





# RLVR-World: Training World Models with Reinforcement Learning

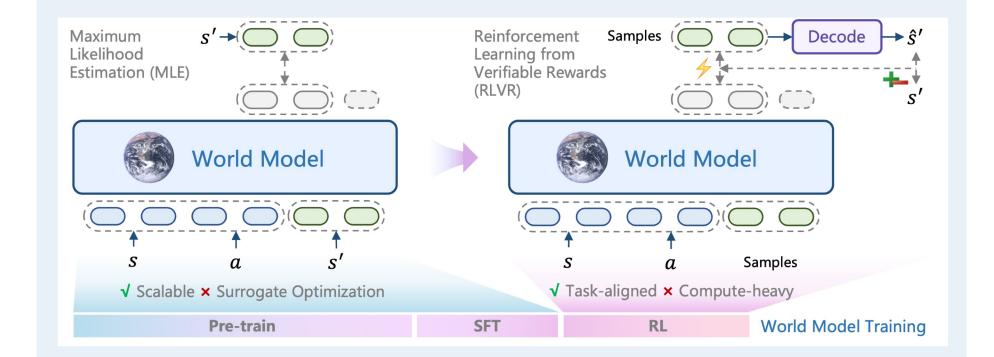


Jialong Wu, Shaofeng Yin, Ningya Feng, Mingsheng Long

#### TL;DR

#### **Key Insight**

As world models are built with more advanced non-end-toend models (discrete autoregressive models or diffusion models), their training methods typically optimizing for likelihood are misaligned with the actual task objective of world modeling.



#### **Contributions**

RLVR as a direct and effective post-training method: Task-specific prediction metrics serve as verifiable rewards

World models across modalities (language and video): Unified under a sequence modeling formulation

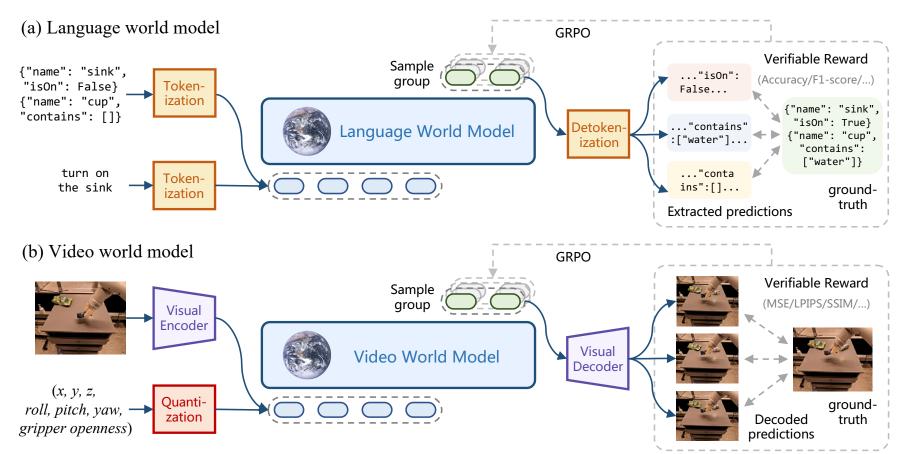
RL-trained world models enhance decision-making: Policy evaluation and model predictive control

Project Website: thuml.github.io/RLVR-World

Code: github.com/thuml/RLVR-World

Datasets & Models: hf.co/collections/thuml/rlvr-world

#### **RLVR-Framework**



Unified Sequence Modeling Formulation (Diffusion in future work!)

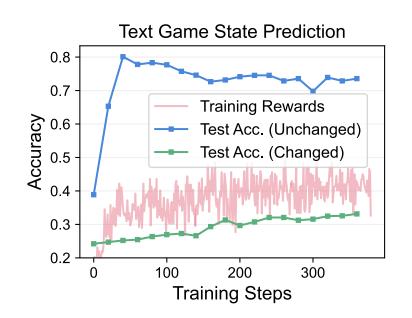
$$\mathcal{J}_{ ext{MLE}}( heta) = \log p_{ heta}(o(s') \mid q(s, a)) = \sum_{t=1}^{|o(s')|} \log p_{ heta}(o_t(s') \mid q(s, a), o_{< t}(s'))$$

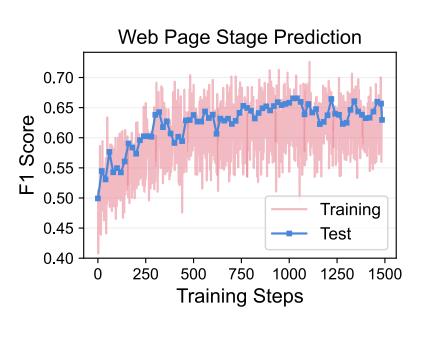
Task-Specific Prediction Metrics as Rewards

sign(D) = -1If Lower D is better GRPO with  $R_i = \text{sign}(D) \cdot D(\hat{s}_i', s')$ sign(D) = 1Otherwise

## Language World Models

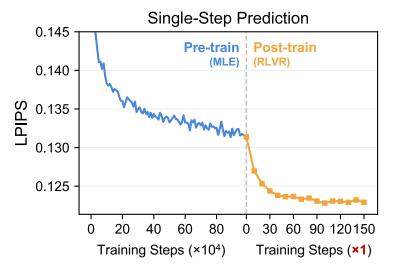
Beyond its success in math and coding, RLVR can also improve LLMs' performance on world modeling tasks involving verbal state transitions.

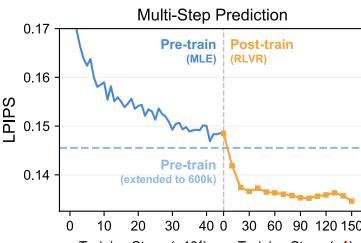




#### **Video World Models**

RLVR bridges the gap between pre-training objectives and visual prediction metrics, leading to more accurate predictions, improved training efficiency, and reduced artifacts such as repetition.





# 1000× Training Efficiency! Ground-truth Base model

### **World Model Applications**

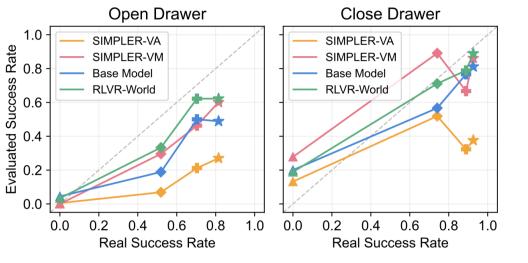
RLVR-trained world models can improve downstream tasks, including policy evaluation and model predictive control.



Success Rate: SFT: 12.06% RLVR: 14.29% (relatively +18.4%)



Real2Sim Evaluation of RT-1 Checkpoints



First work on world model-based robot policy evaluation!