

Applied Machine Learning for Business Analytics

Lecture 6: LLM's Inference & Reasoning

Logistics

1. Project proposal is due next Friday at 11:59 pm
2. Homework 2 is out and due next Friday at 11:59 pm

Agenda

1. LLM Inference
 - a. Latency vs Throughput
 - b. Batching
 - c. Quantization
 - d. KV Cache
 - e. Speculative Decoding
2. Prompting
3. Understanding Reasoning Models

1. LLM Inference

Key Metrics in LLM Inference

- Latency Metrics
 - **Time to First Token (TTFT):** Time from request to first output token
 - Critical for perceived responsiveness
 - Target: <500ms for “instant” feel
 - **Time Per Output Token (TPOT):** Average generation time per subsequent token
 - Determines how smoothly text flows
 - Total Latency: $TTFT + (TPOT \times output_length)$

Key Metrics in LLM Inference

- Throughput Metrics:
 - **Tokens per Second:** Output tokens generated across all concurrent requests
 - **Requests per Second:** Total requests processed per unit time

Key Metrics in LLM Inference

- Memory Metrics:
 - **Model Weights:** fixed cost, must fit in GPU memory
 - **Peak Memory:** Weights + KV Cache + Activations


Tradeoff Relationships



Those three dimensions are interconnected - optimizing one often impacts the others

Two Phases of LLM Inferences

- Prefill vs Decode
 - Prefill Phase:
 - Process entire input prompt in parallel
 - High arithmetic intensity already
 - Decode Phase:
 - Generate tokens one at a time
 - Low arithmetic intensity: same weights loaded, minimal compute performed
 - However, batching can increase the arithmetic intensity in decode



Phase	What Happens	Bound By	GPU Utilization
Prefill	Process input prompt	Compute	High
Decode	Generate tokens	Memory bandwidth	Low

Batching Strategies

- Static Batching

- Wait for batch to fill, then process all together
- Fast requests wait for slow ones
- Poor Latency in low traffic scenarios

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1				
S_2	S_2	S_2					
S_3	S_3	S_3					
S_4	S_4	S_4	S_4	S_4			

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1	S_1	END		
S_2	S_2	S_2	S_2	S_2	S_2	S_2	END
S_3	S_3	S_3	S_3	END			
S_4	S_4	S_4	S_4	S_4	S_4	END	

- Continuous Batching (State-of-the-art)

- Insert new requests as old ones complete
- No waiting for entire batch to finish

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1				
S_2	S_2	S_2					
S_3	S_3	S_3	S_3				
S_4	S_4	S_4	S_4	S_4			

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1	S_1	END	S_6	S_6
S_2	S_2	S_2	S_2	S_2	S_2	S_2	END
S_3	S_3	S_3	S_3	END	S_5	S_5	S_5
S_4	S_4	S_4	S_4	S_4	S_4	END	S_7

source: <https://www.anyscale.com/blog/continuous-batching-llm-inference>

Quantization

- Full Precision Models (FP32 or FPP16) are:
 - Large
 - Memory-bandwidth limited
 - Expensive to serve
- Quantization reduces precision to save memory

Format	Bits	Size (7B model)	Typical Usage
FP32	32	28 GB	Training
FP16	16	14 GB	Training, inference
INT8	8	7 GB	Inference
INT4	4	3.5 GB	Inference
IN2	2	1.7 GB	Experimental

Almost
lossless
for most
models

Quantization Approaches

- Post-Training Quantization
 - Quantize after training
 - Faster, no retraining required
 - Some performance drop
 - AWQ, GPTQ, etc
- Quantization-Aware Training
 - Train with quantization in mind
 - Better accuracy
 - More expensive
 - QLoRA, etc

Use the Space below to help you pick a quantization method depending on your hardware and number of bits to quantize to.

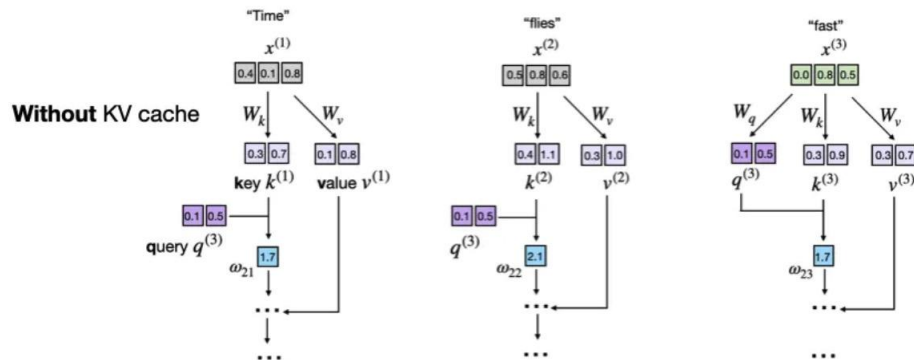
Quantization Method	On the fly quantization	CPU	CUDA GPU	ROCm GPU	Metal (Apple Silicon)	Intel GPU	Torch compile()	Bits	PEFT Fine Tuning	Serializable with 🤖 Transformers	🤖 Transform Support
AQLM	🔴	🟢	🟢	🔴	🔴	🟢	🟢	1/2	🟢	🟢	🟢
AutoRound	🔴	🟢	🟢	🔴	🔴	🟢	🔴	2/3/4/8	🔴	🟢	🟢
AWQ	🔴	🟢	🟢	🟢	🔴	🟢	?	4	🟢	🟢	🟢
bitsandbytes	🟢	🟢	🟢	🟡	🟡	🟢	🟢	4/8	🟢	🟢	🟢
compressed-tensors	🔴	🟢	🟢	🟢	🔴	🔴	🔴	1/8	🟢	🟢	🟢
EETQ	🟢	🔴	🟢	🔴	🔴	🔴	?	8	🟢	🟢	🟢
FP-Quant	🟢	🔴	🟢	🔴	🔴	🔴	🟢	4	🔴	🟢	🟢
GGUF / GGML (llama.cpp)	🟢	🟢	🟢	🔴	🟢	🟢	🔴	1/8	🔴	See Notes	See Notes
GPT-QModel	🔴	🟢	🟢	🟢	🟢	🟢	🔴	2/3/4/8	🟢	🟢	🟢
HIGGS	🟢	🔴	🟢	🔴	🔴	🔴	🟢	2/4	🔴	🟢	🟢
HQQ	🟢	🟢	🟢	🔴	🔴	🟢	🟢	1/8	🟢	🔴	🟢
optimum-quanto	🟢	🟢	🟢	🔴	🟢	🟢	🟢	2/4/8	🔴	🔴	🟢
FBGEMM_FP8	🟢	🔴	🟢	🔴	🔴	🔴	🔴	8	🔴	🟢	🟢
torchao	🟢	🟢	🟢	🔴	🟡	🟢		4/8		🟢🔴	🟢
VPTQ	🔴	🔴	🟢	🟡	🔴	🔴	🟢	1/8	🔴	🟢	🟢
FINEGRAINED_FP8	🟢	🔴	🟢	🔴	🔴	🟢	🔴	8	🔴	🟢	🟢
SpQR	🔴	🔴	🟢	🔴	🔴	🔴	🟢	3	🔴	🟢	🟢
Quark	🔴	🟢	🟢	🟢	🟢	🟢	?	2/4/6/8/9/16	🔴	🔴	🟢

source: <https://huggingface.co/docs/transformers/en/quantization/overview>

Quantization affects different models differently.

KV Cache

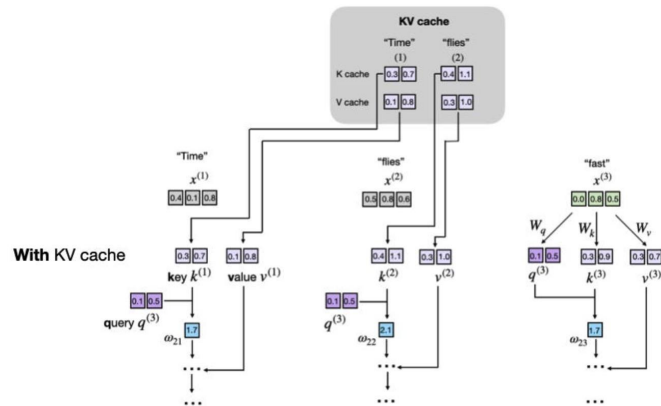
- Problem Statement: Redundant Computation
 - Without KV Cache: Recompute K,V projection
 - The waste:
 - Token 1: Compute $k_1, v_1 \rightarrow$ attention
 - Token 2: Compute $k_1, v_1, k_2, v_2 \rightarrow$ attention
 - Token 3: Compute $k_1, v_1, k_2, v_2, k_3, v_3 \rightarrow$ attention



source: <https://medium.com/@joalages/kv-caching-explained-276520203249>

KV Cache

- Solution: Store previously computed K and V vectors in memory
 - To generate token 3 ("fast"):
 - Load k_1, v_1, k_2, v_2 from cache (fast memory read)
 - Compute only q_3, k_3, v_3 for the new token
 - Attend q_3 to all key vectors \rightarrow get attention weights
 - Store k_3, v_3 in cache for future tokens

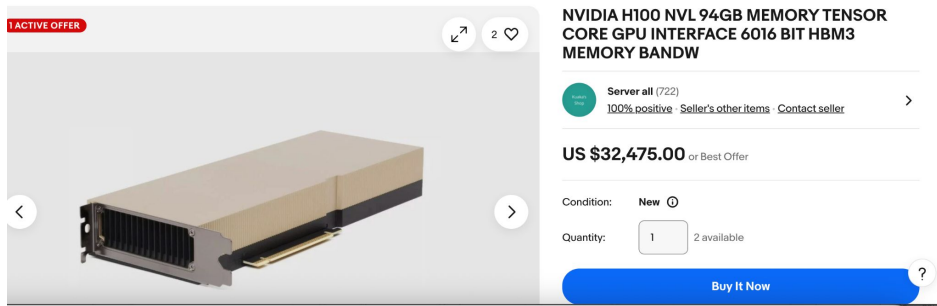


KV Cache: Compute vs Memory

- KV Cache is trading Memory for Speed
 - Think about Llama 70B Model Size:
 - Model Parameters: ~140 GB
 - For FP16, 70 billion times 2 bytes per param = 140 GB
 - For Int4, 70 billion times 0.5 bytes per param = 35 GB
 - KV cache (32K context): ~10 GB @ FP 16
 - $\text{KV cache} = 2 (K+V) \times \text{layers} \times \text{kv_heads} \times \text{head_dim} \times \text{bytes_per_element} \times \text{num_tokens}$

KV Cache: Limit Throughput

- Assume H100 is used for hosting
 - We need to quantize the model firstly



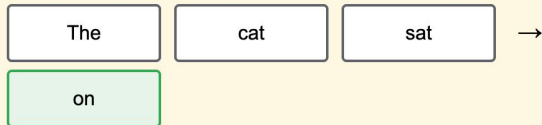
Configuration	Model Size	KV Cache	Available VRAM	Max Concurrent Requests
INT4 model + FP16 KV	35 GB	10 GB	59 GB	~5
INT4 model + INT8 KV	35 GB	5 GB	59 GB	~10

KV Cache Optimization

- Key Problems: for long context and a big batch size, KV cache might dominate memory usage
- Optimization Techniques:
 - Architecture-level:
 - Multi-Query Attention, Group-Query Attention, etc
 - KV Quantization
 - Memory Management:
 - PagedAttention (vLLM), RadixAttention (SGLang), etc
 - Hardware Optimization:
 - Tensor/Sequence Parallelism

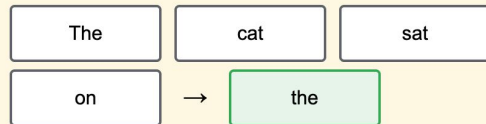
Latency Bottleneck in Decoding

Pass 1 Generate token 1



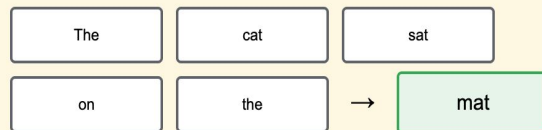
⌚ Must complete before starting Pass 2

Pass 2 Generate token 2



⌚ Must complete before starting Pass 3

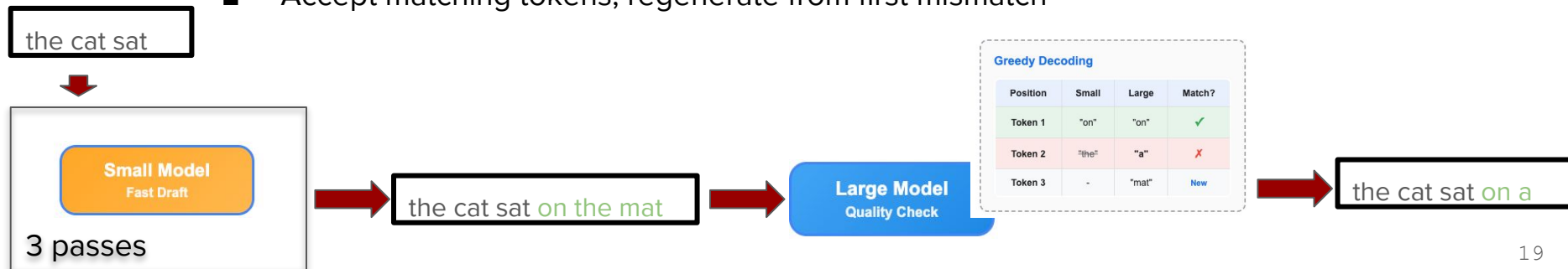
Pass 3 Generate token 3



Three forward pass are required to generate three tokens.

Speculative Decoding: Draft and Verify

- How it works
 - Draft Phase:
 - Small models generate K candidate tokens
 - Verify Phase:
 - Large model process those K tokens in one-forward pass
 - Accept/Reject:
 - Compare draft vs target model probabilities
 - Continue:
 - Accept matching tokens, regenerate from first mismatch



Next Iteration: Input becomes "the cat sat on a -> Small Model generates 3 new candidates

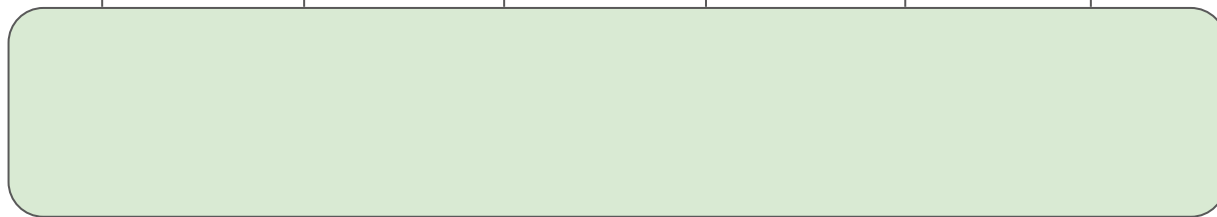
How Large Models Verify in One Pass

Position 0-2 predict tokens we already have (no verification)

Position 3-5
Verification

New
Token

		sat	0.5	an	0.1	the	0.3	cat	0.1	on	0.3
cat	0.4			on	0.7	a	0.5	the	0.2	tee	0.2
				the	0.1	an	0.1	sat	0.3	an	0.6



The cat sat on the mat

Input: Original + candidates

- Verification:

- Position 3: Predicts **on** matches candidates
- Position 4: Predicts **a** doesn't match "the"
- Position 5: Predicts **sat** (not relevant)
- Position 6: Predicts **an** (new token and also irrelevant)

- Result: Verify 3 candidates + generate 1 new token in **one forward pass**

Speculative Decoding: Speed Up Formula

- The Speedup Formula:
 - Sequential (Old Way):
 - K forward pass on Large LM: $K * T_{\text{large}}$
 - Speculative Decoding:
 - K forward pass on small LM: $K * T_{\text{small}}$
 - One forward pass on Large LM: T_{large}
 - With a chance, it would might generate from 1 to K+1 tokens
- When It really speed-up:
 - Draft small model is much faster than target large model ($\gg 5x$)
 - Draft small model has high acceptance rate ($>50\%$)

Inference Optimization

To Improve	Optimize
Max Batch Size	Memory(Quantization, PagedAttention, MQA,etc)
Throughput	Bigger Batch Size, Continuous Batching, etc
Latency	KV Caching, Smaller Batch, Speculative decoding, faster hardware, etc

2. Prompting Techniques

Prompting

- Prompts are the steering wheels to control LLMs to do what we want and also effectively boost their performance.

Example: Sentiment Classification

Software 1.0

```
python @ Casey  
  
def simple_sentiment(review: str) -> str:  
    """Return 'positive' or 'negative' based on a tiny keyword lexicon."""  
    positive = {  
        "good", "great", "excellent", "amazing", "wonderful", "fantastic",  
        "awesome", "loved", "love", "like", "enjoyed", "superb", "delightful"  
    }  
    negative = {  
        "bad", "terrible", "awful", "poor", "boring", "hate", "hated",  
        "dislike", "worst", "dull", "disappointing", "mediocre"  
    }  
  
    score = 0  
    for word in review.lower().split():  
        w = word.strip(",.;:!?") # crude token clean-up  
        if w in positive:  
            score += 1  
        elif w in negative:  
            score -= 1  
  
    return "positive" if score >= 0 else "negative"
```

Software 2.0

10,000 positive examples
10,000 negative examples
encoding (e.g. bag of words)

train binary classifier

parameters

Software 3.0

You are a sentiment classifier. For every review that appears between the tags
<REVIEW> ... </REVIEW>, respond with **exactly one word**, either
POSITIVE or NEGATIVE (all-caps, no punctuation, no extra text).

Example 1
<REVIEW>I absolutely loved this film—the characters were engaging and
the ending was perfect.</REVIEW>
POSITIVE

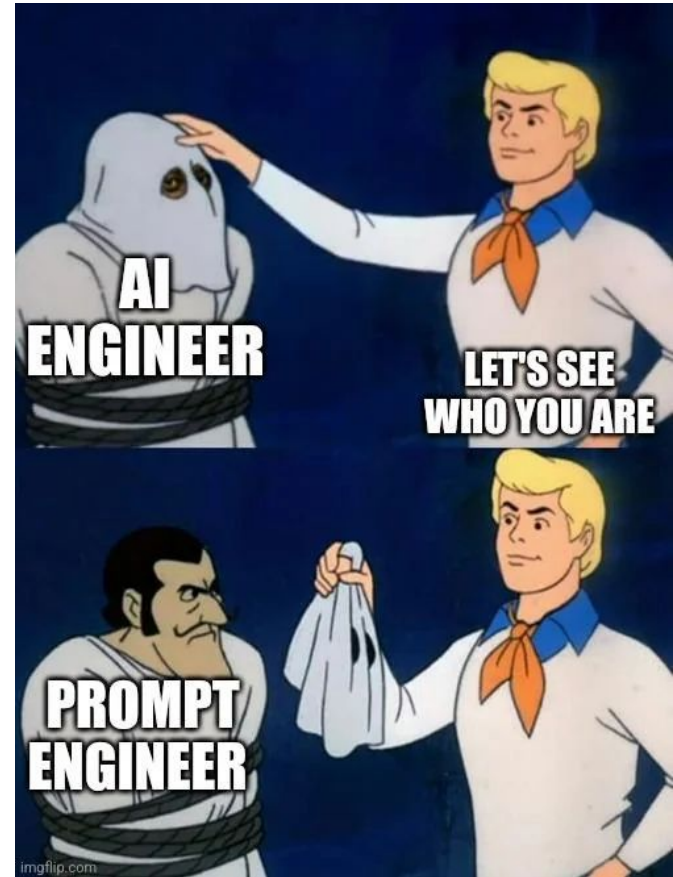
Example 2
<REVIEW>The plot was incoherent and the acting felt forced; I regret
watching it.</REVIEW>
NEGATIVE

Example 3
<REVIEW>An energetic soundtrack and solid visuals almost save it, but
the story drags and the jokes fall flat.</REVIEW>
NEGATIVE

Now classify the next review.

Prompting

Prompting is becoming a must-have skill for AI era



Prompt engineering (Do it first)

- It refers to methods for how to communicate with LLM to steer its behavior for desired outcomes without updating the model
 - More empirical, less scientific
- Advantages
 - Testing and learning early
 - When paired with evaluation it provides your baseline and sets up further optimization
- Disadvantages
 - Introducing new information
 - Reliably replicating a complex style or method
 - Minimizing token usage

Intuition behind Prompt Engineering

- LLMs understand better when you use familiar language and constructs
- LLMs can not know everything. If information is neither in training or in the prompt, they do not know it
- If you look at Prompt and you do not understand it, the prompt can not work

Prompting techniques

- Zero-shot prompting
 - No examples are given in prompt
- Few-shot prompting
 - A few shot examples of tasks are provided
- Chain of Thoughts prompting
 - Examples with the reasoning processes are improved.
 - It could be zero-shot (**Think step by step**) or few-shots
- Self-Consistency Prompting
 - Sample multiple times from output (typically with CoT) and aggregate most common results
- More prompting techniques could be found in this good [survey](#) and this openAI [blog](#).

Prompt Structure

- System prompt:
 - First message provided to LLM (usually hidden towards end user)
 - Provides persona, rules about LLM output, style
 - [Configure](#) your system prompt for Claude

Resources

System Prompts

Copy page

See updates to the core system prompts on [Claude.ai](#) and the Claude iOS and Android apps.

Claude's web interface ([Claude.ai](#)) and mobile apps use a system prompt to provide up-to-date information, such as the current date, to Claude at the start of every conversation. We also use the system prompt to encourage certain behaviors, such as always providing code snippets in Markdown. We periodically update this prompt as we continue to improve Claude's responses. These system prompt updates do not apply to the Anthropic API. Updates between versions are bolded.

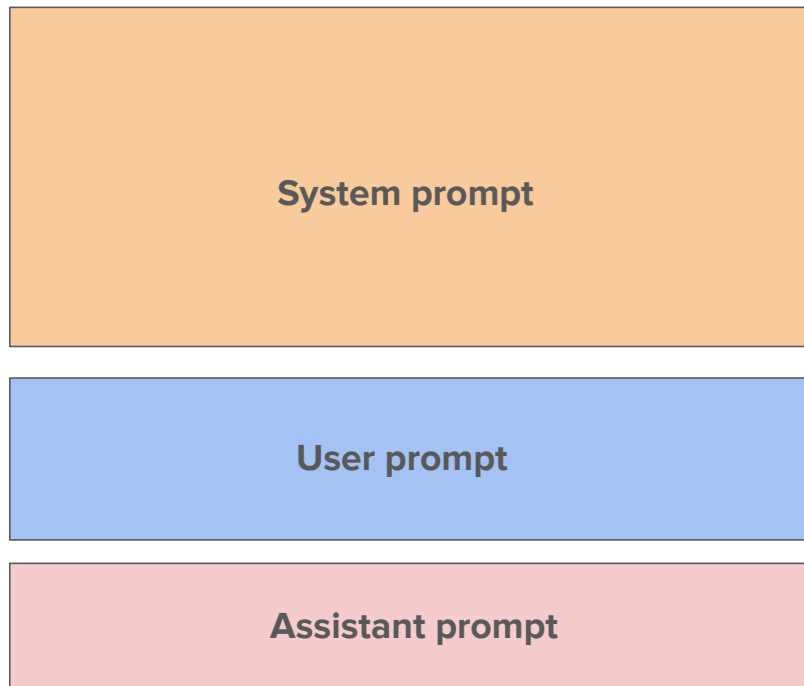
Claude uses a warm tone. Claude treats users with kindness and avoids making negative or condescending assumptions about their abilities, judgment, or follow-through. Claude is still willing to push back on users and be honest, but does so constructively - with kindness, empathy, and the user's best interests in mind. `</tone_and_formatting> <user_wellbeing>` Claude uses accurate medical or psychological information or terminology where relevant.

Claude cares about people's wellbeing and avoids encouraging or facilitating self-destructive behaviors such as addiction, self-harm, disordered or unhealthy approaches to eating or exercise, or highly negative self-talk or self-criticism, and avoids creating content that would support or reinforce self-destructive behavior even if the person requests this. Claude should not suggest techniques that use physical discomfort, pain, or sensory shock as coping strategies for self-harm (e.g. holding ice cubes, snapping rubber bands, cold water exposure), as these reinforce self-destructive behaviors. In ambiguous cases, Claude tries to ensure the person is happy and is approaching things in a healthy way.

source: <https://platform.claude.com/docs/en/release-notes/system-prompts>

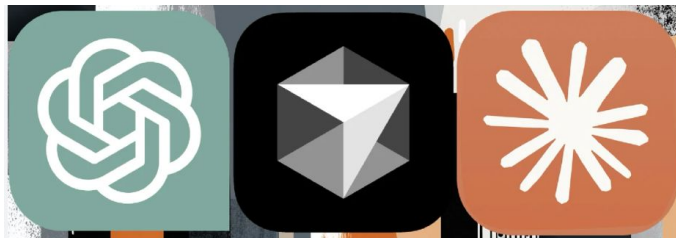
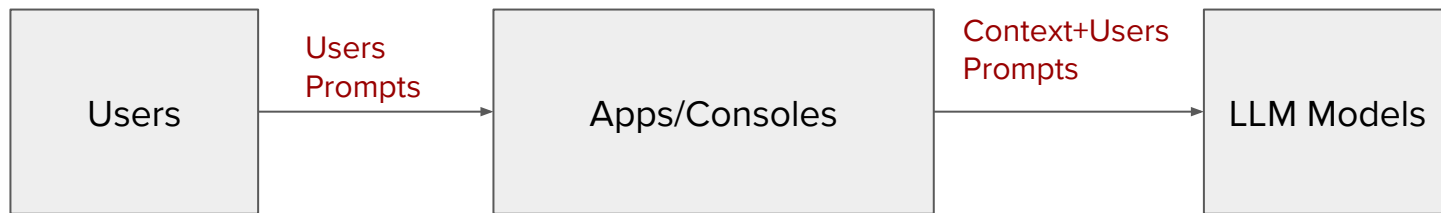
Prompt Structure

- User Prompt:
 - The actual ask or instruction from users
- Assistant Prompt:
 - What the LLM actually generates



What LLM Models actually see?

Claude 3.7 Sonnet <i>64K extended thinking</i>	Claude 3.7 Sonnet <i>No extended thinking</i>	Claude 3.5 Sonnet (new)	OpenAI o1 ¹	OpenAI o3-mini ¹ <i>High</i>	DeepSeek R1 <i>32K extended thinking</i>	Grok 3 Beta <i>Extended thinking</i>
--	---	--------------------------------------	------------------------	---	--	---



How to prepare a good prompt

- Start with:
 - Write clear instructions
 - Split complex tasks into simpler subtasks
 - Give GPTs time to “think”
 - Test changes systematically
- Extend to:
 - Provide reference text
 - Few-shot prompt -> RAG
 - Use external tools
 - Tool-use Prompt
- Prompting Practices from Antropic
 - <https://platform.claude.com/docs/en/build-with-claude/prompt-engineering/overview>
 - Ask LLM to improve your prompts

Best Practices

- Use role prompting to make system prompts more powerful

You are a helpful assistant that loves programming at the level of a senior software developer and is very detailed and pedantic in your answers.



Best Practices

- Prompts should be formatted with structure
 - Multiple components like context, instructions, and examples, XML tags should be used

Example: Generating financial reports

Without XML tags, Claude misunderstands the task and generates a report that doesn't match the required structure or tone. After substitution, there is also a chance that Claude misunderstands where one section (like the Q1 report example) stops and another begins.

Role	No XML Tags	With XML Tags
		<p>You're a financial analyst at AcmeCorp. Generate a Q2 financial report for our investors.</p> <p>AcmeCorp is a B2B SaaS company. Our investors value transparency and actionable insights.</p> <p>Use this data for your report:</p> <pre><data>{{SPREADSHEET_DATA}}</data></pre> <p><instructions></p> <ol style="list-style-type: none">1. Include sections: Revenue Growth, Profit Margins, Cash Flow.2. Highlight strengths and areas for improvement. <p></instructions></p> <p>Make your tone concise and professional. Follow this structure:</p> <pre><formatting_example> {{Q1_REPORT}} </formatting_example></pre>
User	<p>You're a financial analyst at AcmeCorp. Generate a Q2 financial report for our investors. Include sections on Revenue Growth, Profit Margins, and Cash Flow, like with this example from last year: {{Q1_REPORT}}. Use data points from this spreadsheet: {{SPREADSHEET_DATA}}. The report should be extremely concise, to the point, professional, and in list format. It should and highlight both strengths and areas for improvement.</p>	

Source:

<https://platform.claude.com/docs/en/build-with-claude/prompt-engineering/use-xml-tags>

3. Reasoning Models

Source & Acknowledgements

The following pages especially lots of visual guide take reference from the following blog:

[A Visual Guide to Reasoning LLMs](#)

Developing LLM

Stage 1: Building

Building an LLM

Stage 2:
Pre-training

Foundation model

Stage 3: Preference
Tuning (SFT+RLHF)

General-Purpose
Model

Stage 4: Model
Specialization

MultiModal
LLMs

Distilled
Model

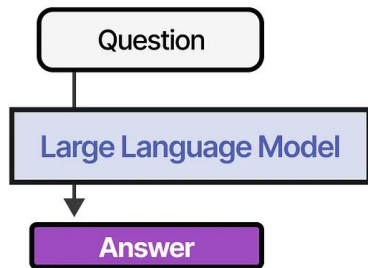
Reasoning
Model

Edge Model

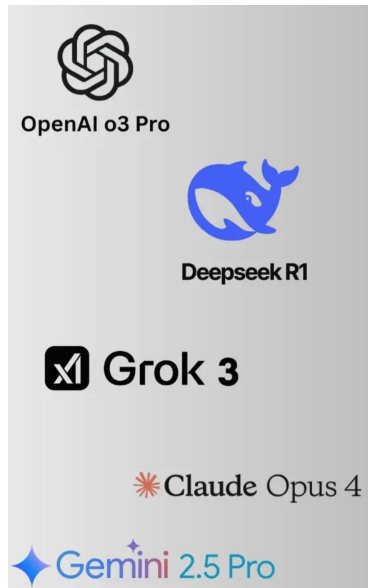
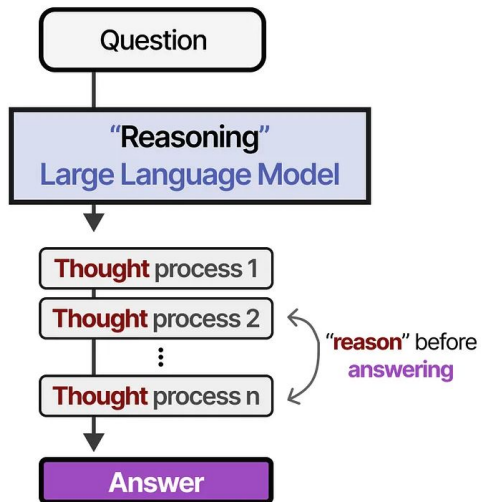
Application
Model

What is a Reasoning Model?

“Regular” LLMs



“Reasoning” LLMs



Popular reasoning models in 2025

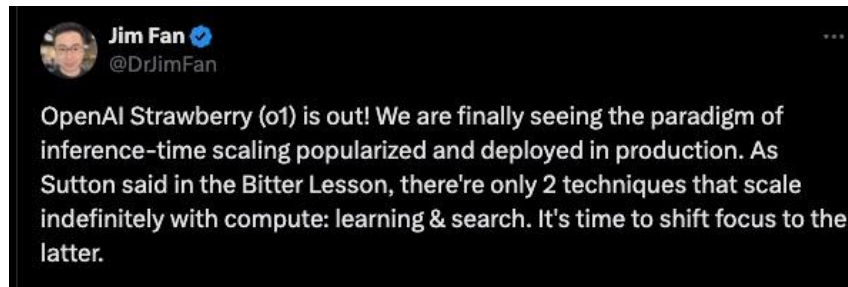
Instead of learning “What” to answer, reasoning models learn “How” to answer through structured intermediate steps called Chain-of-Thought (COT)

A Trivial Example

- Q: What is $3+2$
 - General Purpose LLM: 5
 - Reasoning LLM:
 - `<think>` Adding 3 and 2 gives `</think>` 5
 - `<think>` If $3+1 = 4$, $4+1=5$, the total is `</think>` 5

OpenAI o1

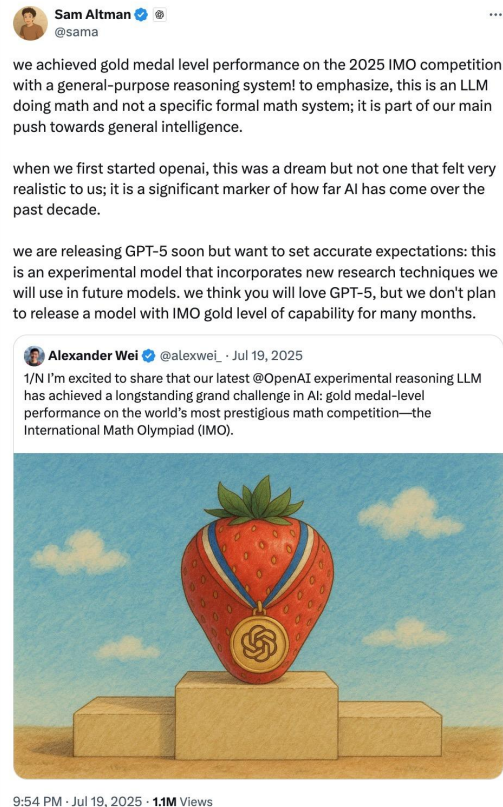
- First Large-scale Reasoning Model
- The [one line](#) from OpenAI:
 - AI models designed to spend more time thinking before they respond



@DrJimFan

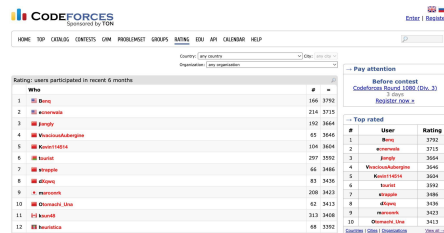
Dominating Math Competition

- OpenAI o1 was the first AI to win IMO Gold Medal
- All reasoning models beat non-reasoning models significantly in AIME 2024
 - Leaderboard is [here](#)



Crushing Coding Competition

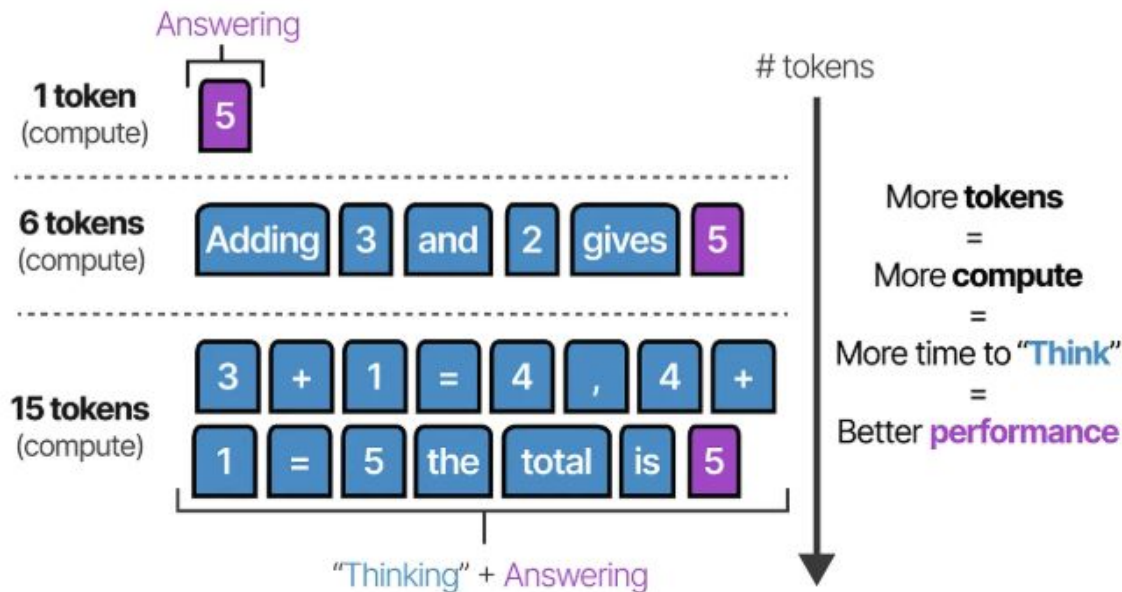
- [Codeforces](#)-based Elo is used to rank competitive programmers
- LLM Models' Performances



Rank	User	Rating	Problems Solved
1	Benq	3792	166
2	antonios	3715	214
3	Benq	3664	162
4	antonios	3564	65
5	antonios	3564	104
6	antonios	3502	297
7	antonios	3446	44
8	antonios	3423	208
9	antonios	3423	42
10	antonios	3413	212
11	antonios	3408	212
12	antonios	3392	68

LLM Models	Elo rating	Beat x% of Coders
DeepSeek-R1	2029	96.3
OpenAI o1	1807	93
o1-mini	1605	86
GPT-4o	668	22

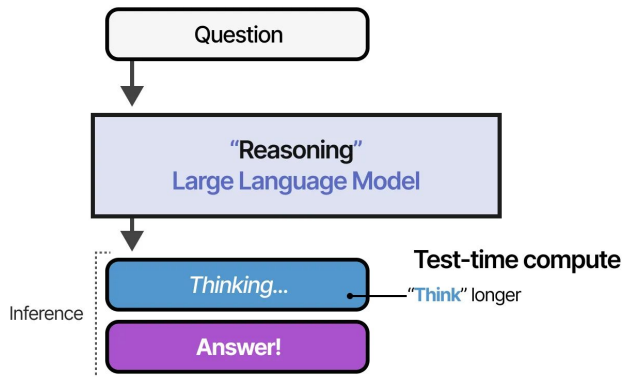
Test-time Compute



Test-time Compute: More compute (e.g., tokens) is spend generating the answer

The Paradigm Shift

- From Train-Time Compute to Test-Time compute



Train-Time Compute

- Model Size: # Parameters
- Dataset Size: # Tokens
- Compute: # FLOPs

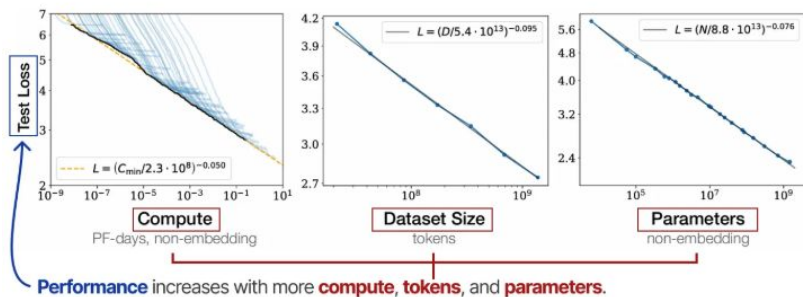
Hitting the wall in 2024

Test-Time Compute

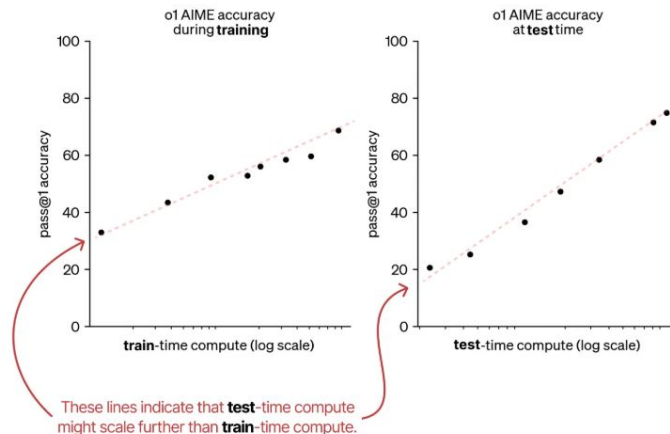
- Reasoning tokens: More thinking steps
- Search/Sampling: Multiple attempts
- Verification: self-checking

Scales with inference budget

The Paradigm Shift



Annotated figure of the “[Scaling laws for neural language models](#)” paper. It shows how performance may increase with different measures of compute (longer training, dataset size, and parameter size).



Annotated figure from “[Learning to reason with LLMs](#)”. The red dotted line was added to demonstrate how OpenAI suggests the new paradigm might be test-time compute.

What is “Reasoning” Really

- It is not human-like thinking - it is intermediate tokens that unlock capabilities
- Why more tokens work?
 - **Hidden working memory:**
 - intermediate tokens provide “working memory” - each step conditions the next prediction, breaking complex problems into simpler sub-problem
 - **Search and path exploration:**
 - instead of just taking the first most likely token, reasoning allows the model to check many more candidates
 - **Self-Correction and “Verifiers”:**
 - Those tokens can act as its own critic. The model can pivot in the middle of process if a mistake is spotted.

Last Letter Concatenating Problems

What is the output when concatenating the last letter of each word in “artificial intelligence”?

No reasoning

The answer is “le”.

Reasoning

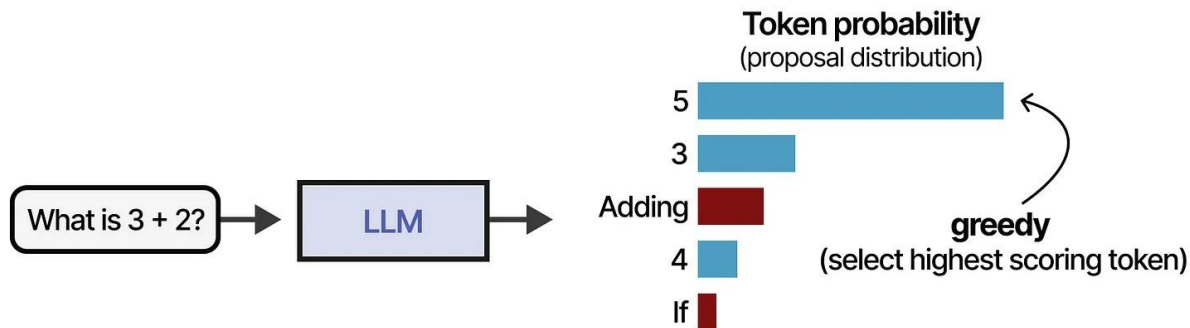


The last letter of “artificial” is “l”. The last letter of “intelligence” is “e”. Concatenating “l” and “e” leads to “le”. So the answer is “le”.

source: <https://dennyzhou.github.io/LLM-Reasoning-Stanford-CS-25.pdf>

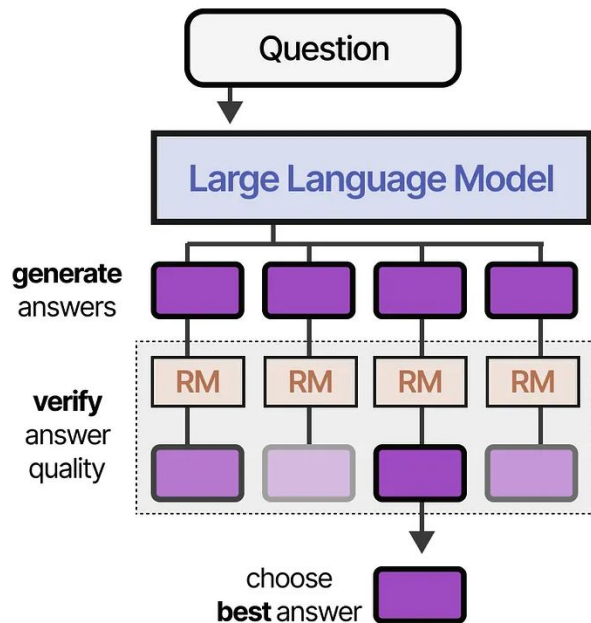
The Model Already Knows

- LLM already has reasoning capabilities embedded in its weights from pre-training
- Challenge: Greedy decoding picks the highest-probability:**5**, skipping **reasoning tokens**
- Solution: Modify the decoding process to favor **reasoning tokens**

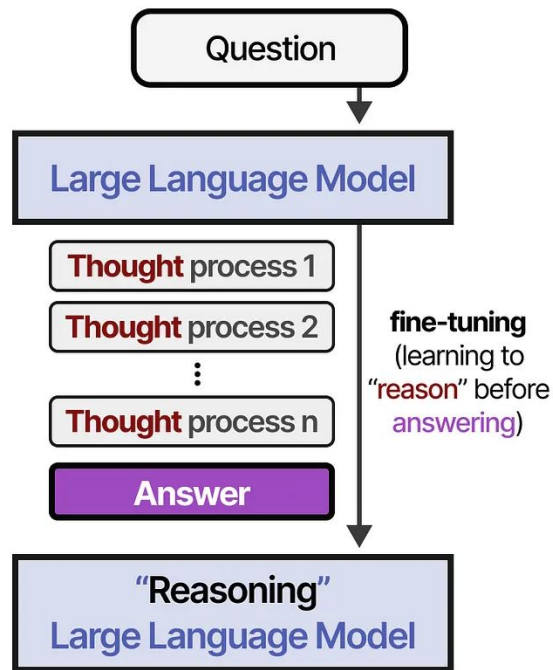


Two Types of Test-Time Compute

Search against Verifiers



Modifying Proposal Distribution



Search Against Verifiers

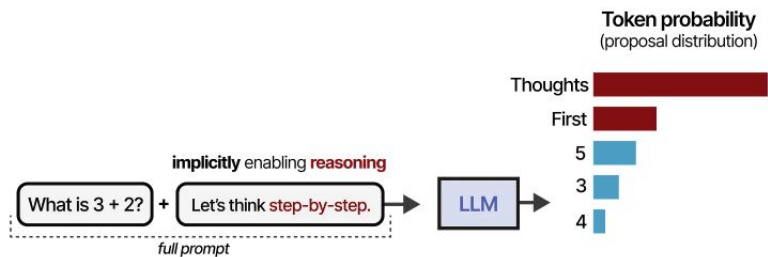
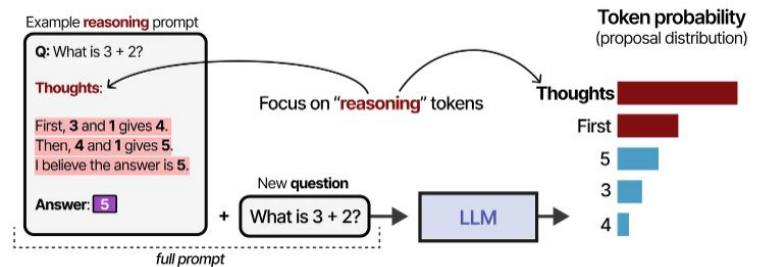
- Output-focused:
 - Generate multiple answers, then score and select the best one
- Two-steps Process:
 - Multiple samples of reasoning processes and answers are created
 - A verifier (reward model) scores the generated output
 - Two types of RM: outcome and process
- Techniques:
 - Majority Voting (self-consistency)
 - Best-of-N Sampling
 - Beam Search with Process RM
 - Monte Carlo Tree Search

Modifying Proposal Distribution

- Input-focused:
 - Tune/train the model to naturally produce better reasoning tokens
- Techniques:
 - Prompting: CoT Prompting
 - Training: SFT on CoT, Self-Taught Reasoner, Reinforcement Learning (RLVR) and etc

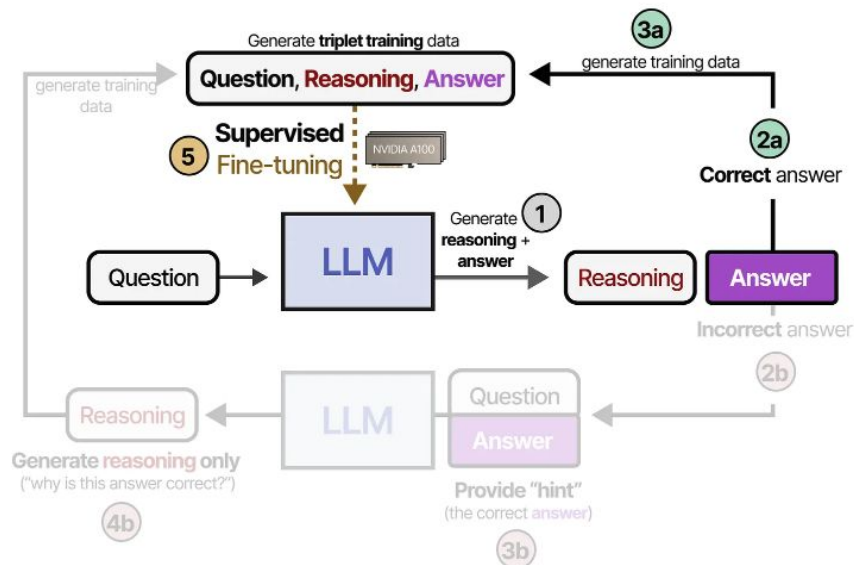
CoT Prompting

- It can be zero-shot or few-shot



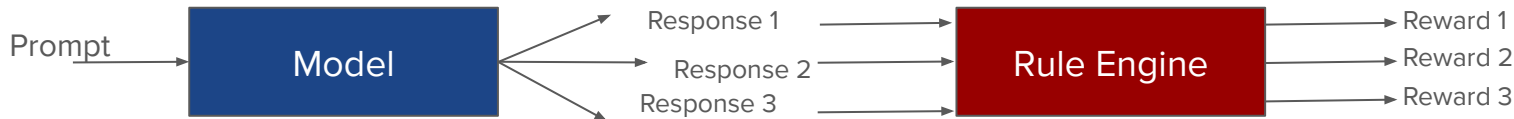
SFT

- The SFT data-triplet training data: question, reasoning, answers
 - Human annotations
 - Synthetic data



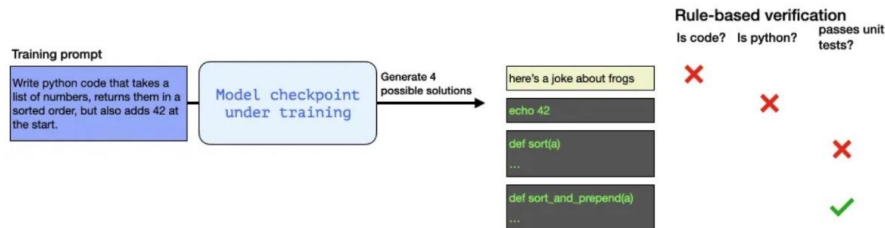
Reinforcement Learning with Verifiable Rewards

- It is used by **DeepSeek-R1**
- The reward is automatically calculated by rule-based systems



Verifiable Rewards

- Accuracy:
 - Math: validate answer
 - Coding: auto-validation
- Format: validate output format
 - Using <thinking> and <answer> tags



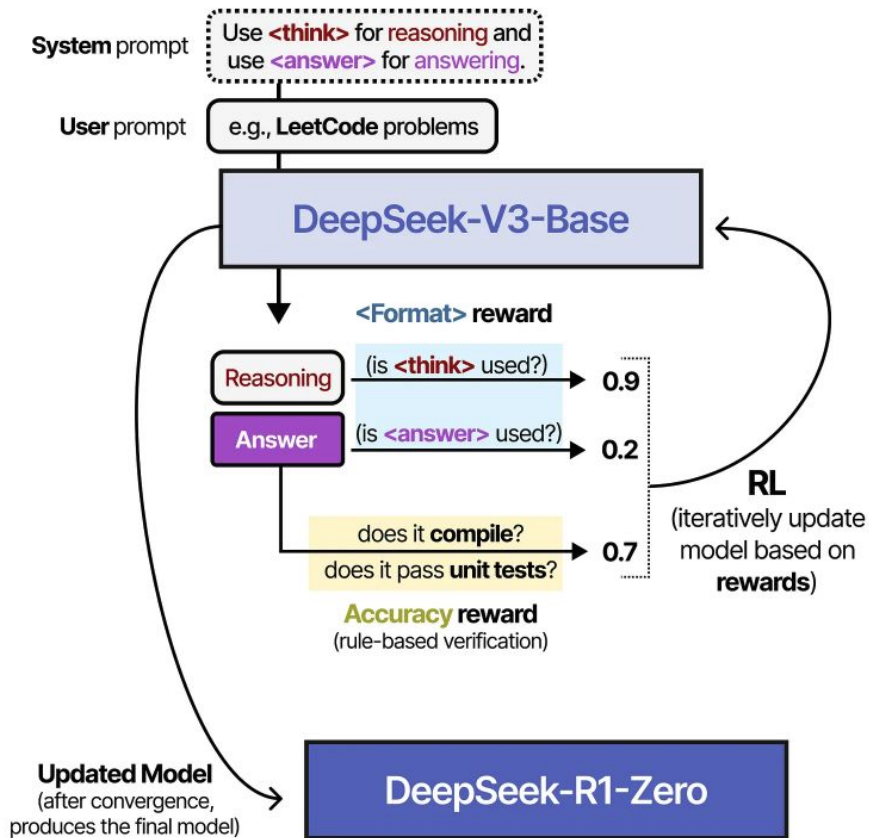
Source: <https://newsletter.languagemodels.co/p/the-illustrated-deepseek-r1>

A conversation between User and Assistant. The user asks a question, and the Assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within <think> </think> and <answer> </answer> tags, respectively, i.e., <think> reasoning process here </think> <answer> answer here </answer>. User: **prompt**. Assistant:

Table 1 | Template for DeepSeek-R1-Zero. **prompt** will be replaced with the specific reasoning question during training.

Source: <https://arxiv.org/pdf/2501.12948>

RLVR for DeepSeek-R1



4 Ways to Build Reasoning Models

Test-time Scaling

No training required. Use prompting (CoT), voting, and search strategies at runtime

SFT+RL

Best approach. Used by DeepSeek-R1 like CoT SFT + RLHF + RLVR.

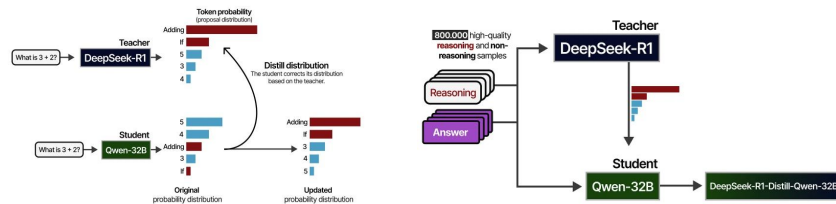
[The illustrated DeepSeek-R1](#)

Test-time Scaling

Train with RL only (no SFT). DeepSeek-R1-Zero showed reasoning can emerge from pure RL

Distillation (Pure SFT)

Fine-tune smaller models on CoT data generated by larger models. Cheap but effective. Like Qwen on R1 outputs



When to Use Reasoning Models

USE FOR

- Complex multi-step math problems
- Advanced coding challenges
- Tasks requiring verification
- Problems with verifiable answers

DON'T USE FOR

- Simple factual questions
- Summarization, translation, creative writing
- Knowledge-based Q&A
- Tasks where speed matters

Next Class: RAG