

13th November, 2025

Date:

Assumption →

Alignment Problem:

Suppose we have a Pre-trained LLM " w_p ", we fine tune it for a downstream application like solving math problem. math problem are not straight-forward to solve.

Issues: ↴

→ we want our LLMs to learn what actions not to take.

$$\begin{array}{c} n \rightarrow y^+ \\ \searrow \\ \rightarrow y^- \end{array}$$

1) Making Preference data

a) ways of making preference data,

1) Point-wise

$$n \rightarrow \text{LLM} \xrightarrow{\quad} O_1 \boxed{95\%} (u_1, o_1)$$
$$\xrightarrow{\quad} O_2 \boxed{40\%} (u_2, o_2)$$
$$\xrightarrow{\quad} O_3 \boxed{60\%} (u_3, o_3)$$

2) Pair-wise

$$O_1 > O_2, O_1 > O_3, O_2 < O_3$$

3) List-wise

Rank wise Rank 1, 2, 3

⇒ Procedure of Generating Pair wise Data:

$$n \rightarrow \begin{bmatrix} f(n) \xrightarrow{\quad} y^+ \\ \oplus \\ \xrightarrow{\quad} y^- \end{bmatrix}$$

1) take LLM $\stackrel{?}{\sim}$ generate responses

2) Rate Relative performance. (compare them)

$$(u, y_1) \quad (u, y_2)$$

• User select his preference

$$n, y_1, \text{winner} \stackrel{?}{\sim} n, y_2$$

- * in training LLM we try to get correct output
- * we don't teach which other solutions are wrong
- * to be able to follow our preferences.

Date: π is policy \rightarrow LLM

\Rightarrow RLHF :

- Reinforcement learning with Human feedback
- aligning LLM with human feedback

\Rightarrow 2 Models:

Reward & Reinforcement.

\Rightarrow Reward Model:

Quantifying quality of responses

(reward by optimizing responses)

based on human feedback, human

feedback is training label.)

$$n \rightarrow f \rightarrow y \rightarrow RM$$

③ BackProp to update Reward model

④ Reinforcement Learning =

$$\Delta(\alpha) = -[\delta(n, y) - \pi_{KL}(\pi_{RL}(y/n))]$$

$\pi_{KL} \rightarrow$ Non-symmetric

distribution b/w 2

distortion.

$$\pi_{RL}(y/n)$$

h Controls change

\Rightarrow Training a Reward Model. $r_{hf} \leftrightarrow$ human feedback out

on human rated preference data:

$$\textcircled{1} (n, y_w) \quad (n, y_L)$$

② Pairwise Loss:

$$\frac{n, y_w}{n, y_L} \xrightarrow{RM} \delta(n, y_w) \quad \delta(n, y_L) \quad \left. \begin{array}{l} \text{Breadly} \\ \text{Model} \end{array} \right]$$

③ Compute Loss: $\sum_{n=1}^N \text{sigmoid}$

$$L(\alpha) = -E \left[\log \left(\frac{1}{2} (\delta(n, y_w) - \delta(n, y_L)) \right) \right]$$

average reward for winner reward for loser

• contrastive loss

• we want distinction b/w winning & losing so that it always choose winner one.

RL for deepSeq Math - Instruct 78

- ① Initialize policy $\pi_0 \leftarrow \pi_{old}$
- ② Update the weights of Policy
 - ↳ generate response $n \xrightarrow{\pi_0} y$
 - ↳ Assign Quality Score. $y \rightarrow RM$
 - ↳ update Policy weights: use generated responses \hat{y} along with its associated estimated $\delta(n, \hat{y})$ to update π_0 using PPO Loss.

$$L^{DPO} = -E[\log(\pi_0(\hat{y}_n/n)) - \beta(\pi_0(\hat{y}_n/n) - \pi_{old}(\hat{y}_n/n))]$$

- here we are calculating ratio.
- β controls divergence
 - less divergence
 - advantage π^*
 - value model.

→ Short Coming of Reward Model:

- 1) Too many steps
- 2) Imperfect proxies \leftarrow avg reward model
- 3) Many hyperparameters to update
 - use Sampling model [Monte Carlo tree search] Thinking model.
 - another is DPO

→ Reward • how correct

→ Value • go to

→ Reference, reward-based line

$n \rightarrow w_p \rightarrow GPO \rightarrow$

Direct Preference Optimization

↳ Steps:

- 1) $\pi_0 \leftarrow w_p$ ($LM \rightarrow$ policy) $\rightarrow w_p$
- 2) for a given input prompt " n " Sample observation from π_{old} . $\rightarrow RM \{o_i\}$
- 3) Rate responses such that we know which one should be the winning & losing observation \rightarrow Reward model \rightarrow PPO / DPO, GPO to backprop
- 4) Minimize DPO loss.

Summary:

\rightarrow pref. data

\rightarrow RM $\{o_i\}$

GPO to backprop

17th November, 2025

Date: Pg. 11

DeepSeek

Reasoning model:

get good, clean, reliable, relevant task specific data

→ Math seed, math examples

→ it is easy to search

domain wise in common crawl

- 2) Validate by using smaller (1-3B) codebase LLM.

→ Near Deduplication, Simple Duplication

→ On 4th epoch, 3rd iteration they found 98% unique math related paper.] scraping better approach

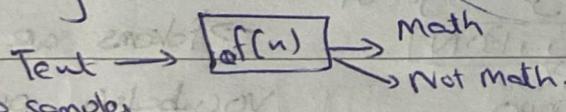
- How to get more math related data:

- 1) Go to commoncrawl

- 2) Seed Math Corpus (already established math datasets)

such as OpenWebMath & others

- they created a model



→ they got 40B HTML pages
→ out of which 120B tokens
→ they used small model just to verify it.

⇒ Supervised fine tuning & Chain of thought.

→ main reason for deepseek performance is deepseek's dataset & SFT + chain of thought.

- model was created using fasttext to classify webpages as math or non-math.

- they created confusion matrix.

20th November, 2025

Tech talk Date: 2025

SAT Benchmark
Score = 84.4%

DeepSeek - Math 7-B Base

Contributions

New Math Benchmarkship dataset

- New State of the Art (SOTA) on MATH

- CnRPO (Thinking model)

Better performance than 10 times Bigger Models on English & Chinese benchmarks.

- Math } logically multistep
- Code } solutions

- code is verified using Test case verification.

- for math we can verify using solution vs label

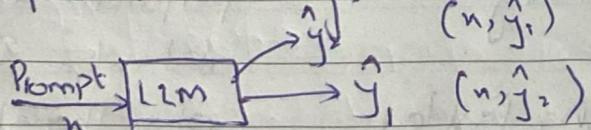
=> Thinking Models: IPPO, BON chain of thought

- Preference Tuning DPO GRPO

↳ with chain of thought

→ Proximal Policy Optimization:

- preference training.



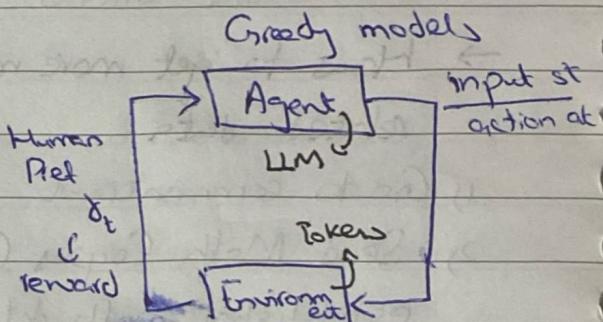
- Change the temp of LLM.
- with new temperature, a new output was calculated.
- we only made data for no, no label (preferred ones) right, wrong

→ For labelling:

- Human Pref
- LLM as judge
- Test case verification.

→ RLHF (Reinforcement

Learning with human feedback



No (at 1st)

→ tokens are set of

vocabulary, out of which LLM will generate based on rule (policy).

→ action is choosing token.

* In RL we want

to update weight of LLM based on human preference.

- slight update to meet human preference.

BLHF:

- 2 components: Reward model, ^①R_L

→ Objective of Reward model:

distinguish good from bad

- RM → Input (Prompt, HF)

- Output → $\gamma(n, y)$ a number 0-1

Issue →

we cannot use these loss functions.

Cross entropy, RMSE, MSE -- etc.

- they only recognise the winner solution, they cannot distinguish b/w slight differences

Breit Bradley Terry formulation:

- Prob that y_i is better than y_j is given as:

gives clear distinction

② Reinforcement Learning:

- The idea of RL step is to use RM to promote good answer.

$$P(y_i > y_j) = \frac{e^{\delta_i}}{e^{\delta_i} + e^{\delta_j}}$$

$$= \sigma(\delta_i - \delta_j)$$

- Sigmoid gives value between 0 & 1.

→ Objective:

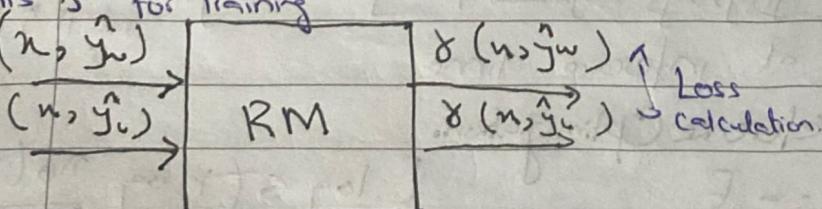
$$L(a) = \text{maximize Reward} +$$

don't deviate too much from

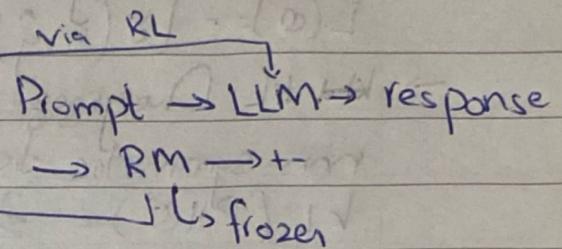
the base / Ref / old model.

- Reward hacking: too much deviation → learning instability

This is for training



$$\delta(a) = -E[\log(\sigma(\gamma(n, \hat{y}_w) - \gamma(n, \hat{y}_b)))]$$



Proximal Policy Optimization:

β trainable controlling parameter

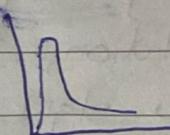
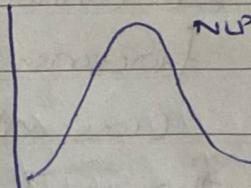
$$L(\alpha) = - [\delta(a_t, \hat{y}) - \beta KL(\pi_\theta(g/a) || \pi_{\theta_{ref}}(g/a))]$$

→ KL Divergence

soft of non-symmetric distribution b/w 2 dist

$$KL(P||Q) = \sum_{i=1}^n P_i \log \frac{P_i}{Q_i} \rightarrow P_i(P_i - Q_i)$$

- We need to keep check b/w change in older & current model by finding distance b/w 2 distributions.



- Issues with PPOs.

- We are calculating reward for each prompt.

Advantage \equiv reward - baseline.

$$\text{PPO-KL} = \underbrace{\mathbb{E}_{\pi_\theta(a_t | s_t)} \left[\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{ref}}(a_t | s_t)} A_t \right]}_{\text{Advantage}} - \beta \underbrace{KL(\pi_{\theta_{ref}}(\cdot | s_t), \pi_\theta(\cdot | s_t))}_{\text{KL Divergence}}$$

- 4 models → issue.

Solution \Rightarrow

Best of N.

$$n \rightarrow \boxed{\text{SFT}} \rightarrow y_1, y_2, y_3, \dots, y_n \rightarrow \boxed{\text{IRM}} \rightarrow k \underset{i \in [1, n]}{\operatorname{argmax}} \delta(n, y_i)$$

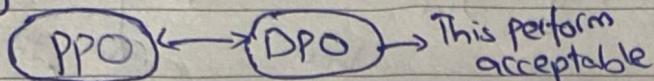
- Super complex at inference time

- DPO: Direct Preference Optimization.

$$L(\pi_\theta; \pi_{\theta_{ref}}) = - \mathbb{E}_{(y_1, y_2, y_3) \sim D} \left[\log \sigma \left(\frac{\pi_\theta(y_1 | x_1)}{\pi_{\theta_{ref}}(y_1 | x_1)} \right) - \beta \frac{\pi_\theta(y_2 | x_2)}{\pi_{\theta_{ref}}(y_2 | x_2)} \right]$$

Date: ..

usually perform better



computationally expensive

- distribution shift due to RM less examples.

\Rightarrow DeepSeek R1 \rightarrow

they changed responses
with thinking steps \rightarrow
solution

Order:

- Preference Training R^{PF}
- Proximal Policy Optimization
- Bag of N $\xrightarrow{\text{factor-critic algo}}$
- Direct preference Optimization
- DeepSeek Math.

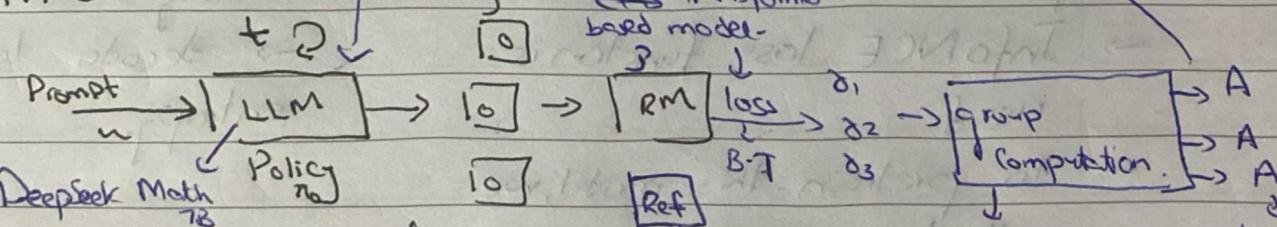
Group Relative Policy Opt

KL divergence

$$L(Q) = \text{maximize}_{\text{reward}} + \text{don't deviate too much from base model}$$

1st Dec. 2025:

GRPO



- * RM: Quantification of the input observation ($\alpha-1$)

any scoring adapts
KL divergence with actual LLM (with advantage)
criteria \rightarrow average, map to low

- * if we have multiple output for the same source:

$$n \rightarrow f \rightarrow \begin{matrix} o_1 \\ o_2 \\ o_3 \end{matrix} \rightarrow P(o_1 > o_2) \text{ BT}$$

$$\ell(\alpha) = G(\delta_1 - \delta_2)$$

- * classification only 0 or 1
cat or dog
- * Quantification

$$L(\theta) = E \left[\underset{\text{prompt}}{\underset{\text{GRPO}}{\sum_i}} P(\theta), \underset{\text{out}}{\underset{\text{Data}}{\sum_i}} \in \pi_\theta(O|q) \right]$$

$\hat{A} = \frac{1}{G} \sum_{i=1}^G \frac{1}{\text{avg size of prompt output}} \sum_{i=1}^{\text{avg size of prompt output}} \left\{ \min_{\text{min}} \left\{ \pi_\theta(O|q) \right\}, \text{clip} \left(\frac{\pi_\theta(O|q)}{\pi_{\theta_{\text{old}}}(O|q)}, 1 - \varepsilon, 1 + \varepsilon \right) \right\}$

or this
 or this
 select among min, the changes

$$\left. \hat{A} \right) - \beta \text{KL} \left(\pi_\theta / \pi_{\theta_{\text{old}}} \right)$$

→ Thinking models =

→ 3rd December, 2026:

ViT, CNN

Encoder/Decoder base LLM

→ Contrastive Learning:

- Lear positive & negative pairs

→ InfoNCE loss function is used.

→ to pretrain this model, we will have input pair like this (text, img)

→ CLIP: dual encoder model.
flaence., SAM

Meme related
→ CLIP 2

VLM
vision Language models.

→ MFCC → Cepstral coefficients

→ Vision Language Model:

images with text

→ decoder loss is cross entropy.

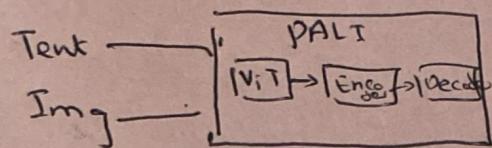
→ WAN - Animate.

→ Generative Adversarial Networks

→ Variational autoencoders

→ Diffusion networks.

⇒ PALI:



⇒ mixture of object.

- Scaling 2 loss functions

⇒ VQAS

- Visual, Questions & Answers
Image text output