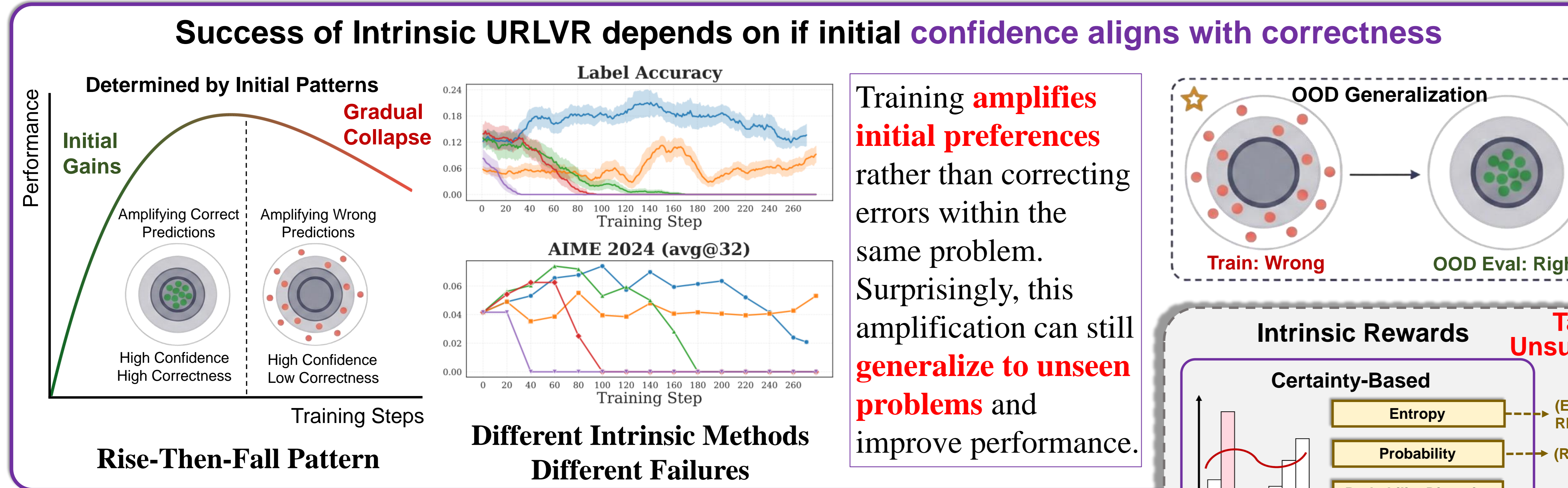


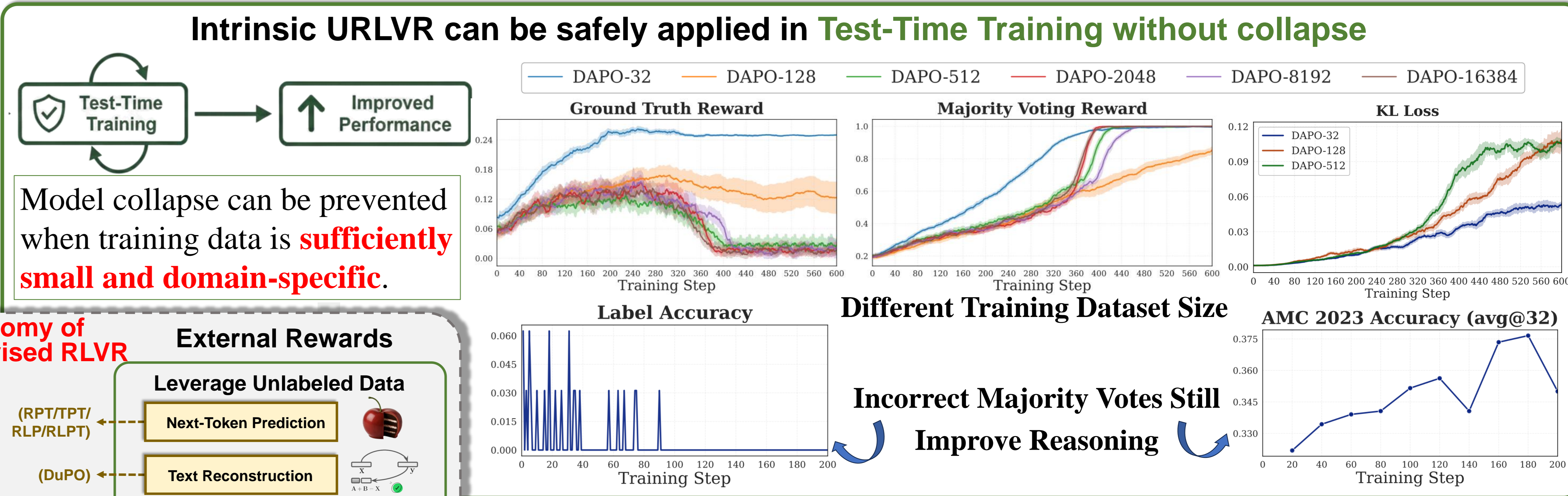
## ❖ When Does Intrinsic Unsupervised RLVR Work?

Intrinsic URLVR universally follows a **rise-then-fall** pattern across all methods. Early gains reflect **confidence-correctness alignment** in the model's prior, while eventual collapse is inevitable when this alignment breaks down.

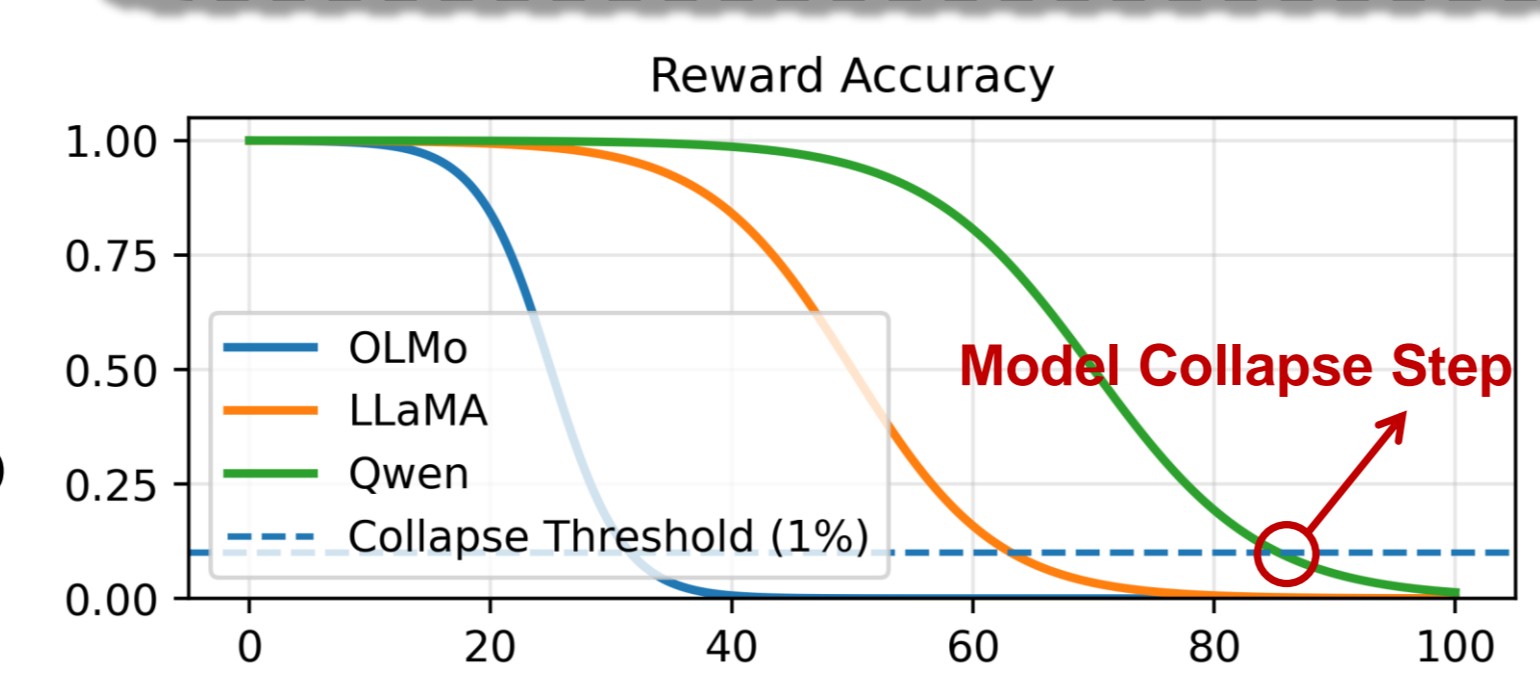
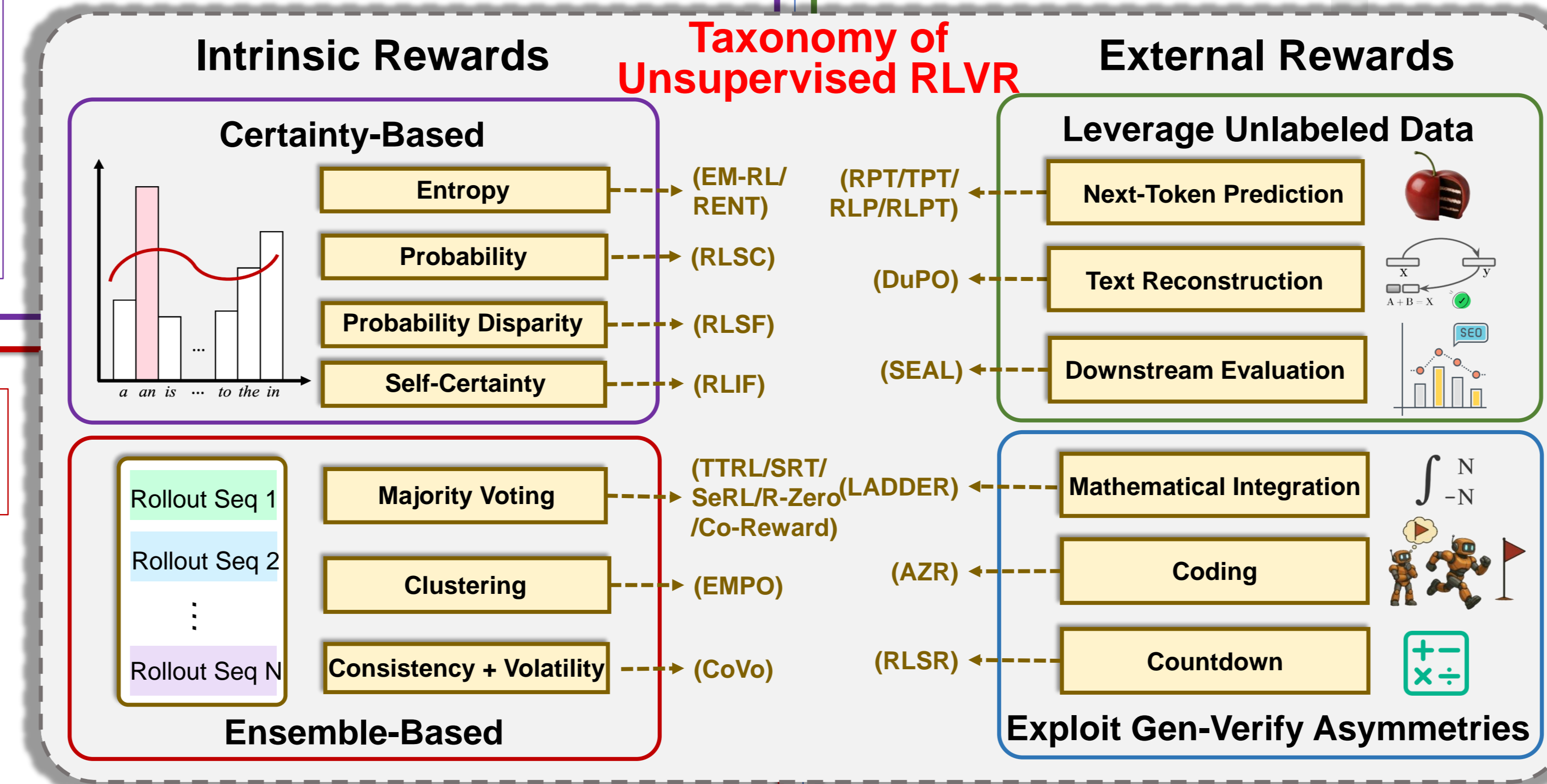
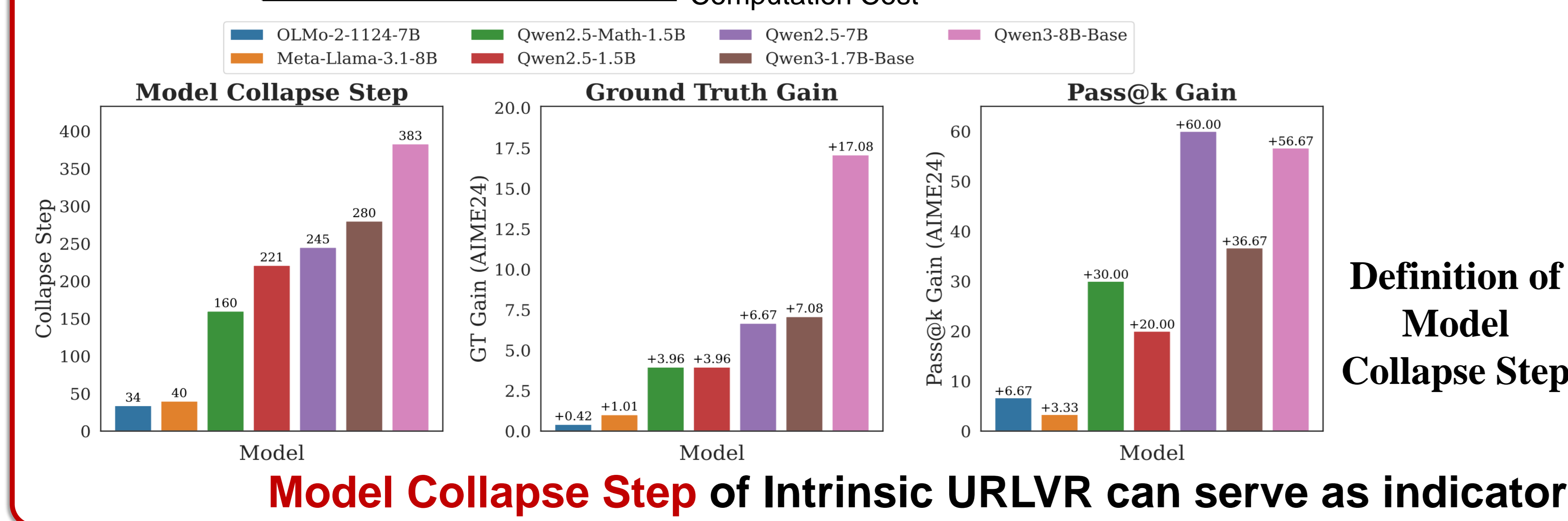


## ❖ How Can Sharpening from Intrinsic URLVR Be Applied Safely?

Small datasets induce localized rather than systematic policy shift, even training on wrong problems can yield gains, making **test-time training a safe and practical application**.

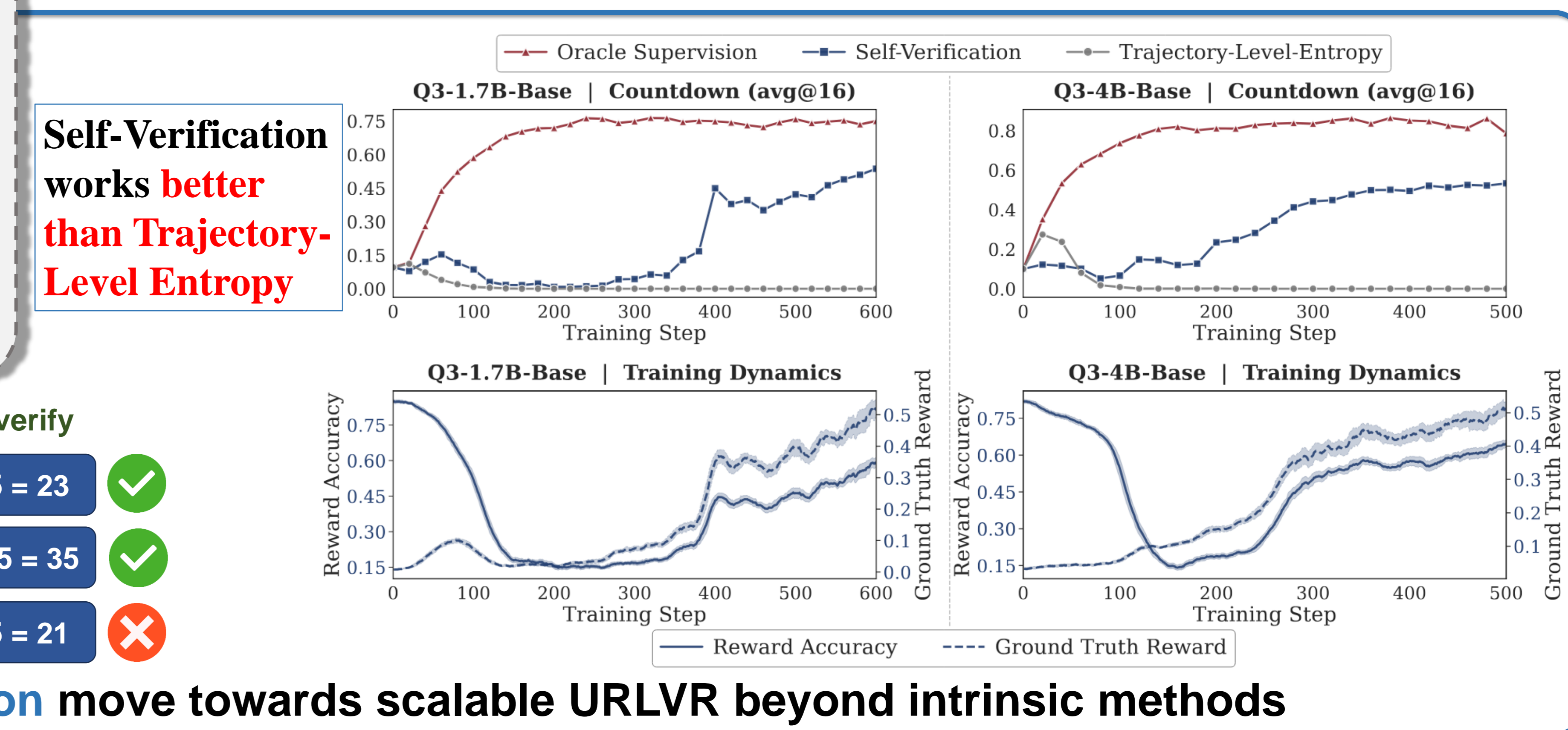


Model Collapse Step correlates strongly with GT Gain, better than Pass@k Gain as indicator. 5.6x Faster.



Hard to generate vs Easy to verify. Asymmetry of Verification. Examples: 3 + 4 x 5 = ? vs 3 + 4 x 5 = 23.

Self-Verification move towards scalable URLVR beyond intrinsic methods



We propose the **Model Collapse Step** as a novel indicator of model priors, which measures standard **RL trainability** by tracking reward accuracy collapses during intrinsic URLVR. This indicator achieves accuracy in assessing trainability on par with running standard RL itself, but with **higher efficiency and outperforming pass@k**.

Intrinsic rewards are fundamentally bounded by what the model already knows. **External rewards grounded in unlabeled data or generation-verification asymmetry** provide signals that **scale with data and computation** rather than saturating with model capacity, offering a more promising path towards scalable URLVR.