

KISA: A Unified Keyframe Identifier and Skill Annotator for Long-Horizon Robotics Demonstrations

Longxin Kou^{*1}, Fei Ni^{*1}, Yan Zheng¹, Jinyi Liu¹, Yifu Yuan¹, Zibin Dong¹, Jianye Hao¹

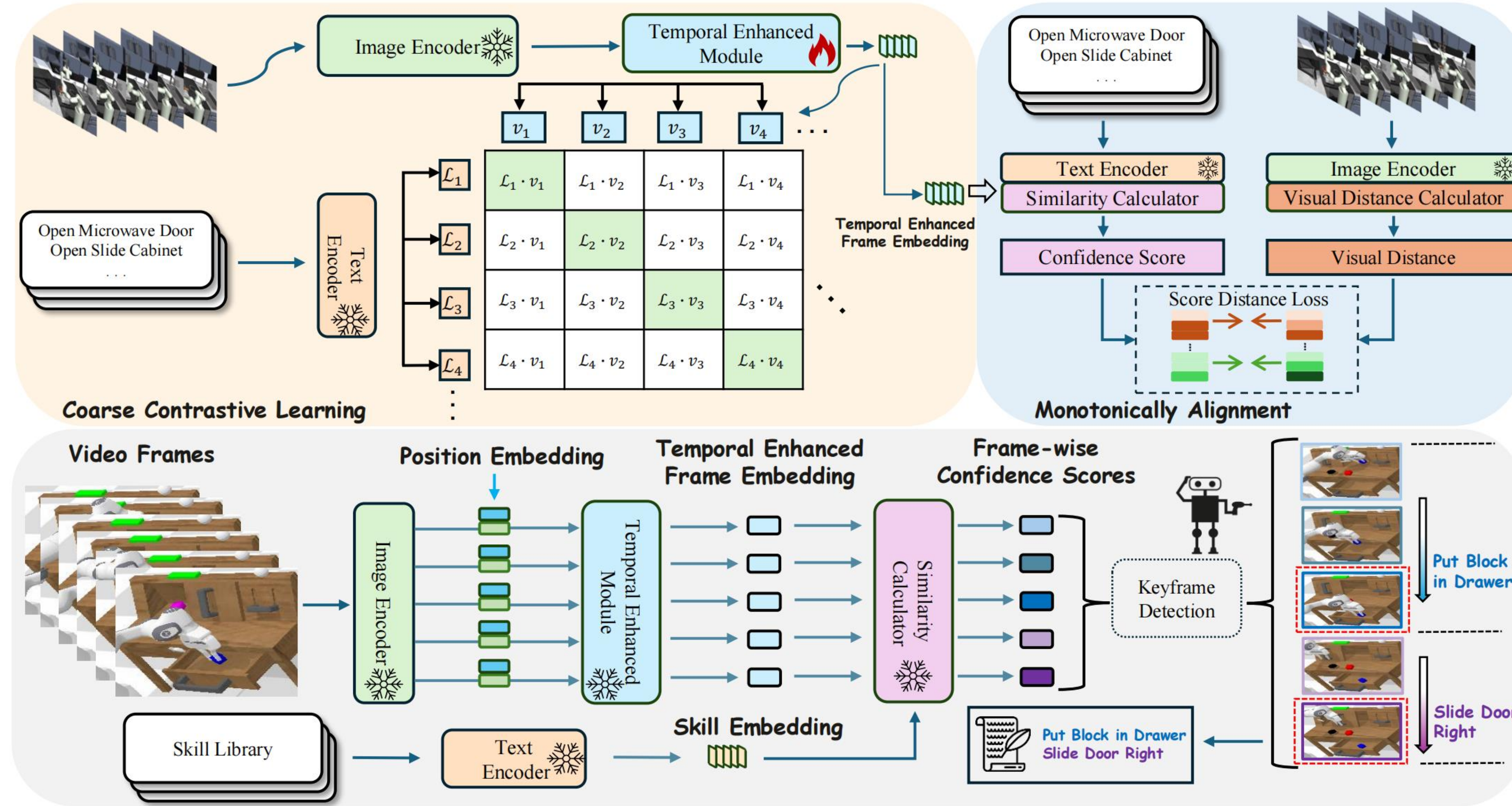
¹Tianjin University



Motivation

- **Complex robotics manipulation** tasks such as desktop tidying often span over long horizons and encapsulate multiple sub-tasks separated by keyframes. Directly learning from long-horizon demonstrations in an end-to-end manner is challenging.
- **Hierarchical policy** learning, by decomposing a complex demonstration into several shorter subtasks to facilitate the reusable skills and further enable modular skill composition for generalization. However, obtaining demonstrations with explicit keyframe boundaries and skill annotations is difficult, especially for real-world human videos.
- For this, we study the following open question - can we develop a framework that enables **automatic**, **scalable**, and **semantically meaningful** keyframe identification and skill annotation from unlabeled demonstrations?

Our Approach



- **Temporal enhancement module**
 - Relying solely on frame-level visual representations would induce training confusion for aligning them to distinct skills. So we propose a simple yet effective temporal-enhanced module on top of pretrained visual representations.
- **History-aware Contrastive Training**
 - we design coarse history-aware contrastive learning via constructing hard negative samples with mismatched historical contexts and incorrect skills.
- **Fine-grained Monotonic Alignment**
 - we additionally fine-grained monotonic alignment to encourage the capture of skill-aware progress within the sub-task, and prevent representation collapse to highly similarity within the same skill.

$$\mathcal{L}_{\text{contrastive}} = -\log \frac{e^{C(o^+, h^+, \ell^+)}}{\sum_{j=1}^k e^{C(o^+, h^+, \ell_j^+)}} - \log \frac{e^{C(o^+, h^+, \ell^+)}}{\sum_{z=1}^k e^{C(o^+, h_z^+, \ell^+)}} - \log \frac{e^{C(o^+, h^+, \ell^+)}}{\sum_{w=1}^k e^{C(\{(o^+, h^+)^{w, ev}, l^+\})}}$$

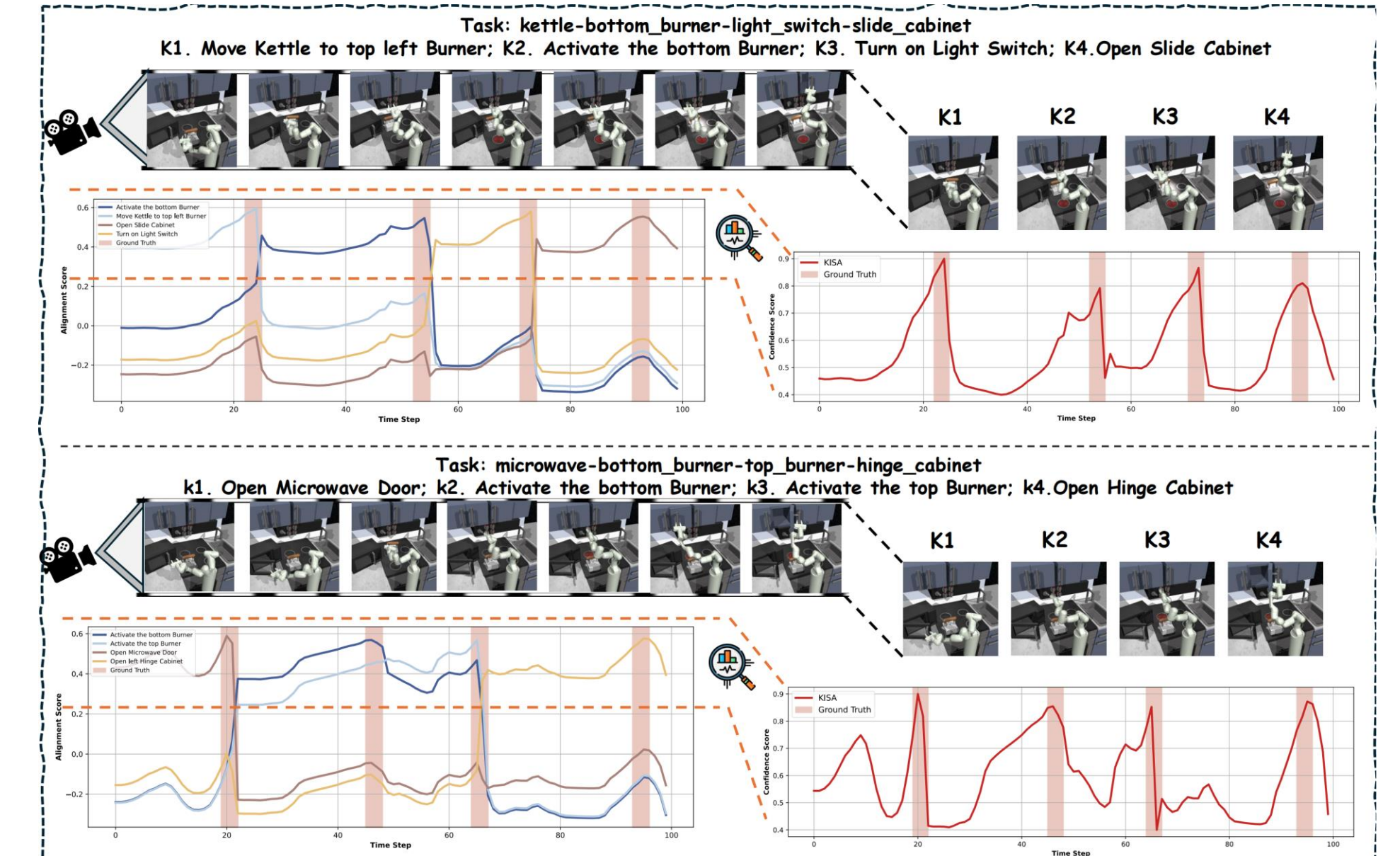
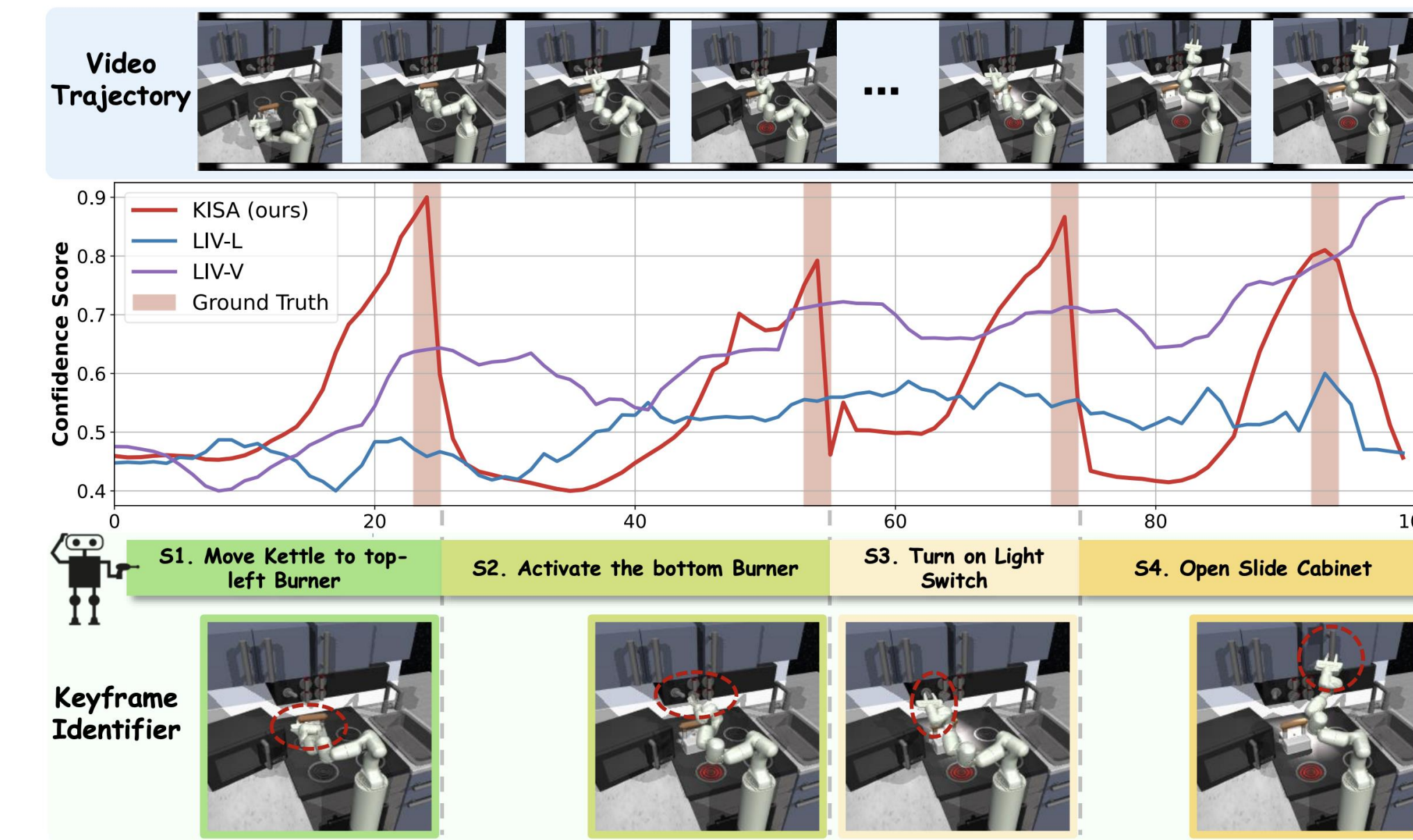
Incorrect Skill Mismatches Disjoint Frame-History Compositions Semantic Reversals via Video Inversion

Experiments

- We conduct experiments on various benchmarks to answer the following questions: 1) Can KISA achieve better **accuracy** and **interpretable** skill alignment compared to other competitive baselines? 2) Does KISA exhibit robust zero-shot generalization across objects, compositional tasks, and cross-embodiments? 3) Is KISA a flexible framework for incorporating pretrained robotics representations?

The Accuracy of Keyframes and Skills Annotation

Model	Maniskill			CALVIN			FrankaKitchen		
	Number Error↓	F1 Score↑	MAE↓	Number Error↓	F1 Score↑	MAE↓	Number Error↓	F1 Score↑	MAE↓
VideoRLCS	10.1 ± 2.1	15.2 ± 0.3%	34.4 ± 0.4	9.5 ± 2.4	15.4 ± 0.5%	54.8 ± 1.0	0.6 ± 0.0	5.5 ± 0.8%	39.3 ± 0.7
KTS	5.1 ± 0.0	15.7 ± 4.0%	24.2 ± 6.8	0.9 ± 0.6	20.2 ± 0.7%	50.9 ± 7.4	0.5 ± 0.2	13.8 ± 3.4%	35.6 ± 6.9
R3M	5.0 ± 0.2	17.1 ± 0.8%	38.2 ± 3.0	5.4 ± 0.2	21.1 ± 1.3%	63.2 ± 1.1	1.0 ± 0.1	53.7 ± 0.8%	30.4 ± 0.1
VIP	4.0 ± 0.2	31.7 ± 2.5%	24.8 ± 1.4	4.6 ± 0.2	24.3 ± 1.6%	63.4 ± 1.2	0.9 ± 0.1	57.2 ± 1.1%	31.4 ± 0.1
LIV	3.9 ± 0.4	30.3 ± 2.2%	23.9 ± 1.4	5.6 ± 0.1	25.9 ± 1.1%	61.7 ± 1.5	1.4 ± 0.0	64.2 ± 0.8%	30.7 ± 0.1
UVD	0.7 ± 0.1	40.2 ± 1.1%	20.3 ± 0.1	0.8 ± 0.1	36.9 ± 0.5%	40.6 ± 1.7	0.6 ± 0.1	64.8 ± 2.4%	31.1 ± 0.2
KISA	0.0 ± 0.0	99.7 ± 0.2%	0.2 ± 0.1	0.1 ± 0.1	85.2 ± 0.9%	11.2 ± 2.4	0.0 ± 0.0	98.7 ± 0.6%	0.4 ± 0.0



The Flexibility for Pre-trained Representations

Methodology	Maniskill2	CALVIN	FrankaKitchen
KISA-R3M	71.8 ± 3.9%	53.6 ± 1.1%	88.9 ± 0.6%
- w/o monotonic align	63.0 ± 3.2% ↓	50.1 ± 0.8% ↓	81.5 ± 1.0% ↓
- w/o historical contrastive	41.3 ± 2.2% ↓	45.7 ± 2.3% ↓	76.6 ± 0.5% ↓
- w/o temporal enhance	23.1 ± 0.5% ↓	24.0 ± 2.2% ↓	21.9 ± 0.8% ↓
KISA-VIP	99.6 ± 0.1%	70.9 ± 2.7%	96.4 ± 0.3%
- w/o monotonic align	88.9 ± 0.5% ↓	64.0 ± 0.4% ↓	90.1 ± 0.3% ↓
- w/o historical contrastive	58.9 ± 1.5% ↓	53.1 ± 1.5% ↓	81.8 ± 0.9% ↓
- w/o temporal enhance	24.0 ± 0.3% ↓	23.4 ± 0.8% ↓	21.4 ± 0.8% ↓
KISA-LIV	99.2 ± 0.1%	94.7 ± 1.1%	96.1 ± 0.2%
- w/o monotonic align	90.2 ± 0.3% ↓	82.1 ± 0.4% ↓	89.1 ± 0.7% ↓
- w/o historical contrastive	59.1 ± 1.9% ↓	58.1 ± 1.3% ↓	73.7 ± 0.4% ↓
- w/o temporal enhance	22.0 ± 1.6% ↓	25.0 ± 1.6% ↓	19.0 ± 0.3% ↓

The Zero-shot Generalization Ability.

Model	Maniskill12 (L1)		CALVIN (L2)		RealKitchen (L3)	
	F1 Score↑	MAE↓	F1 Score↑	MAE↓	F1 Score↑	MAE↓
KTS	12.3 ± 2.7%	25.2 ± 5.3	20.2 ± 0.7%	54.9 ± 7.4	11.9 ± 5.4%	44.9 ± 21.4
R3M	16.8 ± 1.4%	38.8 ± 2.5	20.9 ± 1.1%	63.4 ± 1.2	20.7 ± 0.5%	44.9 ± 21.4
VIP	30.6 ± 1.8%	23.3 ± 1.8	23.6 ± 1.9%	63.2 ± 1.1	20.5 ± 17.8%	34.8 ± 11.0
LIV	30.2 ± 1.7%	23.6 ± 1.0	25.4 ± 1.3%	63.2 ± 1.1	21.7 ± 20.9%	44.9 ± 21.4
UVD	39.2 ± 1.1%	21.5 ± 0.2	36.9 ± 0.5%	40.6 ± 1.7	26.3 ± 11.9%	30.2 ± 6.2
KISA	80.7 ± 0.9%	6.4 ± 0.7	89.4 ± 1.8%	14.2 ± 0.9	40.7 ± 14.8%	27.8 ± 5.0

Visualization of Keyframe Identification

