# Trajectory World Models for Heterogeneous Environments

Shaofeng Yin*, Jialong Wu*, Siqiao Huang, Xingjian Su, Xu He, Jianye Hao, Mingsheng Long#

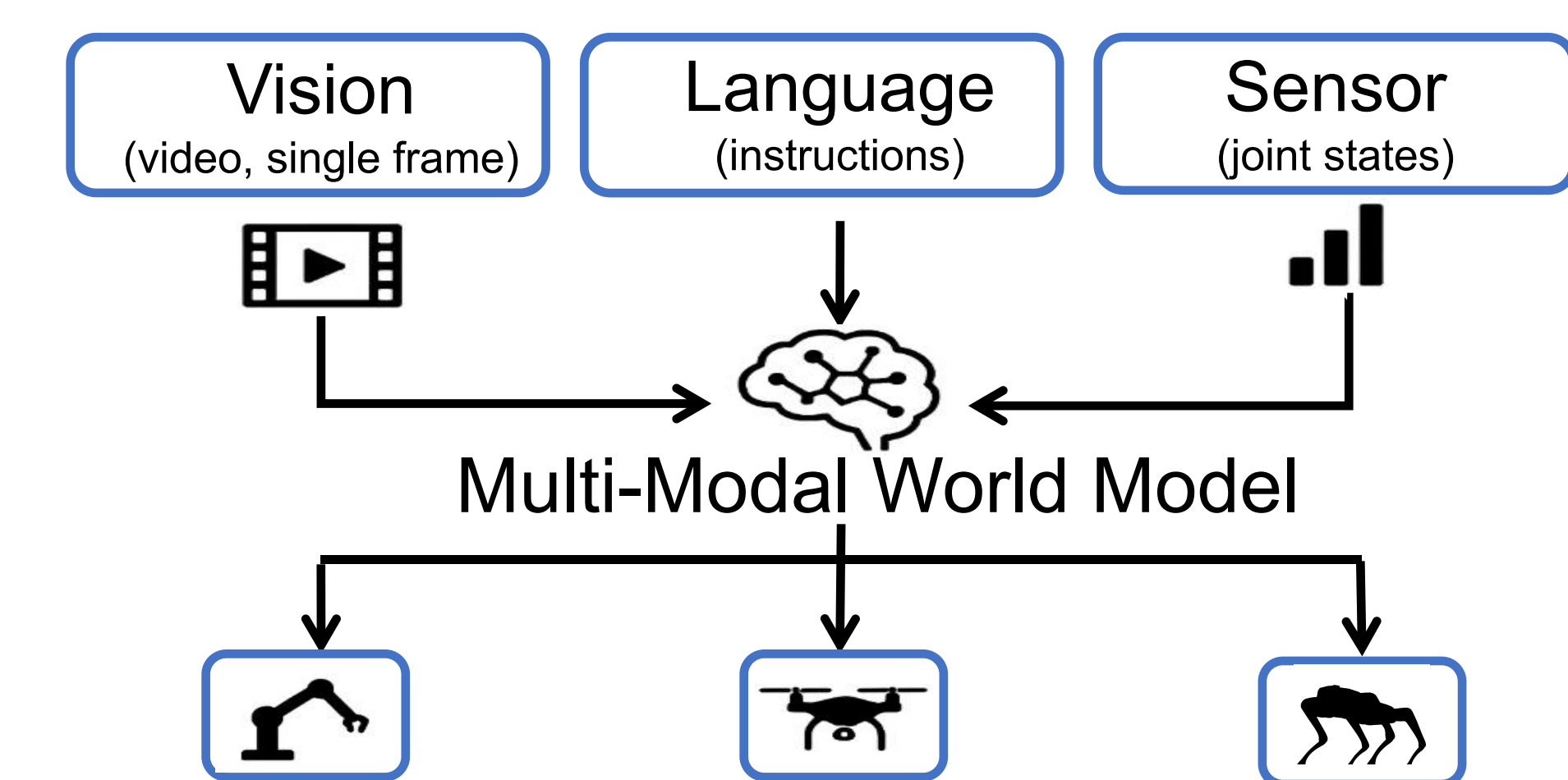## Motivation: Heterogeneity inherent in sensor and actuator information

### Motivation

World models are all with videos or language?

No modality in world models should be left behind, including essential sensor information represented as low-dimensional vectors!

How can we pre-train a world model to extract shared knowledge from trajectories across heterogeneous environments?

### Vision of the future

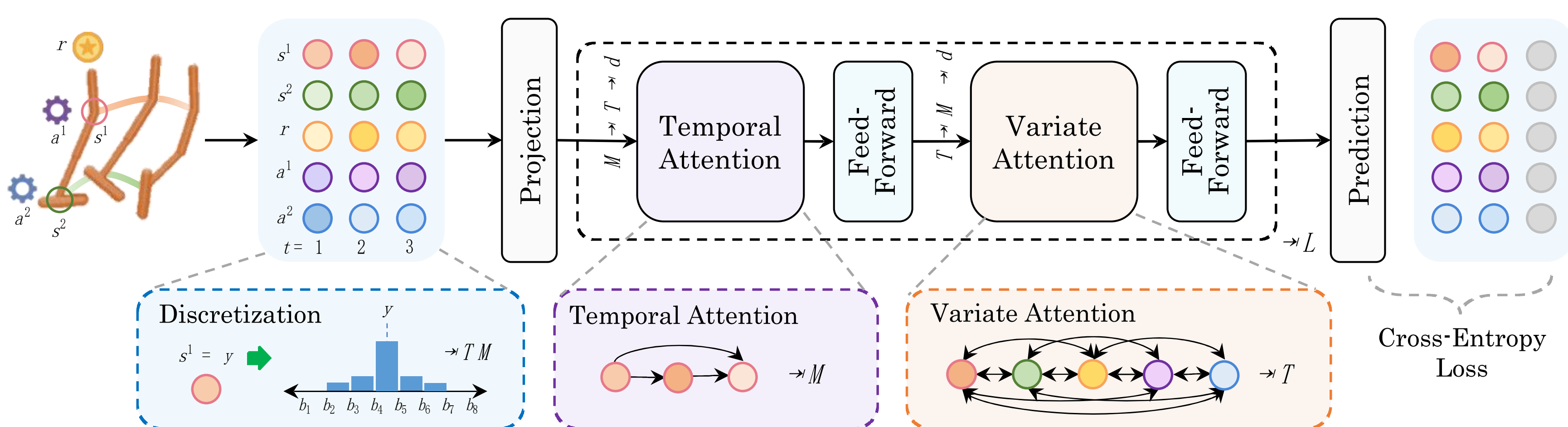Towards multi-modal world models incorporating proprioceptive, visual and linguistical observations



Vision (video, single frame) / Language (instructions) / Sensor (joint states) → Multi-Modal World Model

## Method: TrajWorld (Trajectory World Models)

### Overview:

TrajWorld, designed for flexibility in handling divergent state and action definitions, is capable of flexibly handling varying sensor and actuator information and capturing environment dynamics in-context.



### Intuition:

1. Rediscovering homogeneity in scalars.
2. Identifying environment through historical context.
3. Inductive bias for two-dimensional representations.

### Interleaved temporal-variate attentions:

1. temporal attention
$$U_{1:T,j}^l = \text{CausalAttention}(Z_{1:T,j}^{l-1}), \quad \forall j \in [M],$$
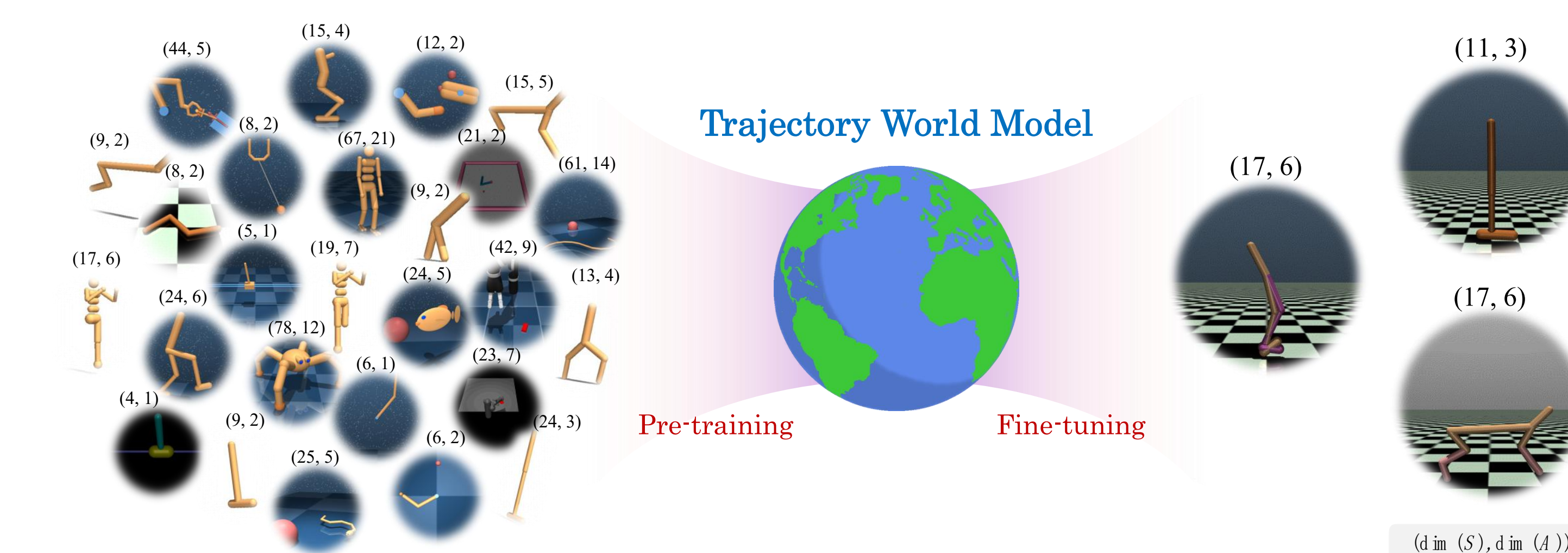2. variate attention
$$V_{i,1:M}^l = \text{Attention}(\hat{U}_{i,1:M}^l), \quad \forall i \in [T].$$
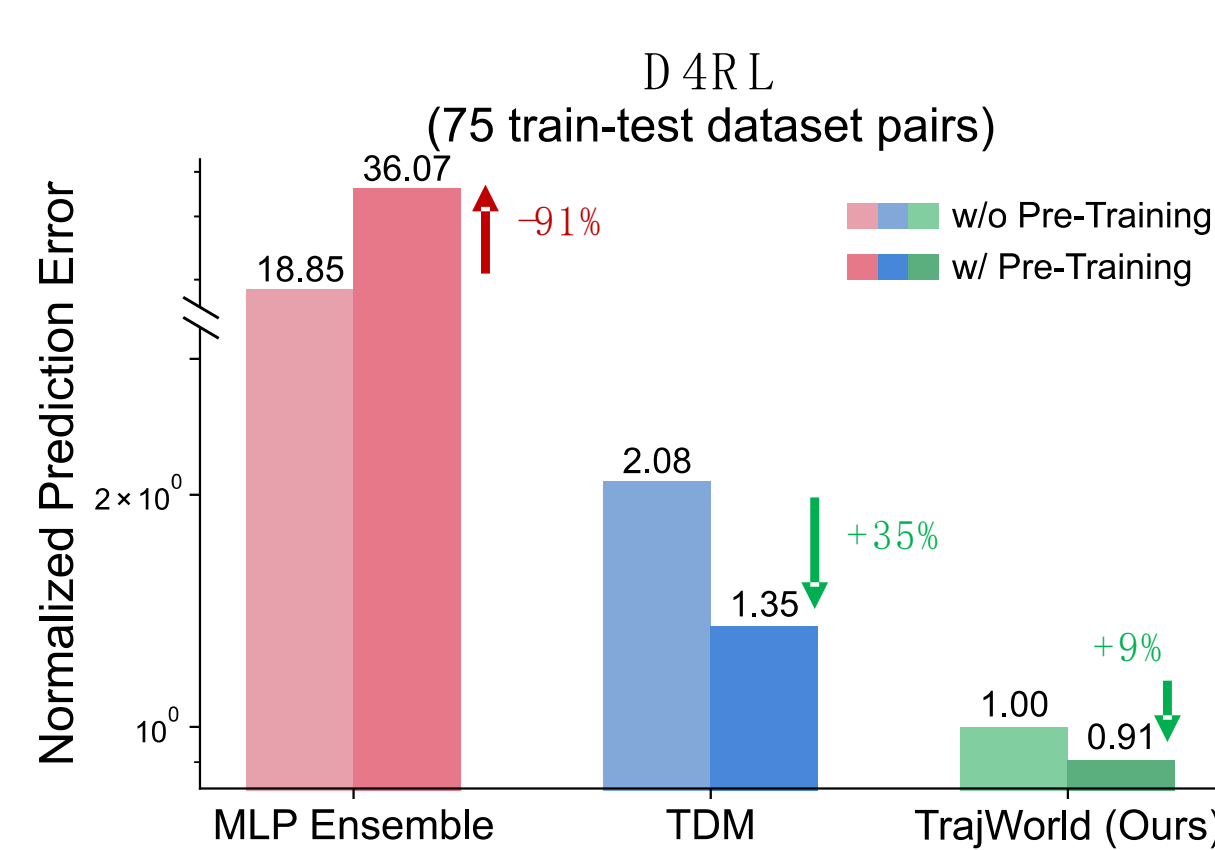
### Pre-training & Finetuning

80 Unique Envs

719M Environment Steps

1.4M Episodes

merge five atasets with different characteristics self-collection

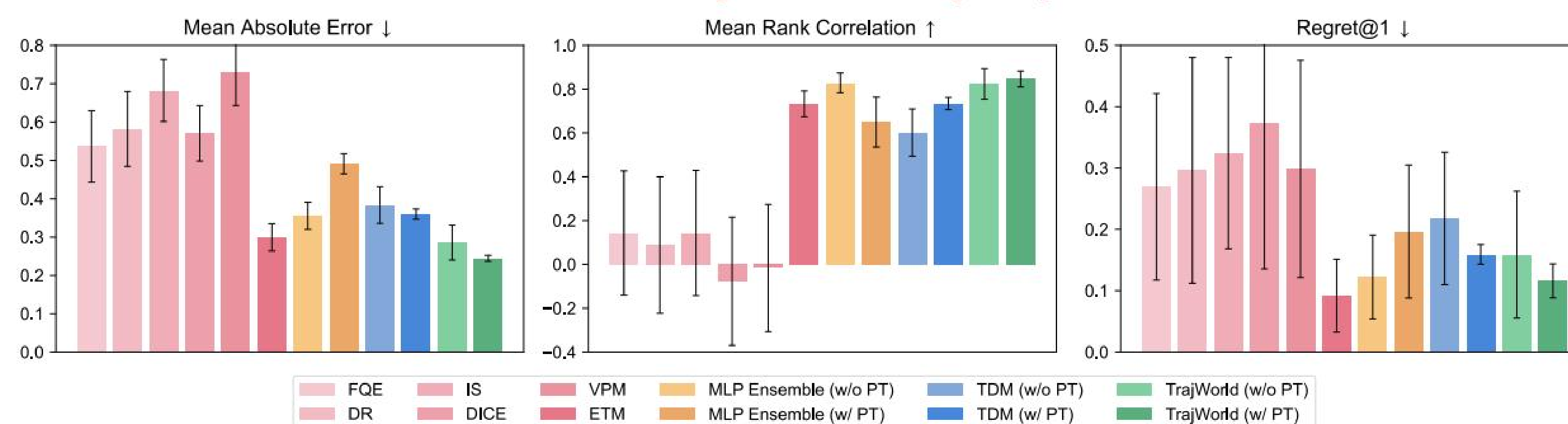environment + distribution diversity
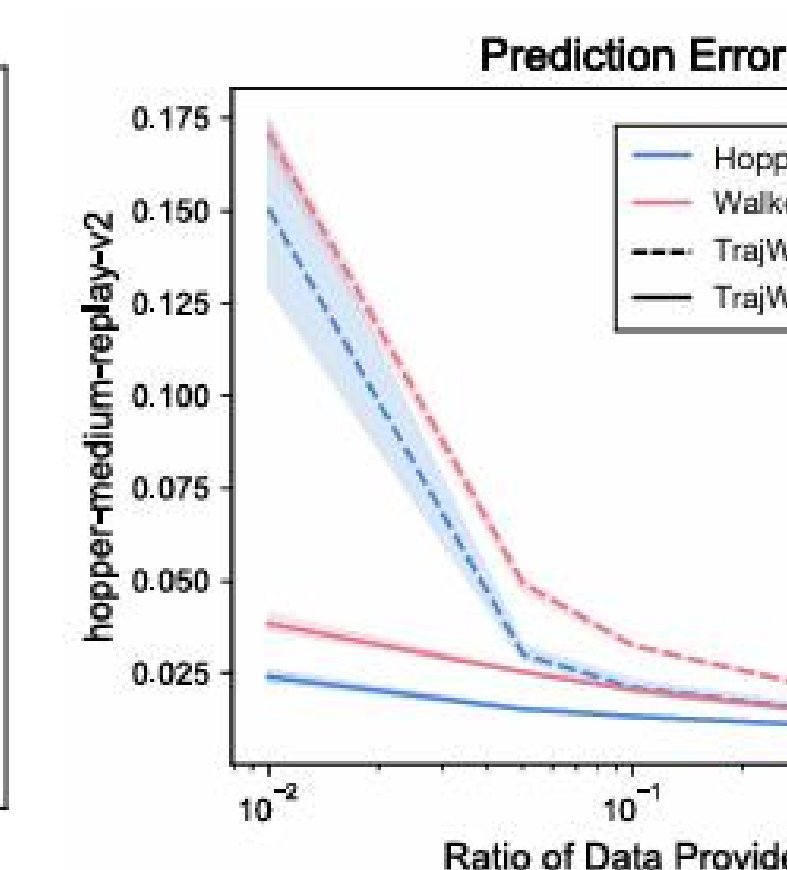


## Experiments

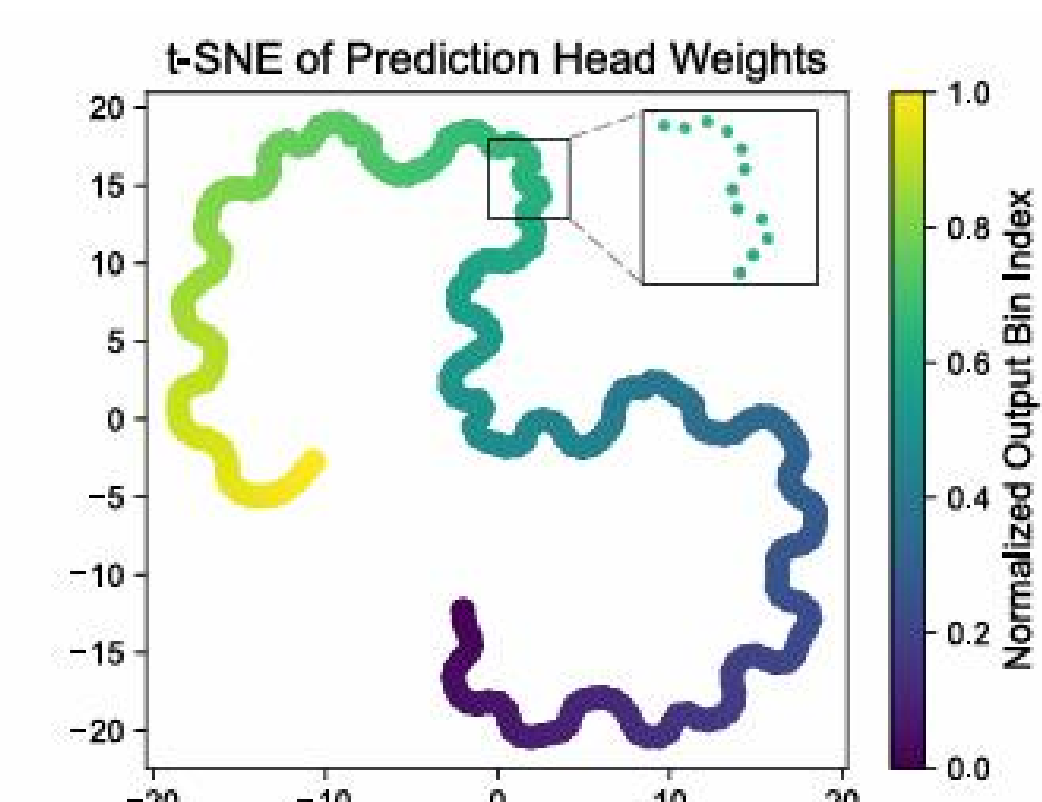### Transition Prediction



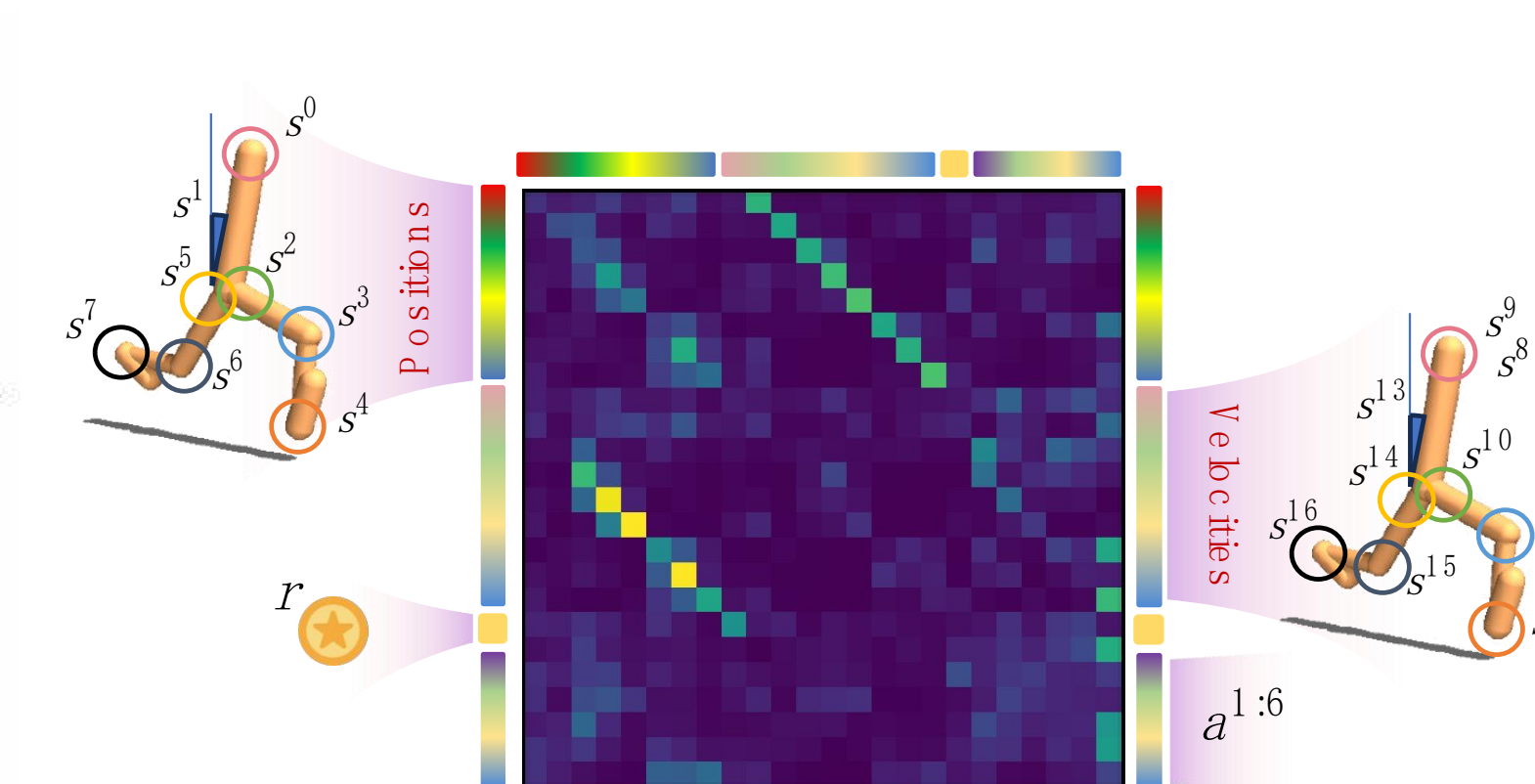### Off-Policy Evaluation (OPE)



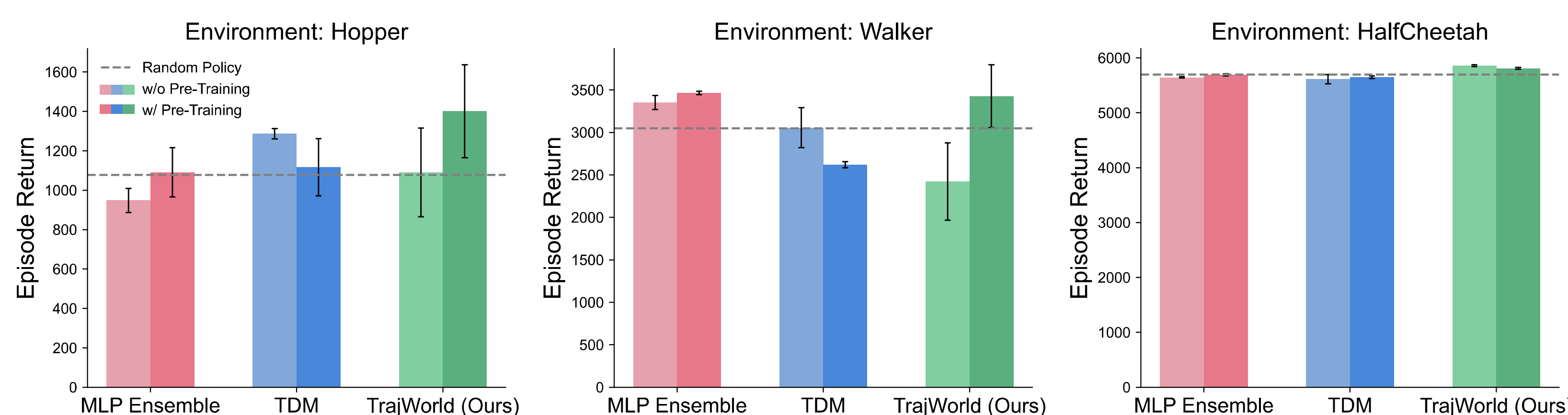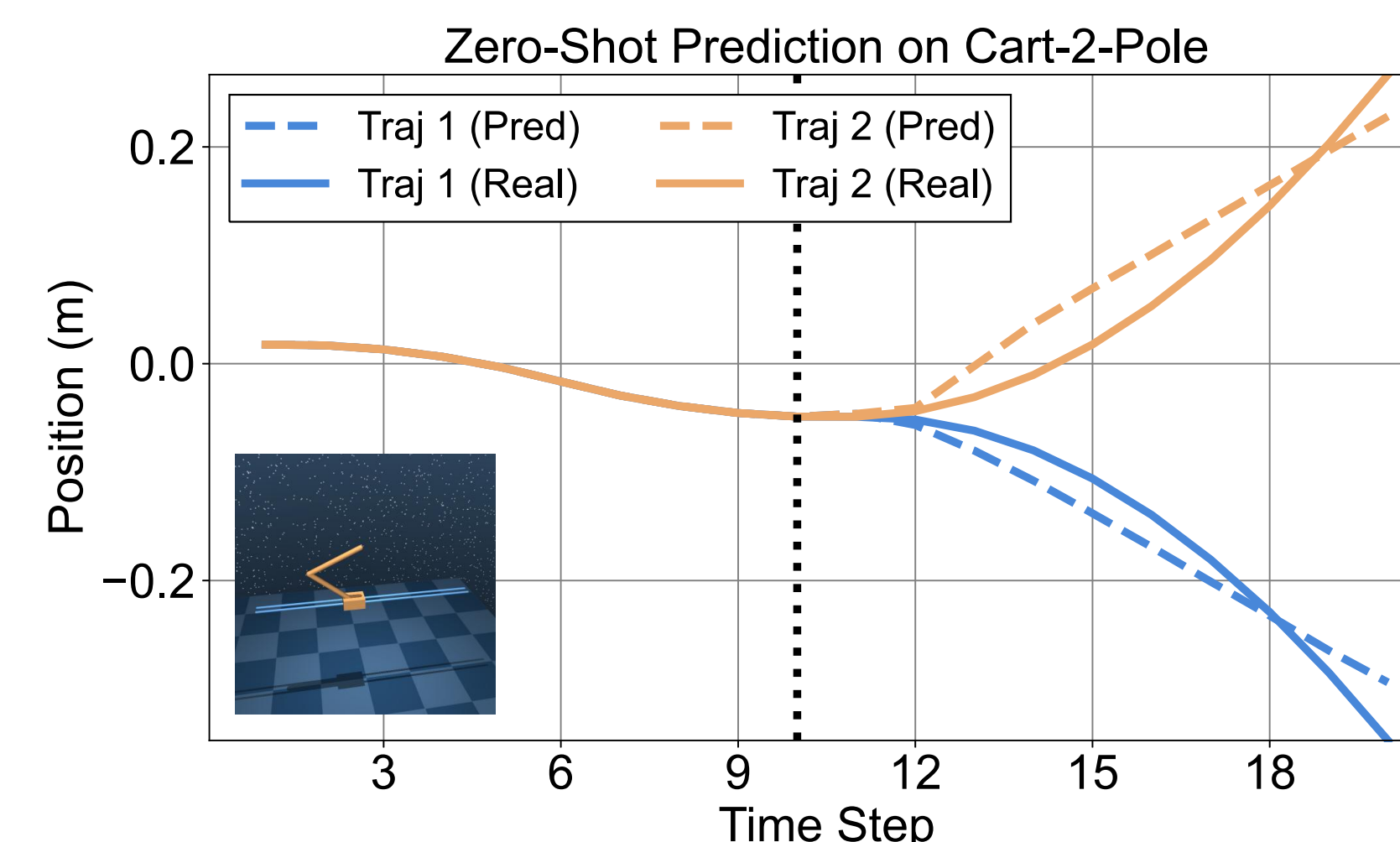### Few-shot adaptation



### Discretization



### Attention Analysis



### Model Predictive Control (MPC)



### Zero-Shot Transfer



### Scaling Trend