



山东大学  
SHANDONG UNIVERSITY

# TSTM: Temporal Segmentation for Task-relevant Mask in Visual Reinforcement Learning Generalization

Weicheng Du<sup>1</sup>, Wenjia Meng<sup>1\*</sup>, Zhengzhe Zhang<sup>1</sup>, Yilong Yin<sup>1</sup>, Xiankai Lu<sup>1,2</sup>

<sup>1</sup>School of Software, Shandong University    <sup>2</sup>School of Mathematics and Computer Science, Quanzhou Normal University

Time: Jun 6th Afternoon    Poster Session: 4

Session Order: 613    Poster #30654

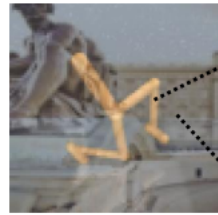
By: Weicheng Du

CVPR  
JUNE 3-7, 2026



DENVER  
COLORADO

- **Motivation**
- **TSTM Framework**
- **Temporal Segmentation Network**
- **Optimization Objectives**
- **Experiments**



Agent with **yellow color**

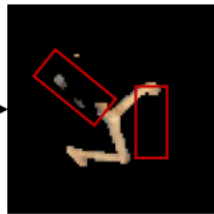
Background with **yellow color**

Similar color makes segmentation difficult.

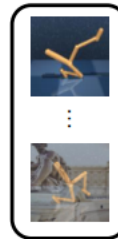
(a) A task-relevant mask segmentation example.



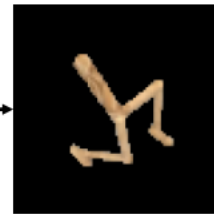
Segmentation Network



(b) Segmentation w/o temporal.



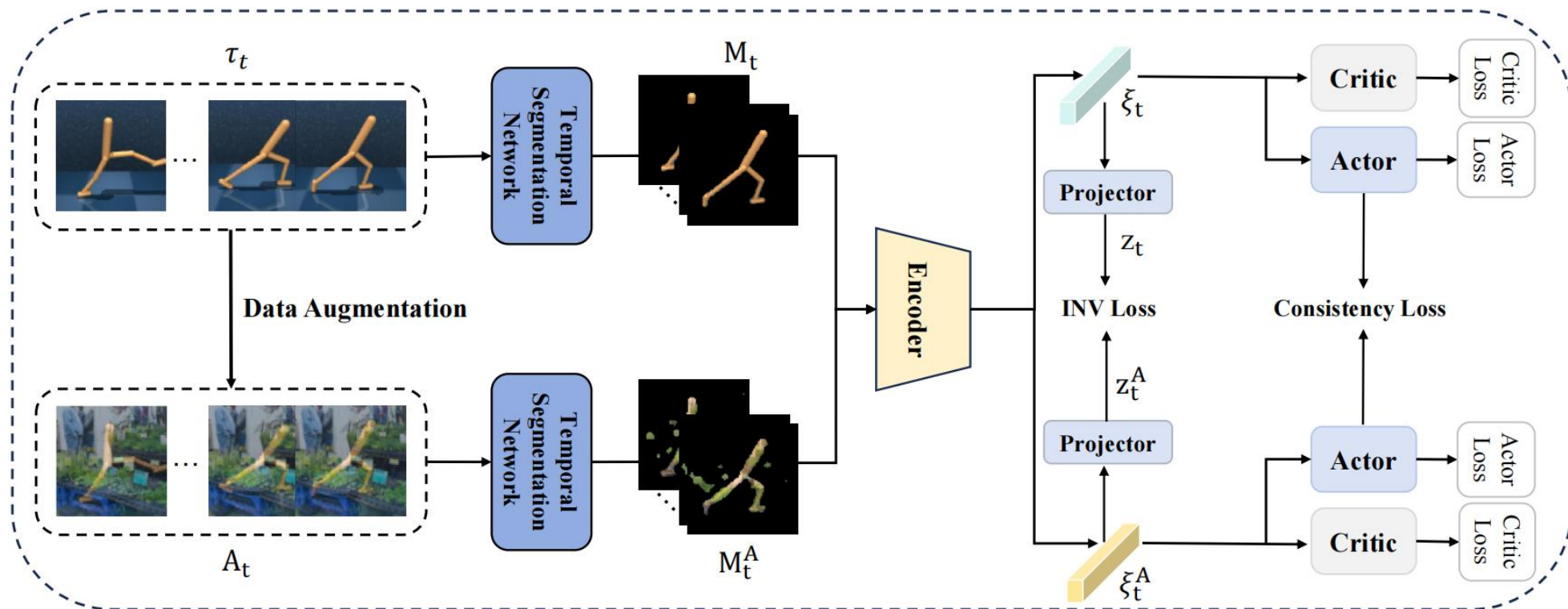
Temporal Segmentation Network



(c) Segmentation w/ temporal.

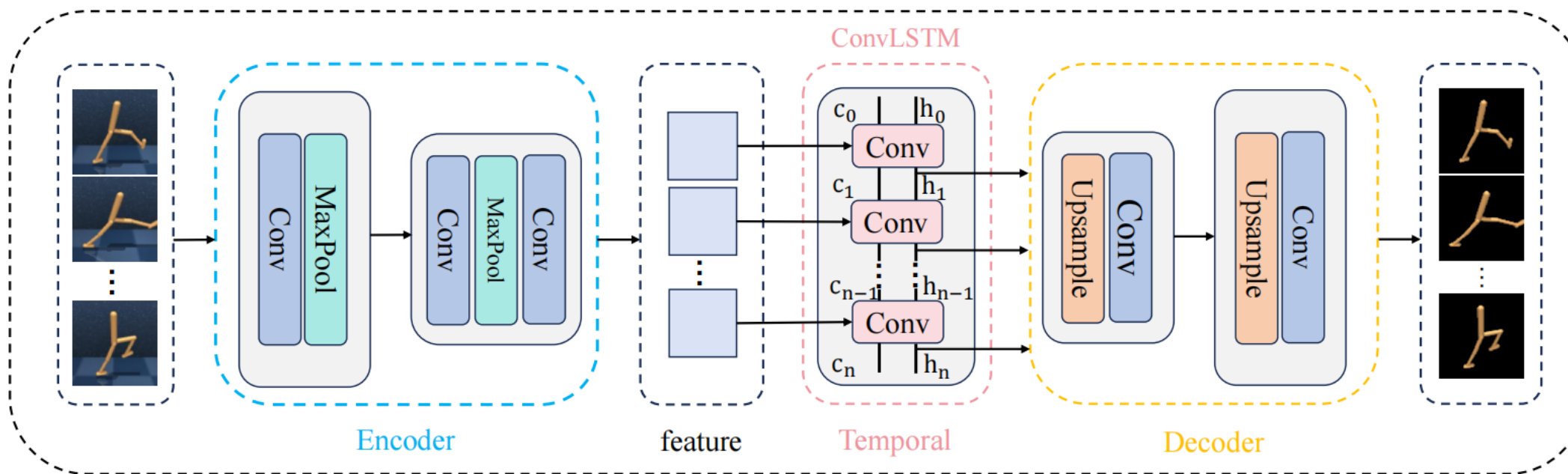
1. **Single-frame masking is unreliable in visual RL.**
2. **Ambiguous backgrounds can confuse task-relevant segmentation.**
3. **Temporal cues help distinguish agents from distracting visual content.**

**Goal: learn more reliable task-relevant masks from observation sequences.**



## key components:

- Temporal segmentation
- Invariant learning
- Policy consistency



## key components:

- **Encoder:** extracts spatial features from sequential observations
- **ConvLSTM:** captures temporal dependencies across frames
- **Decoder:** predicts more reliable task-relevant masks

**Temporal modeling improves segmentation reliability under visual ambiguity.**

## Teacher / Student Segmentation

- We first train a teacher segmentation network and distill it into a lightweight student for efficient inference.

$$\mathcal{L}_{\text{BCE}} = -\frac{1}{H \times W} (y_i \odot \log u_i + (1 - y_i) \odot \log(1 - u_i))$$
$$\mathcal{L}_{\text{Dice}} = 1 - \frac{2\langle y_i, u_i \rangle + \epsilon}{\|y_i\|_1 + \|u_i\|_1 + \epsilon}$$

## Invariant Representation Learning

- Align masked representations of original and augmented observation sequences.

$$\mathcal{L}_{\text{mse}} = \|z_t - z_t^A\|_2^2$$
$$\mathcal{L}_{\text{var}} = \frac{1}{d} \sum_{i=1}^a \left[ \max(0, \Gamma - \sqrt{\text{Var}(z_{t,i}) + \eta}) \right. \\ \left. + \max(0, \Gamma - \sqrt{\text{Var}(z_{t,i}^A) + \eta}) \right]$$
$$\mathcal{L}_{\text{cov}} = \frac{1}{d} \sum_{i \neq j} [C(z_t)_{ij}^2 + C(z_t^A)_{ij}^2]$$
$$\mathcal{L}_{\text{INV}} = \lambda \mathcal{L}_{\text{mse}} + \mu \mathcal{L}_{\text{var}} + \rho \mathcal{L}_{\text{cov}}$$

## Actor Update with Policy Consistency

- The actor is optimized with standard SAC objective plus a policy consistency term between original and augmented masked observations.

$$\mathcal{L}_{\pi} = \mathbb{E}_{s_t} [\alpha \log \pi_{\psi}(a_t | s_t) - Q_{\theta}(s_t, a_t)]$$
$$\mathcal{L}_{\text{PC}} = \mathbb{E}_{\xi_t \sim \mathcal{B}} \left[ D_{\text{KL}} \left( \text{sg}(\pi_{\psi}(\cdot | \xi_t)) \parallel \pi_{\psi}(\cdot | \xi_t^A) \right) \right]$$
$$\mathcal{L}_A = \mathcal{L}_{\pi} + \varsigma \mathcal{L}_{\text{PC}}$$

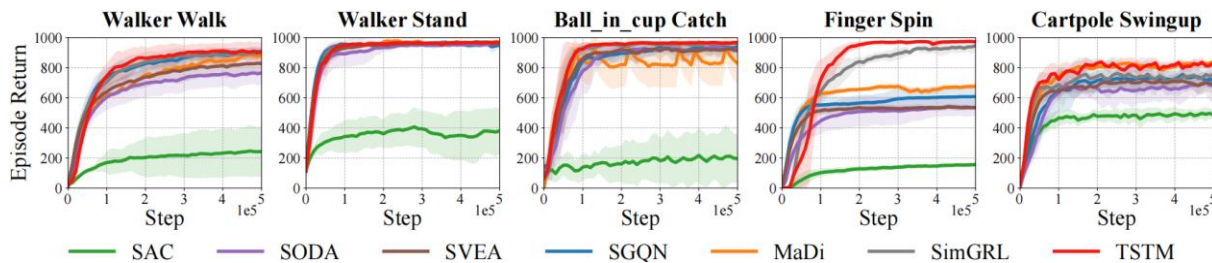
**INV loss acts in representation space, while policy consistency is added to the actor objective.**

## TSTM achieves the best average return on DMC-GB benchmark.

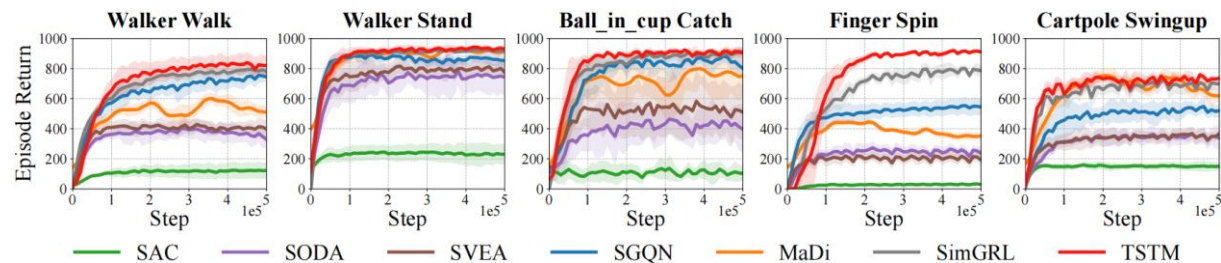
	DMC-GB [21]	SAC [74]	SODA [70]	SVEA [21]	SGQN [25]	MaDi [26]	SimGRL [75]	TSTM (Ours)
Video easy	Walker Walk	245±165	771±66	828±66	910±24	895±24	910±21	<b>912±42</b>
	Walker Stand	389±131	965±7	966±5	955±12	967±3	<b>973±4</b>	<u>968±3</u>
	Ball_in_cup Catch	192±157	939±10	908±55	950±24	807±144	<u>964±7</u>	<b>969±2</b>
	Finger Spin	157±8	535±53	537±11	610±61	679±17	<u>957±16</u>	<b>971±12</b>
	Cartpole Swingup	474±26	678±120	684±74	717±35	<u>848±6</u>	<u>775±60</u>	<b>851±19</b>
average	291	778	785	828	839	<u>916</u>	<b>934</b>	
Video hard	Walker Walk	122±47	312±32	385±63	739±21	504±33	<u>773±31</u>	<b>821±36</b>
	Walker Stand	231±57	736±132	747±43	851±24	920±14	<b>932±17</b>	<u>923±25</u>
	Ball_in_cup Catch	101±37	381±163	498±174	782±57	758±135	<u>902±19</u>	<b>903±27</b>
	Finger Spin	26±6	222±48	174±39	541±53	358±25	<u>784±10</u>	<b>906±10</b>
	Cartpole Swingup	153±22	328±79	381±28	526±42	619±24	<u>696±39</u>	<b>741±22</b>
average	127	396	437	688	632	<u>817</u>	<b>859</b>	

TSTM consistently outperforms prior visual RL baselines on most tasks.

## Learning curves on DMC-GB under the video easy setting



## Learning curves on DMC-GB under the video hard setting

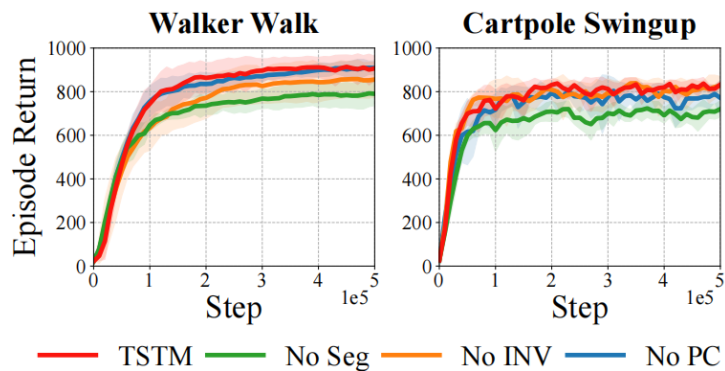


## ► Ablation Study

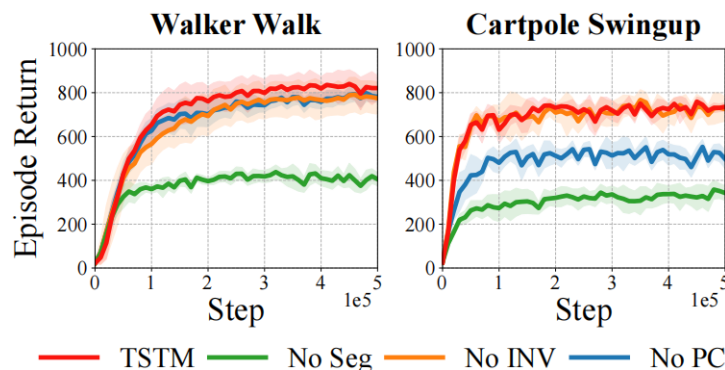
### ● Ablation results on DMC-GB.

	Method	Walker Walk	Cartpole Swingup
Video Easy	TSTM	<b>912±42</b>	<b>851±19</b>
	No Seg	788±47	733±42
	No INV	859±78	834±29
	No PC	908±34	752±120
Video Hard	TSTM	<b>821±36</b>	<b>741±22</b>
	No Seg	391±25	333±33
	No INV	770±73	720±44
	No PC	777±22	469±44

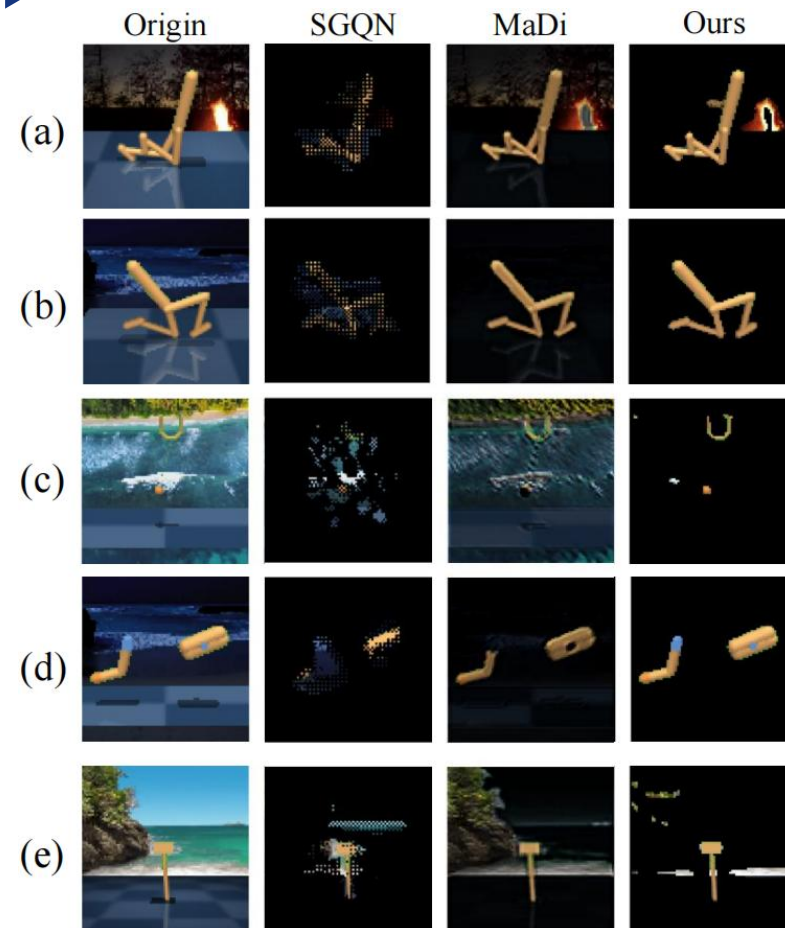
### ● Ablation curves on DMC-GB (video easy).



### ● Ablation curves on DMC-GB (video hard).



## ► Visualization





山东大学  
SHANDONG UNIVERSITY

# Thanks !

---

Time: Jun 6th Afternoon    Poster Session: 4

Session Order: 613            Poster #30654

**CVPR**  
JUNE 3-7, 2026



**DENVER**  
**COLORADO**