# DeepSeek QA

Utrecht University

February 27, 2025
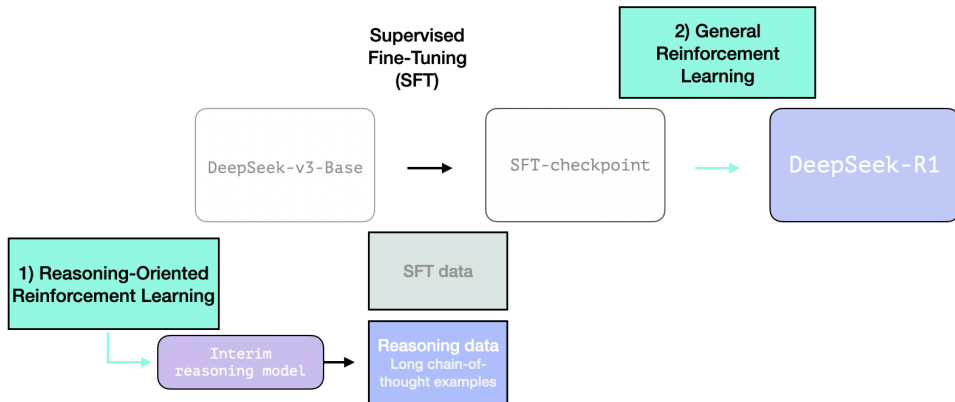
Universiteit Utrecht

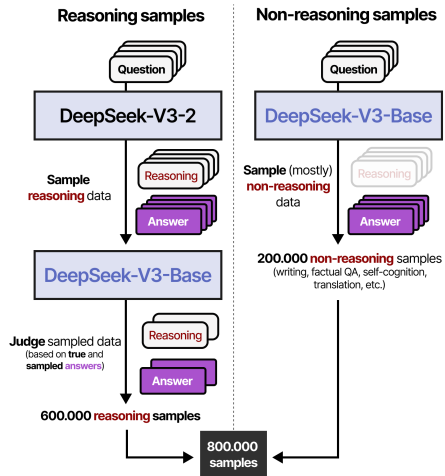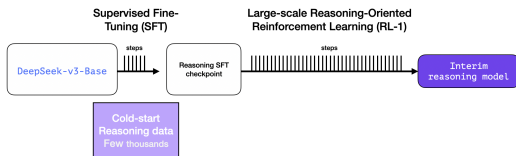## Overview

- Training recipe
- Notable details
- Relevance & Discussion

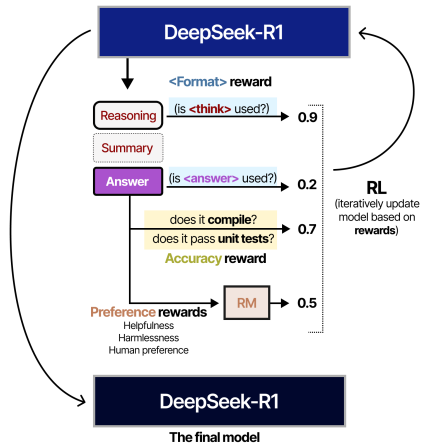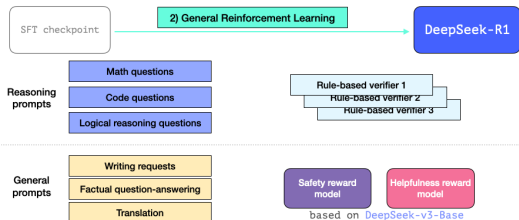Universiteit Utrecht

# DeepSeek Training recipe

# Training: step 1 & 2

# RL process

## Details: Data?

What does the RL data look like?

1. Unreleased by DeepSeek
2. Reproduced by Open-R1 (and openly released).
3. https://huggingface.co/datasets/open-r1/OpenThoughts-114k-math

**Universiteit Utrecht**

# Details: GRPO

1. Avoiding the challenge of learning a value function from a LM backbone, where research hasn't established best practices.
2. Saves memory by not needing to keep another set of model weights in memory.

GRPO does this by simplifying the value estimation and assigning the same value to every token in the episode (i.e. in the completion to a prompt, each token gets assigned the same value rather than discounted rewards in a standard value function) by estimating the advantage or baseline. The estimate is done by collecting multiple completions ($a_i$) and rewards ($r_i$), i.e. a Monte Carlo estimate, from the same initial state / prompt ($s$).

To state this formally, the GRPO objective is very similar to the PPO objective above:

$$J(\theta) = \frac{1}{G} \sum_{i=1}^{G} \left( \min \left( \frac{\pi_\theta(a_i|s)}{\pi_{\theta_{old}}(a_i|s)} A_i, \text{clip} \left( \frac{\pi_\theta(a_i|s)}{\pi_{\theta_{old}}(a_i|s)}, 1 - \varepsilon, 1 + \varepsilon \right) A_i \right) - \beta D_{KL}(\pi_\theta || \pi_{ref}) \right).$$

Note that relative to PPO, the standard implementation of GRPO includes the KL distance in the loss. With the advantage computation for the completion index $i$:

$$A_i = \frac{r_i - \text{mean}(r_1, r_2, \cdots, r_G)}{\text{std}(r_1, r_2, \cdots, r_G)}. \quad (4)$$
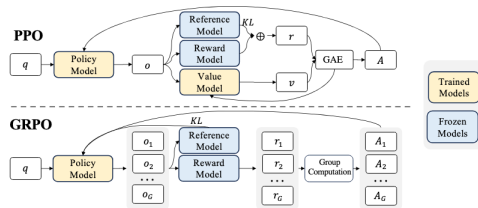


Figure 4 | Demonstration of PPO and our GRPO. GRPO foregoes the value model, instead estimating the baseline from group scores, significantly reducing training resources.

**Universiteit Utrecht**

## Relevance

Can our projects be relevant in this paradigm?

1. Downsampling SFT and reasoning data = data pruning?
2. GRPO: alternative distance functions to KL divergence.
3. Alternatives to coding as reasoning tasks?

**Universiteit Utrecht**

# Discussion

References:

1. https://arxiv.org/abs/2412.19437
2. https://newsletter.languagemodels.co/p/the-illustrated-deepseek-r1
3. https://www.interconnects.ai/p/deepseek-r1-recipe-for-o1
4. https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-reasoning-llms

Universiteit Utrecht