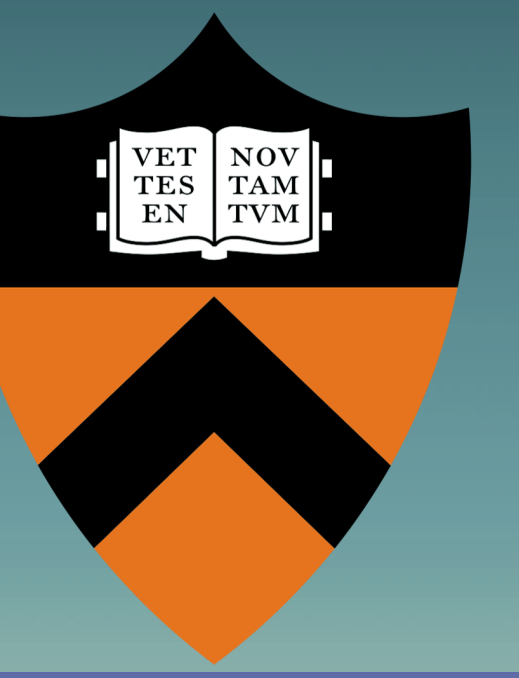# AUGMENTING LARGE REASONING MODELS WITH CONTRASTIVE GOAL-CONDITIONED RL

DEVAN SHAH, KEVIN WANG, DAVID YAN

## MOTIVATION

As the paradigm in artificial intelligence shifts from pre-training scaling laws toward test-time training with RL, reasoning models have emerged as the next frontier. Ever since the release of DeepSeek-R1, a growing trend involves training on complex reasoning tasks with verifiable outcomes (i.e. reward of 1 if the solution is correct and 0 otherwise, with simple format rewards. We make the observation that this outcome-oriented reward paradigm is effectively a goal-conditioned setup.

Meanwhile, in the broader RL community, recent self-supervised RL algorithms have shown strong success on classical goal-conditioned settings, where sparse reward only provides a single bit of reward feedback for each trajectory. A core question thus arises: given this quasi-goal-conditioned paradigm in NLP, can these same goal-conditioned self-supervised RL methods be used to advance LLM reasoning?

## INTRODUCTION

**Contrastive RL:** Contrastive RL [1] is a goal-conditioned method learns representations of state-action pairs ($\varphi(s, a)$) and future states ($\psi(s_f)$) such that the representations of future states are closer than the representations of random states. Formally, we maximize the InfoNCE loss:

$$\mathcal{L}(s, a, s_f^+, s_f^-) \triangleq \log \sigma( \underbrace{f(s, a, s_f^+)}_{\phi(s,a)^T \psi(s_f^+)} ) + \log(1 - \sigma( \underbrace{f(s, a, s_f^-)}_{\phi(s,a)^T \psi(s_f^-)} ))$$

The final reward model captures the expected reward of achieving a goal state at intermediate reasoning steps.

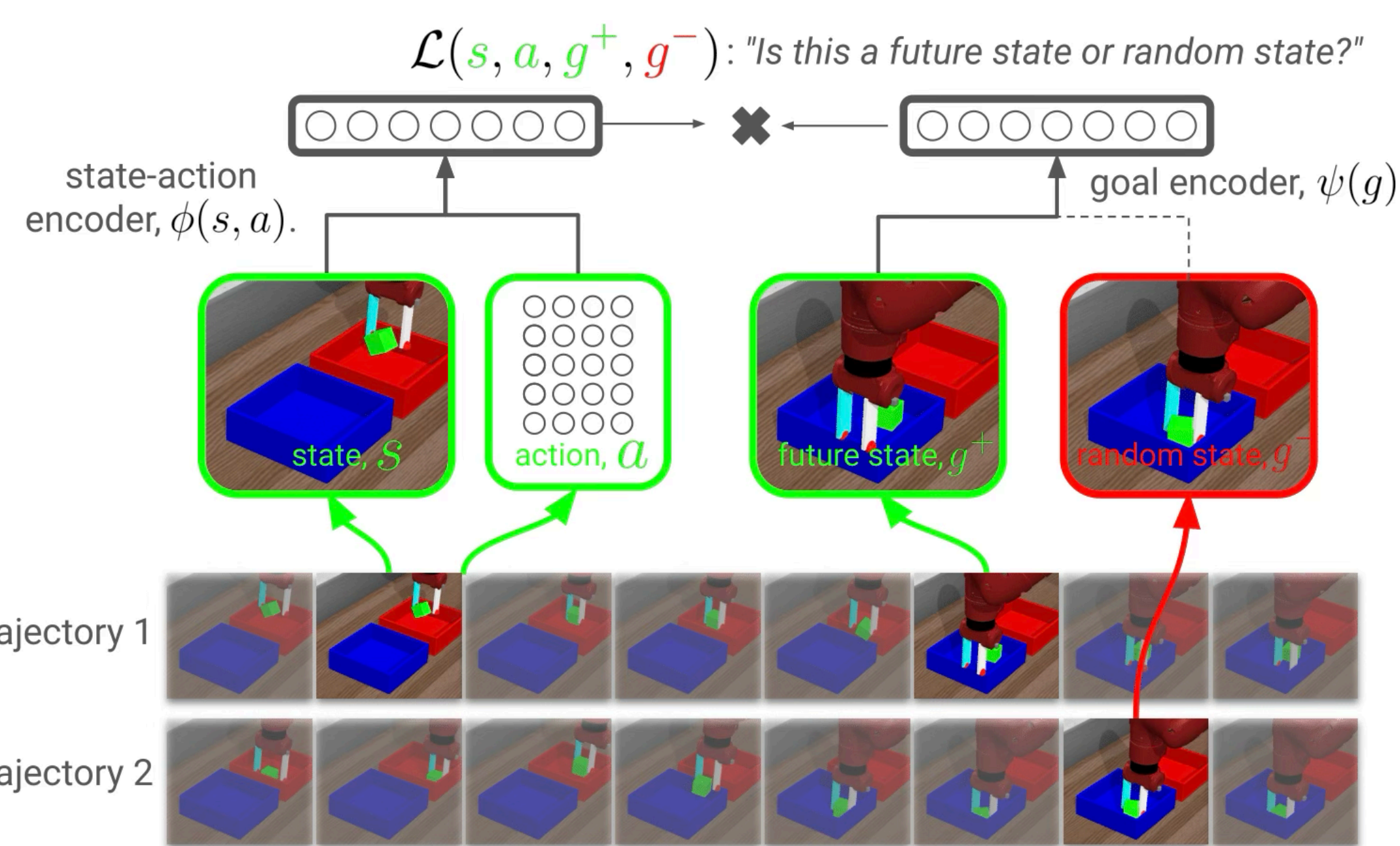$$\mathcal{L}(s, a, g^+, g^-): \text{"Is this a future state or random state?"}$$

state-action encoder, $\phi(s, a)$.     goal encoder, $\psi(g)$

state, $S$    action, $a$    future state, $g$    random state

trajectory 1

trajectory 2



**Figure 1:** Classic example of GCRL for a robotics task [2].

**Related Work:** We build off LongProc's [3] countdown benchmark prompt and validation pipeline to train our reward model. ScaleAI presented a method to train reward models using contrastive learning [4], but their method requires a human-labeled dataset of positive and negative examples, limiting generalizability. Our method requires no labels and works fully unsupervised.

## APPROACH

### State-Action-Goal Sampling

#### Trajectory 1

<Search Procedure>
Initial number set: [21, 16, 17, 26], target: 46

Pick two numbers (21, 16) (numbers left: [17, 26]). Try possible operations. ### Current State: [21, 16, 17, 26]

|- Try 21 + 16 = 37. Add 37 to the number set. Current number set: [37, 17, 26], target: 46. Options for choosing two numbers: [(37, 17), (37, 26), (17, 26)]. ### Current State: [37, 17, 26]

⋮  **State**

|- Try 37 + 17 = 54. Add 54 to the number set. Current number set: [54, 26], target: 46, just two numbers left. ###  **Action**

Current State: [54, 26]

|- Try 54 + 26 = 80. Evaluate 80 != 46, drop this branch. ###
Current State: [80]

⋮

Current State: [20, 26]  **Future Goal**

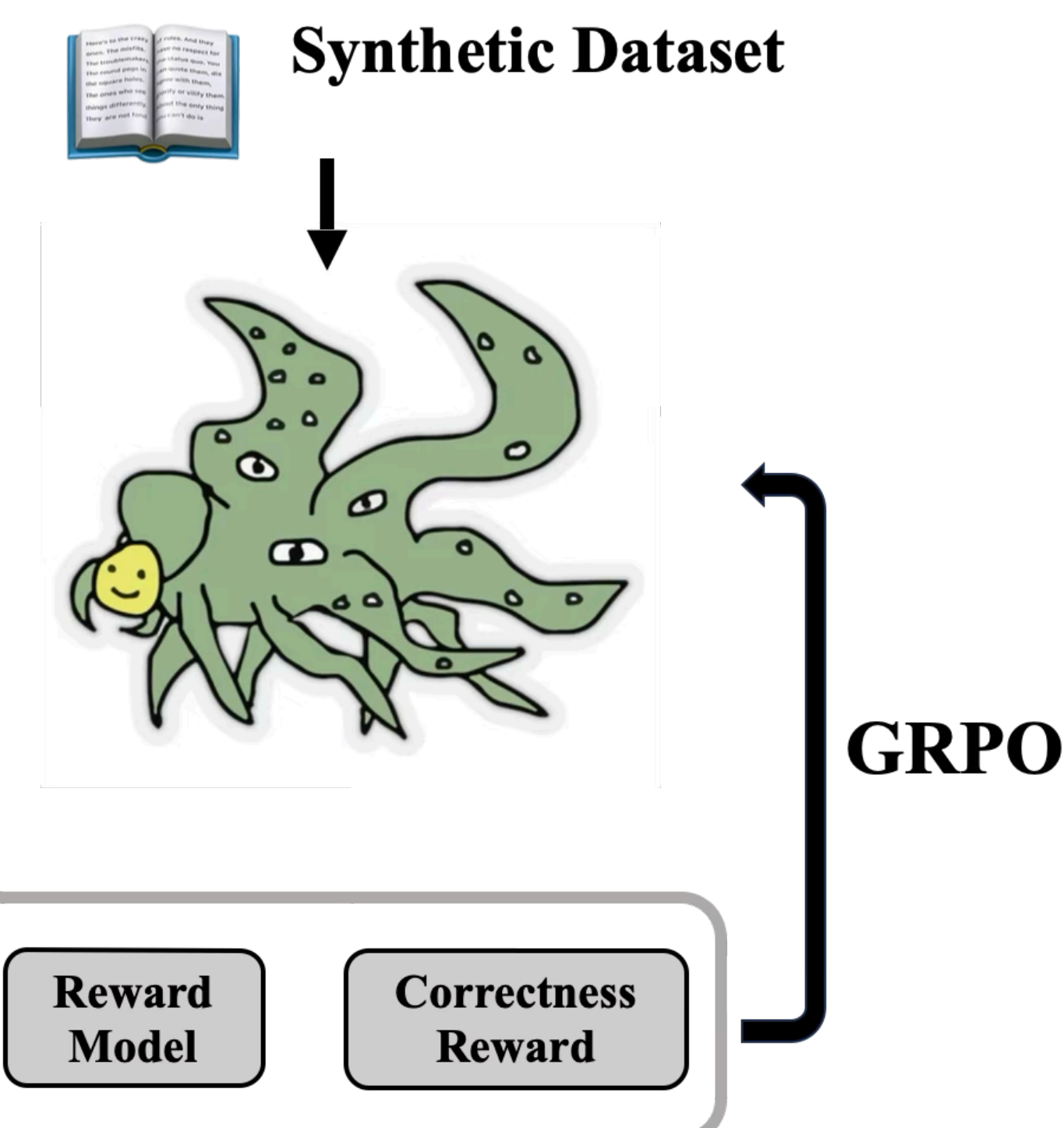|- Try 26 + 20 = 46. Evaluate 46 == 46, target! ### Current State: [46]
</Search Procedure>

#### Trajectory 2

⋮
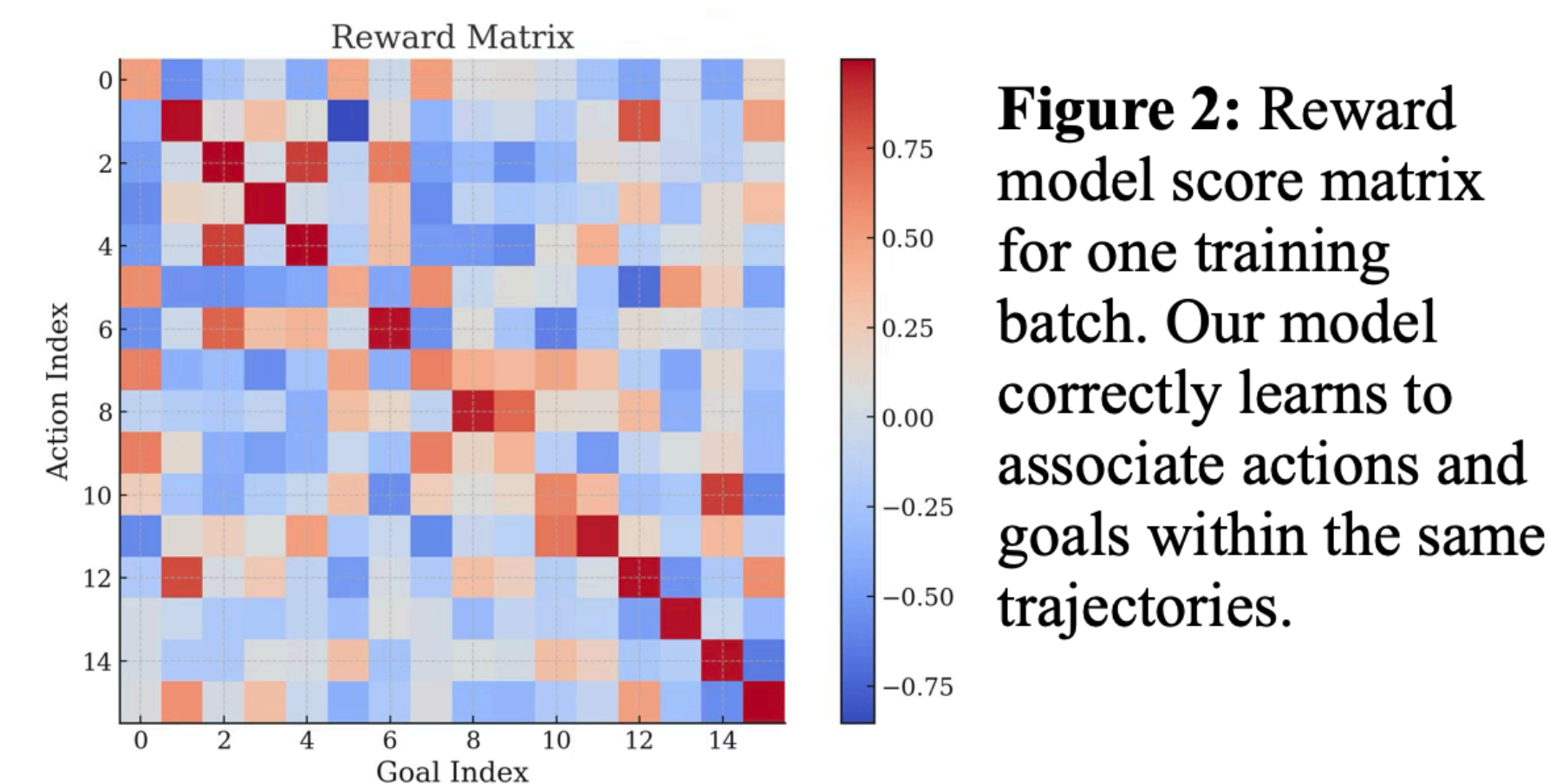
Current State: [336, 442]  **Random Goal**

⋮



**Figure 2:** Reward model score matrix for one training batch. Our model correctly learns to associate actions and goals within the same trajectories.

**Figure 2.** Reasoning traces are split into state–action pairs, each matched with a later goal from the same trace. A pretrained LM is fine-tuned on the resulting triples via GRPO, while a contrastive critic (heat-map) learns dense rewards that align actions with future goals. The improved policy then generates new traces, which are resampled into fresh triples, closing a self-supervised actor-critic loop for goal-conditioned reasoning.

### GRPO Finetuning

**Synthetic Dataset**



GRPO

| Reward Model | Correctness Reward |

### Critic Model Rewards

Reward Matrix

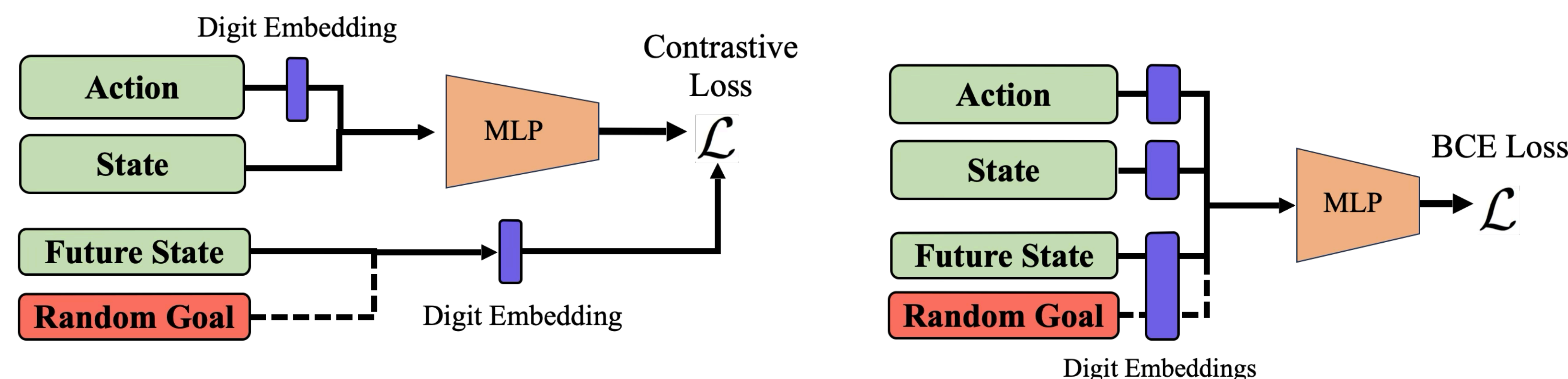### Critic Model Architectures



**Figure 3:** Architectures experimented on for the critic model. We found that the left architecture using a contrastive loss was insufficiently expressive because it relies on dot-product similarity to compare the state-action and goal embeddings.

## CONTRASTIVE CRITIC



InfoNCE Critic Training Loss
Steps: 2000000
Min Loss: 2.0557
Max Loss: 3.1453
— Raw Loss
— 100-Moving Average
--- Average Loss: 2.4823

BCE Critic Training Loss
Steps: 2000000
Min Loss: 0.2268
Max Loss: 0.7099
— Raw Loss
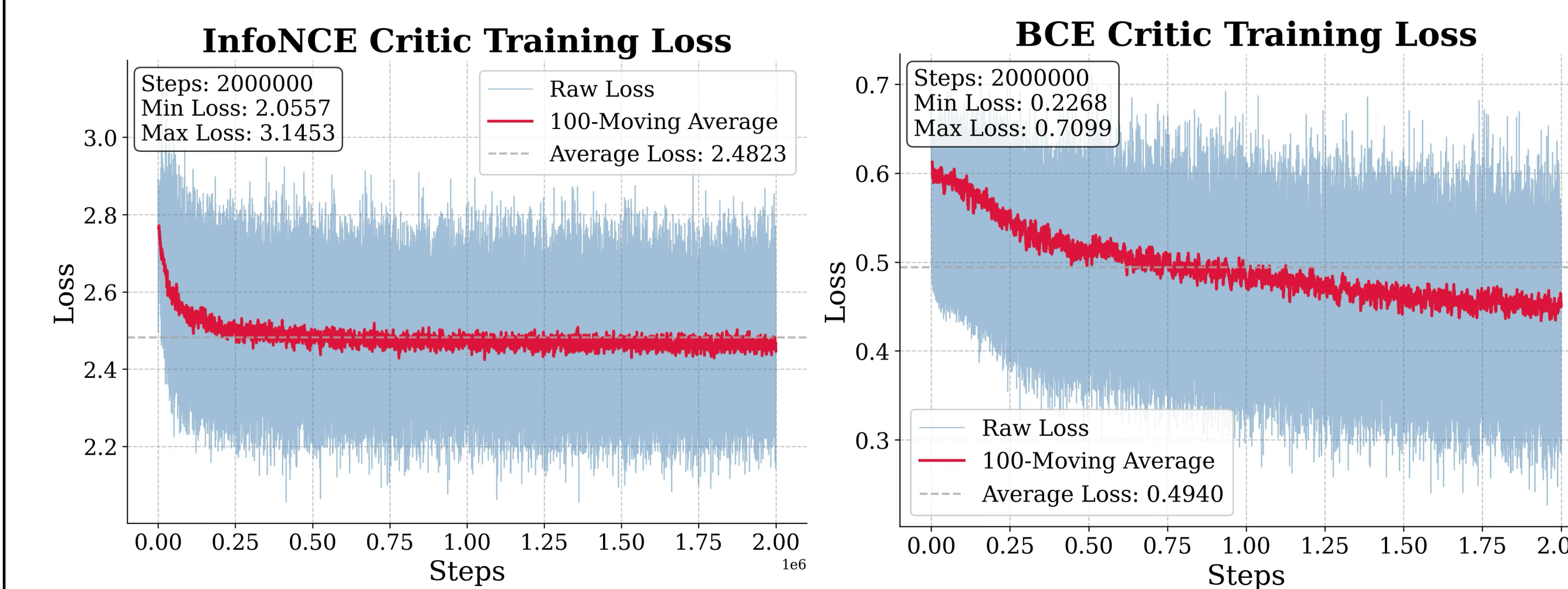— 100-Moving Average
--- Average Loss: 0.4940

**Figure 4:** We train our critic models on 2,000,000 synthetic examples extracted from 100,000 procedurally generated trajectories to associate goals with trajectories.
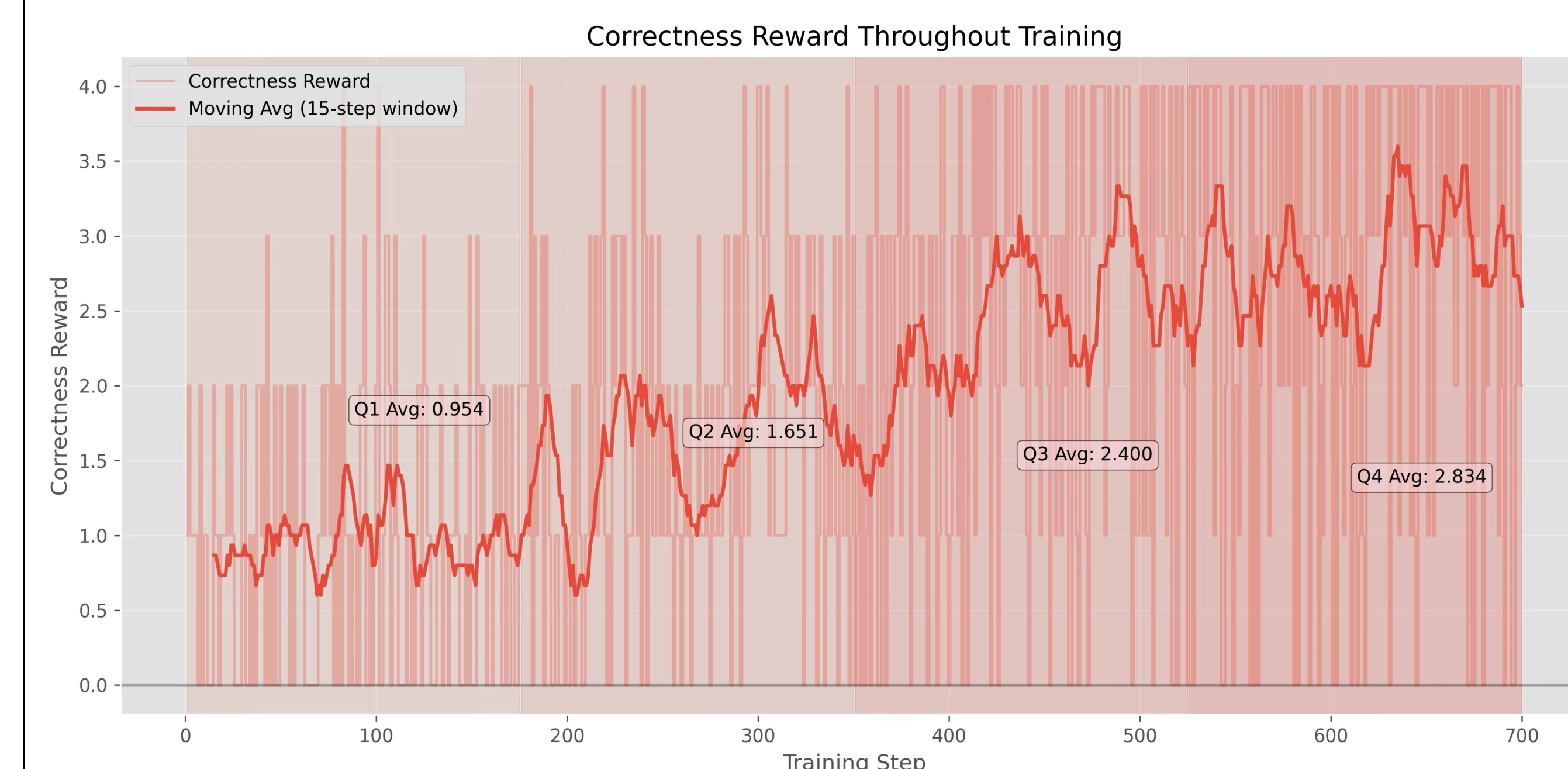
## GRPO RESULTS



Correctness Reward Throughout Training
— Correctness Reward
— Moving Avg (15-step window)
Q1 Avg: 0.954    Q2 Avg: 1.651    Q3 Avg: 2.400    Q4 Avg: 2.834

**Figure 5:** Baseline success rate over time during GRPO without contrastive loss.

## REFERENCES

[1] DeepSeek-AI, et al. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. https://arxiv.org/abs/2501.12948

[2] Benjamin Eysenbach, Tianjun Zhang, Ruslan Salakhutdinov and Sergey Levine. Contrastive Learning As a Reinforcement Learning Algorithm. https://arxiv.org/abs/2206.07568

[3] Xi Ye, Fangcong Yin, Yinghui He, Joie Zhang, Howard Yen, Tianyu Gao, Greg Durrett and Danqi Chen. LongProc: Benchmarking Long-Context Language Models on Long Procedural Generation. https://arxiv.org/abs/2501.05414

[4] Vaskar Nath, Dylan Slack, Jeff Da, Yuntao Ma, Hugh Zhang, Spencer Whitehead, Sean Hendryx. Learning Goal-Conditioned Representations for Language Reward Models. https://arxiv.org/abs/2407.13887