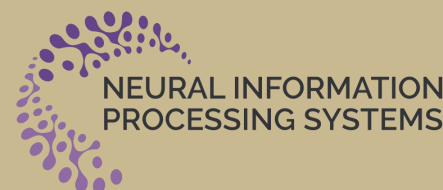


PRIMT: Preference-based Reinforcement Learning with Multimodal Feedback and Trajectory Synthesis from Foundation Models

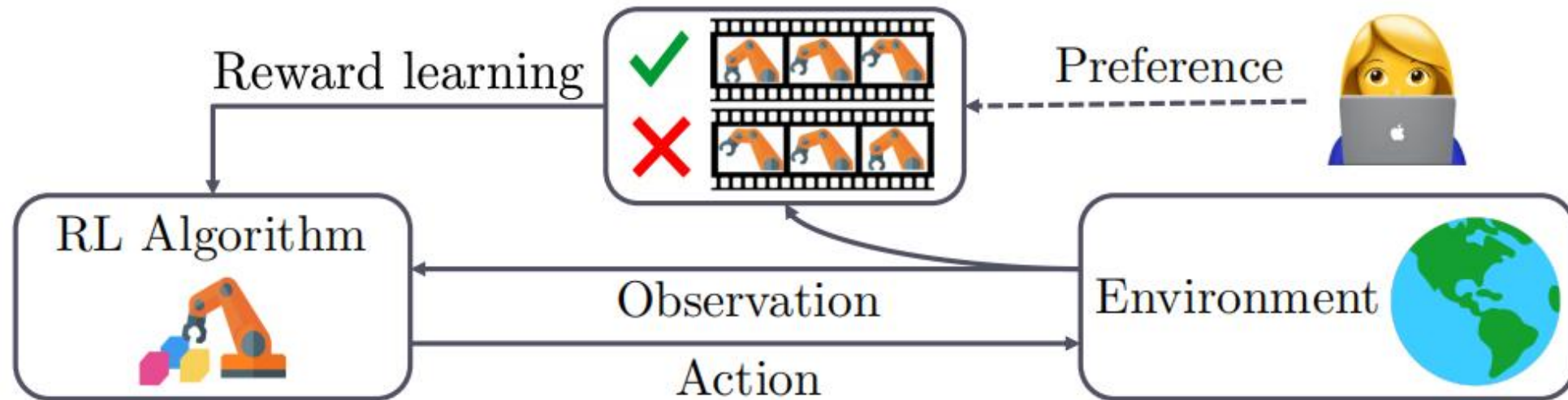
Ruiqi Wang, Dezhong Zhao, Ziqin Yuan, Tianyu Shao,
Guohua Chen, Dominic Kao, Sungeun Hong, Byung-Cheol Min

Purdue University, Indiana University Bloomington,
UIUC, BUCT, Sungkyunkwan University



Preference-based Reinforcement Learning

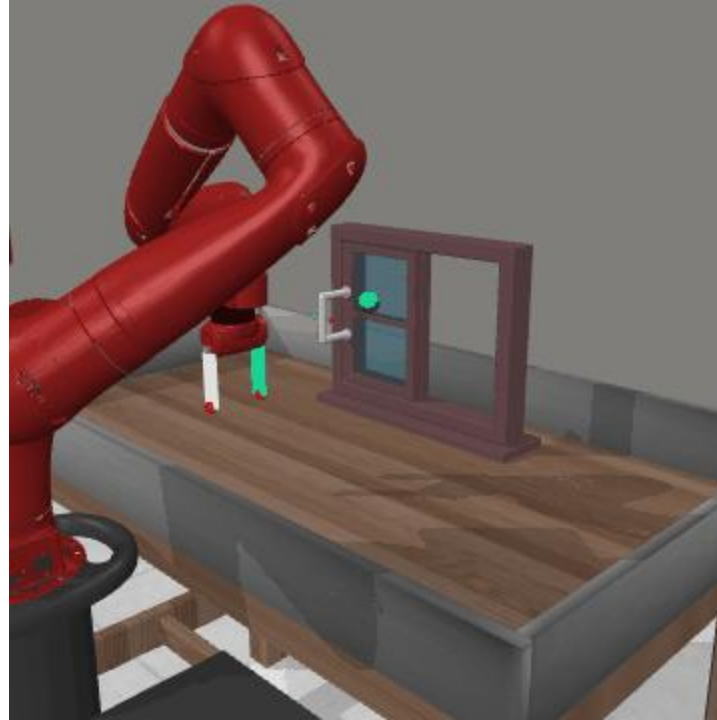
Preference-based RL (PbRL) has emerged as a promising paradigm for teaching robots complex behaviors without reward engineering.



(Lee et al., 2022)

Preference-based Reinforcement Learning

However, the reliance on **extensive human labels** prevents its scalability

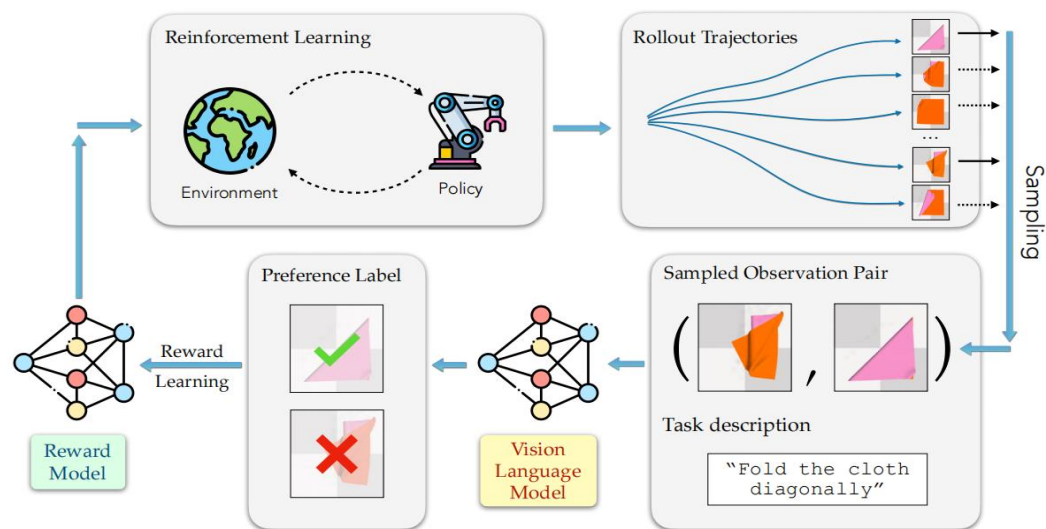


~10000 rounds of human preference are needed!

~84 hours

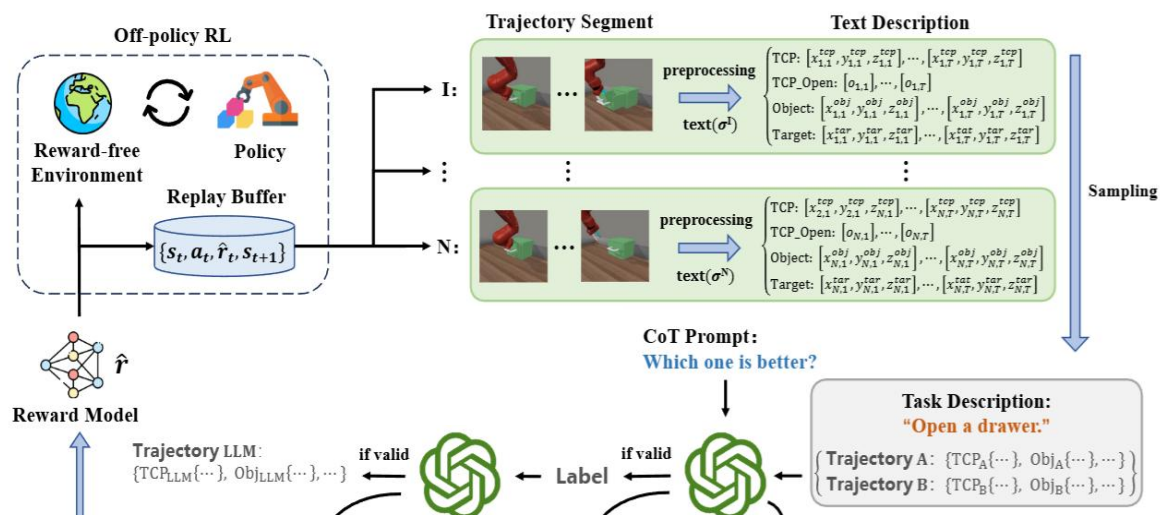
Foundation Models for Synthetic Feedback

VLM agents giving preference labels by analyzing state images



(Wang et al., 2024)

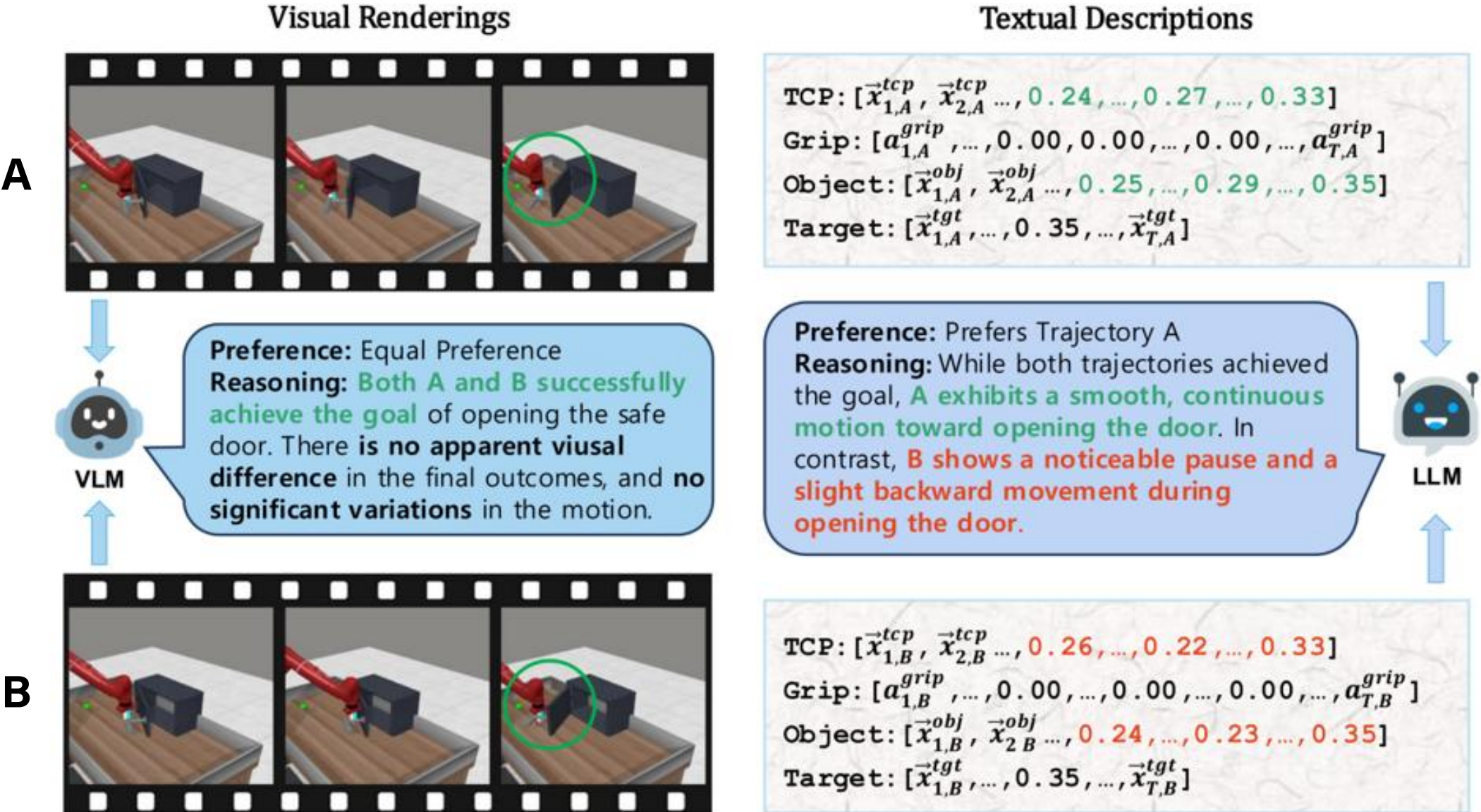
LLM agents giving preference label by analyzing textual trj descriptions



(Tu et al., 2025)

However, existing work still faces challenges in ensuring **reliable synthetic feedback** due to the **single-modal evaluation patterns**.

Limitations of Single-Modal Evaluation



(Answers are from gpt-4o)

Limitations of Single-Modal Evaluation

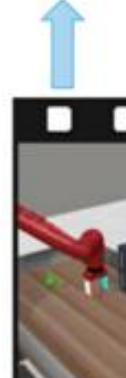
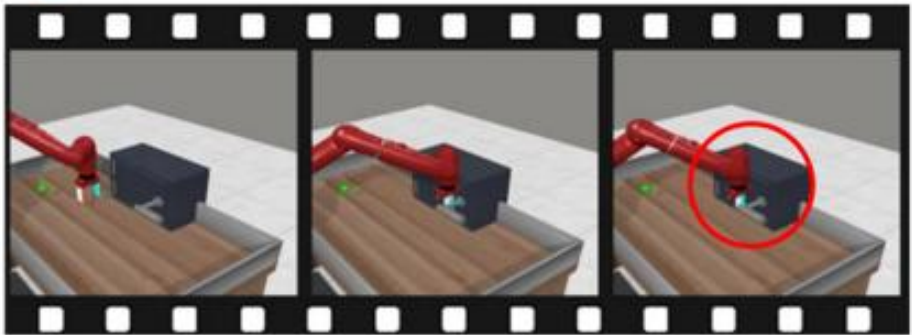
A

Visual Renderings



Preference: Prefers Trajectory A
Reasoning: Trajectory A shows **better task completion with the door slightly opening**. In contrast, Trajectory B **fails to open the door**.

B



Textual Descriptions

TCP: $[\vec{x}_{1,A}^{tcp}, \vec{x}_{2,A}^{tcp}, \dots, 0.23, \dots, 0.01, \dots, 0.08]$
Grip: $[a_{1,A}^{grip}, \dots, 0.95, \dots, 0.00, \dots, 0.00, \dots, a_{T,A}^{grip}]$
Object: $[\vec{x}_{1,A}^{obj}, \vec{x}_{2,A}^{obj}, \dots, 0.00, \dots, 0.00, \dots, 0.08]$
Target: $[\vec{x}_{1,A}^{tgt}, \dots, 0.35, \dots, \vec{x}_{T,A}^{tgt}]$



Preference: Equal Preference
Reasoning: Both Trajectory A and Trajectory B **exhibit clear approaching and opening behaviors, leading to successful task completion**

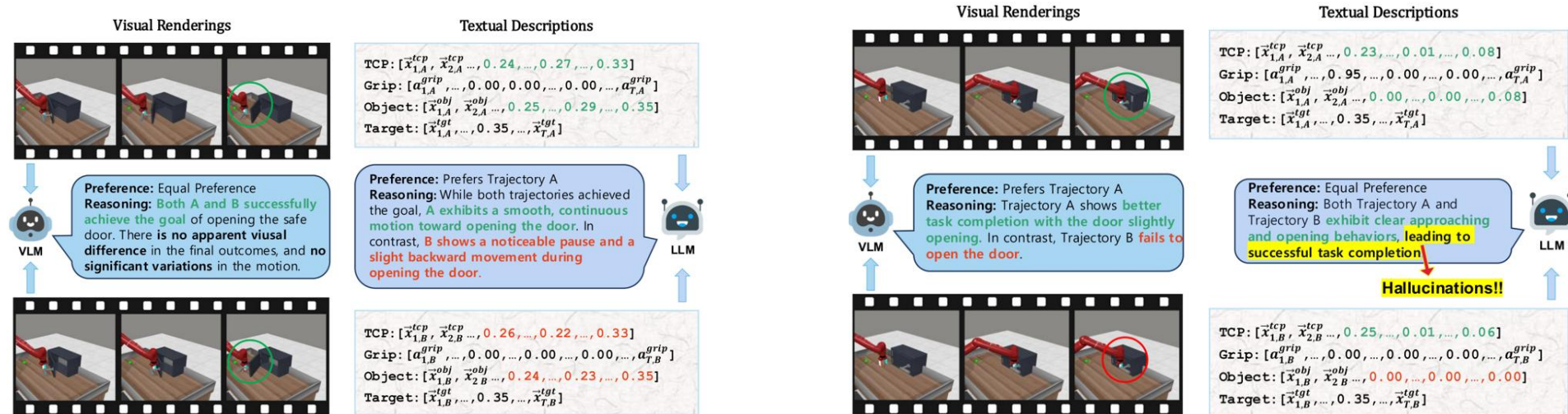
Hallucinations!!



TCP: $[\vec{x}_{1,B}^{tcp}, \vec{x}_{2,B}^{tcp}, \dots, 0.25, \dots, 0.01, \dots, 0.06]$
Grip: $[a_{1,B}^{grip}, \dots, 0.00, \dots, 0.00, \dots, 0.00, \dots, a_{T,B}^{grip}]$
Object: $[\vec{x}_{1,B}^{obj}, \vec{x}_{2,B}^{obj}, \dots, 0.00, \dots, 0.00, \dots, 0.00]$
Target: $[\vec{x}_{1,B}^{tgt}, \dots, 0.35, \dots, \vec{x}_{T,B}^{tgt}]$

(Answers are from gpt-4o)

Limitations of Single-Modal Evaluation



Visual-modal-only (VLMs) → reliable **spatial grounding and key events assessment**, but limited ability to interpret temporal progression or subtle motion dynamics.

Textual-modal-only (LLMs) → good **temporal and logical reasoning**, but often hallucinate or miss fine-grained spatial interactions and key events.

Calls for **multimodal evaluation** methods for reliable synthetic feedback!

PbRL Inherent Challenges

Even if feedback from FMs reaches human-expert-level quality, PbRL still faces two inherent challenges in reward learning.

- Query Ambiguity
- Preference Credit Assignment Uncertainty

PbRL Inherent Challenges - Query Ambiguity

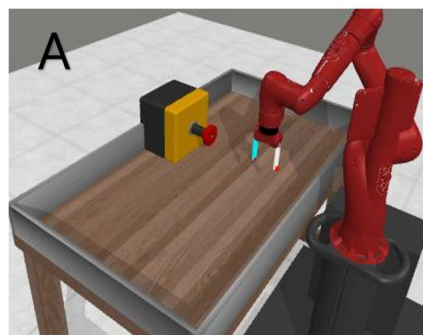
Early trajectories from randomly initialized policies are uniformly low quality
lacking meaningful task differences

→ **cannot provide informative comparisons**

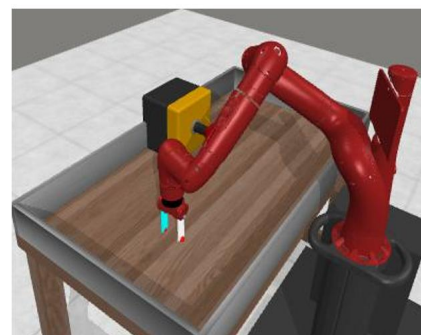
"Which is better?"



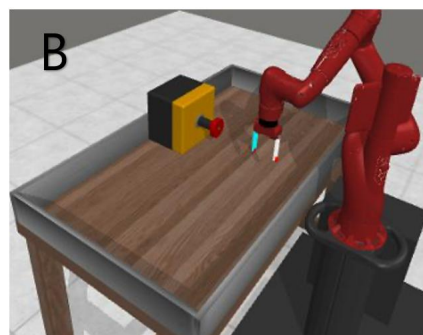
✗ Both show minimal progress
✗ No meaningful differences
→ Hard to decide



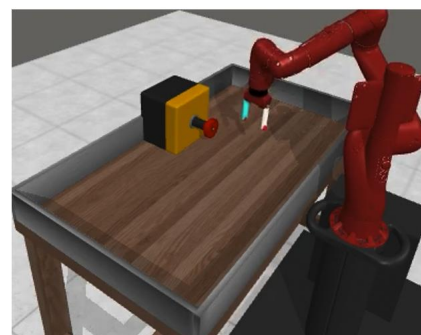
...



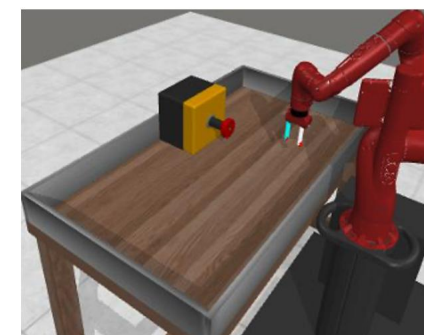
...



...



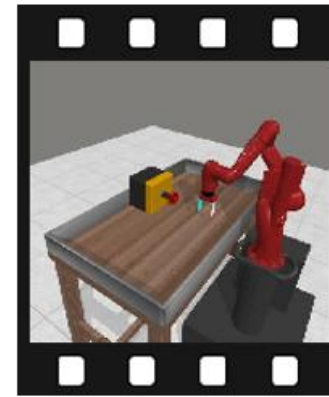
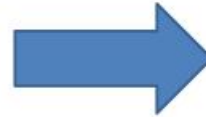
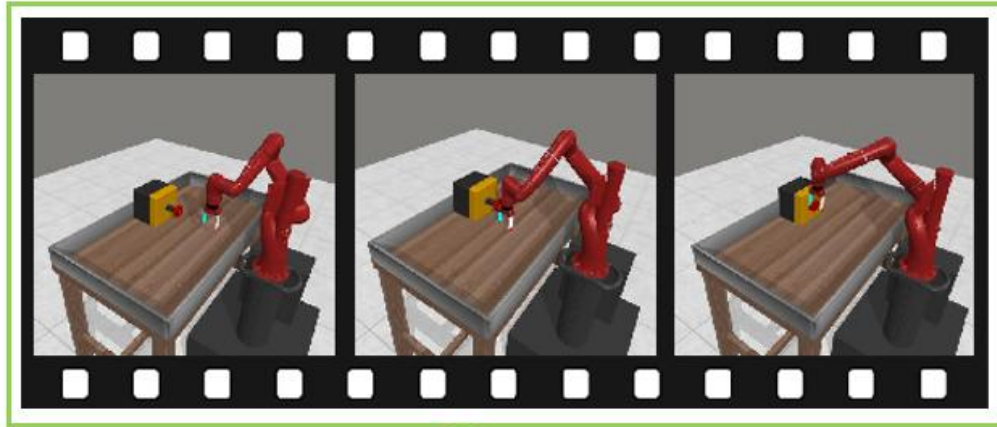
...



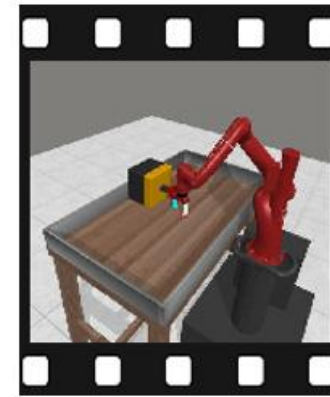
PbRL Inherent Challenges - Preference Credit Assignment Uncertainty

Preferences are given at **trajectory level**, but desired reward models operate at **state-action level**

→ **uncertainty in attributing preference credit to specific steps in the trajectory, impairing the alignment of the learned reward model**



$r_1?$



$r_2?$



$r_3?$

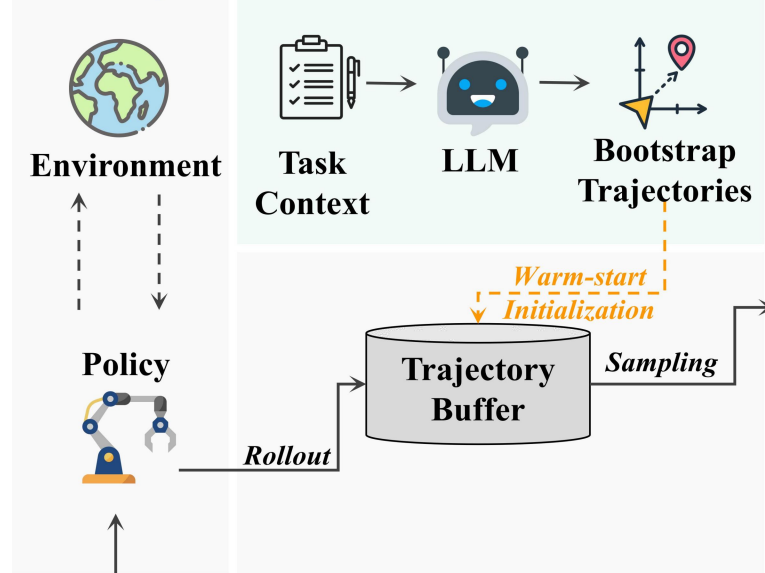
Our Goal

Enable scalable, efficient, and zero-shot PbRL with FMs

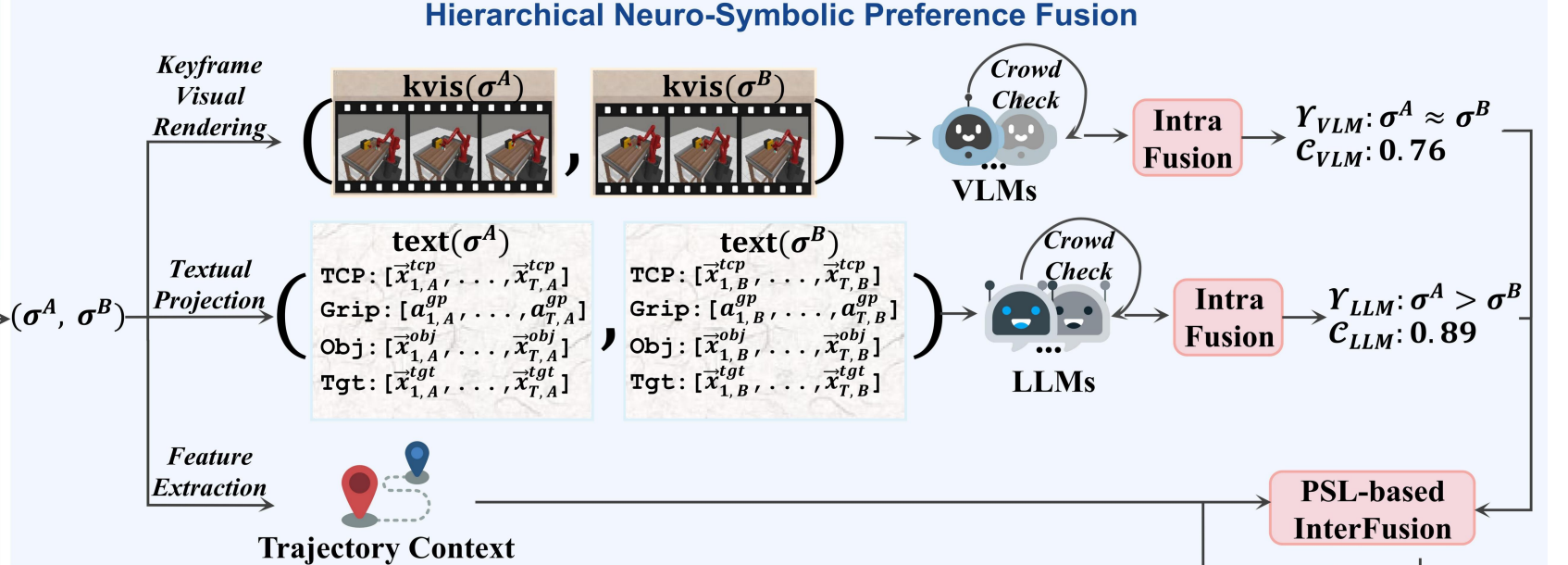
- Improve the quality and reliability of synthetic feedback
- Reduce query ambiguity
- Enhance preference credit assignment

Method - PRIMT

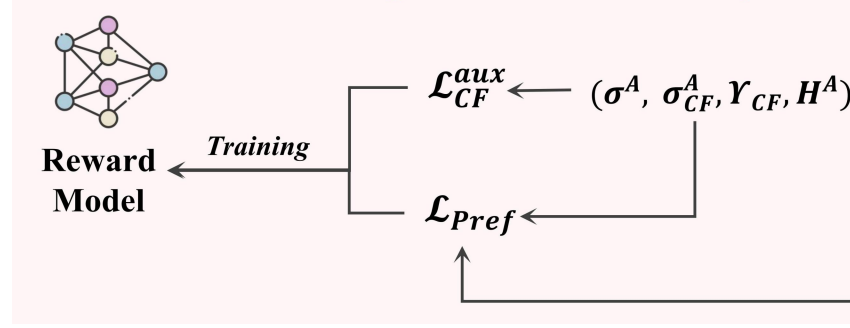
Off-Policy RL Foresight Trajectory Generation



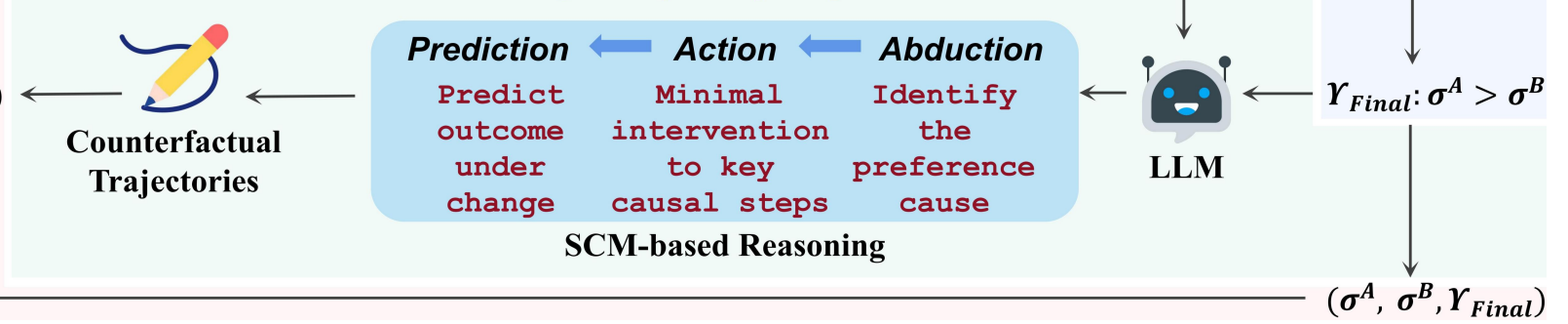
Hierarchical Neuro-Symbolic Preference Fusion



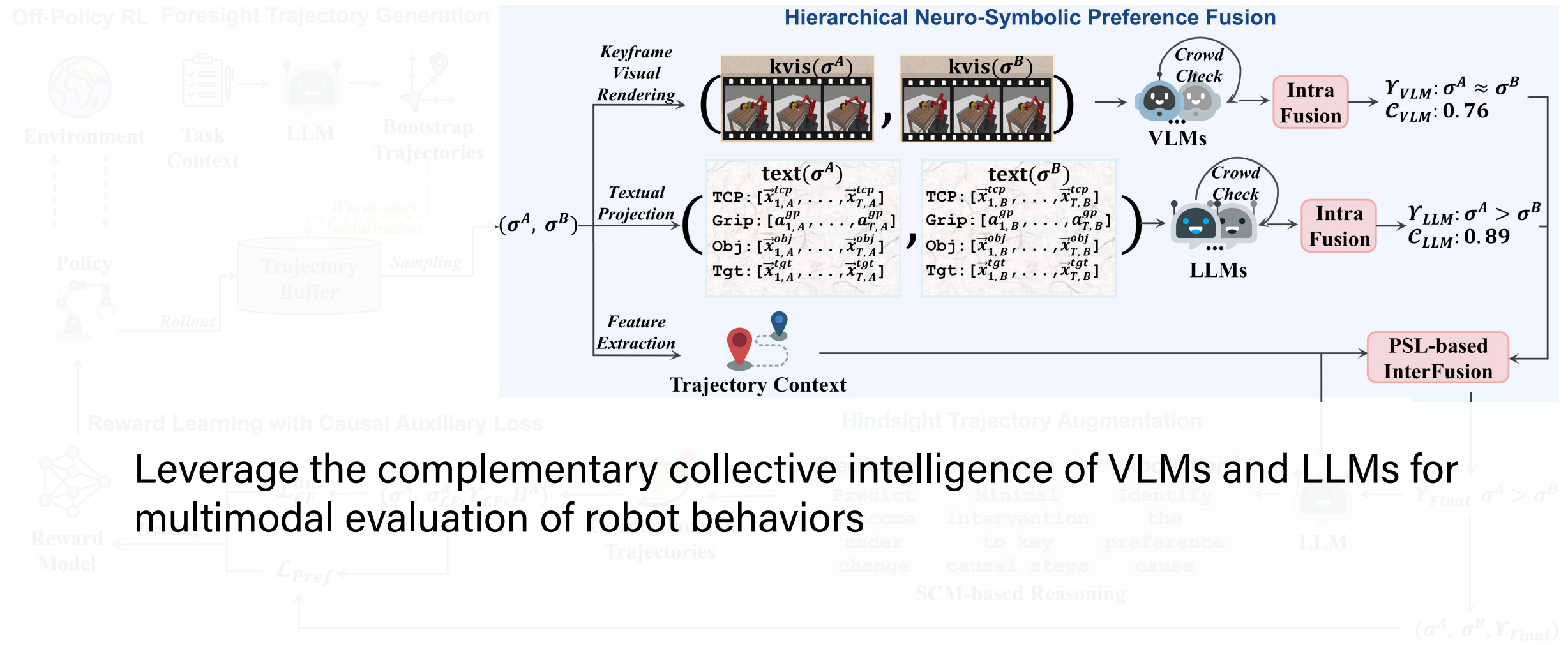
Reward Learning with Causal Auxiliary Loss



Hindsight Trajectory Augmentation

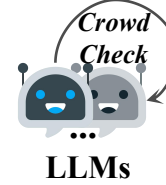
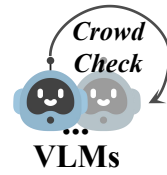


PRIMT - Multimodal Preference Generation and Fusion



PRIMT - Hierarchical Preference Fusion

Intra-modal Fusion



- i) query LLM or VLM multiple times with randomly permuted trajectory orderings
- ii) get final intra-modal labels via major voting

$$\Upsilon_M = \operatorname{argmax}_{l \in \{-1, 0, 1\}} \sum_{k=1}^K \mathbb{I}(\Upsilon_M^{(k)} = l)$$

- iii) calculate label confidence

$$\bar{c}_M = \frac{1}{N} \sum_{k=1}^K c_M^{(k)} \cdot \mathbb{I}(\Upsilon_M^{(k)} = \Upsilon_M); \quad \dot{c}_M = \frac{1}{K} \sum_{k=1}^K \mathbb{I}(\Upsilon_M^{(k)} = \Upsilon_M)$$

$$c_M = \alpha \cdot \bar{c}_M + (1 - \alpha) \cdot \dot{c}_M$$

PRIMT - Hierarchical Preference Fusion

Inter-modal Fusion

Consider **intra-modal uncertainty**, **inter-modal conflicts**, and **trajectory context** that reflects the relative difficulty of visual v.s. textual evaluation.

To model latent dependencies across these factors, we employ the Probabilistic Soft Logic (PSL) framework

PRIMT - Hierarchical Preference Fusion

Defined Rules in PSL:

i) Agreement Rule:

$$\forall \Upsilon, M : \text{IsAgree}(\Upsilon) \wedge \text{ConfHigh}(M) \rightarrow \text{FinalLabel}(\Upsilon)$$

ii) Conflict Resolution Rules:

$$\forall \Upsilon : \neg \text{IsAgree}(\Upsilon) \wedge \text{VLMLabel}(\Upsilon) \wedge \text{ConfHigh}(\text{VLM}) \wedge \text{VDHigh} \rightarrow \text{FinalLabel}(\Upsilon)$$

$$\forall \Upsilon : \neg \text{IsAgree}(\Upsilon) \wedge \text{LLMLabel}(\Upsilon) \wedge \text{ConfHigh}(\text{LLM}) \wedge \text{TDHigh} \rightarrow \text{FinalLabel}(\Upsilon)$$

iii) Indecision Rule:

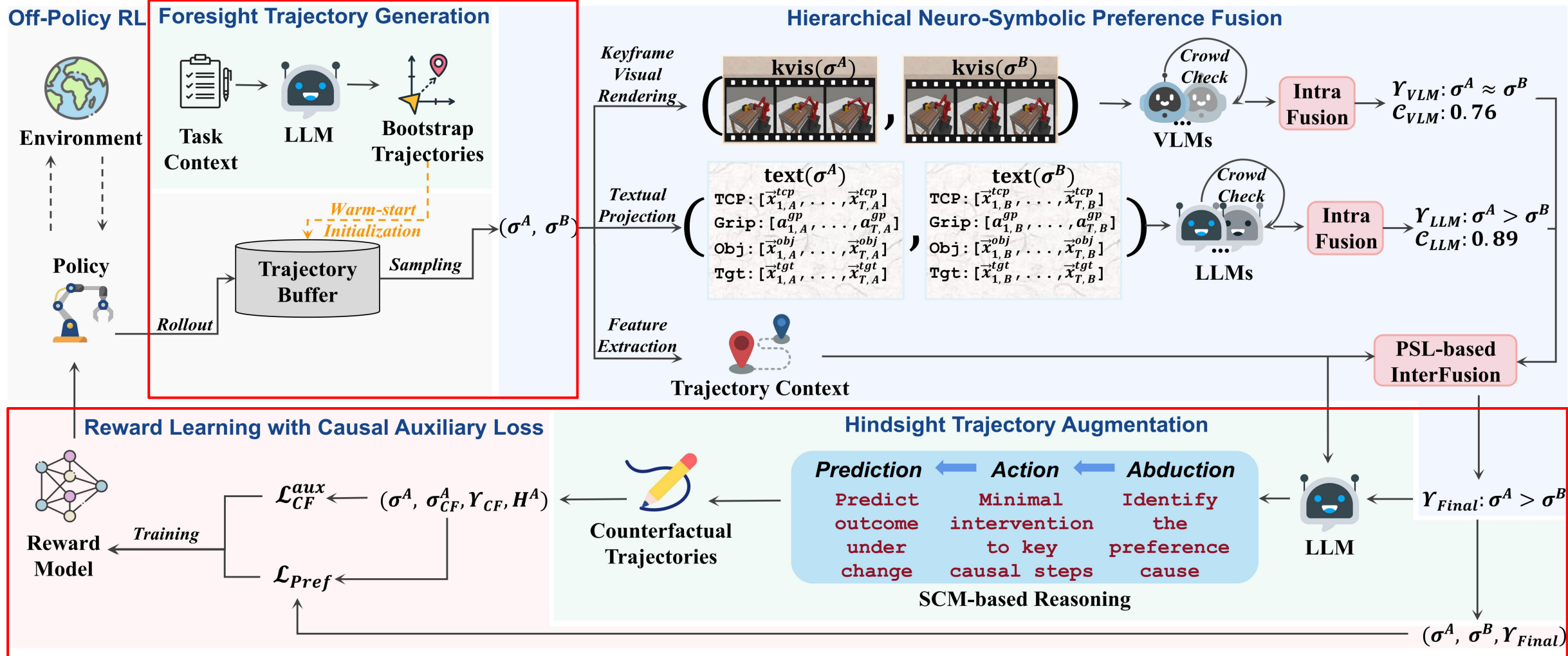
$$\neg \text{ConfHigh}(\text{VLM}) \wedge \neg \text{ConfHigh}(\text{LLM}) \rightarrow \text{FinalLabel}(-1)$$

PRIMT - Hierarchical Preference Fusion

PSL turns these rules into one or one more hinge-loss penalties and solves an online convex optimization to find the label that best fits all rules

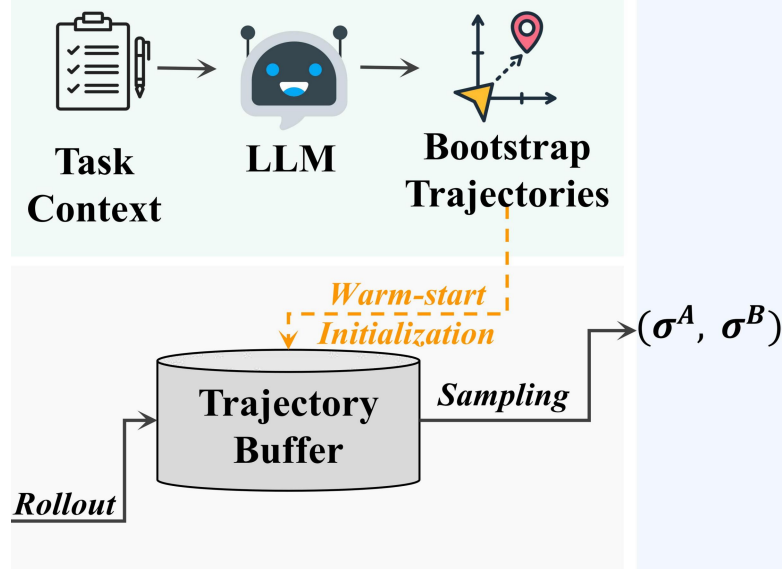
$$Y^* = \arg \min_Y \sum_{i=1}^m w_i \phi_i(Y, X) \quad \text{s.t.} \quad \sum_{\Upsilon \in \{-1, 0, 1\}} \text{FinalLabel}(\Upsilon) = 1$$

PRIMT - Bidirectional Trajectory Synthesis



PRIMT - Foresight Trajectory Generation

Foresight Trajectory Generation

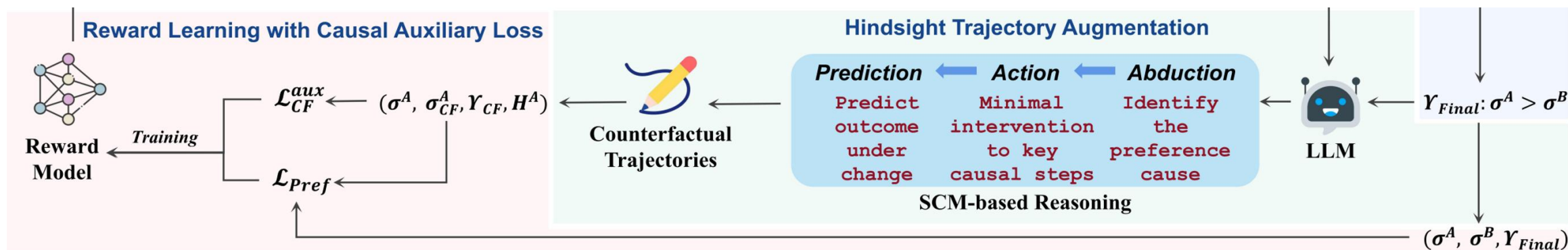


Before training, we proactively generate diverse, task-aligned trajectory samples using LLMs:

high-level multi-step plan \rightarrow *code for each step*
 \rightarrow *rollout trajectories under varied conditions*

serve as ***informative preference anchors*** rather than optimal demonstrations, making early comparisons easier when combined with query sampling strategies

PRIMT - Hindsight Trajectory Augmentation



Once we have a clear preference for a trajectory pair, we perform the **hindsight counterfactual reasoning** by asking:

"What minimal change would make the preferred trajectory less preferred?"

identify key steps in the preferred trj \rightarrow minimal edits \rightarrow create counterfactuals

PRIMT - Hindsight Trajectory Augmentation

Preferred Trajectory



Counterfactual Trajectory



Hold the cube for a longer period of time before releasing it to the target area

PRIMT - Hindsight Trajectory Augmentation

Preferred Trajectory



Counterfactual Trajectory



Add a short hesitation before grasping the cube

PRIMT - Auxiliary Causal Loss

Preferred Trajectory



Counterfactual Trajectory



Add a short hesitation before grasping the cube

$$\mathcal{L}_{\text{cf}}^{\text{aux}} = \underbrace{\sum_{t=1}^T H_t \cdot \log \left(1 + \exp \left(r_{\psi}(s_t^{\text{cf}}) - r_{\psi}(s_t^*) \right) \right)}_{\text{i) causal contrast loss}} + \underbrace{\sum_{t=1}^T (1 - H_t) \cdot \left\| r_{\psi}(s_t^*) - r_{\psi}(s_t^{\text{cf}}) \right\|_2^2}_{\text{ii) reward consistency loss}}$$

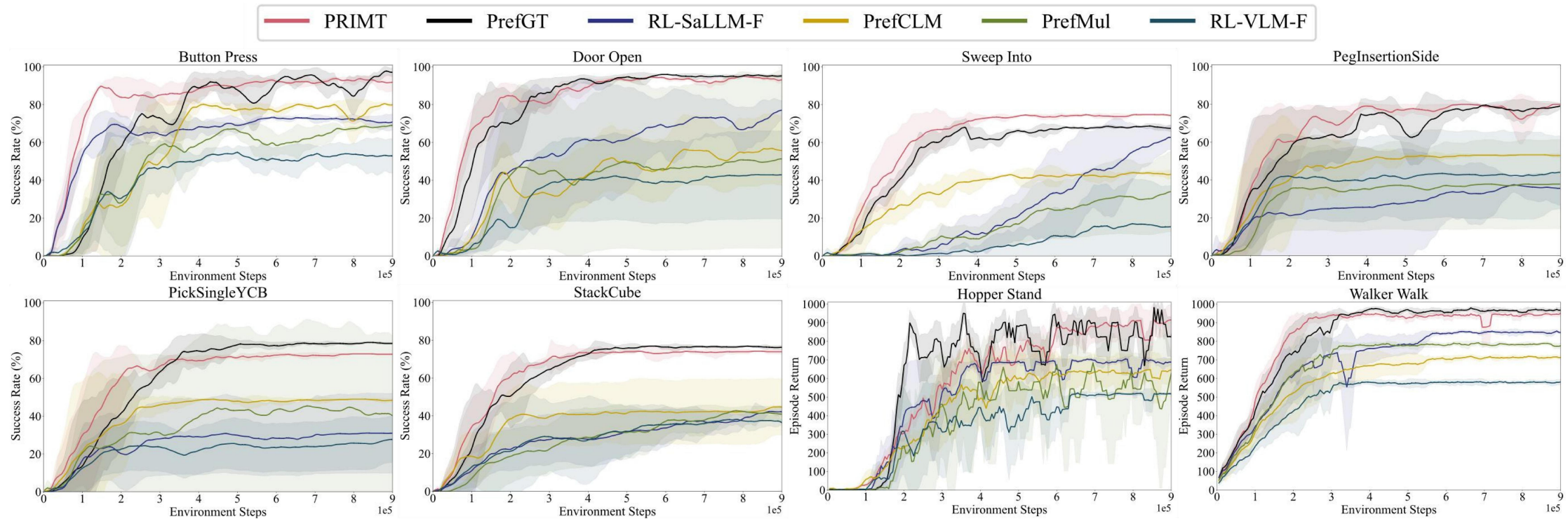
i) causal contrast loss

i) contrast rewards at edited (causal) steps

ii) reward consistency loss

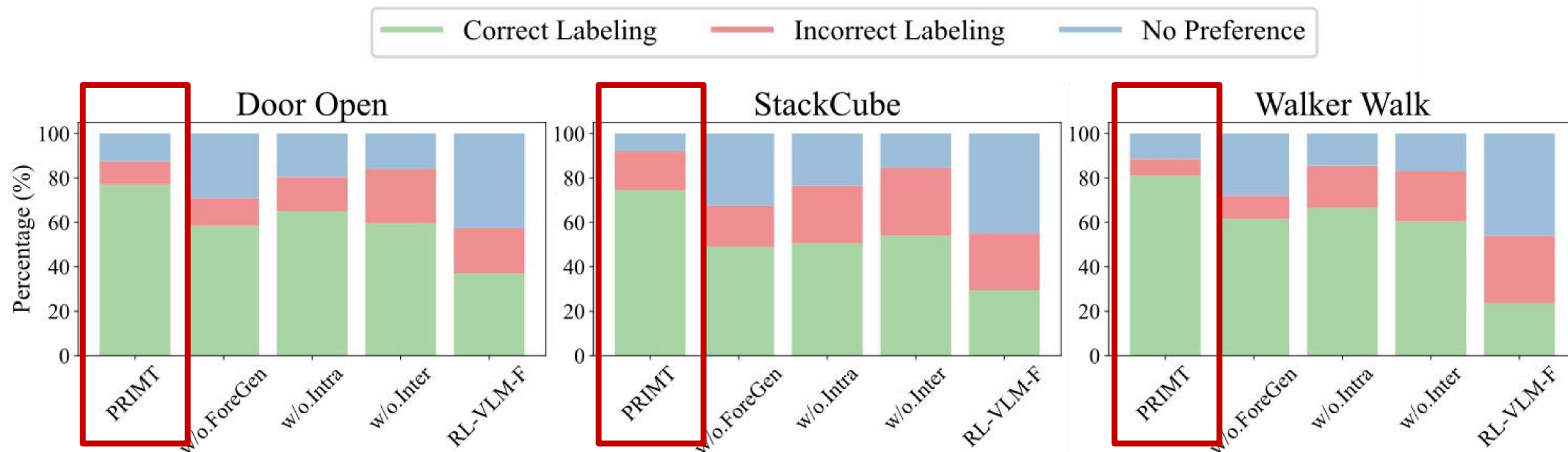
ii) enforce consistency on unedited regions

Evaluation - Task Performance on Manipulation and Locomotion Tasks



- PRIMT (red) outperforms all FM-based baselines and the naïve multimodal fusion
- PRIMT matches the oracle PrefGT (black, with ground-truth preference labels) and even surpasses it on 2 of 8 tasks
- PRIMT learns more efficiently in early training -> less query ambiguity

Evaluation - Synthetic Feedback Quality



- Improved label accuracy (green)
- Reduced indecision and query ambiguity (blue)

Evaluation - Reward Alignment

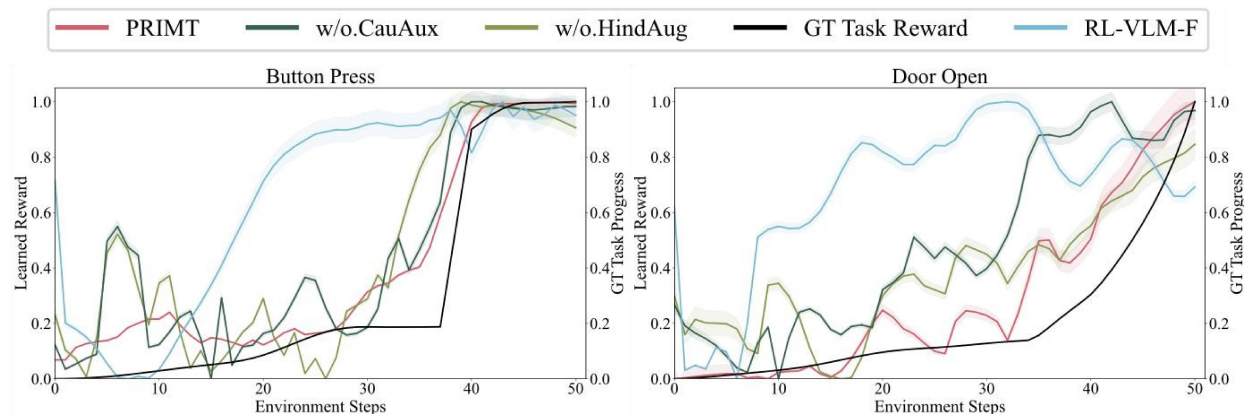


Table 5: R^2 Coefficient Analysis (Reward Alignment with Ground Truth).

Task	PRMT	w/o CauAux	w/o HindAug	RL-VLM-F
PegInsertionSide	0.56	0.28	0.23	0.37
PickSingleYCB	0.84	0.01	0.34	-0.05
StackCube	0.78	-1.31	-2.28	-1.50
ButtonPress	0.87	0.68	0.53	-0.61
DoorOpen	0.64	-1.19	0.15	-4.72
SweepInto	0.88	0.83	0.73	-0.27
WalkerWalk	0.33	0.19	0.02	-2.29

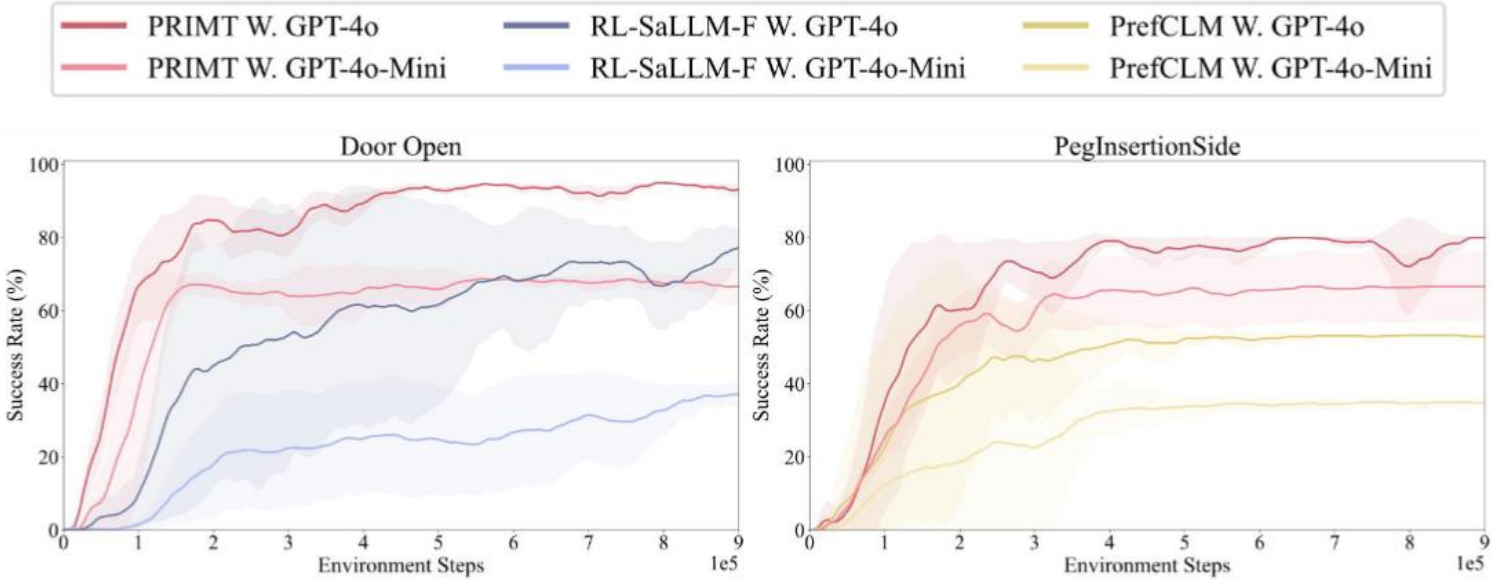
- Aligned more closely with ground-truth rewards and task progress

-> Enhanced Preference Credit Assignment

Evaluation - Cost-Performance Efficiency

Baseline	Cost	Time <i>MetaWorld / ManiSkill / DMC</i>	Performance	Efficiency [†]
vs RL-VLM-F	+43%/+39%/+47%	+58%/+43%/+69%	+95%/+117%/+68%	2.0×
vs Sa-LLM-F	+45%/+44%/+38%	+31%/+30%/+46%	+32%/+109%/+19%	1.4×
vs Human (PrefGT)	−92%/−95%/−92%	—/—/—	−1%/−%/−%	47×

- [†] *Efficiency = Average performance gain / Average resource increase*
- *Performance is measured using the final return from the learning curves presented in Figure 2*
- *Values shown are relative changes across MetaWorld / ManiSkill / DMC environments, respectively*



Improved cost–performance efficiency and increased robustness to weaker FMs

Evaluation - Real-World Experiments

PRIMT

Block Lifting



Block Stacking



PrefCLM



Thank you !

Poster Section

Dec 3rd, Wednesday

4:30 p.m. - 7:30 p.m. PST

In Room: Exhibit Hall C,D,E

Poster Location: **#2209**

Ruiqi Wang

wang5357@purdue.edu

Purdue University



Project Website



Contact Information