





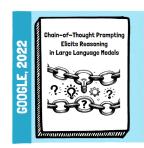




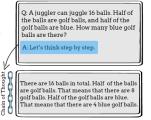
DISTILLATION

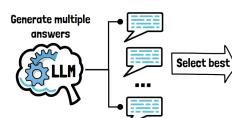














MAJORITY VOTING

Most frequent / consensus answer chosen



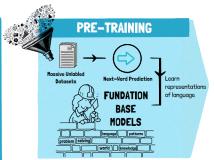
BEAM SEARCH

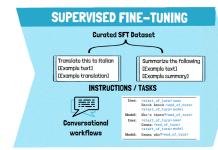
Maintain multiple candidate sequences and select most promising iteratively

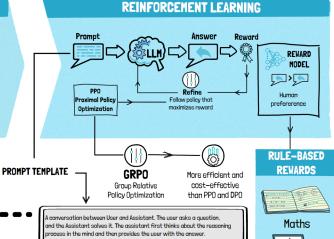


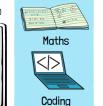
DEEPSEEK R1-ZERO

DISTILLATION





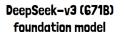














000

Cold-start





ACCURACY Does code complie? Is the maths solution correct?













humans prefer

Predict which responses



R1-Zero Drawbacks Poor readability DEEPSEEK R1 Language mixing "To solve this, we need to 'consider' the '关键因素'

which is the key factor



poorly formatted or

unstructured responses







Cleaned-up CoT examples

generated with

DeepSeek R1-Zero

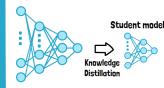
REINFORCEMENT LEARNING

1ST ROUND OF

FINAL ROUND OF REINFORCEMENT LEARNING

Multiple reinforcement strategies to target both reasoning and alignement with human preferences

Teacher model







| DeepSeek-R1-Distill-Llama 8B |
|-------------------------------|
| DeepSeek-R1-Distill-Llama 70B |
| DeepSeek-R1-Distill-Qwen 1.5B |
| DeepSeek-R1-Distill-Qwen 7B |
| DeepSeek-R1-Distill-Qwen 14B |
| DeepSeek-R1-Distill-Owen 32B |

| 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 |
| OpenAl-o1-mini | 63.6 | 80.0 | 90.0 | 60.0 | 53.8 | 1820 |
| OwQ-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-Owen-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 |

| | Based on: |
|---|---------------------------------------|
| - | https://magazine.sebastianraschka.com |
| | /n/understanding-reasoning-llms |



