

Simpler is Better: Finding the Best Reward Function in Long **Chain-of-Thought Reinforcement Learning for Small Language Models**

Junkuan Liu, Zichen Zhang, Luning Wang {junkuan, zhangzzc, lnwang}@umich.edu

Cosine Reward rewards decreasing from +2.0 (shortest) to +1.0 (longest). with rewards increasing from -10.0 (shortest) to 0.0 b) Cosine Reward a) Classic Reward 1.0 0.8 -6.0 **ard** vard **6** 0.4**a** −4 correct 0.2 wrong 0.0 8,000 16,000 8,000 16,000 **Generation Length Generation Length Dynamic Reward**

- a stronger repetition penalty.
- becomes less likely, requiring a weaker repetition penalty.



on completion length.

show an explicit relationship with reward functions.

- samples for deeper understanding of the phenomena.

