

## Pre-training Contextualized World Models with In-the-wild Videos for Reinforcement Learning

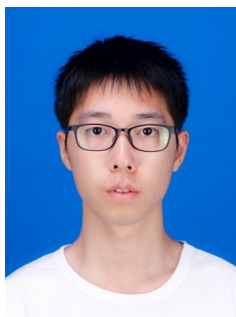
Code Available: <https://github.com/thuml/ContextWM>

**Jialong Wu\***, **Haoyu Ma\***, **Chaoyi Deng**, **Mingsheng Long**✉

School of Software, BNRist, Tsinghua University, China

wujialong0229@gmail.com, {mhy22, dengcy23}@mails.tsinghua.edu.cn

mingsheng@tsinghua.edu.cn



Jialong Wu



Haoyu Ma



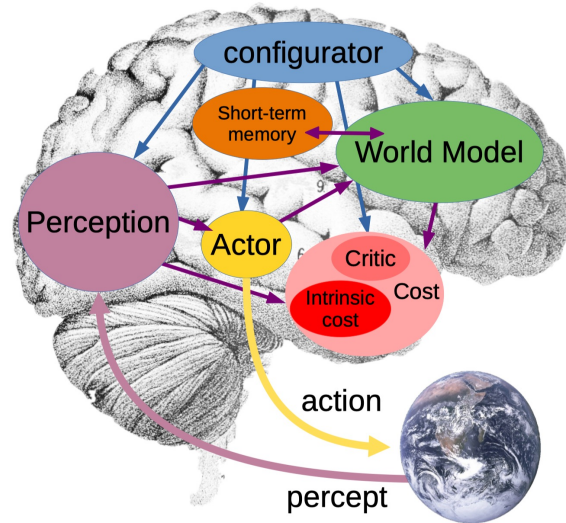
Chaoyi Deng



Mingsheng Long

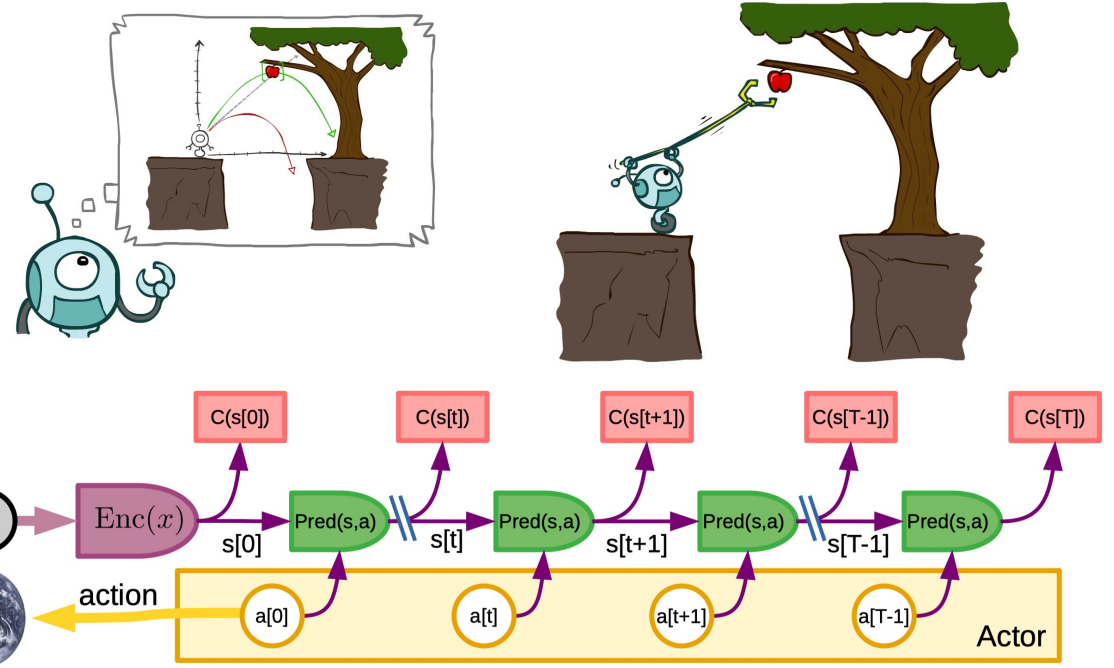


# World Models



## World Models:

internal models of how the world works



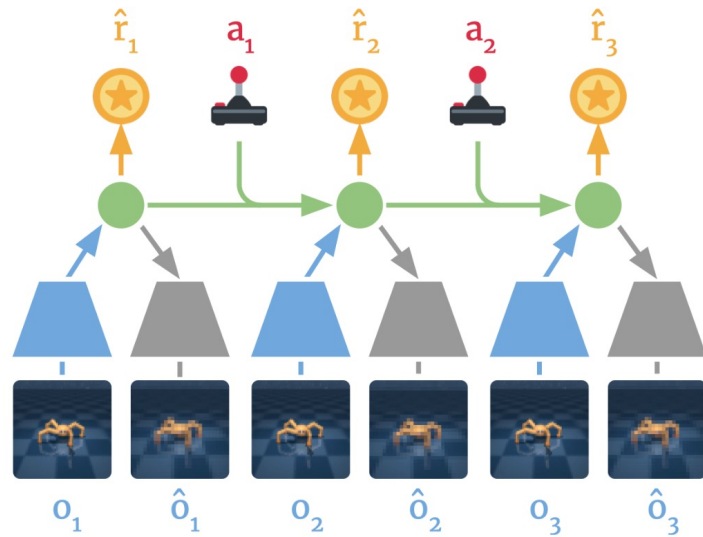
## Model-based Agents:

Act through an optimization procedure (**planning**) running the **world model**.

Yann LeCun. A path towards autonomous machine intelligence. 2022.

Dan Klein and Pieter Abbeel. Introduction to Artificial Intelligence.

# Dreamer: An Instantiation of World Models



Representation model:  $z_t \sim q_\theta(z_t | z_{t-1}, a_{t-1}, o_t)$

Transition model:  $\hat{z}_t \sim p_\theta(\hat{z}_t | z_{t-1}, a_{t-1})$

Image decoder:  $\hat{o}_t \sim p_\theta(\hat{o}_t | z_t)$

Reward predictor:  $\hat{r}_t \sim p_\theta(\hat{r}_t | z_t)$

Model Learning  
with **Sequential**  
**Variational Inference**

$$\mathcal{L}(\theta) \doteq \mathbb{E}_{q_\theta(z_{1:T} | a_{1:T}, o_{1:T})} \left[ \sum_{t=1}^T \left( \underbrace{-\ln p_\theta(o_t | z_t) - \ln p_\theta(r_t | z_t)}_{\text{reconstruction loss}} + \underbrace{\beta_z \text{KL} [q_\theta(z_t | z_{t-1}, a_{t-1}, o_t) \parallel p_\theta(\hat{z}_t | z_{t-1}, a_{t-1})]}_{\text{KL loss between prior and posterior}} \right) \right].$$

Behavior Learning: Purely on **imaginary latent trajectories**

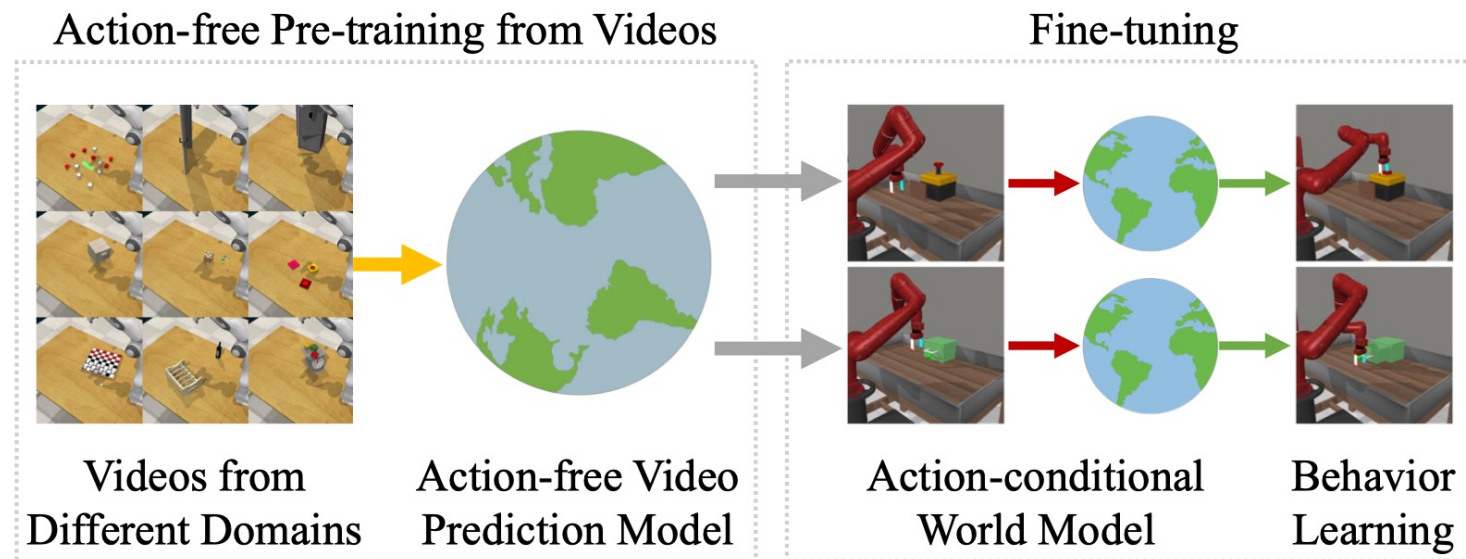
Hafner, Danijar, et al. Dream to control: Learning behaviors by latent imagination. ICLR 2020.

Hafner, Danijar, et al. Mastering atari with discrete world models. ICLR 2021.

# APV: Action-free Pre-training from Videos

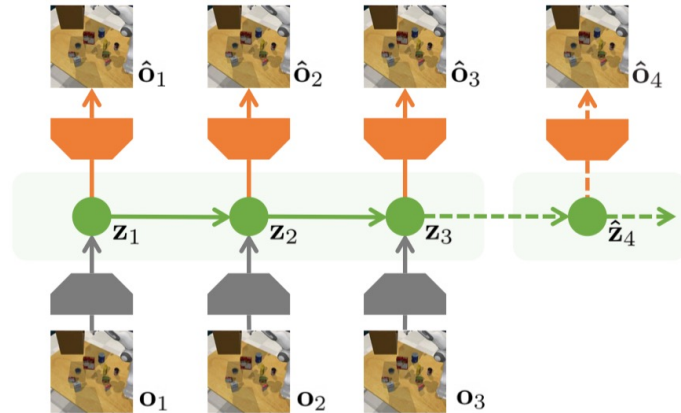
## How to represent and acquire prior knowledge for RL?

Learning **representations** useful for understanding the **dynamics**  
via **generative pretraining on videos**





# APV: Action-free Pre-training from Videos



## Stacked Latent Prediction Model

Action-free

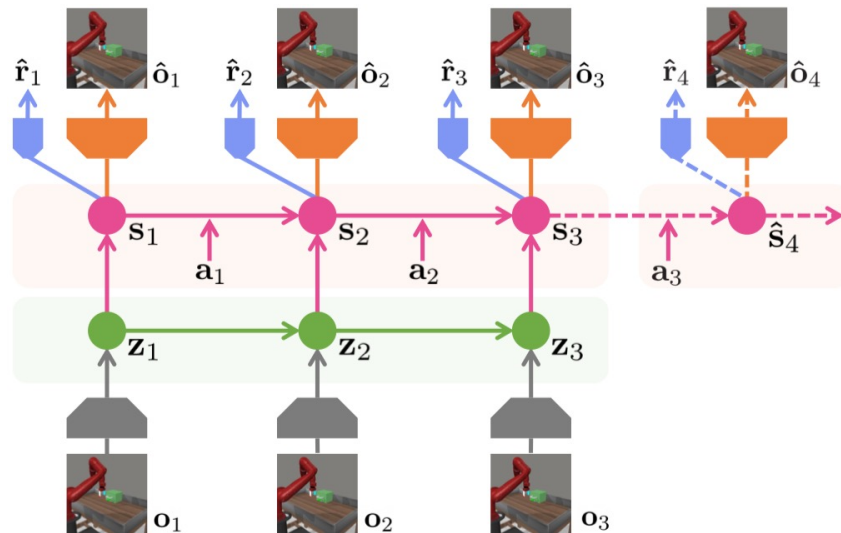
Representation:  $q_{\theta}(z_t | z_{t-1}, o_t)$   
 Transition:  $p_{\theta}(\hat{z}_t | z_{t-1})$

Image decoder:  $p_{\theta}(\hat{o}_t | s_t)$

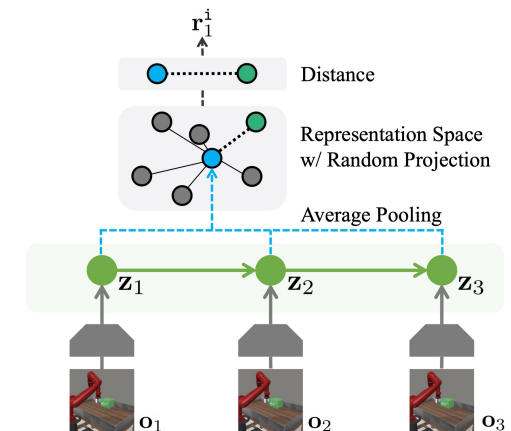
Action-conditional

Representation:  $q_{\phi}(s_t | s_{t-1}, a_{t-1}, z_t)$   
 Transition:  $p_{\phi}(\hat{s}_t | s_{t-1}, a_{t-1})$

Reward predictor:  $p_{\theta}(\hat{r}_t | z_t)$

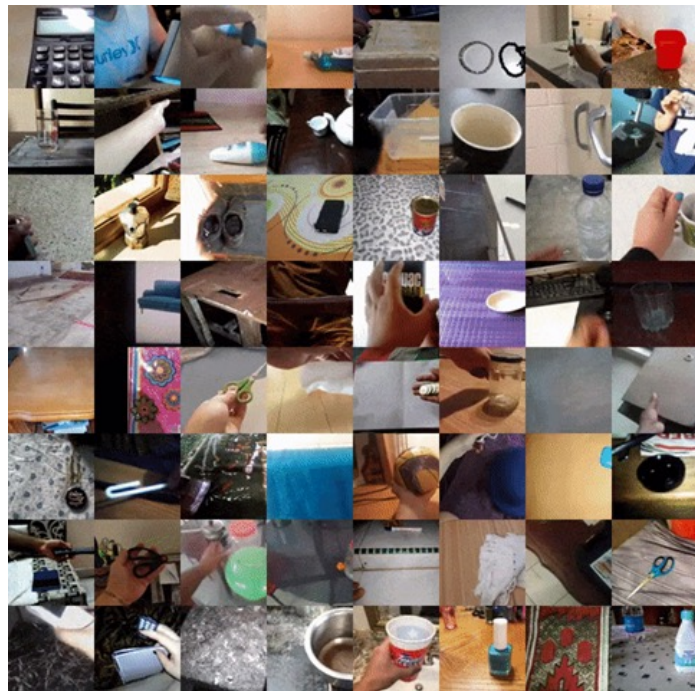


1. Pre-train an action-free latent video prediction model
2. Stack an action-conditional model when fine-tuned for MBRL
3. Video-based intrinsic bonus for better exploration



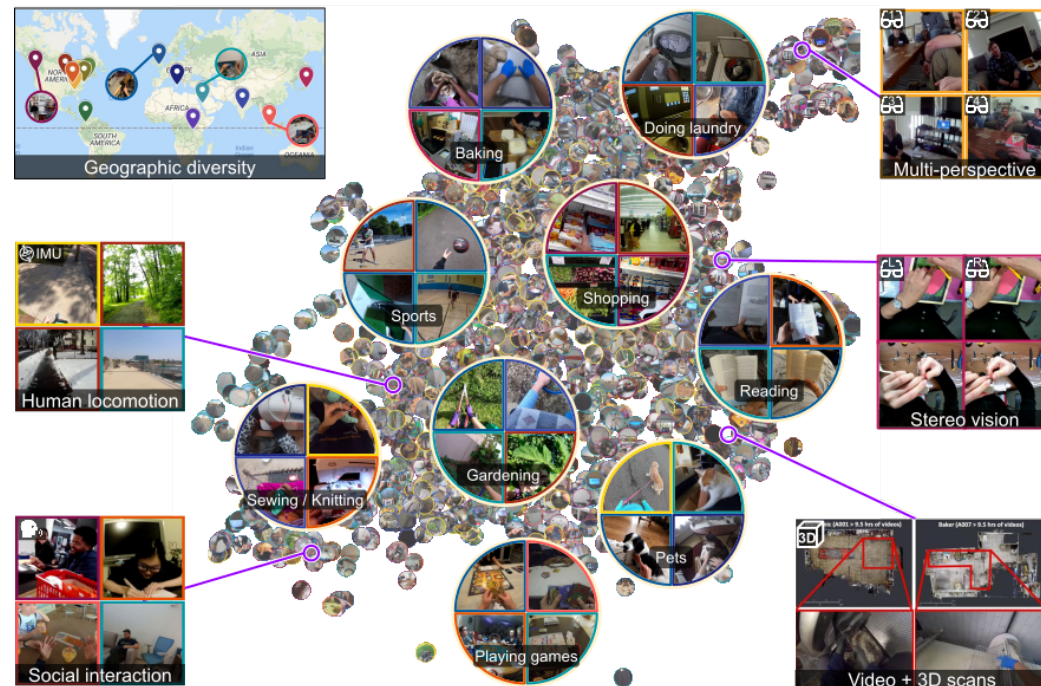
# Our Work: Towards a **General** World Model

**General world knowledge** for a variety of downstream tasks  
from **abundant in-the-wild videos** on the Internet



Something-Something V2

Goyal et al. ICCV 2017



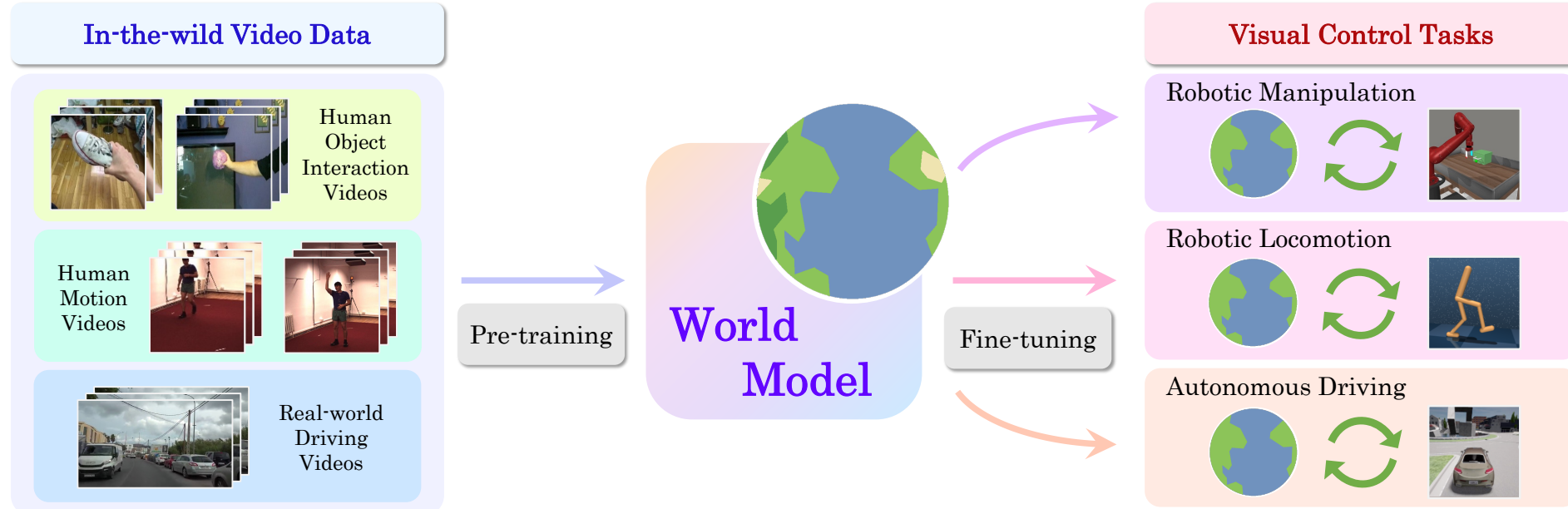
Ego4D

Grauman et al., Facebook AI. CVPR 2022

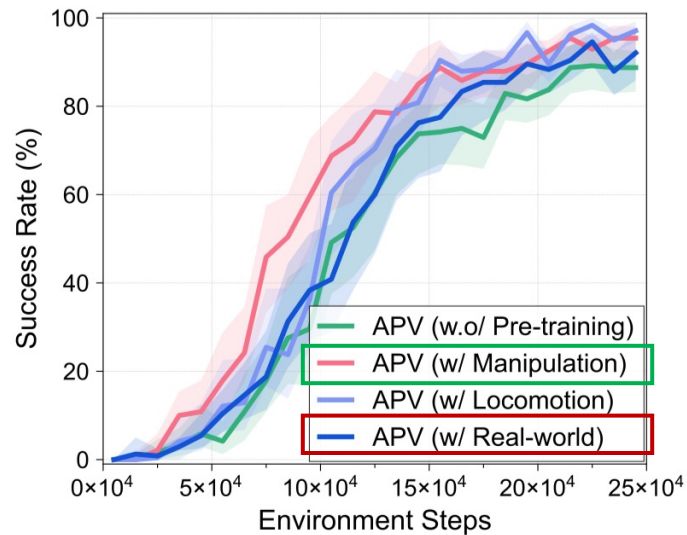
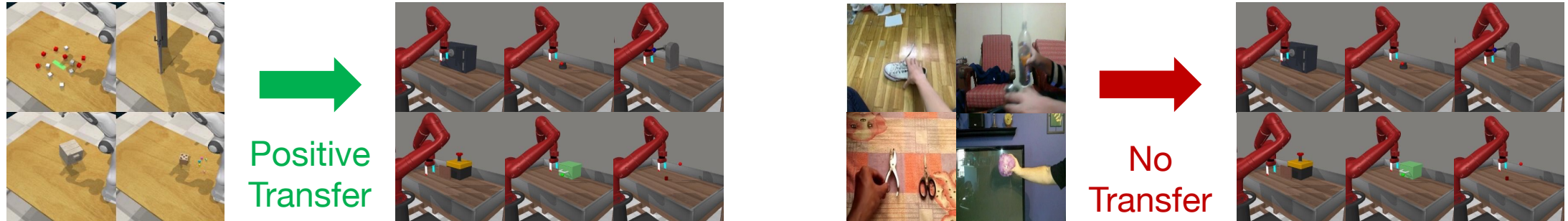
# IPV: In-the-wild Pre-training from Videos

## Towards a **general world model**:

- How to overcome the **visual complexity** and diversity?
- What is the **shared knowledge** transferable from in-the-wild video domain to visual control tasks?



# Failure of **Plain** World Models on In-the-wild Videos



## Why pre-training fails?

Seo et al.: Video prediction model suffers from **severe underfitting**

Wasting model capacity on modeling low-level **contextual** information!

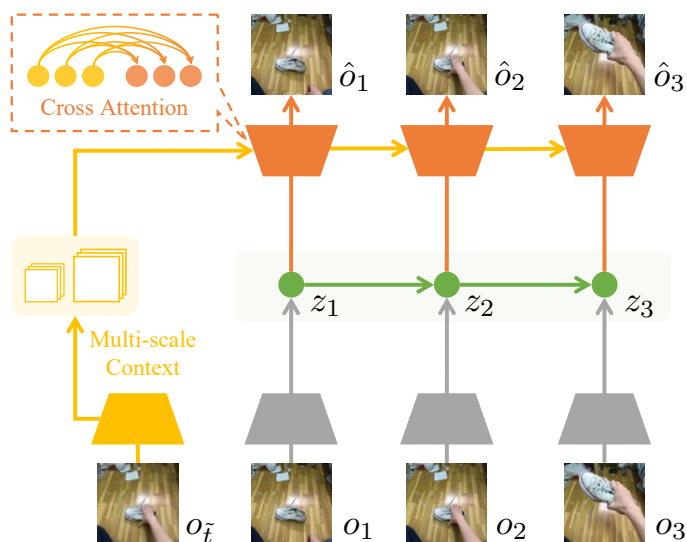


# Contextualized World Models (ContextWM)

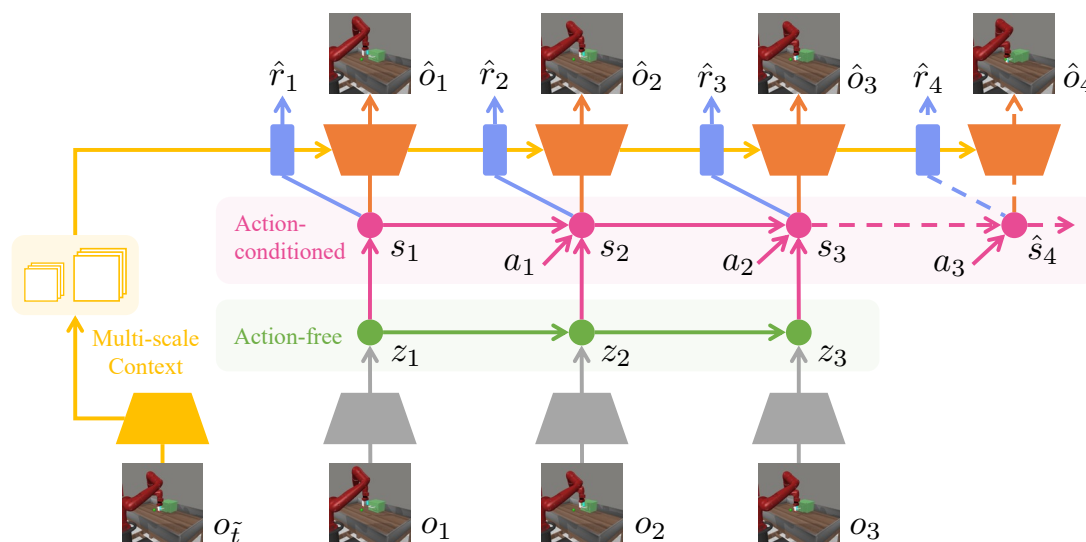
## Overview:

ContextWM empowers the **image decoder** by incorporating a **context encoder** that operates in parallel with the **latent dynamics model**

- ✓ Less inductive bias
- ✓ Diverse datasets & tasks



Step 1. Pre-training with in-the-wild videos by action-free video prediction



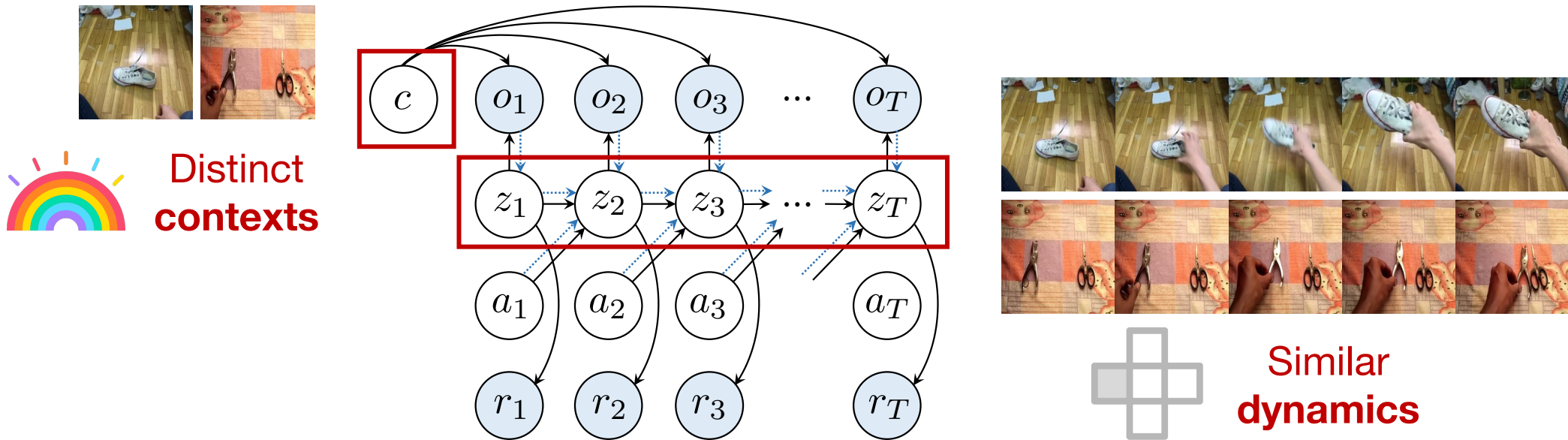
Step 2. Fine-tuning on downstream visual control tasks with MBRL



# Contextualized Latent Dynamics Models

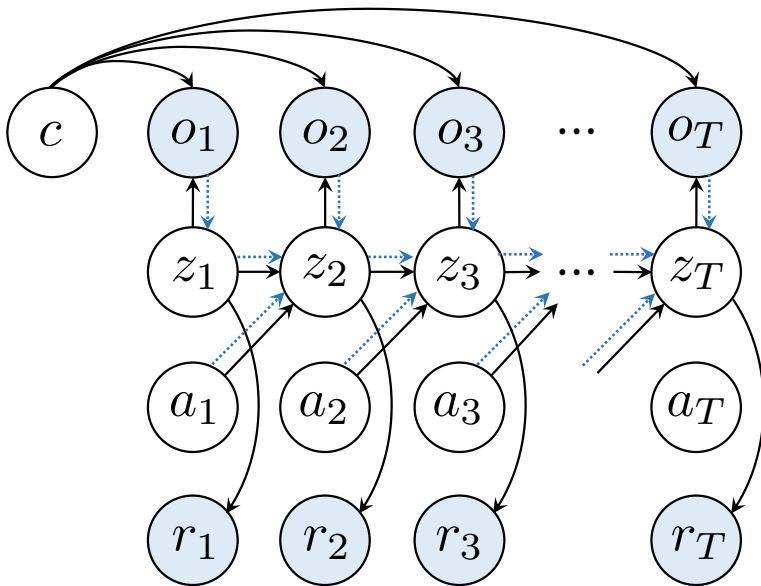
## Our insight:

Even across **distinct** scenes (**contexts**), the environment **dynamics** and physics **share a similar structure**.



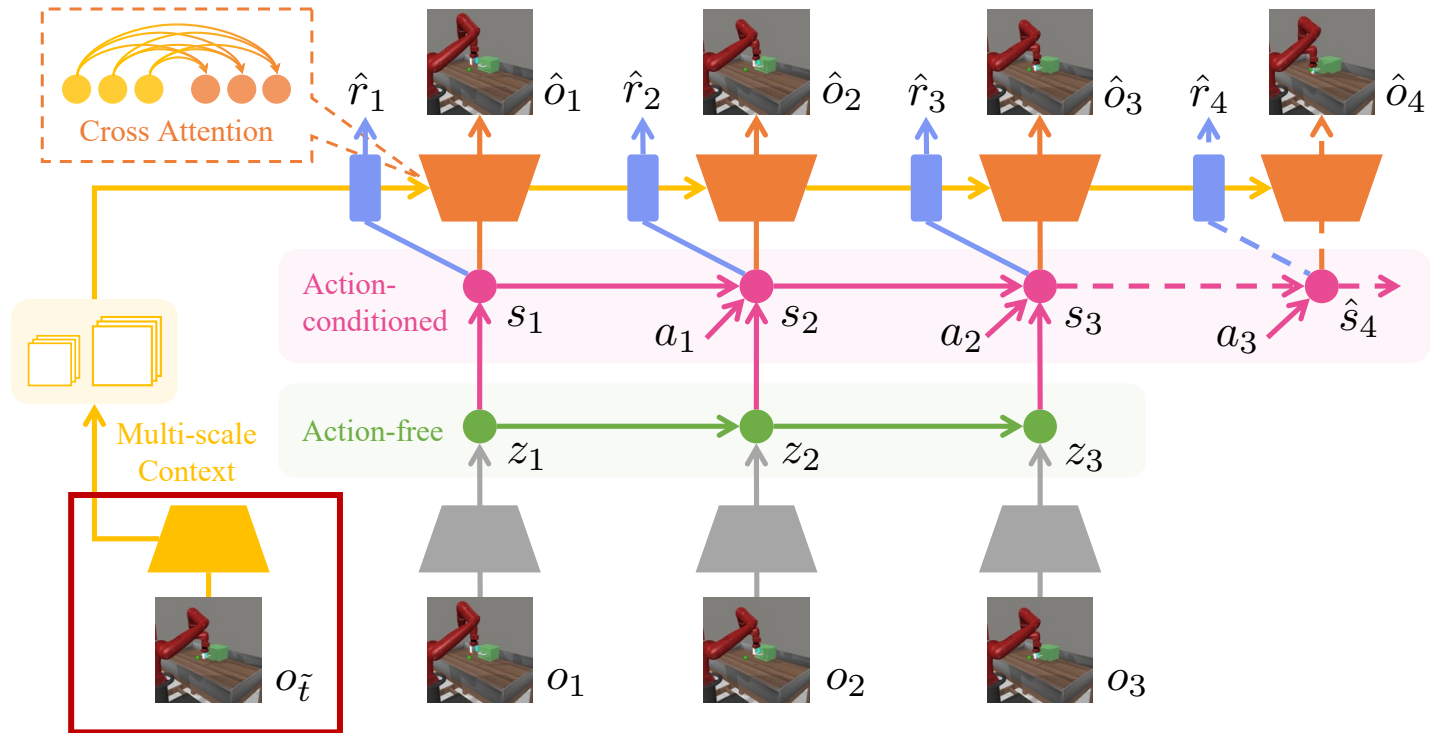
# Contextualized Latent Dynamics Models

$$\mathcal{L}(\theta) \doteq \underbrace{\mathbb{E}_{q_\theta(z_{1:T}|a_{1:T},o_{1:T})}}_{\substack{\text{context-unaware} \\ \text{latent inference}}} \left[ \sum_{t=1}^T \left( \underbrace{-\ln p_\theta(o_t | z_t, c)}_{\text{contextualized image loss}} - \ln p_\theta(r_t | z_t) + \beta_z \text{KL} [q_\theta(z_t | z_{t-1}, a_{t-1}, o_t) \| p_\theta(\hat{z}_t | z_{t-1}, a_{t-1})] \right) \right]$$



- Learn with ELBO of **conditional**  $\ln p_\theta(o_{1:T}, r_{1:T} | a_{1:T}, c)$  without the need to model the context distribution
- **Contextualized** image decoders with **rich information** beyond the expressiveness of latent variables
- Latent **dynamics** inference concentrates on **essential temporal variations**

# Contextualized World Models: An Implementation



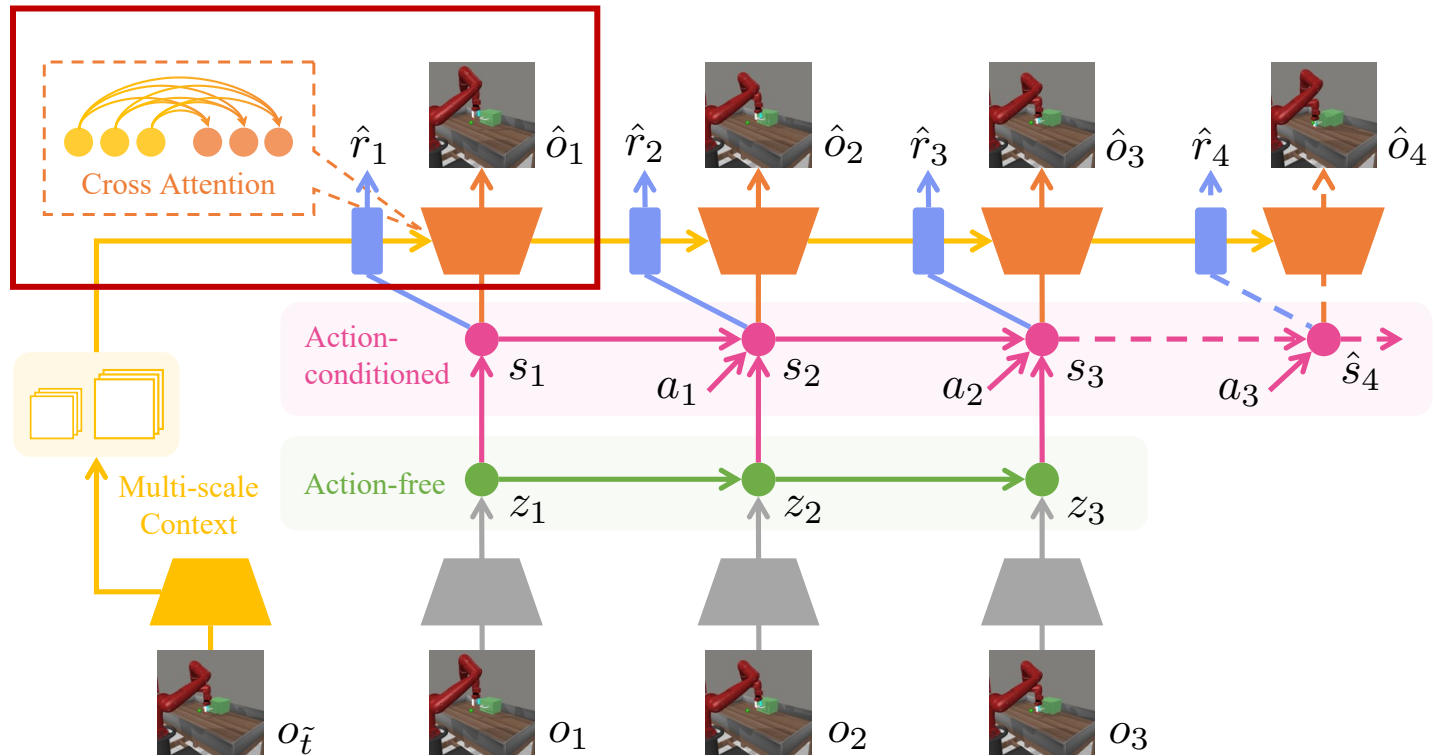
## Context formulation:

A random single frame from the trajectory segment

$$c \doteq o_{\tilde{t}}, \tilde{t} \sim \text{Uniform}\{T\}$$

By random selection, the context encoder learns to be **robust to temporal variations**

# Contextualized World Models: An Implementation



## Multi-scale cross-attention:

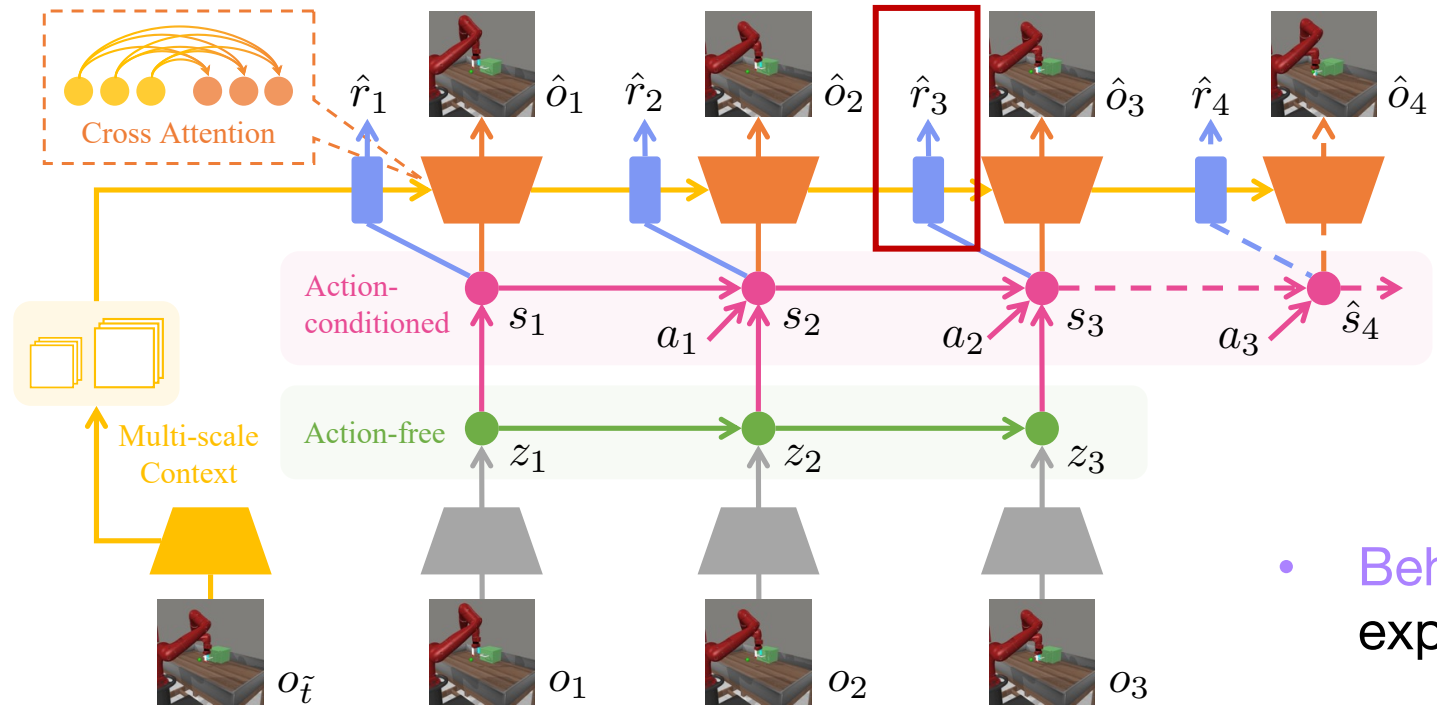
1. U-Net-style multi-scale feature shortcuts
2. Instead of naive concatenation forcing a spatial alignment, adaptive cross-attention mechanism is utilized

Flatten:  $Q = \text{Reshape}(X) \in \mathbb{R}^{hw \times c}$ ,  $K = V = \text{Reshape}(Z) \in \mathbb{R}^{hw \times c}$

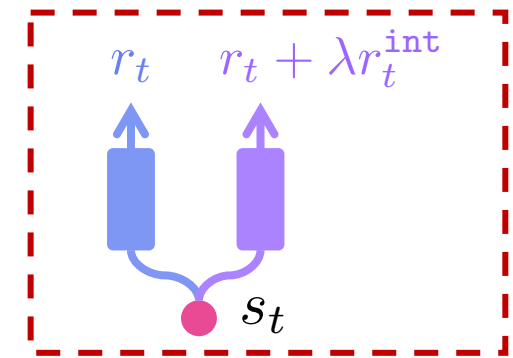
Cross-Attention:  $R = \text{Attention}(QW^Q, KW^K, VW^V) \in \mathbb{R}^{hw \times c}$

Residual-Connection:  $X = \text{ReLU}(X + \text{BatchNorm}(\text{Reshape}(R))) \in \mathbb{R}^{c \times h \times w}$ .

# Contextualized World Models: An Implementation



## Dual reward predictors:



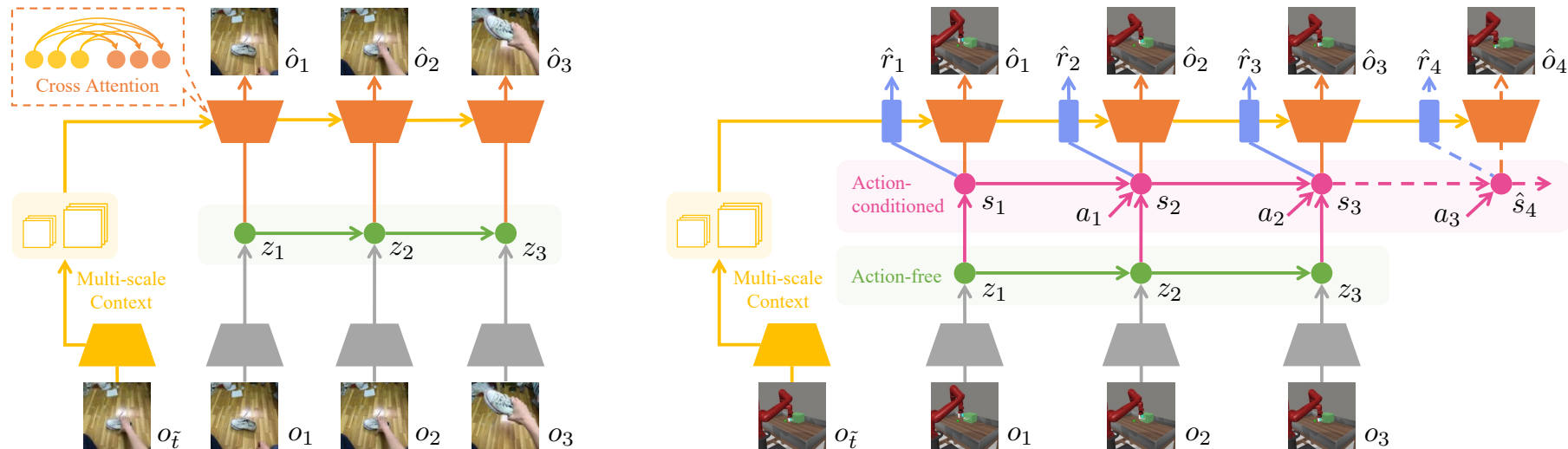
- Behavioral reward predictor: exploratory reward for behavior learning
- Representative reward predictor: pure task reward for task-relevant representation learning



# Contextualized World Models: An Implementation

## Overall objective:

$$\begin{aligned}
 \mathcal{L}^{\text{CWM}}(\phi, \varphi, \theta) \doteq & \underbrace{\mathbb{E}_{q_\phi(s_{1:T}|a_{1:T}, z_{1:T}), q_\theta(z_{1:T}|o_{1:T})}}_{\text{context-unaware latent inference}} \left[ \sum_{t=1}^T \left( \underbrace{-\ln p_\theta(o_t|s_t, c)}_{\text{contextualized image loss}} \right. \right. \\
 & \underbrace{-\ln p_\phi(r_t + \lambda r_t^{\text{int}}|s_t)}_{\text{behavioral reward loss}} \underbrace{-\beta_r \ln p_\varphi(r_t|s_t)}_{\text{representative reward loss}} \underbrace{+\beta_z \text{KL}[q_\theta(z_t|z_{t-1}, o_t) \| p_\theta(\hat{z}_t|z_{t-1})]}_{\text{action-free KL loss}} \\
 & \left. \left. \underbrace{+\beta_s \text{KL}[q_\phi(s_t|s_{t-1}, a_{t-1}, z_t) \| p_\phi(\hat{s}_t|s_{t-1}, a_{t-1})]}_{\text{action-conditional KL loss}} \right) \right].
 \end{aligned}$$



# Experiments: Diverse Datasets & Tasks



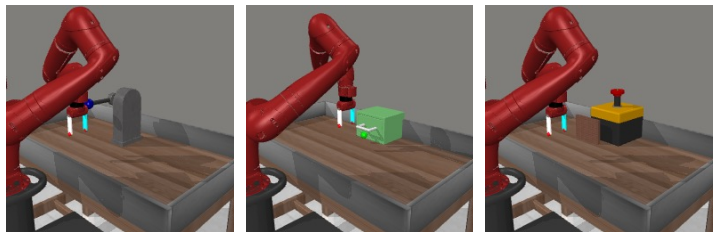
Something-Something V2  
Goyal et al. ICCV 2017



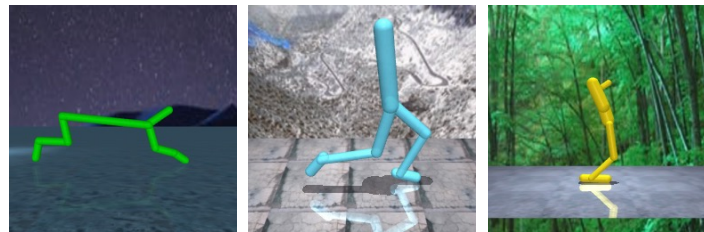
Human3.6M  
Ionescu et al. TPAMI 2014



YouTube Driving  
Zhang et al. ECCV 2022



Meta-World  
Yu et al. CoRL 2020

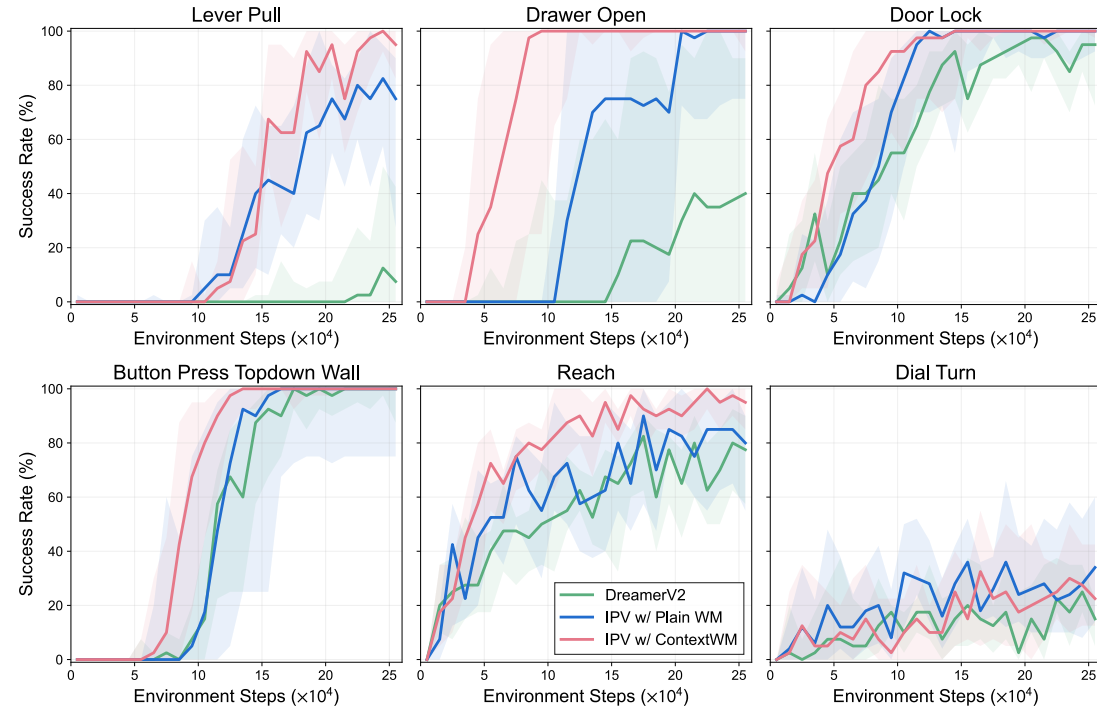
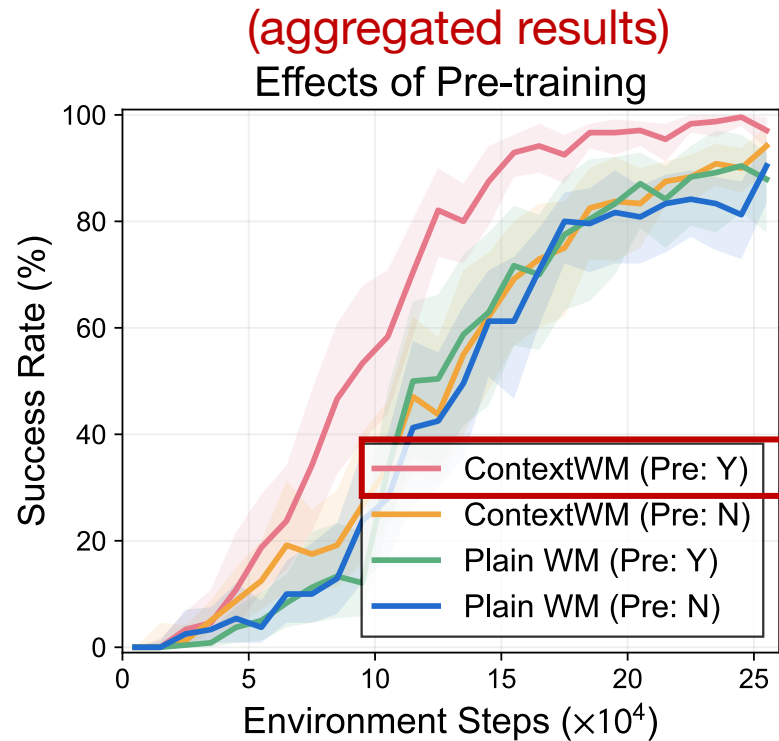


DMC Remastered  
Grigsby et al. 2020



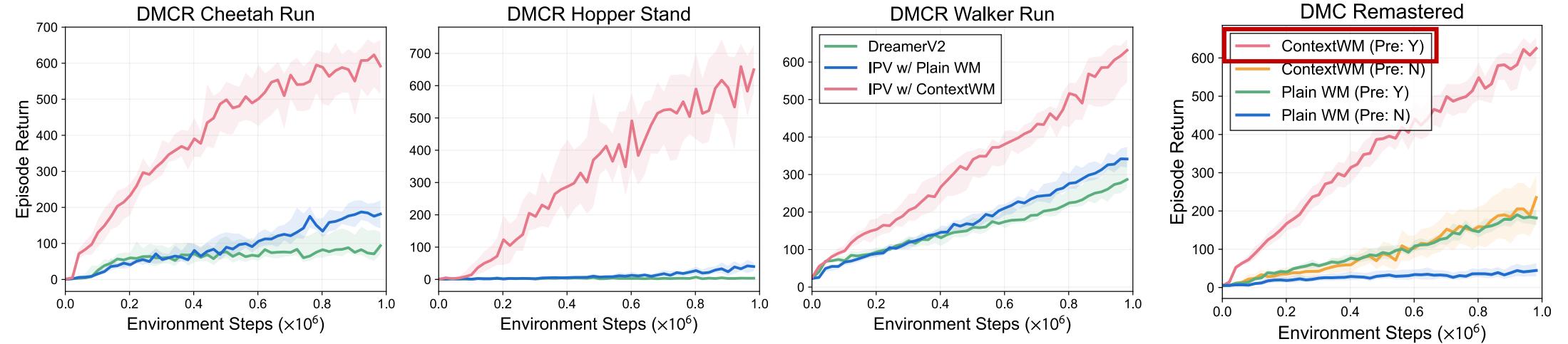
CARLA  
Dosovitskiy et al. CoRL 2017

# Main Results: Meta-world

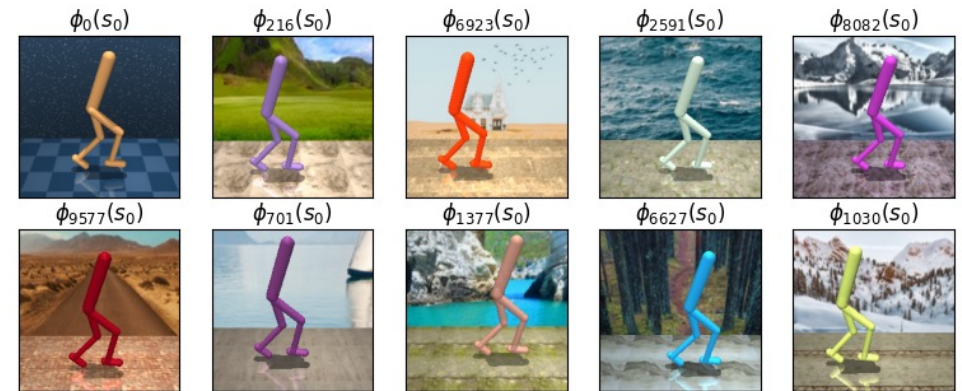


**On six Meta-world tasks, ContextWM achieves significant positive transfer (from SSv2) in terms of sample efficiency, while a plain WM fails.**

# Main Results: DMC Remastered

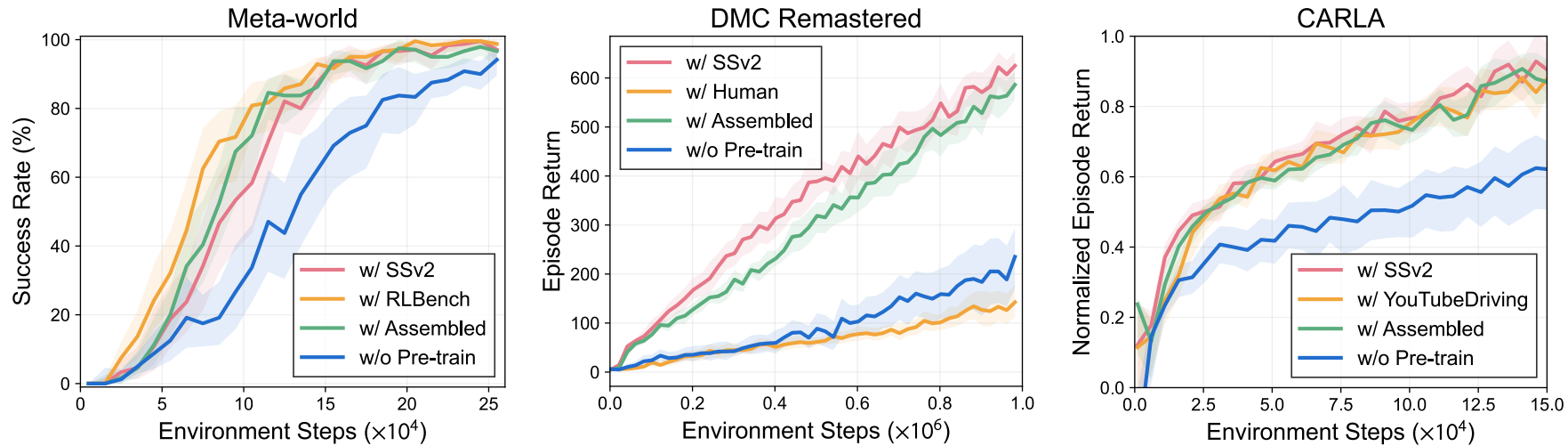


**On visual generalization benchmark,  
pre-training from in-the-wild videos (SSv2)  
incurs significant performance boost,  
which is further unleashed by ContextWM.**



Visual generalization benchmark: Seven visual factors randomly initialized on each episode

# Effects of Pre-training Dataset Domain

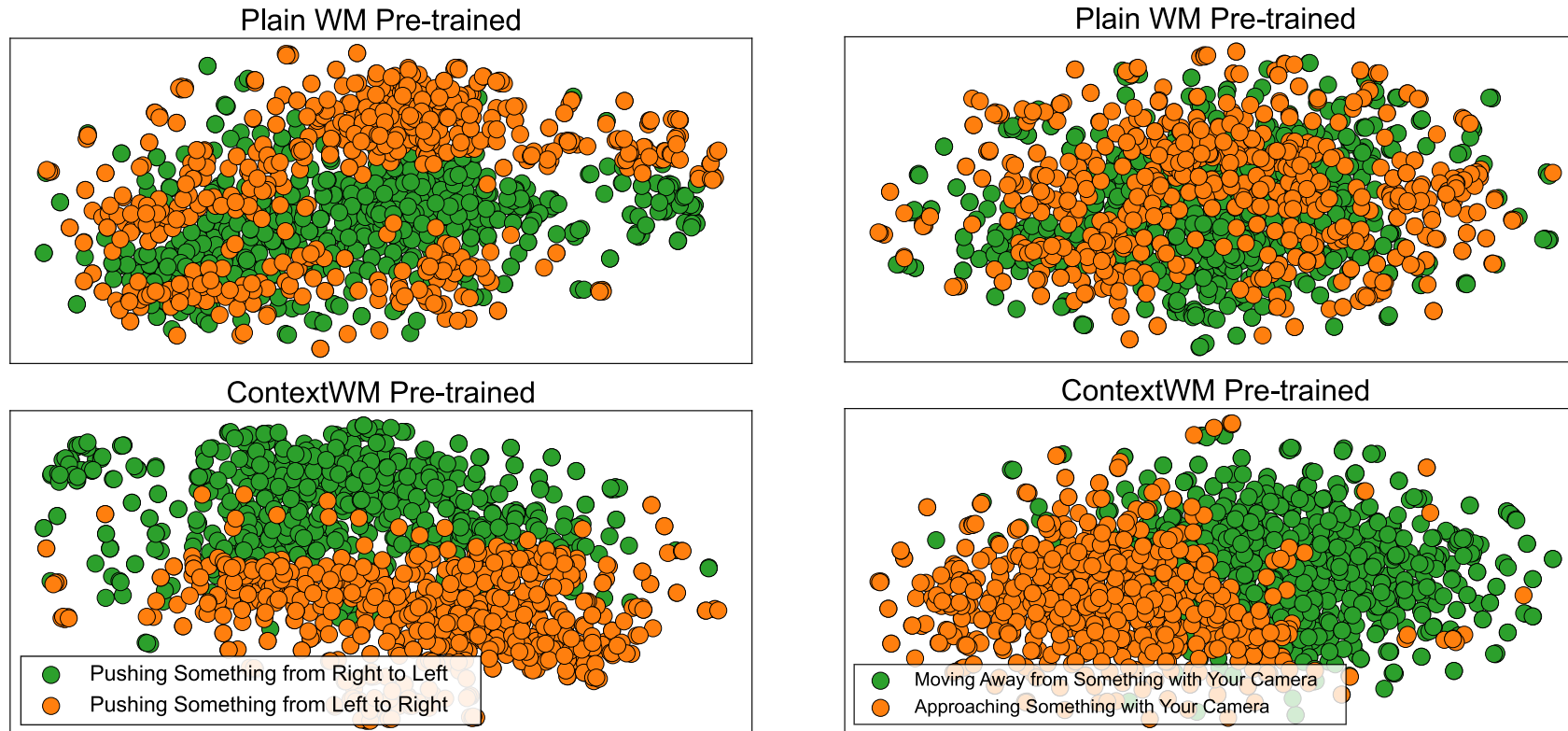


## Takeaways:

1. Human-object interaction data (SSv2) are generally beneficial.
2. A more similar domain (e.g. RL/Bench) is more useful, but more diverse datasets can serve as promising scalable alternatives.
3. Pre-training data lack of diversity (Human3.6M) can even be harmful.

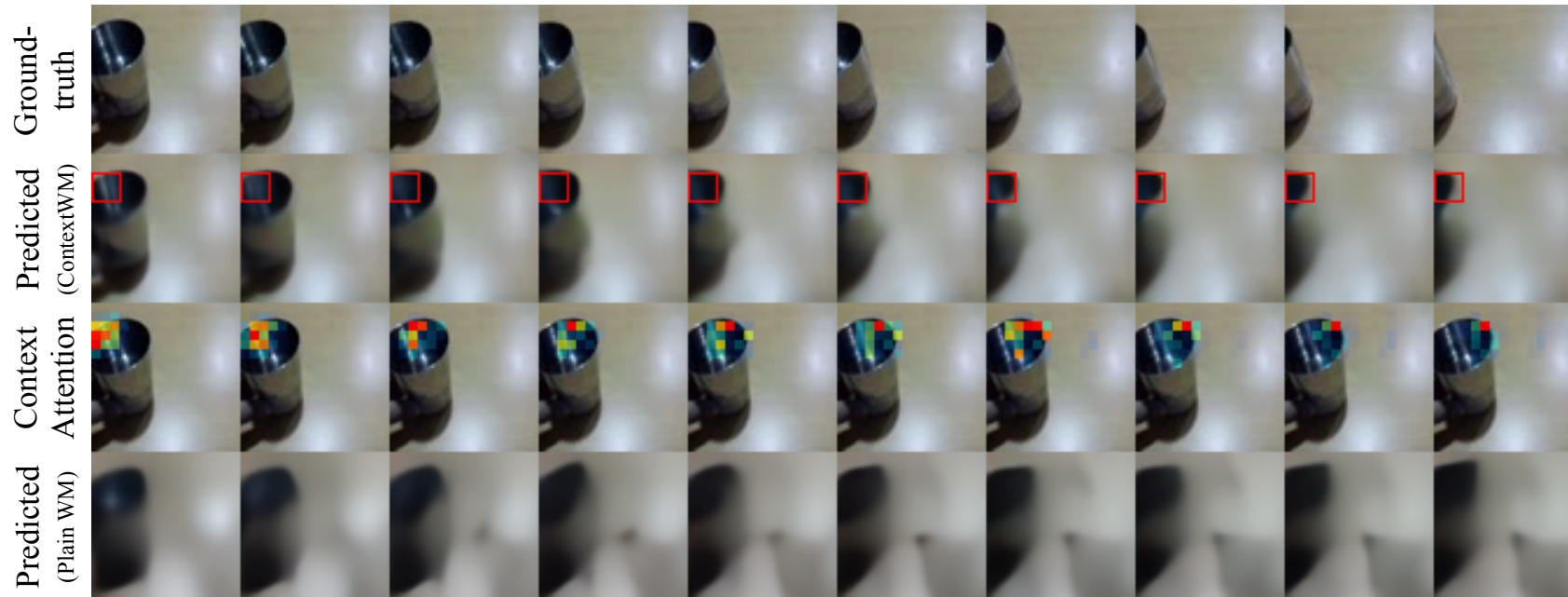


# Qualitative Evaluation: Video Representations



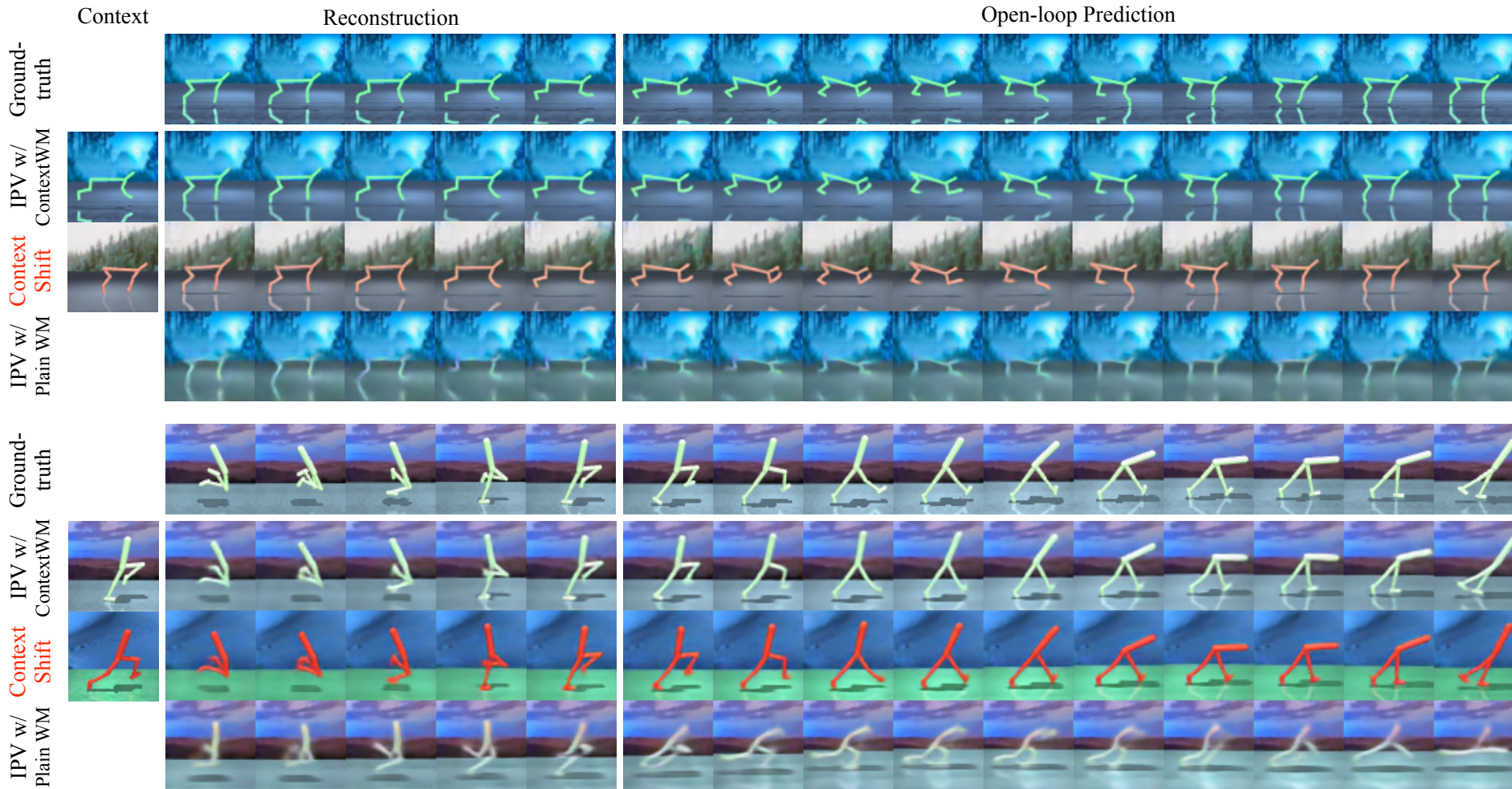
ContextWM learns representations well distributed according to **different directions of motion**, while **not utilizing any labels** of the videos in pre-training

# Qualitative Evaluation: Video Prediction



1. Predictions from ContextWM **well capture the shape and motion** of the water cup.
2. Cross-attentions from different frames **successfully attend to varying spatial positions** of the context frame.

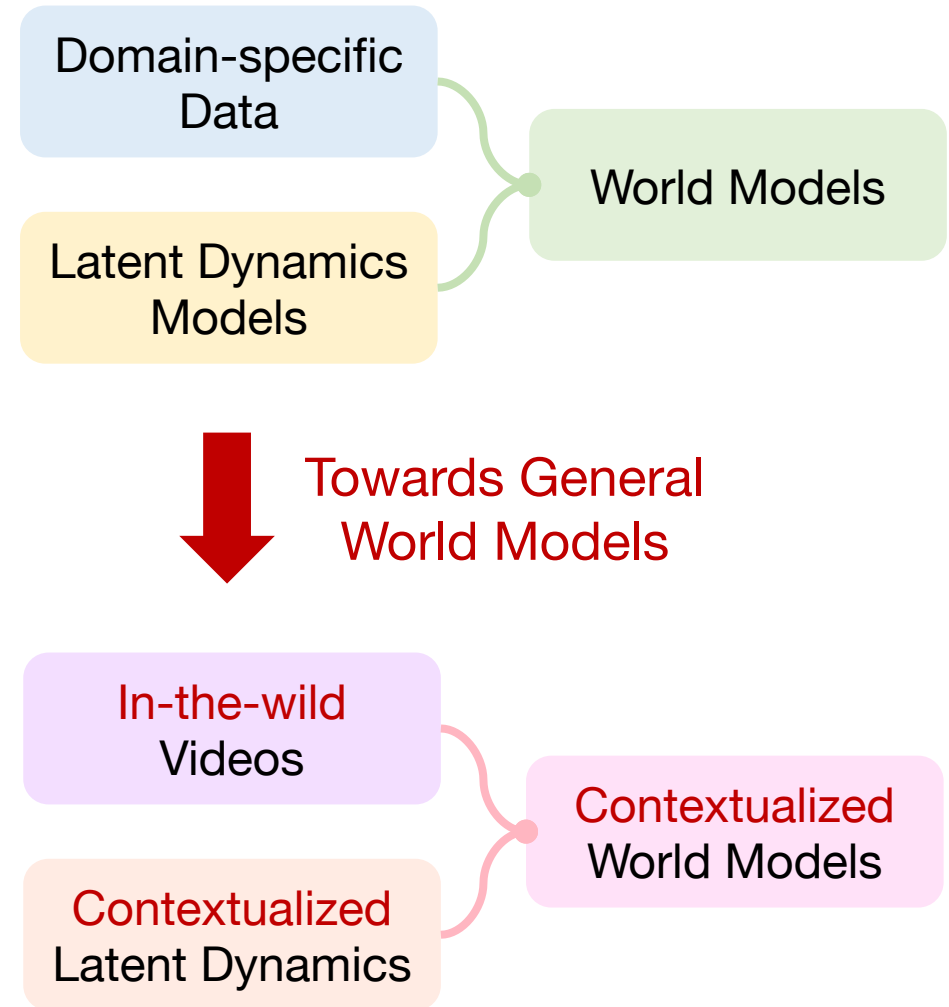
# Qualitative Evaluation: Compositional Decoding



**Excellent**  
**compositionality**  
to combine new  
contexts with the  
original dynamics  
by **disentangled**  
**representations**

# Summary

- A **world model** is an internal model of how the world works, which empowers the agent to **reason and plan**.
- Our work:
  - Introduces **Contextualized World Models (ContextWM)**
  - Applies it to the paradigm of **In-the-wild Pre-training from Videos (IPV)**
  - Followed by fine-tuning on downstream tasks to **boost learning efficiency of MBRL**





# Thank You!

Code Available: <https://github.com/thuml/ContextWM>

Contact: [wujialong0229@gmail.com](mailto:wujialong0229@gmail.com)

Machine Learning Group, School of Software, Tsinghua University

<http://ise.thss.tsinghua.edu.cn/~mlong/>

