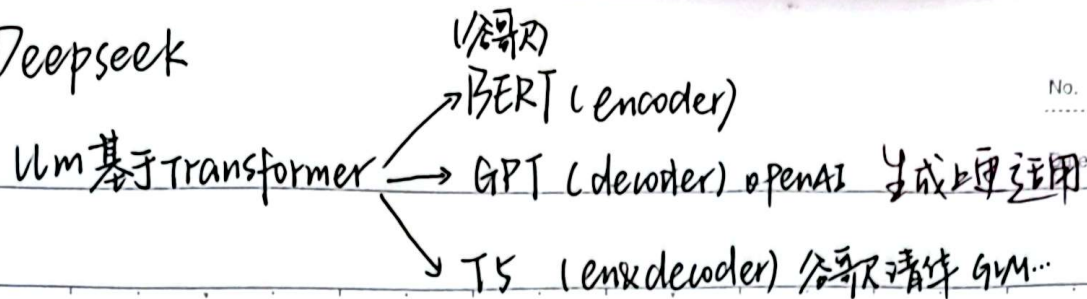


Deepseek



早期囿于“预训练+微调”范式 (pretraining+finetune)

pretraining 缺点 (高) 在非常大的数据集上从头训练模型需要大量计算资源和时间。

1. 数据集要求高
2. 模型参数规模要求高
3. 算力要求高

Deepseek 预训练 V3 模型。

{ MLA (Multi-Head Latent Attention) 降维。

Deepseek MoE 混合专家模型。

★ (R1) 提升推理性能 → post-training 的一种。

后训练机制 {

- SFT supervised fine-tune “小修小改”对下游任务加训
- RLHF (GPT) 引入强化学习
- R1 这不是强化学习更强。
- KD (Knowledge Distillation) 知识蒸馏

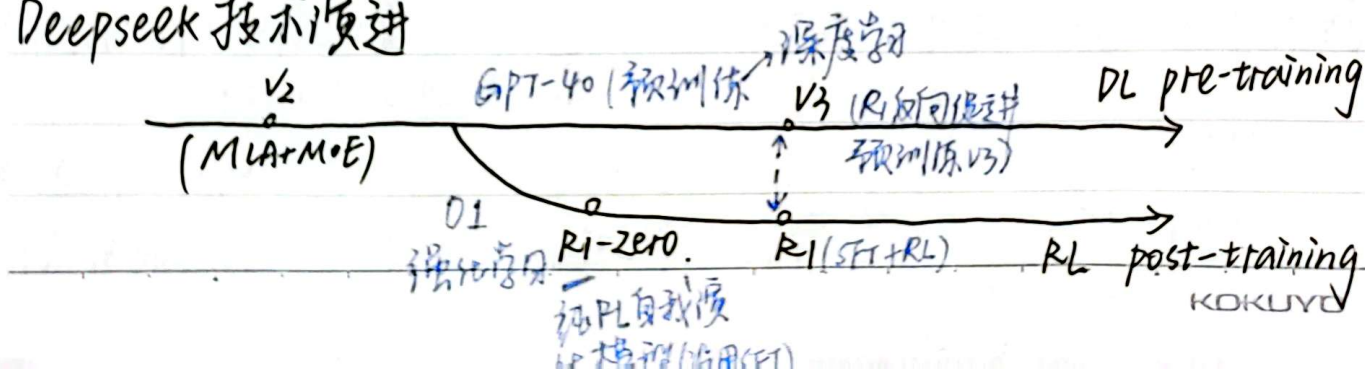
区分 fine-tune / KD

微调的来源是 额外被标注的数据集

而 KD 是用大模型提供的知识, 让小模型拟合大模型。

训练的是小模型 (A new model!)

Deepseek 技术演进





R1-zero.

## 损失函数

损失函数: GPRO (Group Relative Policy Optimization) + 修改损失函数

$$J_{\text{GPRO}}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|q)]$$

$$\mathbb{E}_{q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|q)} \left\{ \frac{1}{G} \sum_{i=1}^G \left( \min \left( \frac{\pi_{\theta}(O_i|q)}{\pi_{\theta_{\text{old}}}(O_i|q)} A_i, \text{clip} \left( \frac{\pi_{\theta}(O_i|q)}{\pi_{\theta_{\text{old}}}(O_i|q)}, 1-\epsilon, 1+\epsilon \right) A_i \right) - \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) \right) \right\}$$

两个随机变量.

 $q = \text{query}$  输入固定查询

$$D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) = \frac{\pi_{\text{ref}}(O_i|q)}{\pi_{\theta}(O_i|q)} - \log \frac{\pi_{\text{ref}}(O_i|q)}{\pi_{\theta}(O_i|q)} - 1$$

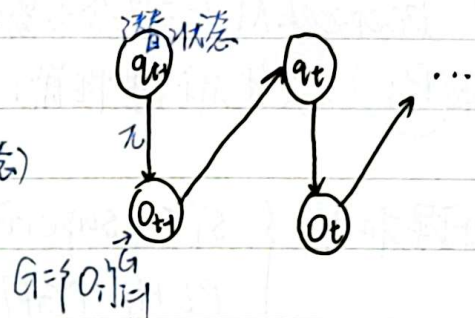
$G$  所有回答构成的样本空间  
 $O_i$  每个生成的回答



$$\pi_{\theta_{\text{old}}}(O_i|q)$$

策略模型  
策略输出 ↑ 输入 (替换状态)

old: 原/前/老策略.



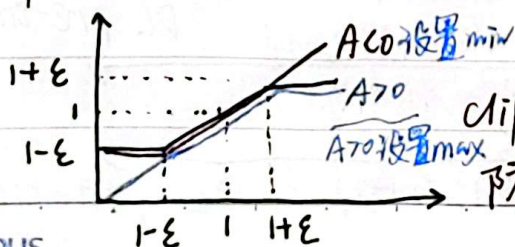
$$\frac{1}{G} \sum_{i=1}^G \left( \min \left( \frac{\pi_{\theta}(O_i|q)}{\pi_{\theta_{\text{old}}}(O_i|q)} A_i, \text{clip} \left( \frac{\pi_{\theta}(O_i|q)}{\pi_{\theta_{\text{old}}}(O_i|q)}, 1-\epsilon, 1+\epsilon \right) A_i \right) - \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) \right)$$

重要性采样  
权重  $w$   
“VBI” 优势 Advantage  
价值函数  $Q-V$  “刹车”  
“教练” 理想值

$\pi_{\theta}(O_i|q)$  以  $q$  生成  $O_i$  的概率 (在现在策略  $\pi_{\theta}$  下)  
 $\pi_{\theta_{\text{old}}}(O_i|q)$  以  $q$  生成  $O_i$  的概率 (在原策略  $\pi_{\theta_{\text{old}}}$  下)  
 $w \geq 1$  倾向于  $A_i$

 $A_i$ : 相当于给  $O_1, O_2, \dots, O_{t+1}, O_t, \dots$  生成序列打分函数

函数中  $\min$ 、 $\text{clip}$ 、 $\text{KL}$  散度均源于 GPRO 近端优化策略算法 (PPO 改进)  
 用  $\text{KL}$  项进行正则化约束



clip 划分区间 → 充当“安全带”/“限速 80/120” 高速公路

防止在优化过程中变化过大, 导致模型崩溃/剧烈波动

Campus

clip: 限制最小值最大值为  $[1-\epsilon, 1+\epsilon]$



损失函数: (未体现方向感  
未体现各种情况下场景)

GPRO 损失函数目的: 防止策略更新过大不稳定导致崩溃。  
稳定策略优化过程。

KL 散度:

标准 KL 散度计算公式:  $D_{KL}(\pi_\theta | \pi_{ref}) = \sum_i \pi_\theta(o_i | q) \log \frac{\pi_\theta(o_i | q)}{\pi_{ref}(o_i | q)}$

变体公式:  $D_{KL}(\pi_\theta | \pi_{ref}) = \frac{\pi_{ref}(o_i | q)}{\pi_\theta(o_i | q)} \left[ \log \frac{\pi_{ref}(o_i | q)}{\pi_\theta(o_i | q)} - 1 \right]$

此处可改进

本质: 实现非线性制约。  
“反向阻尼运动”



奖励函数  $A$  (reward)  
TII 值函数

$A_i = \frac{r_i - \text{mean}\{r_1, r_2, \dots, r_n\}}{\text{std}\{r_1, r_2, \dots, r_n\}}$

完全取决于  $r_i$  (方向)

商业机密未披露, 但仍是基于规则的设计

基于推理链条 CoT 设计奖励函数

e.g. 解答题与题

按步骤聚 两步的结果准确性格式  
得到相应的奖励) 步步给分

奖励汇总

得到最终得分。

数据来源不同  
不同后处理方式: 加权方式不同

SFT 受制于数据质量  
RLHF 受制于打分数据  
RL-zero 纯 RL 规则奖励 × 不标数据

解题思路模板化 “白箱”

区别 RLHF 黑箱深度学习

Deepseek 在坚定走 RL 路线的同时

比较好的解决了降低搜索难度、提升

收敛效率的问题。

“基于思维链的奖励编程解决了 RL 搜索路径问题”



# 放空

P1-Zero 存在的问题:

可读性和语言混合方面存在困难 → 贴近人类数据 (SFT)

① 生成结果符合人类预期, 但人看不懂

② 中英混杂语无伦次

Deepseek-R1: SFT+RL

流程:

-RL  
SFT 1. 冷启动: 用 SFT 优化强化学习的起始点 (可控)

人工“过滤+后处理”收集一批数据, 定义并输出数据格式样例

用 SFT 训练得到能生成清晰思维链的初始策略

-RL  
PL 2. 面向推理的 RL: 准确性奖励 + 语言一致性奖励

针对地给出解题思维链 → 奖励函数

-RL  
SFT 3. 拒绝采样与监督微调

用 RL 训练出的策略生成数据, 人工筛选、过滤、优化 SFT.

-RL  
PL 4. 全场景 RL