

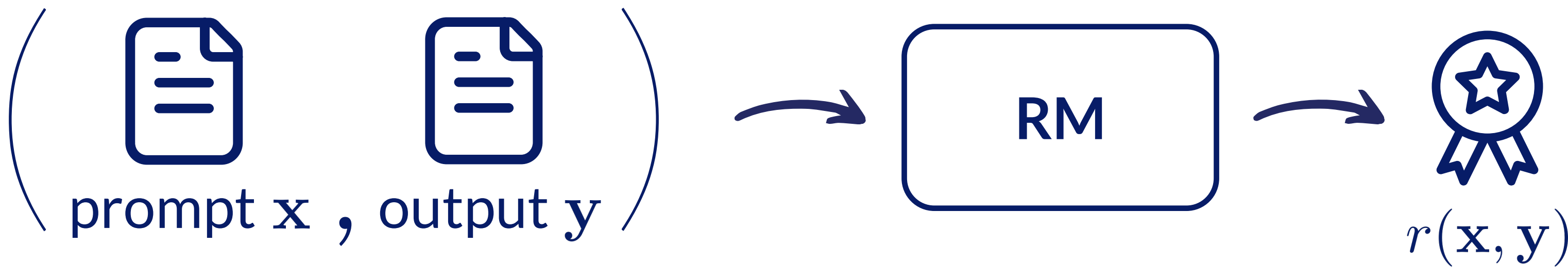
# Why is Your Language Model a Poor Implicit Reward Model?

Noam Razin, Yong Lin, Jiarui Yao, Sanjeev Arora

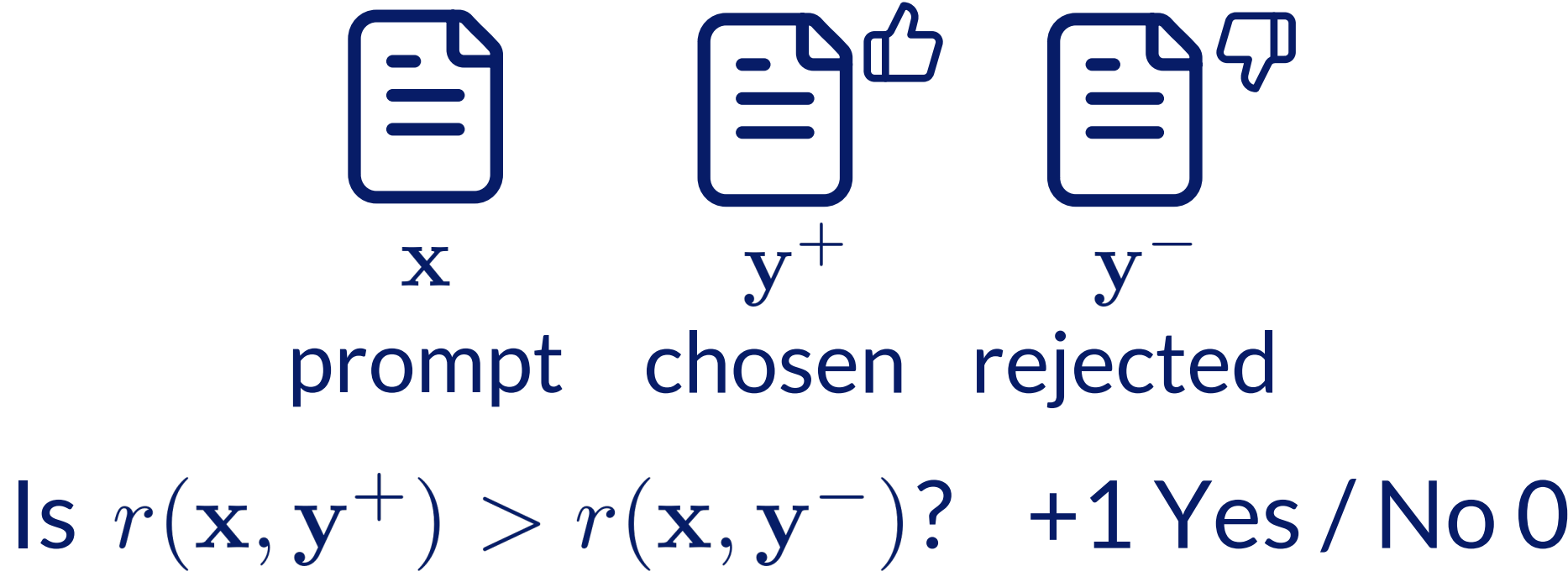


Princeton Language & Intelligence, Princeton University  
University of Illinois Urbana-Champaign

Reward Models (RMs) are widely used in language model (LM) post-training and inference

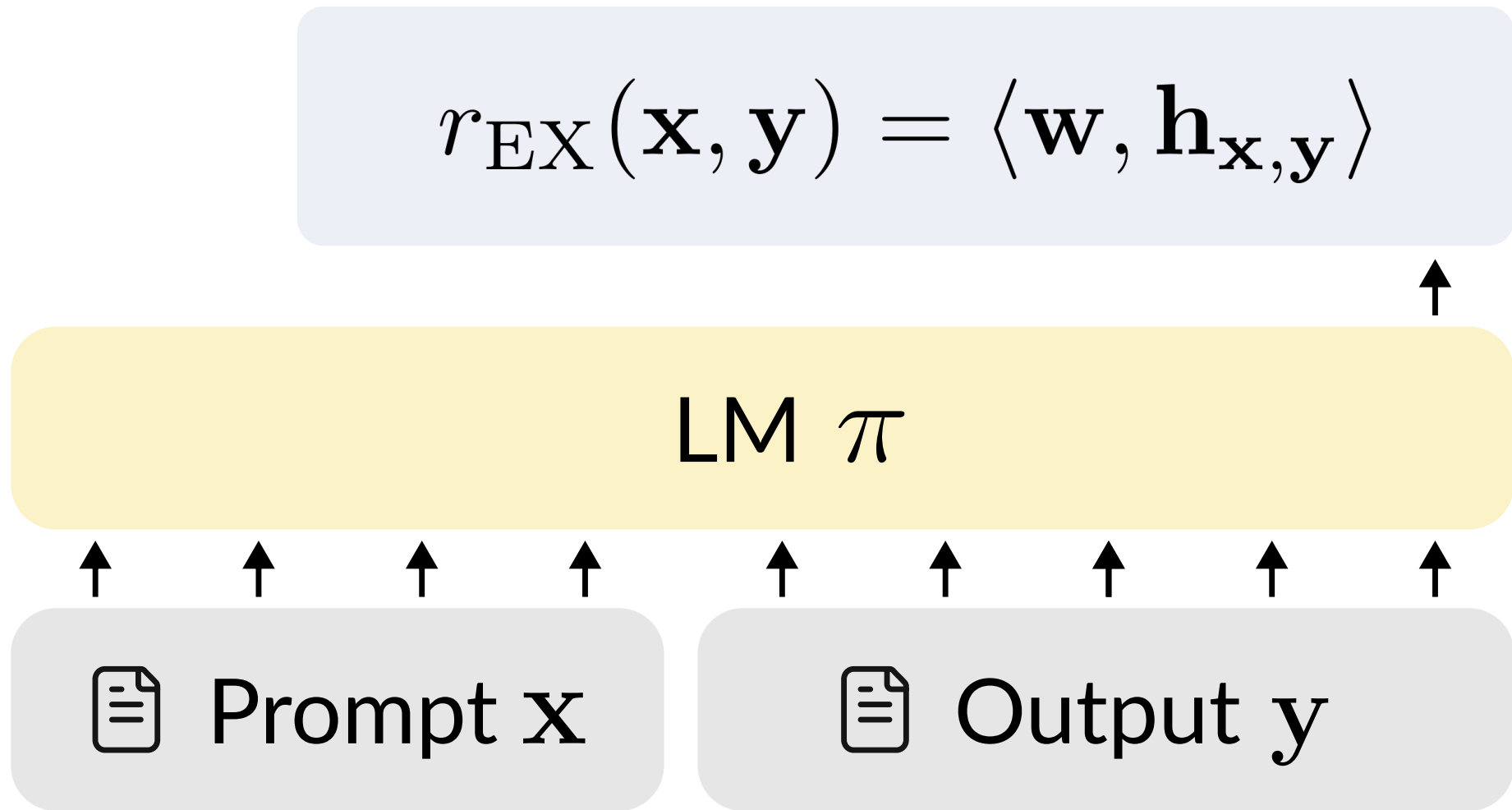


Ability of RMs to evaluate quality of outputs is measured via **accuracy**



## Explicit Reward Models (EX-RMs) vs Implicit Reward Models (IM-RMs)

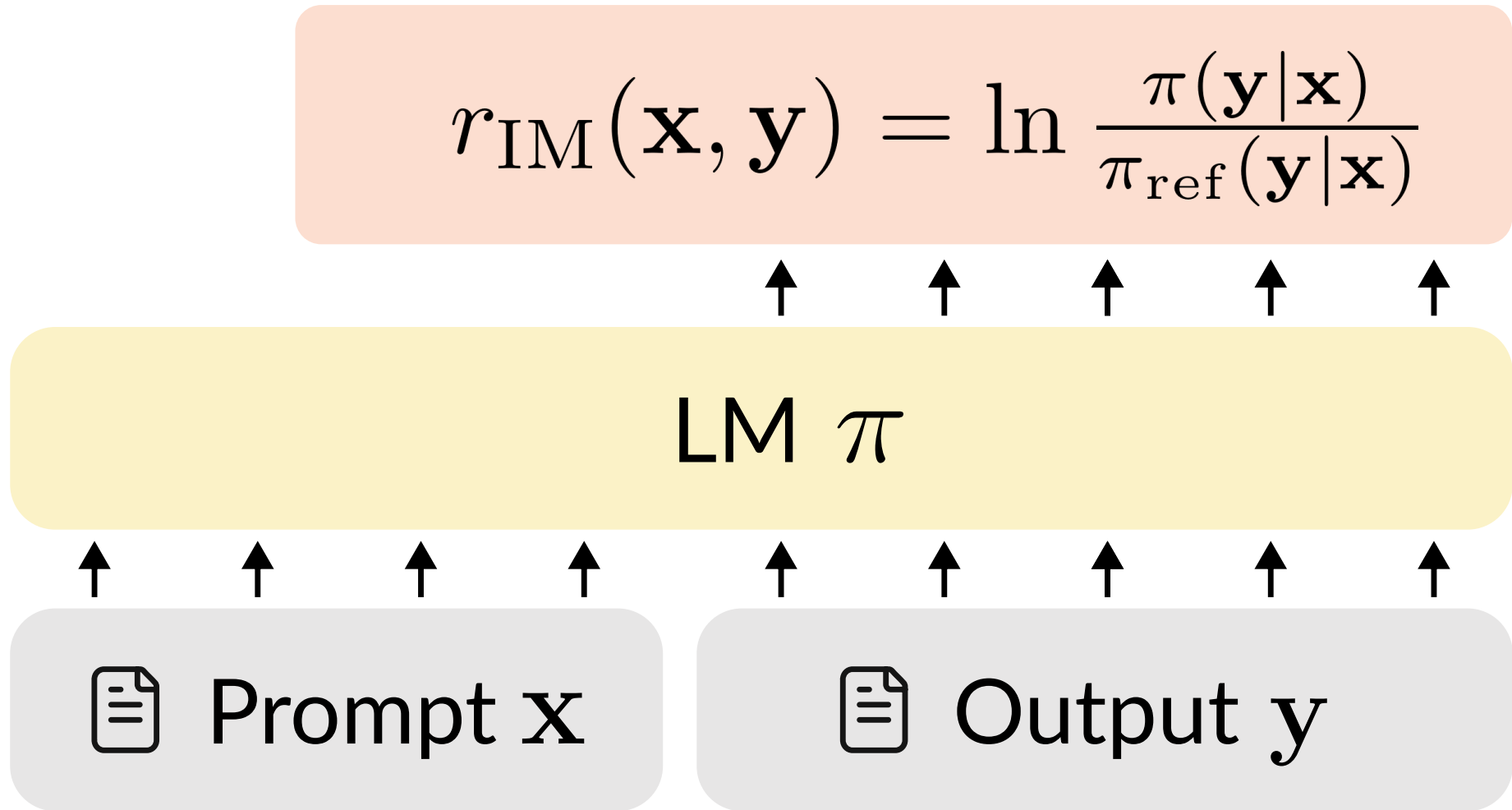
**EX-RM:** Apply a linear head over hidden representation of LM



**Similarities:** Trained using the same data, loss, and LM

**Difference:** How reward is computed based on the LM

**IM-RM:** Every LM defines an RM via its log probabilities (Rafailov et al. 2023)



**Prior Work:** EX-RMs often generalize better than IM-RMs (e.g., Lin et al. 2024, Lambert et al. 2024, Swamy et al. 2025)

**Q:** Why is there a *generalization gap* between EX-RMs and IM-RMs despite their similarity?

### I) Challenge Existing Hypothesis

Do IM-RMs struggle in tasks where generation is harder than verification? **X**

### II) Identify Cause for the Gap

IM-RMs rely more heavily than EX-RMs on superficial token-level cues **✓**

**Existing Hypothesis:** If generation is harder than verification, IM-RM should be harder to learn than EX-RM (e.g., Dong et al. 2024, Singhal et al. 2024)

Trained to:	Verify	Generate
EX-RM	✓	✗
IM-RM	✓	✓

We provide evidence against this hypothesis

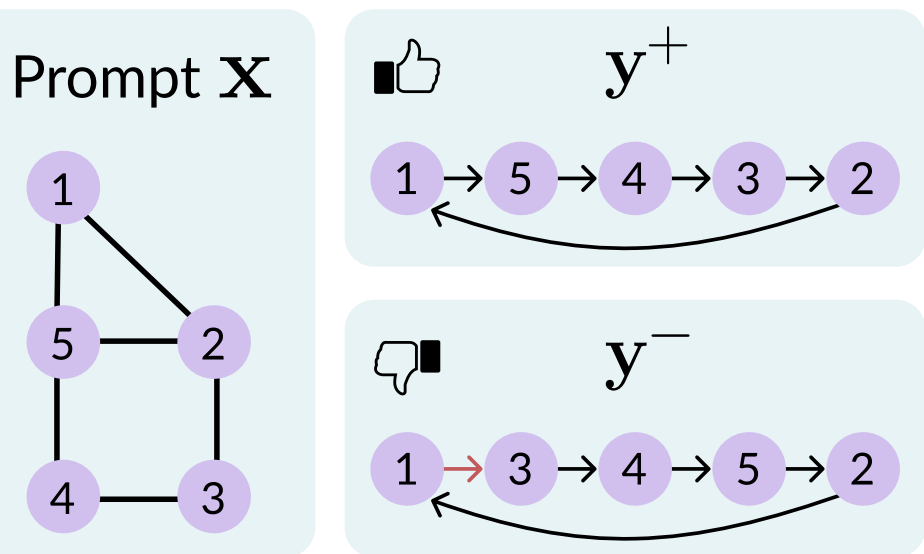
#### Theory

**Theorem:** IM-RM can learn to verify without learning to generate

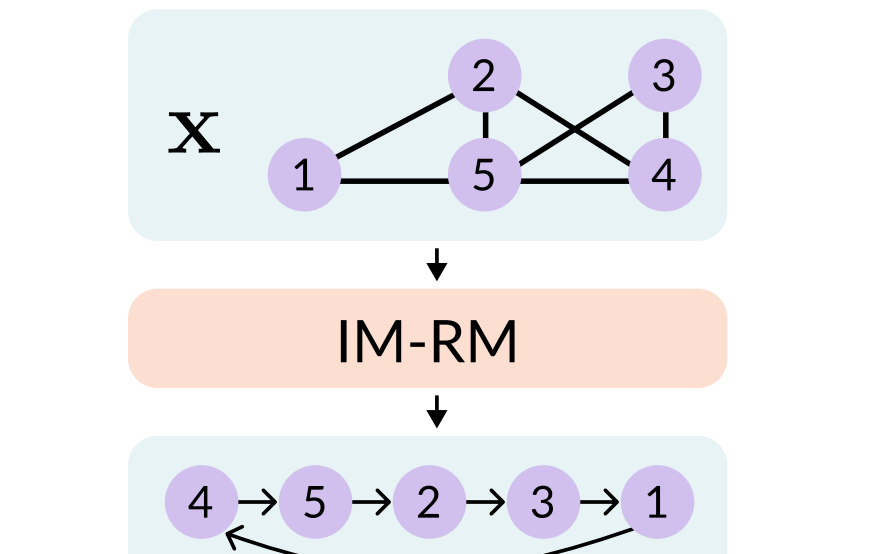
reward of "good" outputs > reward of "bad" outputs + const

#### Experiments

Hamiltonian Cycle **Verification**



Hamiltonian Cycle **Generation**



	EX-RM	IM-RM
Test Accuracy	0.980	0.993
Correct Generations	-	0

LM: Pythia-1B

**Result:** Despite the generation-verification gap, the IM-RM accurately verifies outputs (without being able to generate)

### Theory: Learning Dynamics

#### Approach

Characterize how a gradient update on  $(x, y^+, y^-)$  affects reward assigned to unseen prompt-output pair  $(\bar{x}, \bar{y})$

$$\Delta r(\bar{x}, \bar{y}) \approx \langle -\nabla \text{loss}(x, y^+, y^-), \nabla r(\bar{x}, \bar{y}) \rangle$$

#### Results: EX-RM

$\Delta r_{\text{EX}}(\bar{x}, \bar{y})$  depends on outputs through hidden representations

The reward increases when  $h_{\bar{x}, \bar{y}}$  is more aligned with  $h_{x, y^+}$  than with  $h_{x, y^-}$

#### Results: IM-RM

$\Delta r_{\text{IM}}(\bar{x}, \bar{y})$  depends directly on tokens in the outputs

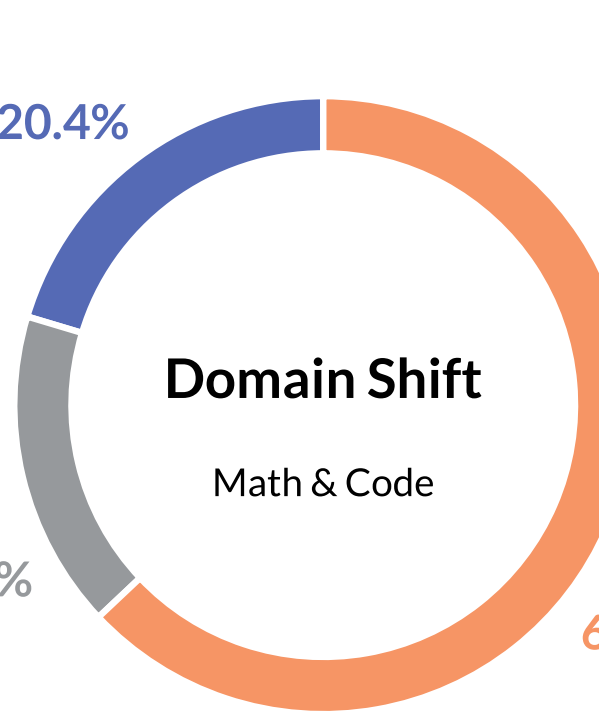
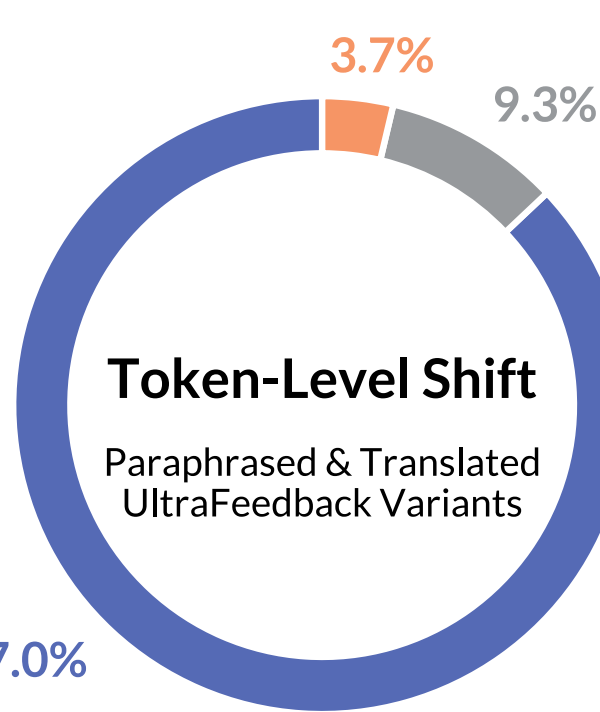
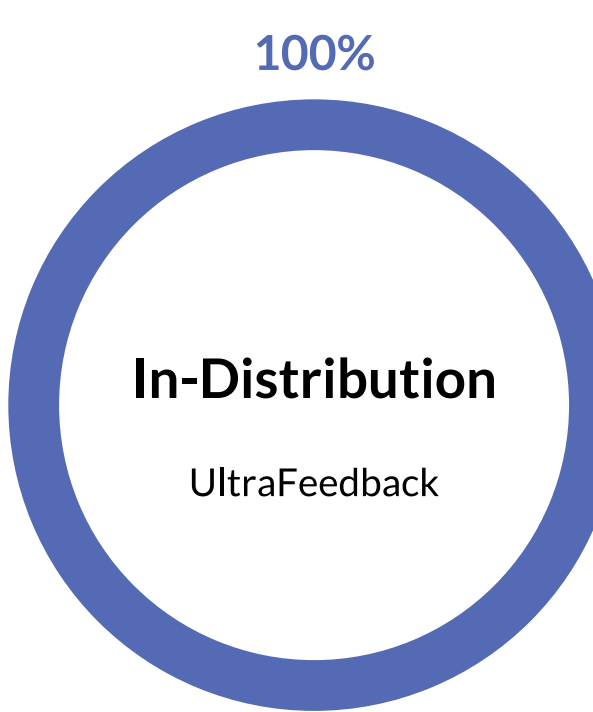
Tokens of  $\bar{y}, y^+$  overlap? Effect similar to EX-RM  
Tokens of  $\bar{y}, y^+$  are distinct? Effect opposite to EX-RM

IM-RM may decrease rewards of outputs semantically similar to  $y^+$  if their tokens have little overlap!

### Experiments

Training Data: UltraFeedback

EX-RM Win  
Tie  
IM-RM Win



Based on LMs of up to 8B size from the Gemma, Qwen, and Llama families

**Result:** In line with our theory, IM-RMs are less robust to token-level shifts, but perform comparably or better under domain shifts