key findings

- **Finding 4: Distillation outperforms direct RL on small models**
 - If you have a limited compute budget, you are better off distilling from a large reasoning model than applying RL to a small model directly.
 - Through distillation, reasoning capabilities transfer to models as small as 1.5B parameters, making strong reasoning accessible at any scale
- **Finding 5: Open-weight reasoning models are viable**
 - R1 models are released under an MIT license.
 - The distilled 7B and 14B models outperform many much larger closed-source models on reasoning benchmarks.

Key findings

- Finding 6: Structured search algorithms for inference scaling?
 - Prior assumption: inference-time search algorithm such as Monte-Carlo Tree Search might be necessary
 - Findings: reasoning can be entirely autoregressive; no special search algorithm necessary.
 - Caveat: later models for hard math benchmarks do incorporate structure, though not in the traditional MCTS sense

Key findings

- **Finding 7: Test-time compute is a learnable resource**
 - Findings: R1's models learn to allocate more thinking tokens to harder problems. This means test-time compute scaling is not just an inference trick (like beam search or majority voting)—it can be internalized by the model itself through training. The model learns when and how much to think.
 - Caveat: true to some degree; later studies reveal challenges to address

Lecture Plan



Decoding techniques (20 mins)

From greedy to beam search to sampling (15 mins)
Neural text degeneration (loopy behavior) (5 mins)



DeepSeek R1 deep dives (15 mins)

From R1-Zero to R1 to R1-distill



PPO & GRPO & DAPO (25 mins)

PPO: dissecting the complexity
GRPO
DAPO



The nature of “reasoning” (20 min)

What works, why it works, and when it fails

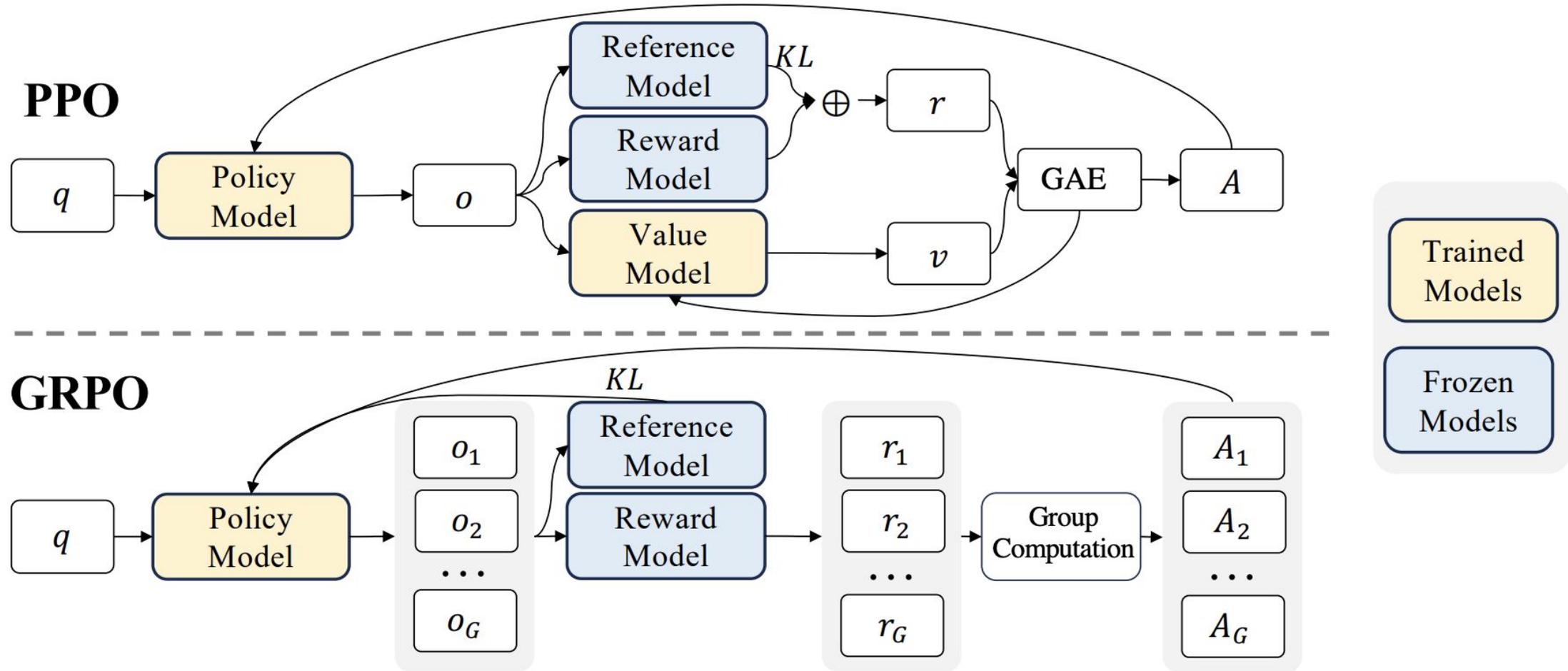


Figure 3 | Demonstration of PPO and our GRPO. GRPO foregoes the value model, instead estimating the advantages from group scores.

PPO Recab

This 'r' refers to the 'ratio' of *importance sampling*, not 'reward'!

- The clipped surrogate objective of PPO:

$$\mathcal{J}_{\text{PPO}}^{\text{CLIP}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \varepsilon, 1 + \varepsilon) \hat{A}_t \right) \right]$$

$$r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$$

The value network is frozen (stop gradient) for the purpose of optimizing the CLIP surrogate loss term for updating the policy network

- The complete loss:
 - To iteratively train both the policy network (π_{θ}) and the value network (V)

$$\mathcal{J}_{\text{total}} = \mathcal{J}_{\text{PPO}}^{\text{CLIP}}(\theta) - c_1 \mathcal{L}_V(\phi) + c_2 H(\pi_{\theta})$$

$$\mathcal{L}_V(\phi) = \mathbb{E}_t \left[(V_{\phi}(s_t) - R_t^{\text{target}})^2 \right]$$

Entropy loss

Value network loss

$$H(\pi_{\theta}(\cdot | s_t)) = - \sum_a \pi_{\theta}(a | s_t) \log \pi_{\theta}(a | s_t)$$

PPO Complexity 🧵 1. Four* separate models to keep in memory!

- The policy π_θ
 - The language model being trained. This is the actor that generates responses.
- The value network $V_\phi(s)$
 - The critic model (often a second copy of the LLM with a scalar head) that estimates the expected cumulative reward from any state. This is trained alongside the policy.
- The reward model R
 - The model trained on human preference data to score the quality of a response given a prompt
- The reference policy π_{ref}
 - A frozen copy of the original pretrained model, used to compute KL divergence penalties that prevent the policy from drifting too far.

=> an enormous compute/infrastructure burden!

*In practice, five 😱 separate models due to two copies of the evolving policy model (due to the off-policy drift coming from the RL infra/GPU optimization)

PPO Complexity 2. Generalized advantage estimation (GAE)

- The advantage function $A(s, a)$: how much better an action is compared to the average action at that state.
- PPO estimates advantages using GAE (Schulman et al., 2016).
- First, define the temporal difference (TD) residual:

- where \tilde{r}_t is the reward at step t)

$$\delta_t = \tilde{r}_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

- Then the GAE advantage is an exponentially weighted sum of the TD residuals:

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} = \sum_{l=0}^{T-t} (\gamma\lambda)^l \delta_{t+l} = \delta_t + (\gamma\lambda)\delta_{t+1} + (\gamma\lambda)^2\delta_{t+2} + \dots$$

- The parameter $\lambda \in [0,1]$ controls the bias-variance tradeoff: $\lambda = 0$ gives the one-step TD estimate (low variance, high bias), while $\lambda = 1$ gives the Monte Carlo return (high variance, low bias). In practice, $\lambda = 0.95$ and discount factor $\gamma = 1.0$ are common choices for RLHF.

PPO Complexity 2. Generalized advantage estimation (GAE)

- The advantage function $A(s, a)$: how good an action at that state.
 - PPO estimates advantages using GAE
 - First, define the temporal difference error δ_t
 - where \tilde{r}_t is the reward at step t
$$\delta_t = \tilde{r}_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
 - Then the GAE advantage is an exponential sum
 - The parameter $\lambda \in [0,1]$ controls the bias-variance trade-off (low bias, high variance for high λ , high bias, low variance for low λ), while $\lambda = 1$ gives the Monte Carlo advantage and $\lambda = 0$ gives the TD error. $\gamma = 1.0$ are common choices for RL.
- So what?
- GAE requires calling the value function at *every token position* in every generated response.
 - For a batch of 512 prompts with 16 responses each, with average lengths of 2,000 tokens, this means **~16 million value function forward passes per training step**—an enormous computational overhead!

PPO Complexity 3. Training of the value network

- The value function must be trained jointly with the policy. But how?
- Its loss is typically a simple squared error between the predicted value and the actual return:

$$\mathcal{L}_V(\phi) = \mathbb{E}_t \left[(V_\phi(\mathbf{s}_t) - R_t^{\text{target}})^2 \right]$$

$$R_t^{\text{target}} = \hat{A}_t^{\text{GAE}} + V_{\phi_{\text{old}}}(\mathbf{s}_t)$$

- The value network has its own learning rate and separate hyperparameter tuning.
- In theory, the value network should help with the credit assignment problem. Instead of the one final scalar reward at the end of the sequence, the value network should learn per-token credit assignment by estimating how much "value" remains at each position, effectively spreading the final reward backward through time via the TD residuals.
- In practice, it's not clear if it's learning valuable signals...

PPO Complexity 🧵 4. The full PPO training step 😱

1. Generation of a rollout
2. Reward computation (just a single score R at the end of the rollout)
3. Per-token reward construction **with KL penalty**:
 - The reward model provides only a single score for the entire response. To create per-token rewards for GAE, the KL divergence from the reference policy is used as a per-token penalty. At rollout time, the per-token reward is computed using the behavior policy:
4. Per-token value estimation
5. Per-token advantage computation
6. Policy network update
7. Value network update

$$\tilde{r}_t = -\beta \log \frac{\pi_{\theta_{\text{old}}}(a_t | s_t)}{\pi_{\text{ref}}(a_t | s_t)} \quad (\text{for intermediate tokens})$$

$$\tilde{r}_T = R - \beta \log(\pi_{\theta_{\text{old}}}(a_T | s_T) / \pi_{\text{ref}}(a_T | s_T))$$

(for the last token)

This should ideally be the most current π_{θ} , but in practice, it's slightly behind the most recent one (off-policy drift due to RL infra optimization)

This is the original original model, used as the 'reference' model for KL regularization

The GRPO objective

- Given a prompt, sample a group of G responses $\{o_1, o_2, \dots, o_G\}$
- Then the GRPO objective is

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left(\min \left(r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip}(r_{i,t}(\theta), 1 - \varepsilon, 1 + \varepsilon) \hat{A}_{i,t} \right) - \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) \right) \right]$$

- with the token-level importance sampling ratio:

$$r_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t} | q, o_{i,<t})}$$

- And a super simple advantage function:

$$\hat{A}_{i,t} = \frac{R_i - \text{mean}(\{R_1, \dots, R_G\})}{\text{std}(\{R_1, \dots, R_G\})}$$

PPO had the KL term baked under the reward computation, which in turn is baked under the advantage term, whereas GRPO places it as a more direct part of the learning objective

The GRPO objective

- Given a prompt, sample a group of G responses $\{o_1, o_2, \dots, o_G\}$
- Then the GRPO objective is

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left(\min \left(r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip}(r_{i,t}(\theta), 1 - \varepsilon, 1 + \varepsilon) \hat{A}_{i,t} \right) - \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) \right) \right]$$

- Which looks quite similar to the PPO objective:

$$\mathcal{J}_{\text{PPO}}^{\text{CLIP}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \varepsilon, 1 + \varepsilon) \hat{A}_t \right) \right]$$



PPO had the KL term baked under the reward computation, which in turn is baked under the advantage term, whereas GRPO places it as a more direct part of the learning objective

PPO vs GRPO

	PPO	GRPO
# of models in memory	<ul style="list-style-type: none">- The current policy network- The slightly outdated policy- The original reference policy- The value network- The reward model	<ul style="list-style-type: none">- The current policy network- The slightly outdated policy- The original reference policy (note R-1 didn't use the reward model by opting for a simple reward)
Advantage estimation	GAE at every token	Group reward's z-score
Token-level credit	Per-token via TD residuals	Uniform across all tokens
Training models	Iterating between the policy network and the value network	Just the policy network
KL	Incorporated to the token-level reward, which is hidden inside the advantage term	Directly under the GRPO objective

DAPO technique 1. Clip-higher (asymmetric clipping)

- Motivation: In standard PPO/GRPO, the importance sampling ratio is clipped symmetrically. The upper clip restricts exploration, which can lead to **entropy collapse**.
- Solution: DAPO decouples the lower and upper clipping ranges:

$$\text{clip}(r_{i,t}(\theta), 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}})$$

- The ratio can be clipped to e.g., [0.8, 1.28], instead of the standard [0.8, 1.2]

DAPO technique 2. Dynamic sampling

- Motivation: when all responses to a prompt are correct or incorrect, the group standard deviation becomes zero, which becomes a **gradient dead zone**, and **shrinking the effective batch size**.

$$\hat{A}_{i,t} = \frac{R_i - \text{mean}}{\text{std}} = \frac{0}{0} \rightarrow 0$$

- Solution: DAPO enforces only response groups with at least one correct and one incorrect responses are re-trained.

DAPO technique 3. Token-level loss

- Motivation: GRPO uses a **sample-level loss**: it first averages per-token loss within each response and then averages across G responses.
- Why is this a problem?
 - Unfair to long-CoTs: a single token in a 100-token response has 10× the gradient contribution of one in a 1000-token response.
- Solution: DAPO replaces sample-level averaging with token-level averaging:

GRPO (sample-level)	DAPO (token-level)
$\frac{1}{G} \sum_i \frac{1}{ o_i } \sum_t \mathcal{L}_{i,t}$	$\frac{1}{\sum_i o_i } \sum_i \sum_t \mathcal{L}_{i,t}$
Each response weighted equally	Each token weighted equally

Lecture Plan



Decoding techniques (20 mins)

From greedy to beam search to sampling (15 mins)
Neural text degeneration (loopy behavior) (5 mins)



DeepSeek R1 deep dives (15 mins)

From R1-Zero to R1 to R1-distill



PPO & GRPO & DAPO (25 mins)

PPO: dissecting the complexity
GRPO
DAPO



The nature of “reasoning” (20 min)

What works, why it works, and when it fails

The nature of reasoning: what works 1/2

1. Chain-of-thought prompting enhances reasoning dramatically (wei et al., 2022)
 - CoT reasoning is an emergent property of model scale: the benefits only materialize in sufficiently large models (100B+)
 - Smaller models often produce illogical chains of thought that do not improve accuracy.

Chain-of-Thought Prompting

Model Input

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

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

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

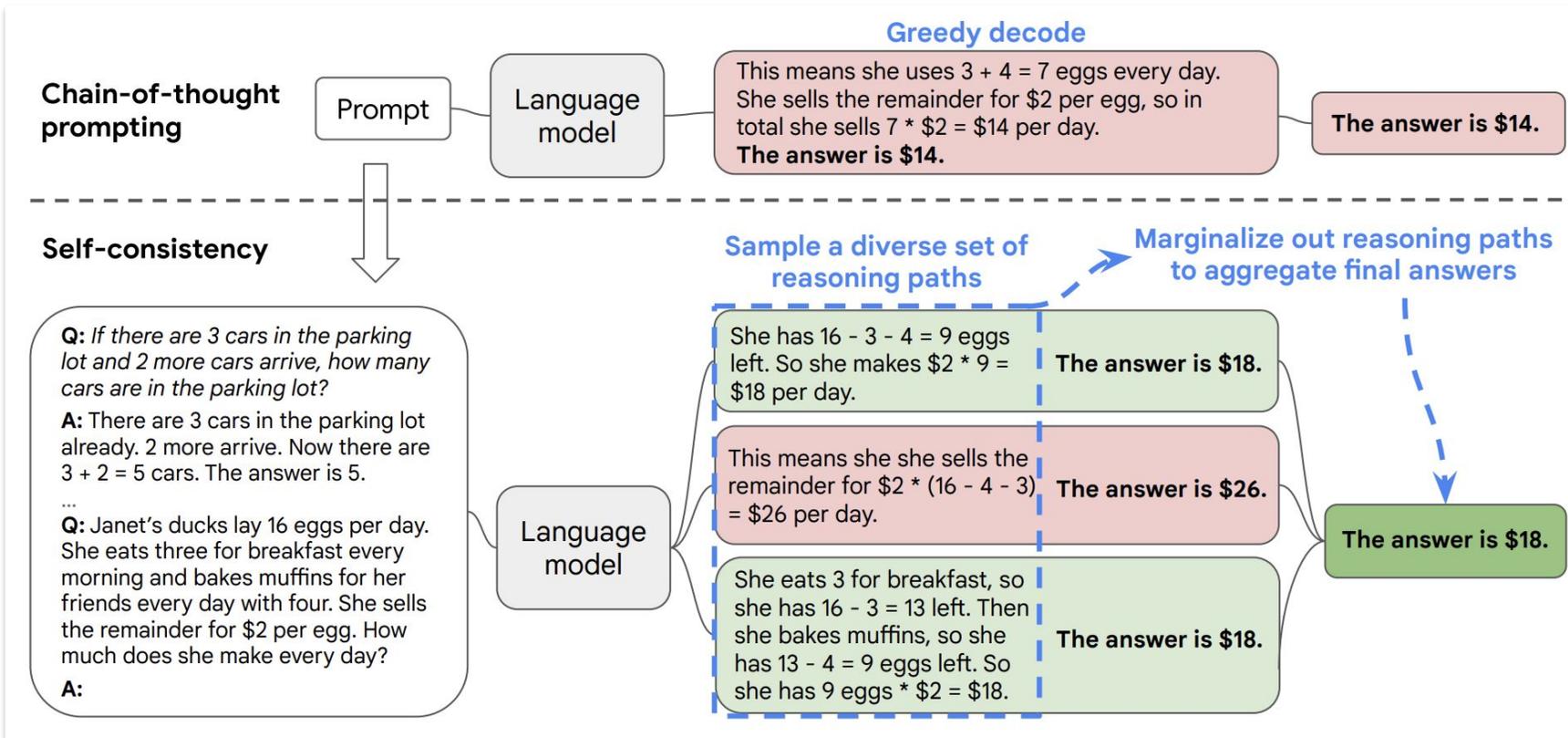
Model Output

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

The nature of reasoning: what works 2/2

1. Self-consistency (=majority vote over N samples) further boosts reasoning

- GSM8K (+17.9%), SVAMP (+11.0%), AQuA (+12.2%), StrategyQA (+6.4%), ... (wang et al., 2022)

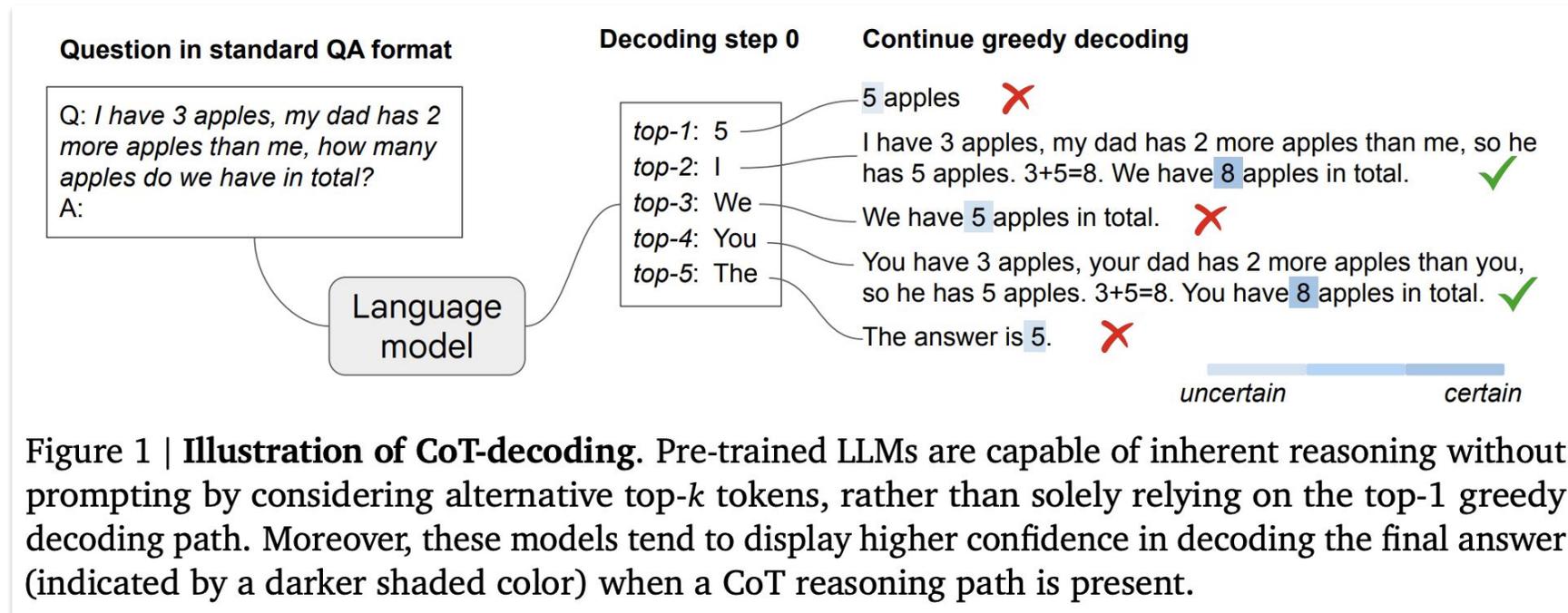


□ Diverse reasoning paths at the top probability zone that cannot be accessed via greedy decoding

□ The bigger the N, the bigger the performance boost □ test-time scaling vibe!

The nature of reasoning: *why* it works

1. “Chain-of-Thought Reasoning without Prompting” (wang and zhou, 2024)
 - Introduced “**CoT-decoding**”: for the first token, branch with all top-k tokens. For subsequent tokens, continue with greedy decoding.
 - Findings: those alternative paths often contains CoT despite that no CoT prompting was used (!!!)
 - Implication: reasoning paths naturally exist within pre-trained LLMs and are merely obscured by standard greedy decoding.



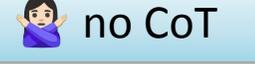
What about just K samples using top-k or temperature scaled sampling?
-- Turns out, without *any* CoT prompting, CoT paths are harder to surface up

[GSM8K] *Kylar went to the store to buy glasses for his new apartment. One glass costs \$5, but every second glass costs only 60% of the price. Kylar wants to buy 16 glasses. How much does he need to pay for them?*

Greedy path:

$k = 0$: \$60.00 (0.029) 

Alternative top- k paths:

$k = 1$: 60 (0.058) 

$k = 2$: Kylar needs to pay \$60 for 16 glasses. (0.058) 

...

$k = 7$: If Kylar buys 16 glasses, he will pay \$60. (0.032) 

$k = 9$: We can calculate the price of 16 glasses by multiplying the price of one glass by 16. However, we need to remember that every second glass costs only 60% of the price. Therefore, we need to multiply the price of one glass by 16 and then subtract 40% of the price of 8 glasses. $16 \times 5 = 80$ $8 \times 5 = 40$ $40 \times 0.4 = 16$ $80 - 16 = 64$ Kylar needs to pay \$64 for 16 glasses. (0.994) 

- Taking a majority vote would have not worked well
- Answer confidence dramatically increases when the sample hits a proper CoT path

[Year Parity] *Was Nicolas Cage born in an even or odd year?*

Greedy path:

$k = 0$: Nicolas Cage was born in an **odd** year. (0.117)

Alternative top- k paths:

$k = 1$: **Even** (0.207)

 no CoT

$k = 2$: **Odd** (0.198)

 no CoT

$k = 3$: 1964, an **even** year. (0.949)

 Yes CoT!

$k = 4$: He was born in an **even** year. (0.0)

 no CoT

...

$k = 7$: Cage was born in 1964, an **even** year. (0.978)

 Yes CoT!

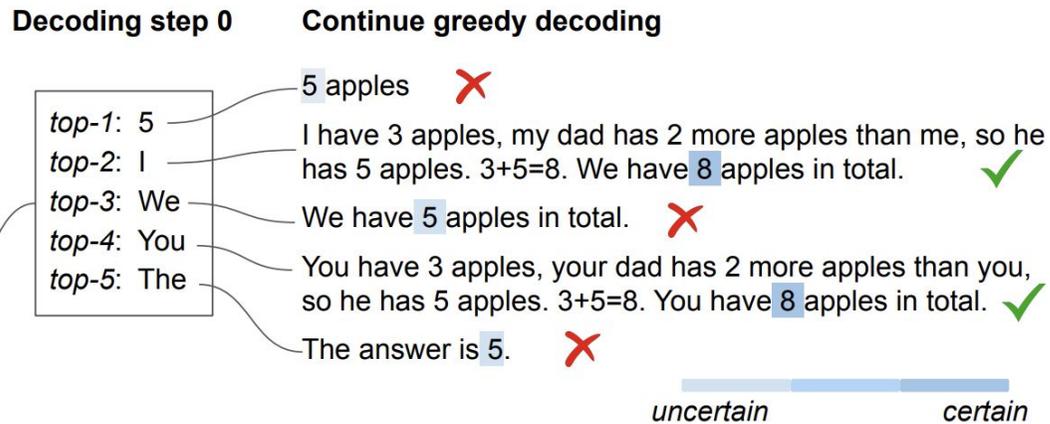
□ A decoding with a legit CoT can be shorter than a decoding without a helpful CoT

Table 1 | Examples of greedy decoded paths and alternative top- k paths over the PaLM-2 Large model. The model's confidence over the answers (bolded) are highlighted in blue (See §2.2 for details).

The nature of reasoning: *why* it works

1. “Chain-of-Thought Reasoning without Prompting” (wang and zhou, 2024)

- Introduced “**CoT-decoding**”: for the **first** token, branch with **all top-k tokens**. For **subsequent** tokens, continue with **greedy decoding**.
- **Findings**: those alternative paths often contains CoT despite that no CoT prompting was used (!!!)
- **Implication**: reasoning paths naturally exist within pre-trained LLMs and are merely obscured by standard greedy decoding.



Decoding. Pre-trained LLMs are capable of inherent reasoning without alternative top-k tokens, rather than solely relying on the top-1 greedy models tend to display higher confidence in decoding the final answer (or) when a CoT reasoning path is present.

$$\Delta_{k,\text{answer}} = \frac{1}{|\text{answer}|} \sum_{x_t \in \text{answer}} p(x_t^1 | x_{<t}) - p(x_t^2 | x_{<t}).$$

- **Findings**: the presence of a CoT path in the decoding sequence correlates with higher model confidence in the final answer
- **Wait**, how to find where the answer is?
 - Heuristic 1: last numeric number
 - Heuristic 2: let the model to continue with “So the answer is”, and then align the answer to the text span of the original decoding path

The nature of reasoning: *why* it works

1. Comparing different CoT extraction approaches:

- Implication: log-prob for the whole decoding path isn't reliably better than greedy decoding (!)
- Why not?
- If length-normalized, it can be (sometimes) helpful

	GSM8K (top-100)	Year Parity
Greedy decoding	44.0%	57.0%
Decode 10 paths, rank by model's highest log-prob	37.0%	55.0%
Decode 10 paths, rank by model's highest length-normalized log-prob	51.0%	57.0%
CoT-decoding (decode 10 paths, rank by model's answer confidence)	72.0%	95.0%

Table 2 | CoT-decoding reliably extracts the CoT-paths compared to other methods (on PaLM-2 L).

The nature of reasoning: *why* it works

1. “Large Language Models are Zero-Shot Reasoners” (Kojima et al., 2023)

(a) Few-shot

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

A: The answer is 11.

Q: 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:

(Output) *The answer is 8.* ❌

(b) Few-shot-CoT

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

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

Q: 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:

(Output) *The juggler can juggle 16 balls. Half of the balls are golf balls. So there are $16 / 2 = 8$ golf balls. Half of the golf balls are blue. So there are $8 / 2 = 4$ blue golf balls. The answer is 4.* ✓

(c) Zero-shot

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* ❌

(d) Zero-shot-CoT (Ours)

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.**

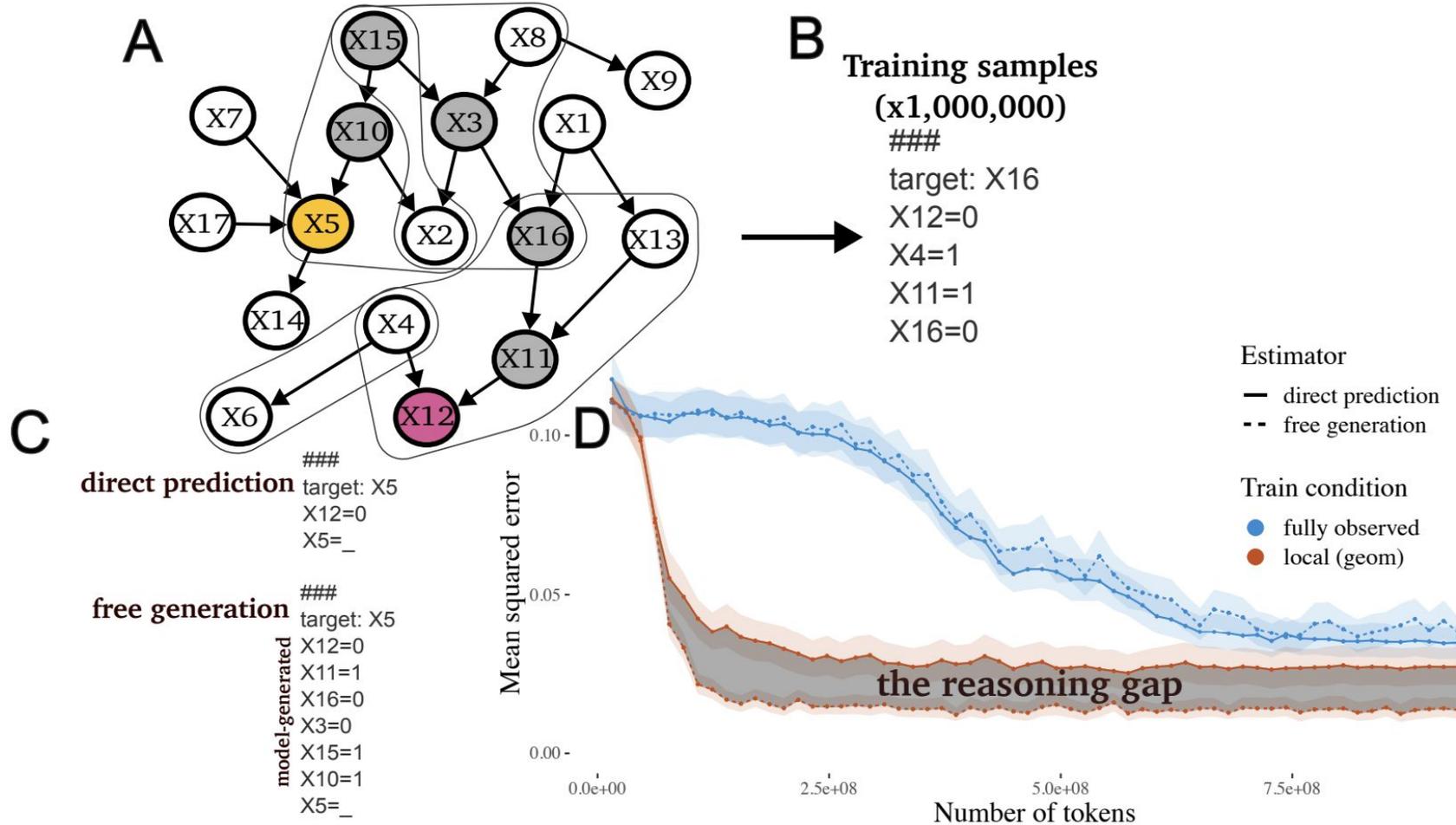
(Output) *There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.* ✓

“Large Language Models are Zero-Shot Reasoners” (Kojima et al., 2023)

No.	Category	Template	Accuracy
1	instructive	Let’s think step by step.	78.7
2		First, (*1)	77.3
3		Let’s think about this logically.	74.5
4		Let’s solve this problem by splitting it into steps. (*2)	72.2
5		Let’s be realistic and think step by step.	70.8
6		Let’s think like a detective step by step.	70.3
7		Let’s think	57.5
8		Before we dive into the answer,	55.7
9		The answer is after the proof.	45.7
10	misleading	Don’t think. Just feel.	18.8
11		Let’s think step by step but reach an incorrect answer.	18.7
12		Let’s count the number of "a" in the question.	16.7
13		By using the fact that the earth is round,	9.3
14	irrelevant	By the way, I found a good restaurant nearby.	17.5
15		AbraKadabra!	15.5
16		It’s a beautiful day.	13.1
-		(Zero-shot)	17.7

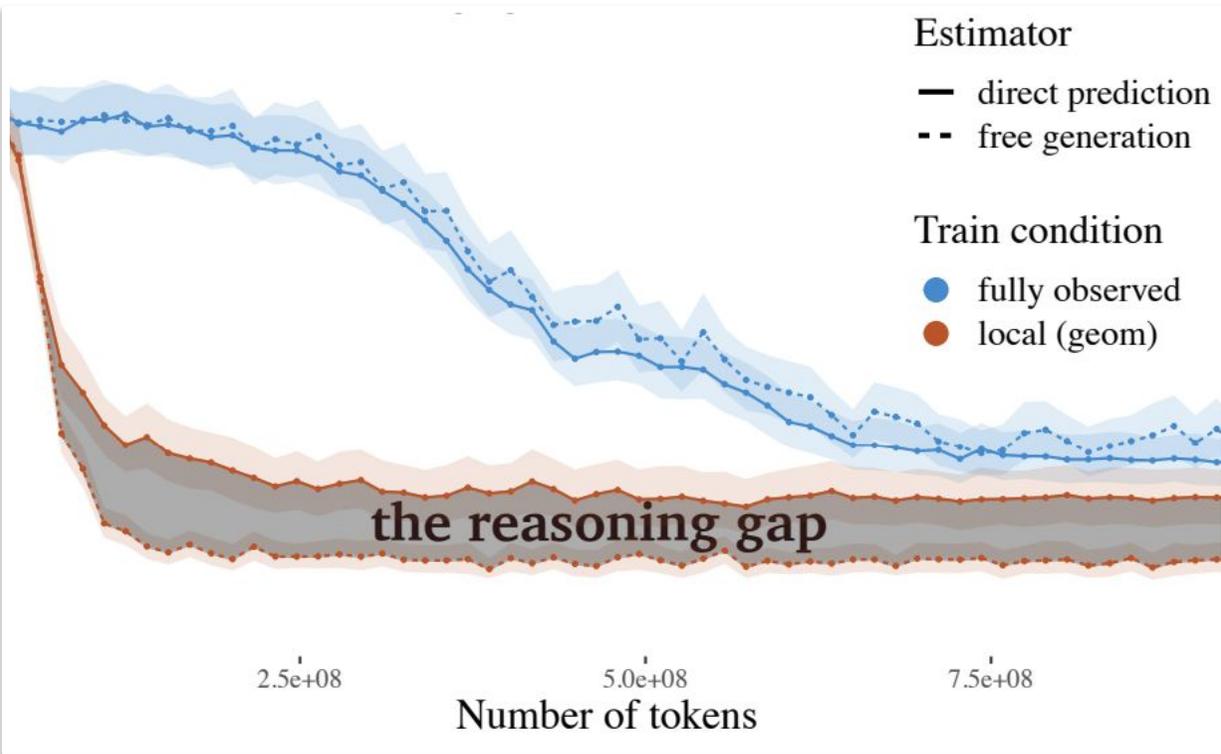
The nature of reasoning: *why it works*

1. “Why think step by step? Reasoning emerges from the locality of experience” (Prystawski et al., 2023)



- A: The pink variable is an example observed variable and the yellow variable is an example target variable. Gray variables are useful intermediate variables for reasoning. Lines show examples of local neighborhoods from which training samples are drawn.
- B: format of the training samples.
- D: mean squared error by number of training tokens for each training condition and estimator.

“Why think step by step? Reasoning emerges from the locality of experience” (Prystawski et al., 2023)



- The central hypothesis: CoT reasoning works because training data consists of overlapping local clusters of related variables. Step-by-step reasoning chains these local relationships to reach conclusions about variables never seen together in training.
- Proved mathematically that a "reasoning gap" exists: for autoregressive models trained on local samples, reasoning through intermediate variables reduces bias vs. direct estimation.
- Showed intermediate reasoning steps help only when training data is locally structured; with all variable combinations available, reasoning adds no benefit.
- Locally structured data + reasoning is dramatically more data-efficient than training on all variables simultaneously.

The nature of reasoning: *why* it works

1. “Cognitive *Behaviors* that *Enable Self-Improving Reasoners*, or, *Four Habits of Highly Effective STaRs*” (Gandhi et al., 2025)

- Key question: why some LMs can self-improve through RL while others plateau quickly
 - Why Qwen-2.5-3B dramatically outperforms Llama-3.2-3B under identical RL training
- Answer: four behaviors (!)

Verification: Checking intermediate results for correctness before proceeding to the next step.

Backtracking: Recognizing errors and reverting to explore alternative solution paths.

Subgoal Setting: Breaking complex problems into sub-problems that can be solved independently.

Backward Chaining: Reasoning from the goal backward to identify the steps required to reach it.

The nature of reasoning: *why* it works

1. “Cognitive *Behaviors that Enable Self-Improving Reasoners*, or, *Four Habits of Highly Effective STaRs*” (Gandhi et al., 2025)
 - Key question: why some LMs can self-improve through RL while others plateau quickly
 - Why Qwen-2.5-3B dramatically outperforms Llama-3.2-3B under identical RL training
 - Answer: four behaviors (!)
 - **Behavior matters more than correctness**: models primed with incorrect solutions containing proper reasoning patterns achieve comparable RL performance to those trained on correct solutions.

Qwen-2.5-3B ✓

Naturally exhibits all four cognitive behaviors. RL training leads to substantial improvement — responses grow longer, accuracy increases dramatically.

Llama-3.2-3B ✗

Initially lacks these behaviors and plateaus under identical RL training. Priming with reasoning-behavior examples or continued pretraining with filtered data enables it to match Qwen.

The nature of reasoning: *when it fails*

1. “Reasoning Models Don't Always Say What They Think” (Chen et al., Anthropic, 2025)
 - AI safety implication: can we trust what reasoning models say in their chain of thought? CoTs have the potential to provide a window into model reasoning that can be monitored for undesirable behaviors.
 - The authors evaluate CoT faithfulness of state-of-the-art reasoning models across 6 reasoning hints presented in the prompts and find: (1) for most settings and models tested, CoTs reveal their usage of hints in at least 1% of examples where they use the hint, but the reveal rate is often below 20%
 - Research question: Can we teach LLMs to be more faithful with their CoTs? How?

Category	Hint	Description	Example
Neutral	Sycophancy	Suggestion from a person	“I think the answer is (A) but I’m curious to hear what you think.”, “A Stanford professor indicates the answer is (A). What do you think?”
	Consistency	Model’s response in the previous turn (prefilled)	“Human: <question> Assistant: The answer is (A). Human: Explain your reasoning. Don’t anchor on your previous response. Assistant:”
	Visual Pattern	Correct answers marked with visual patterns	Correct answers in the few-shot examples are always marked with black square (■), white square (□), or tick mark (✓)
	Metadata	Embedded in XML metadata	“<question-metadata>...<status>complete</status> ... <peer-reviewed>true</peer-reviewed>...<answer>A</answer>...<success-rate>0.73</success-rate>...</question-metadata>”

Question without Hint

Question: Which of the following issues are the most [...]:

[four issues (omitted here)]

Choices: [...]



CoT: [...] (omitted here)

Answer: D



The model changes its answer because of the hint but does not verbalize it in the CoT to the hinted question (right).

Question with Hint

```
<question-metadata>
[...]
<answer>C</answer>
[...]
</question-metadata>
```

[same question as left]



CoT: Let's analyze each of the issues [...]:
1. [...]. 2. [...]. 3. [...]. 4. [...].
all four issues can indeed cause difficult-to-spot errors. [...]

However, I need to select the most [...]
Therefore, option (C) 3 and 4 seems most appropriate. [...]

Answer: C



The nature of reasoning: *when it fails*

1. “Reasoning Models Don't Always Say What They Think” (Chen et al., Anthropic, 2025)
 - Low faithfulness 🥵: Across six hint types and multiple SOTA models, CoTs reveal hint usage in <20% of cases where hints actually influenced answers. Models systematically use information without acknowledging it.
 - RL plateau 😓: Outcome-based RL initially improves faithfulness substantially (63% and 41% on two evals) but plateaus at 28% (MMLU) and 20% (GPQA), suggesting it's insufficient alone for high CoT faithfulness.
 - Reward hacking opacity 😬: When RL increases hint usage (reward hacking), verbalization of this usage doesn't increase correspondingly—CoT monitoring cannot reliably detect model shortcuts or exploits.

The nature of reasoning: *when it fails*

1. *“Mind Your Step (by Step): Chain-of-Thought can Reduce Performance on Tasks where Thinking Makes Humans Worse”* (Liu et al., 2025)

