Background

Paper Discussion

Conclusion O

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Reading is fun

"Yue, Y., Chen, Z., Lu, R., Zhao, A., Wang, Z., Song, S., & Huang, G. (2025). Does reinforcement learning really incentivize reasoning capacity in IIms beyond the base model?"

Chong Cher

July 16, 2025

Background

Paper Discussion

Conclusion 0

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで



Ground Rules SILE Office Access Administrative Stuff

Background

Paper Discussion What is the problem?

Conclusion



Paper Discussion

Conclusion 0

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

SILE Office Access

Thanks to SILE for permission to use the premises for the reading group.

All participants of the reading group are to enter from the main door, and are only allowed to use the designated meeting room.

Please do not access the main office.



Paper Discussion

Conclusion 0

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Havelock 2 Washroom

The washroom is located outside of the office, on the opposite end of the lift lobby.

Please press the doorbell once you are at the main door; someone in the meeting room will open the door for you.

Background

Paper Discussion

Conclusion O

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Chatham House Rule

The reading sessions will be held under Chatham House Rule —

When a meeting, or part thereof, is held under the Chatham House Rule, participants are free to use the information received, but neither the identity nor the affiliation of the speaker(s), nor that of any other participant, may be revealed.

Feel free to share ideas openly, but please respect confidentiality.

Paper Discussion

Conclusion O

Reading Group Expectations

This is a community-led reading group broadly interested in Machine Learning / Artificial Intelligence topics. Some simple guidelines:

- Treat others as you would want to be treated
- Offensive language is prohibited
- Read (or skim) the paper; the discussion is much more productive if we have a common baseline for discussion
- Please volunteer to present; this could be a paper you find interesting, or even a project you would like feedback on or to discuss.

Paper Discussion

Conclusion 0

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

How to Contribute

Please feel free to use the Telegram group chat ("Discussion is fun") to discuss the presented paper, or any interesting AI/ML news you may have come across.

If you would like to present, please contact me (@cheekycheeky on Telegram) with a brief description of the topic.

Once I have approved it, I will announce the next session (typically Thursdays) on the Telegram channel ("Reading is fun"), ideally one month in advance.

Background •000 Paper Discussion

Conclusion O

◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

Transformers



Background

Paper Discussion

Conclusion 0

Language Model

LLM – Next token predictor



◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ● ● ●

Background

Paper Discussion

Conclusion O

◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

Reinforcement Learning



Step 2

Paper Discussion

Step 3

using PPO.

Optimize a policy against

the reward model using reinforcement learning.

Conclusion O

Example: RLHF (GPT-3)

Collect comparison data.

and train a reward model

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is A prompt and A new prompt sampled from our several model is sampled from Explain the moon Explain the moon prompt dataset. inding to a 6 year old outputs are landing to a 6 year old the dataset. sampled. A 8 Explain gravity Explain was The policy A labeler generates demonstrates the C O Acon is nature astolists of an output. desired output hehavior Some people went to the moon.. A labeler ranks the outputs from best to worst. This data is used SET D . C . A . B The reward model to fine-tune GPT-3 calculates a with supervised reward for learning. This data is used the output. BBB to train our reward model. The reward is used to update 0 · O · A - B the policy

ndel

See also: Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. Advances in neural information processing systems,

35, 27730-27744.

 r_{ν}

Background

Paper Discussion

Conclusion 0

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Problems with RLHF?

???

Background

Paper Discussion

Conclusion O

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Problems with RLHF?

???

• What is the limiting factor for RLHF?

Background

Paper Discussion

Conclusion O

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Problems with RLHF?

???

- What is the limiting factor for RLHF?
- How can we scale it up?

Background

Paper Discussion

Conclusion O

Verifiable Rewards



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Paper Discussion

Effect of current RLVR on LLM's reasoning ability



Background

Paper Discussion

Conclusion O

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Paper's hypothesis



Paper Discussion

Conclusion O

Pass@K



Figure 2: Passible curves of base models and their RUNR-trained counterparts across multiple mathematical benchmarks. When h is small, RL-trained models outperform their base versions. However, as h increases to the tens or hundreds, base models consistently catch up and surpass RL-trained models. More results on GSM8K and AMC23 can be found at Figure 3.

Pass@K used by the authors of the paper, instead of more traditional metrics (e.g., Pass@1, Best-of-N, Majority Voting). Why? What are the authors trying to measure?

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─の�?

ackground 000 Paper Discussion

Conclusion O

RLVR (coding benchmarks)



Figure 3: RLVR for Coding.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○三 のへ⊙

Paper Discussion

Conclusion O

Author's summary

- 1. \cdots problems solved by the RLVR model are also solvable by the base model;
 - observed improvement in average scores stems from more efficient sampling on these already solvable problems, rather than learning to solve new problems.
- 2. ··· after RLVR training, the model often exhibits narrower reasoning coverage compared to its base model.
- 3. · · · all the reasoning paths exploited by the RLVR model are already present in the sampling distribution of the base model.
- 4. ··· RLVR does not introduce fundamentally new reasoning capabilities and that the reasoning capacity of the trained model remains bounded by that of its base model.

Background

Paper Discussion

Conclusion O

Distillation versus RLVR

• · · · training data consist of long CoT reasoning traces generated by the teacher model · · ·



Figure 7: Coverage comparison of base, Instruct, RLVR, and distilled models.

Background

Paper Discussion

Conclusion

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

CC's Thoughts

- Interesting paper, experiments and hypothesis are both convincing
- Use of Pass@K instead of other metrics such as Pass@1, Best-of-N, and Majority Voting is interesting for understanding model *capability*
- RLVR is still useful; improves performance for Pass@1
- Distillation is potentially more costly, but would has higher *potential* upside (e.g., smaller model size, improved reasoning capabilities)