

Inside DeepSeek- R1

How Reasoning LLMs Are Born



@shivanivirdi





The Problem with Standard Training

trained models to imitate reasoning, not to execute it.

Supervised learning on millions of examples teaches pattern matching.

- The model's ability is bounded by the quality and scope of human data.
- It appears to reason, but fails when problems deviate from its training examples.



@shivanivirdi





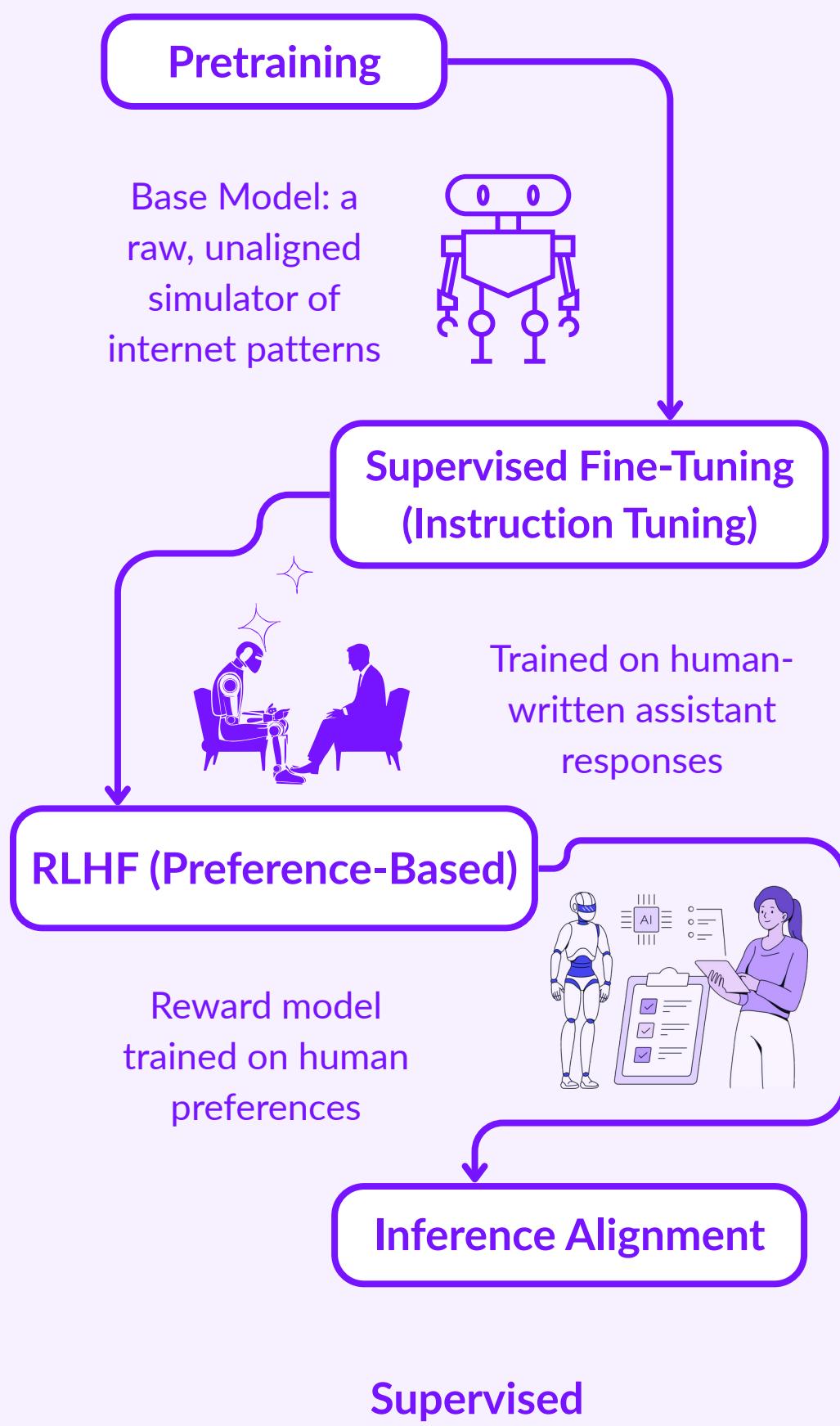
NEOSAGE

DeepSeek-R1

Flipped the Script

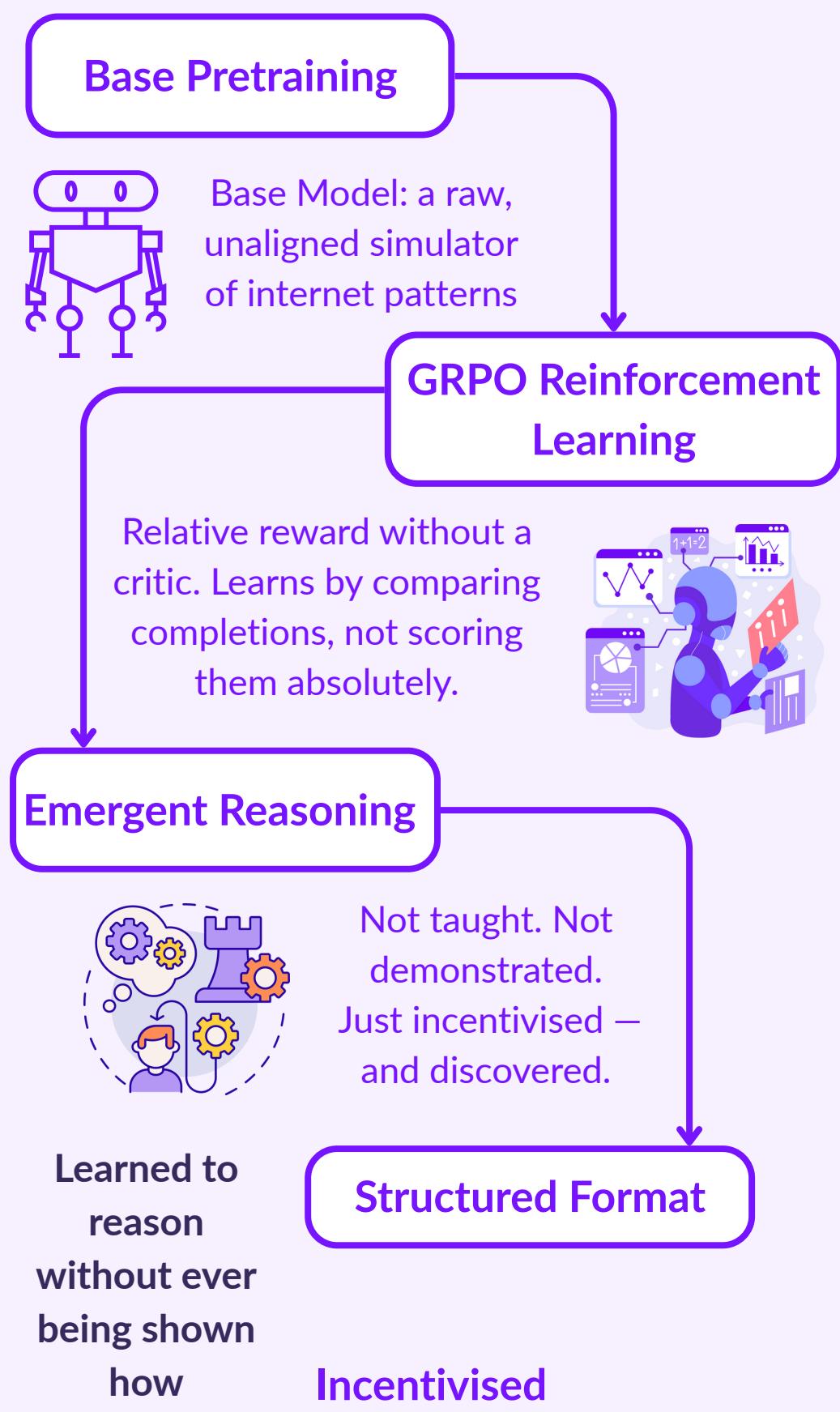
It wasn't taught reasoning with labelled examples.
It was incentivised to discover it via reinforcement.

Typical LLMs



Supervised

DeepSeek-R1-Zero



Incentivised



@shivanivirdi





NEOSAGE

The Engine

Group Relative Policy Optimization (GRPO)

A cheaper, more stable alternative to PPO/RLHF.

- **No Critic Model:** Drops the expensive second network.
- **Group Comparison:** Samples multiple outputs for a prompt and scores them with a programmatic reward (e.g., for accuracy).
- **Advantage Calculation:** Reinforces outputs that score higher than the group average.



@shivanivirdi



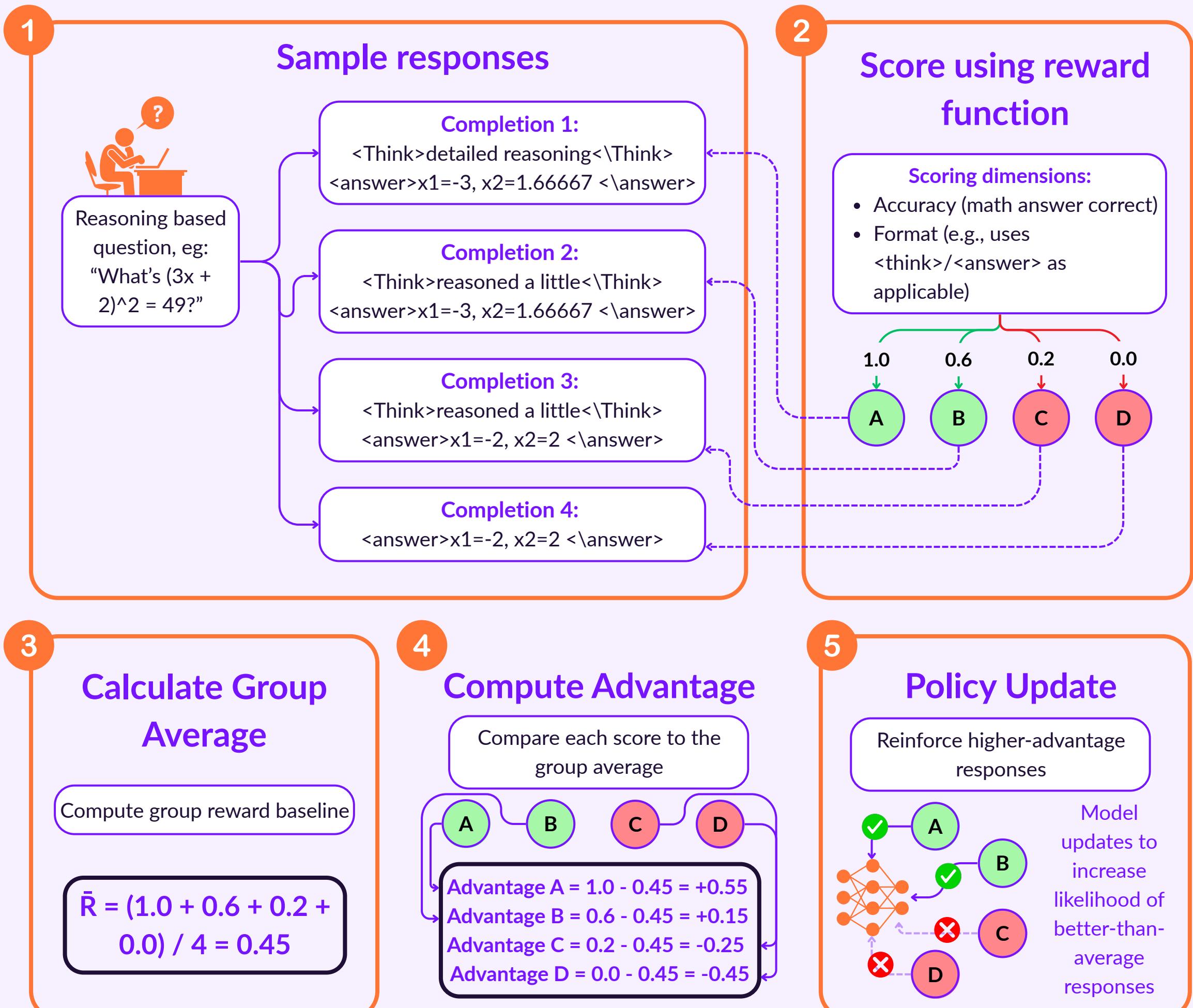


NEOSAGE

The Engine

GRPO: Reinforcement Without a Critic

Trains by comparing completions within a group



@shivanivirdi





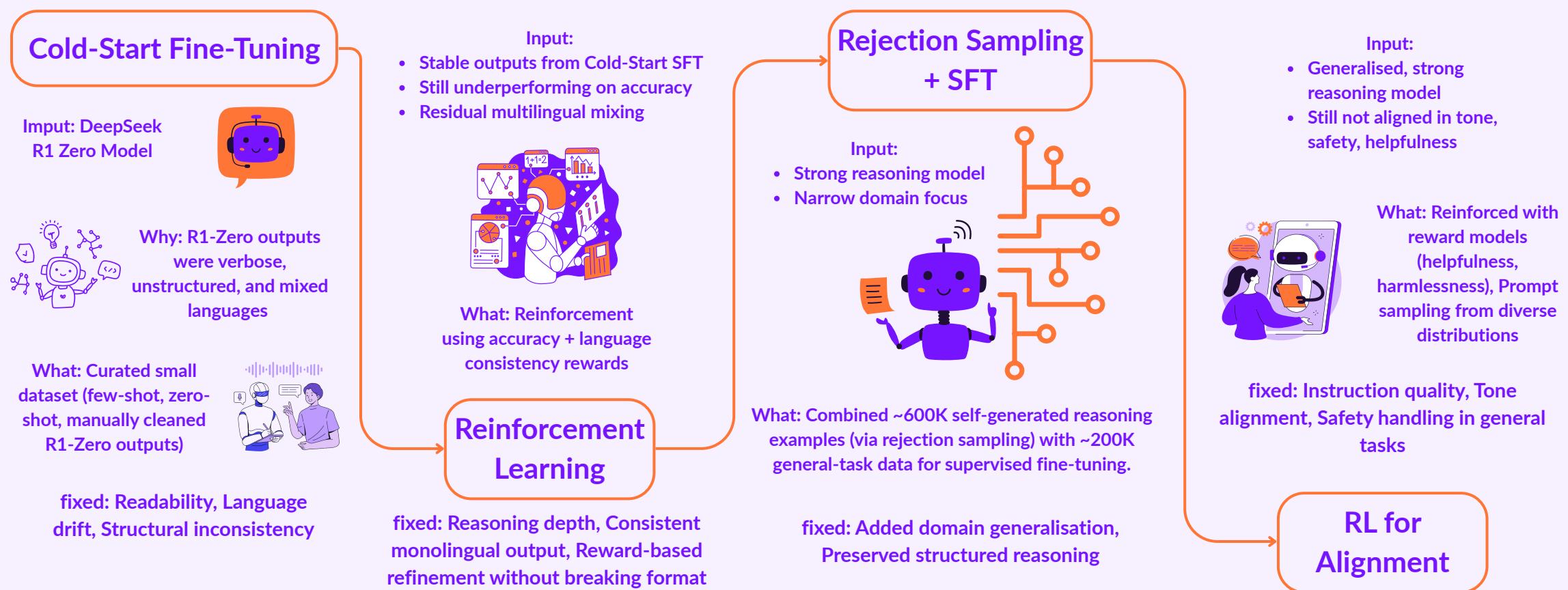
NEOSAGE

From R1-Zero to R1

R1 zero had emergent reasoning, but a raw reasoner isn't a deployable system.

You need a structured, generalized, and aligned model. DeepSeek built it in four sequential stages.

DeepSeek-R1: A Four-Stage System to Sculpt Reasoning into Alignment



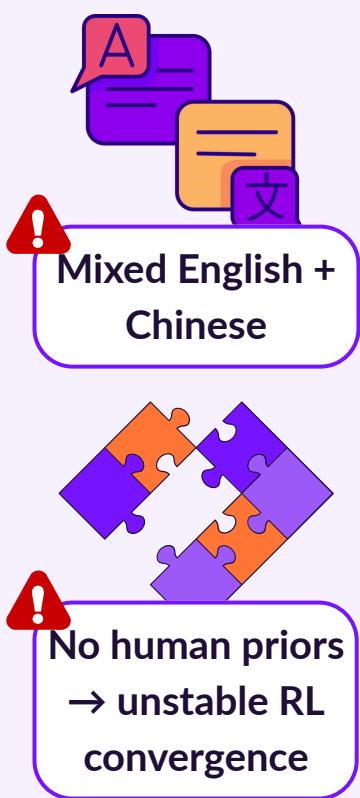
@shivanivirdi





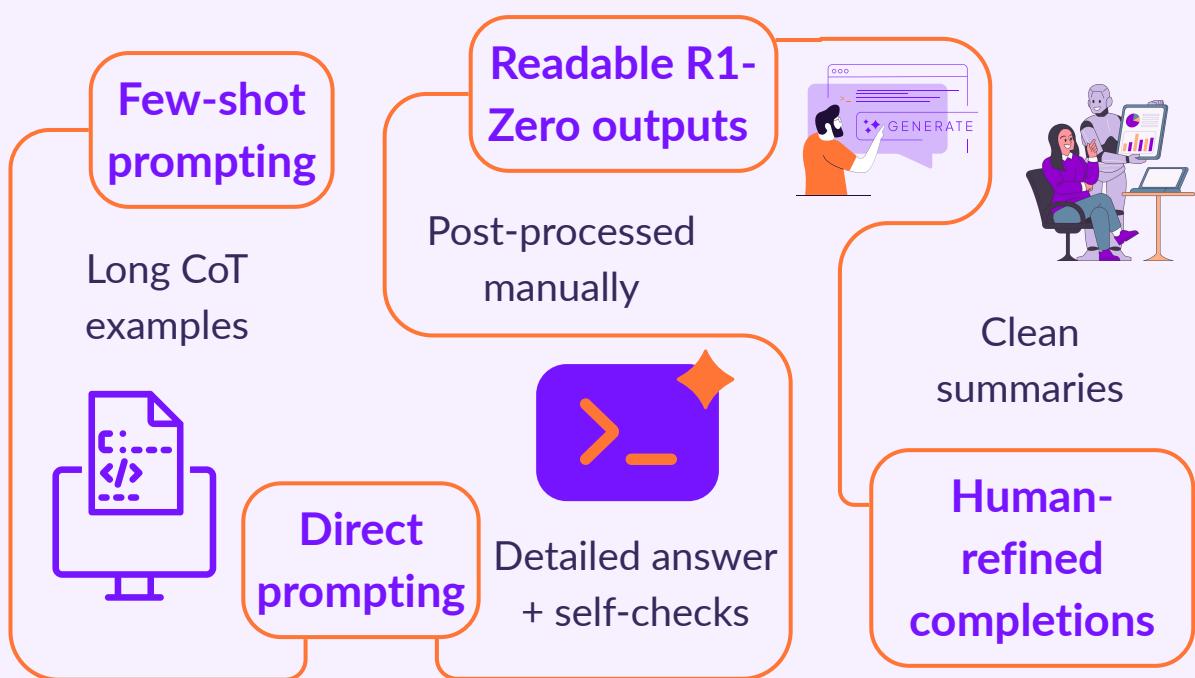
Cold Start: How DeepSeek Engineered a Stable Launch Point for RL

What Was Broken in R1-Zero?



Hard to train from this. Needed a human-prior stabilised base.

How Cold-Start Dataset Was Constructed

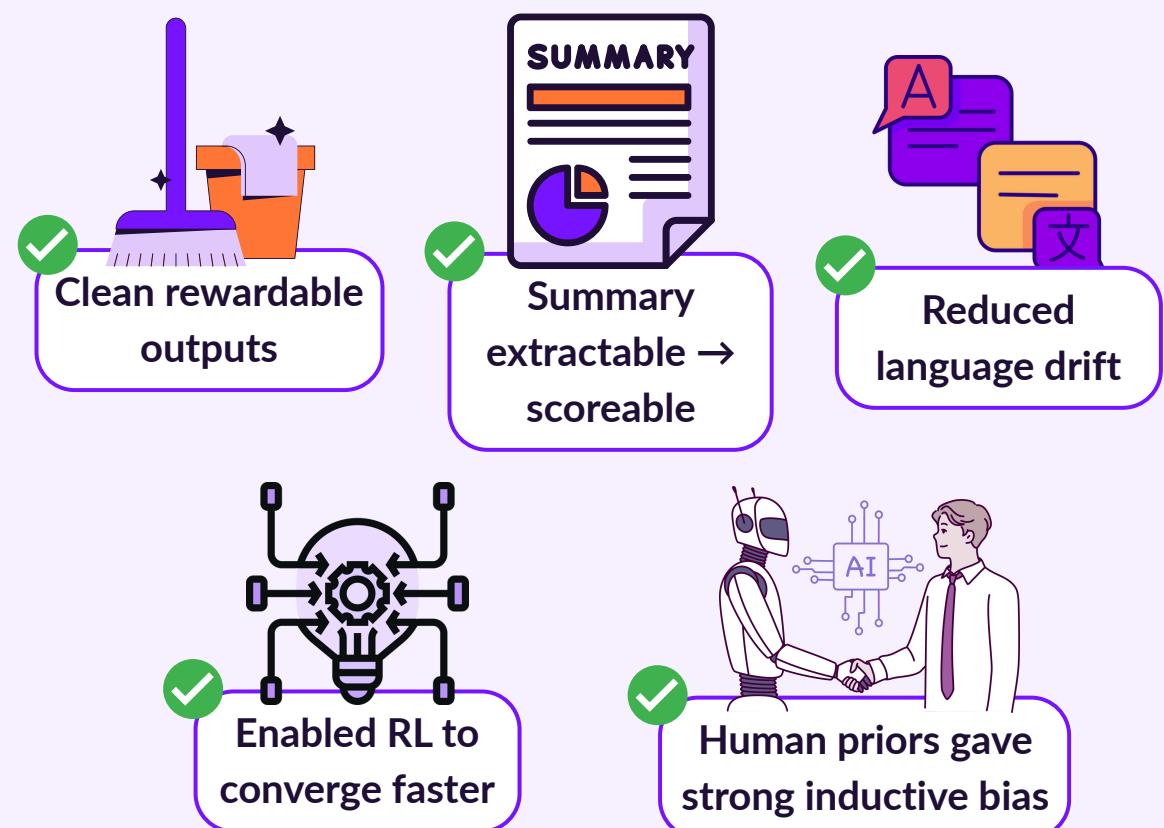


Enforced Output Pattern

Response format
<reasoning_process>detailed CoT style output<\reasoning_process>
<summary>Human readable summary<\summary>

- Reasoning stays structured
- Summary is extractable
- Markdown- and user-friendly

What Cold Start Fixed



@shivanivirdi





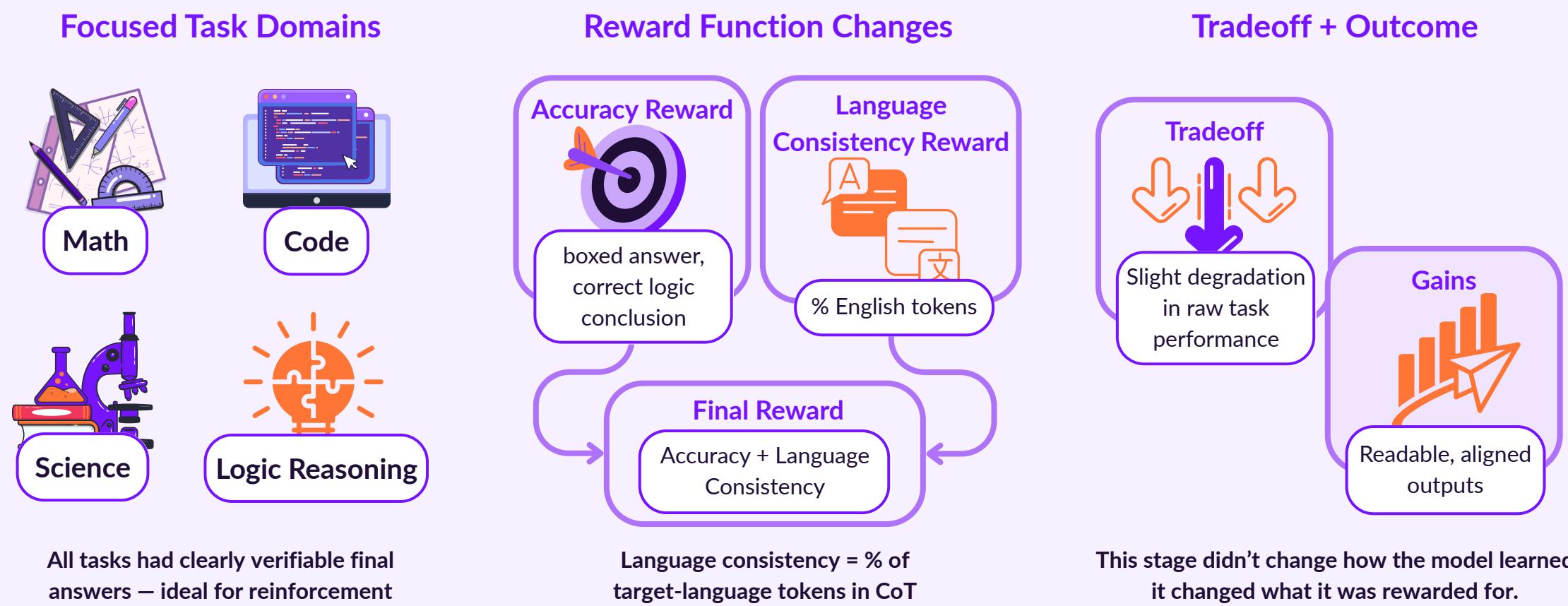
NEOSAGE

Reasoning-Oriented Reinforcement Learning

Problem: The format was stable, but reasoning accuracy could still be improved and language mixing persisted.

Solution: Applied GRPO with a composite reward function, optimizing for both task accuracy and language consistency simultaneously.

Reinforcing Reasoning While Controlling for Language Drift



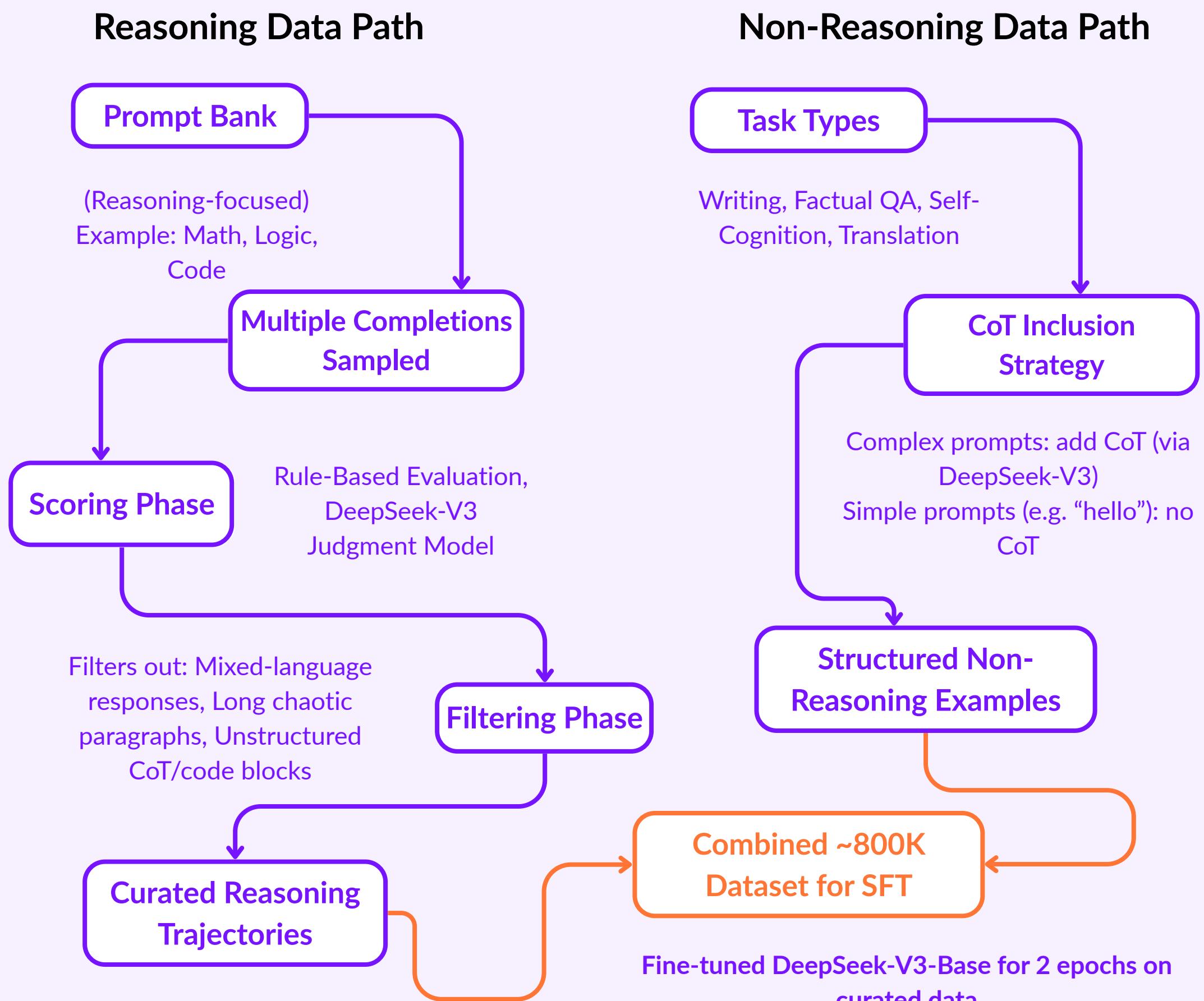
@shivanivirdi





Rejection Sampling + Supervised Fine-Tuning

- The model was a specialist, excelling at STEM.
- It lacked general-purpose abilities like writing, summarization, and chat.



@shivanivirdi

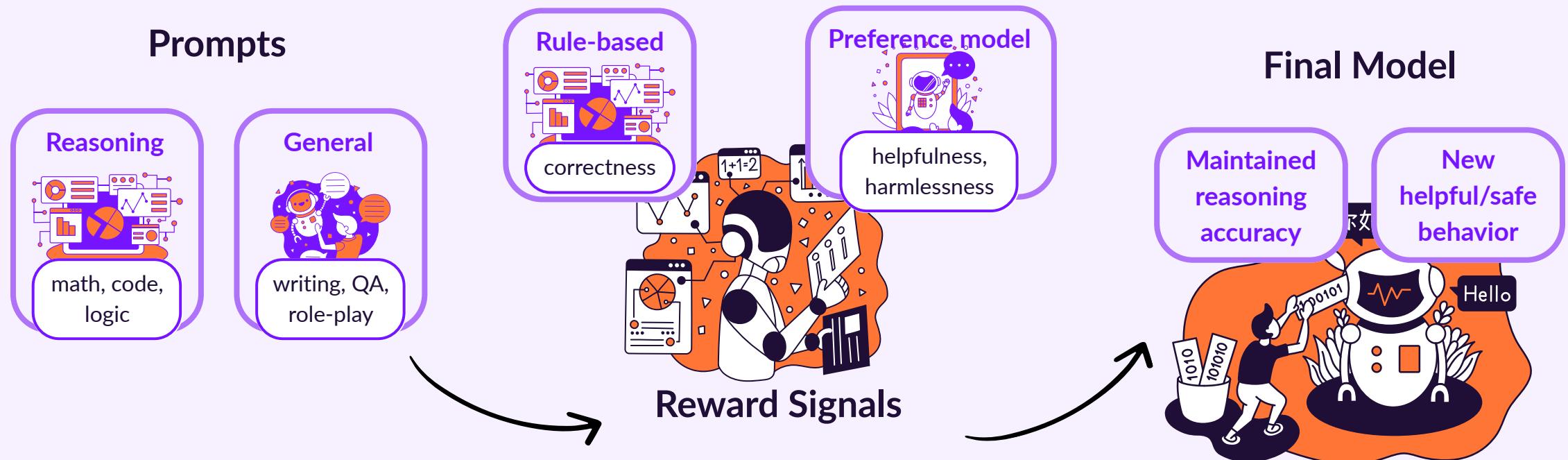




Reinforcement for Alignment and Safety

- **Problem:** The model was capable but not yet aligned on subjective qualities like helpfulness and safety.
- **Solution:** A final RL pass using a hybrid reward system: rule-based rewards for reasoning tasks and preference models for helpfulness/harmlessness on general tasks.

Reinforcement for All Scenarios



@shivanivirdi





NEOSAGE

The Payoff: Distilling Genius

- Used the curated ~800K dataset from Stage 3.
- Fine-tuned smaller, dense models like LLaMA3 & Qwen2.5.
- No RL needed, just standard Supervised Fine-Tuning (SFT).
- Result: The smaller models inherited the emergent reasoning behaviors of the large R1 model, proving the capability is transferable.



@shivanivirdi





Final Note Think Like an Architect

This isn't just about DeepSeek. It's a new paradigm for building advanced AI.

Systematic Refinement: Each stage deliberately fixed a failure from the previous one.

Incentivised Behaviour: Reasoning wasn't labelled in a dataset; the system was designed to make it the optimal strategy.

Distilled Capability: Smaller models don't have to learn from scratch; they can inherit complex behaviours from a well-trained teacher model.

Reasoning is emergent. It cannot be taught through imitation alone

Read the full deep dive:

<https://blog.neosage.io/p/inside-deepseek-r1-a-masterclass>



@shivanivirdi





@shivanivirdi



Liked This?

SAVE

REPOST

FOLLOW

Read full issue on my
newsletter blog.neosage.io



NEOSAGE