# Mixture-of-Experts Meets In-Context Reinforcement Learning

Wenhao Wu (吴文浩)

Department of Control Science and Intelligent Engineering
Nanjing University

Email: wenhaowu@smail.nju.edu.cn

Code: https://github.com/NJU-RL/T2MIR

# Contents

## Algorithm Distillation (AD)

**Data Generation**

Task 1
...
Task $n$

$h_T^{(n)} = (o_0, a_0, r_0, o_1, a_1, r_1, \ldots, o_T, a_T, r_T)_n$

RL algorithm learning histories

*learning progress*

**Model Training**

$o_0$ | $a_0$ | $r_0$ | $\cdots$ | $o_{t-1}$ | $a_{t-1}$ | $r_{t-1}$ | $o_t$

Predict actions using across-episodic contexts

Causal Transformer

$P_\theta(a_t | h_{t-1}, o_t)$

## Decision-Pretrained Transformer (DPT)

$\tau_N$
$\tau_1$

$s_1, a_1, r_1, s_1'$ | $s_2, a_2, r_2, s_2'$ | $\cdots$ | $s_n, a_n, r_n, s_n'$

Transformer

$s_{query}$

$a_{query}^*$

$M_\theta(a^* | s_{query}, D)$

$M_\theta(\cdot | \cdot, D)$

**Online Exploration**

$M_\theta(\cdot | \cdot, D)$

**Offline Learning**

➢ **Algorithm Distillation**

- train a causal transformer to predict actions given preceding learning histories as context
- cross-episodic trajectories
- a dataset of learning histories is generated by a source RL algorithm

$$\mathcal{L}(\theta) = -\sum_{n=1}^{N} \sum_{t=0}^{T-1} \log P_\theta\left(a = a_t^n \middle| \tau_{\mathrm{pro}}^n, s_t^n\right), \tau_{\mathrm{pro}}^n = h_{t-1}^n$$

➢ **Decision-Pretrained Transformer**

- predict the optimal action given a query state and a prompt of interactions
- need expert policy to label actions
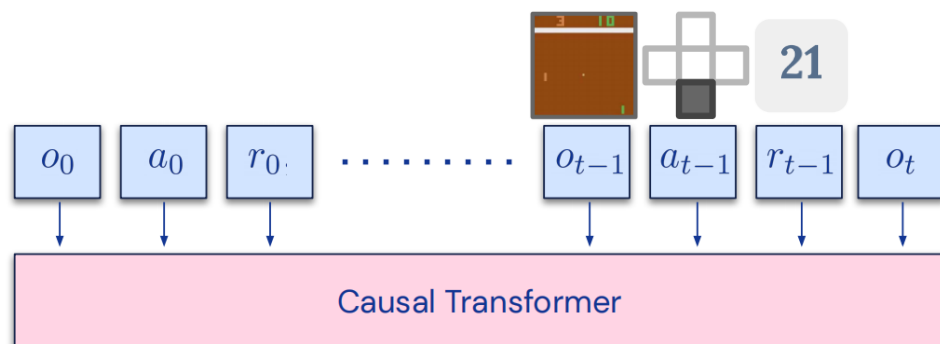- robust to different dataset qualities

$$\mathcal{L}(\theta) = -\sum_{n=1}^{N} \sum_{t=0}^{T} \log P_\theta\left(a = a_t^{n*} \middle| \tau_{\mathrm{pro}}^n, s_t^n\right), \tau_{\mathrm{pro}}^n \sim \mathcal{D}^n$$
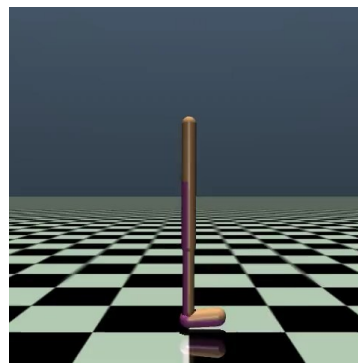
## multi-modality in state-action-reward sequence



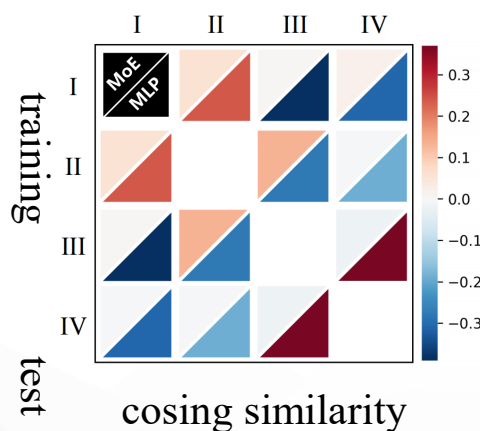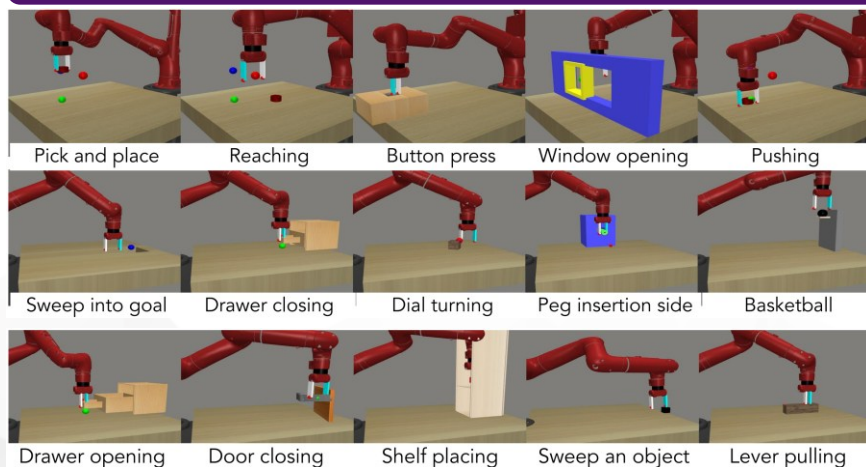$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_T, a_T, r_T)$$

Walker-Param

➢ **multi-modality in state-action-reward sequence**

• states: physical quantities (position, velocity, and acceleration)

• actions: joint torques or discrete commands

• rewards: simple scalars

## task diversity and heterogeneity



Pick and place  Reaching  Button press  Window opening  Pushing

Sweep into goal  Drawer closing  Dial turning  Peg insertion side  Basketball

Drawer opening  Door closing  Shelf placing  Sweep an object  Lever pulling

cosing similarity

➢ **task diversity and heterogeneity**

• some tasks are similar and others differ significantly

• intrinsic gradient conflicts in challenging scenarios with significant task variation
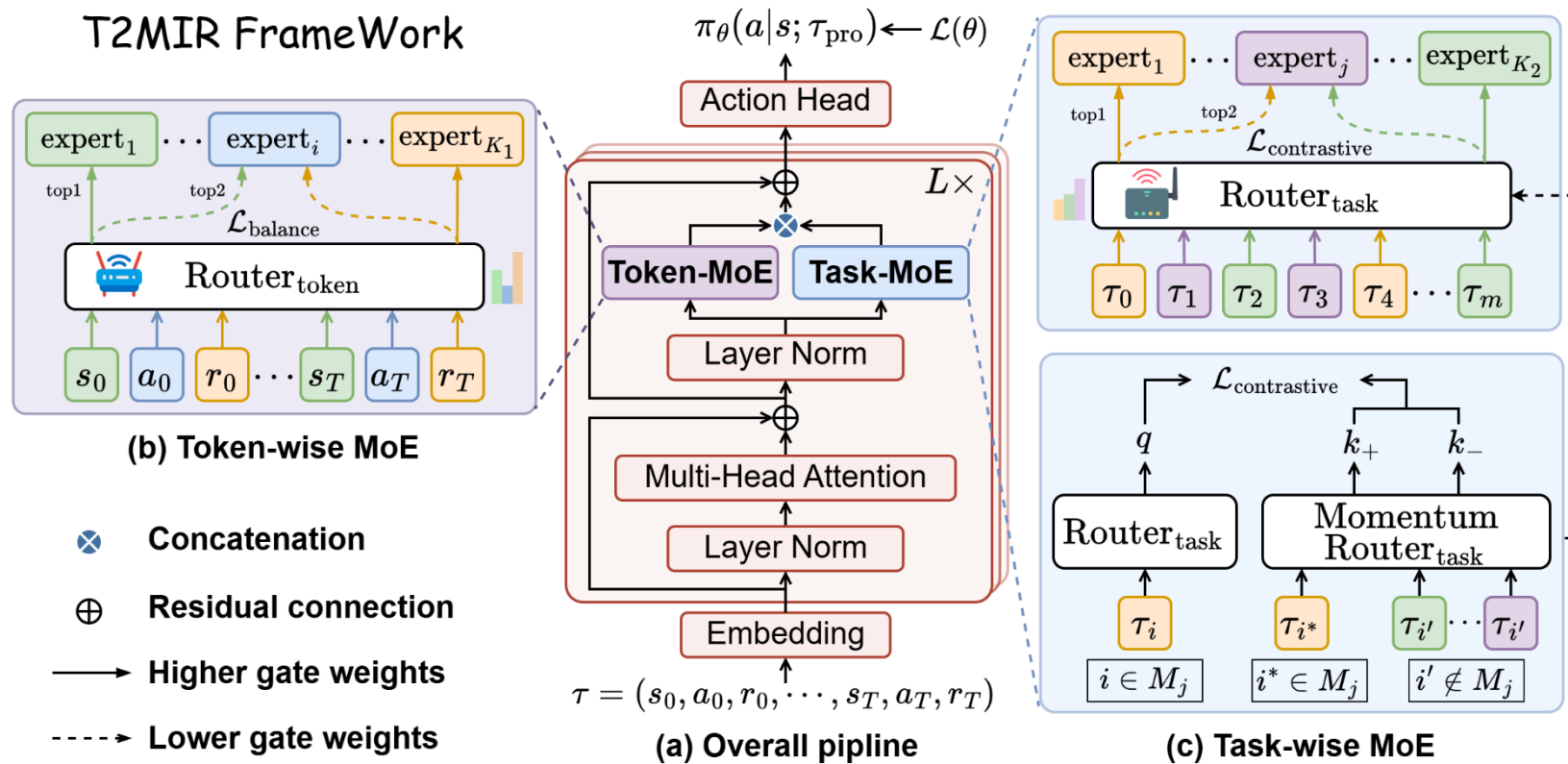
- **two parallel MoE layers**

- Token-wise MoE

✓ balance loss

✓ enables different experts to process tokens with distinct semantics

- Task-wise MoE

✓ InfoNCE loss

✓ effectively manages a broad task distribution

✓ alleviate gradient conflicts



**The overview of T2MIR**

T2MIR FrameWork

(b) Token-wise MoE

$\otimes$ Concatenation

$\oplus$ Residual connection

⟶ Higher gate weights

⤑ Lower gate weights

$\tau = (s_0, a_0, r_0, \cdots, s_T, a_T, r_T)$

(a) Overall pipline

(c) Task-wise MoE

$$\mathcal{L}(\theta) = -\sum_{n=1}^{N} \sum_{t=0}^{T-1} \log P_\theta\left(a = a_t^n \mid \tau_{\mathrm{pro}}^n, s_t^n\right), \tau_{\mathrm{pro}}^n = h_{t-1}^n$$

$$\mathcal{L}(\theta) = -\sum_{n=1}^{N} \sum_{t=0}^{T} \log P_\theta\left(a = a_t^{n*} \mid \tau_{\mathrm{pro}}^n, s_t^n\right), \tau_{\mathrm{pro}}^n \sim \mathcal{D}^n$$

➢ **motivation: intrinsic multi-modality in state-action-reward sequence**

- semantic gap among states, actions and rewards

- router $G_{\text{tok}}$ learns to assign each token to specific experts at the token level

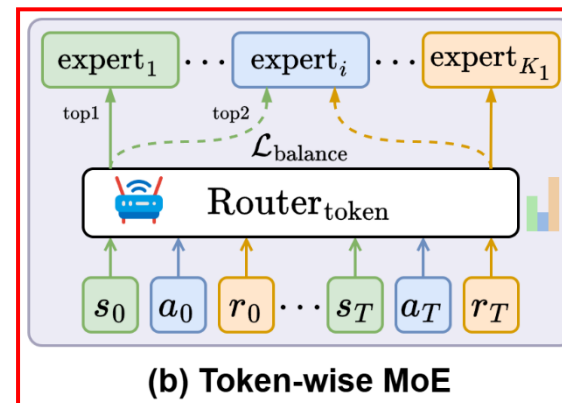$$y_{\text{tok}} = \sum_{i=1}^{K_1} w_{\text{tok}}(i; h) \cdot E_{\text{tok}}(i|h)$$

$$w_{\text{tok}}(i; h) = \text{softmax}\left(\text{topk}(G_{\text{tok}}(i|h))\right)[i]$$
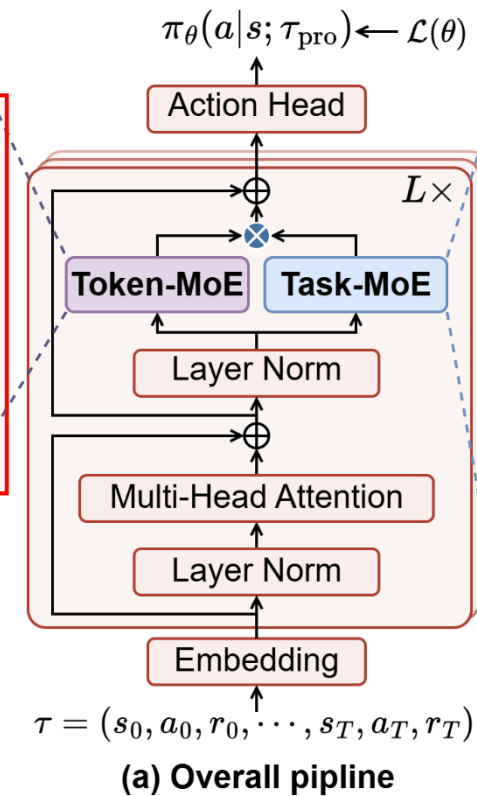
- balance expert utilization with balance loss

$$\mathcal{L}_{\text{balance}} = w_{\text{imp}} \cdot CV\left(\text{Imp}(h)\right)^2 + w_{\text{load}} \cdot CV\left(\text{Load}(h)\right)^2$$

**Token-wise MoE**

T2MIR FrameWork

$\pi_\theta(a|s; \tau_{\text{pro}}) \leftarrow \mathcal{L}(\theta)$



(b) Token-wise MoE

- ⊗ **Concatenation**
- ⊕ **Residual connection**
- ⟶ **Higher gate weights**
- ╌╌▸ **Lower gate weights**

Action Head

$L\times$

| Token-MoE | Task-MoE |

Layer Norm

Multi-Head Attention

Layer Norm

Embedding

$\tau = (s_0, a_0, r_0, \cdots, s_T, a_T, r_T)$

(a) Overall pipline

**Task-wise MoE**

$\pi_\theta(a|s;\tau_{\text{pro}}) \leftarrow \mathcal{L}(\theta)$

Action Head

Token-MoE  Task-MoE

Layer Norm

Multi-Head Attention

Layer Norm

Embedding

$\tau = (s_0, a_0, r_0, \cdots, s_T, a_T, r_T)$

**(a) Overall pipline**

$\text{expert}_1 \cdots \text{expert}_j \cdots \text{expert}_{K_2}$

top1    top2

$\mathcal{L}_{\text{contrastive}}$

Router$_{\text{task}}$

$\tau_0 \; \tau_1 \; \tau_2 \; \tau_3 \; \tau_4 \cdots \tau_m$

$\mathcal{L}_{\text{contrastive}}$

$q$      $k_+$    $k_-$

Router$_{\text{task}}$    Momentum Router$_{\text{task}}$

$\tau_i$     $\tau_{i*}$    $\tau_{i'} \cdots \tau_{i'}$
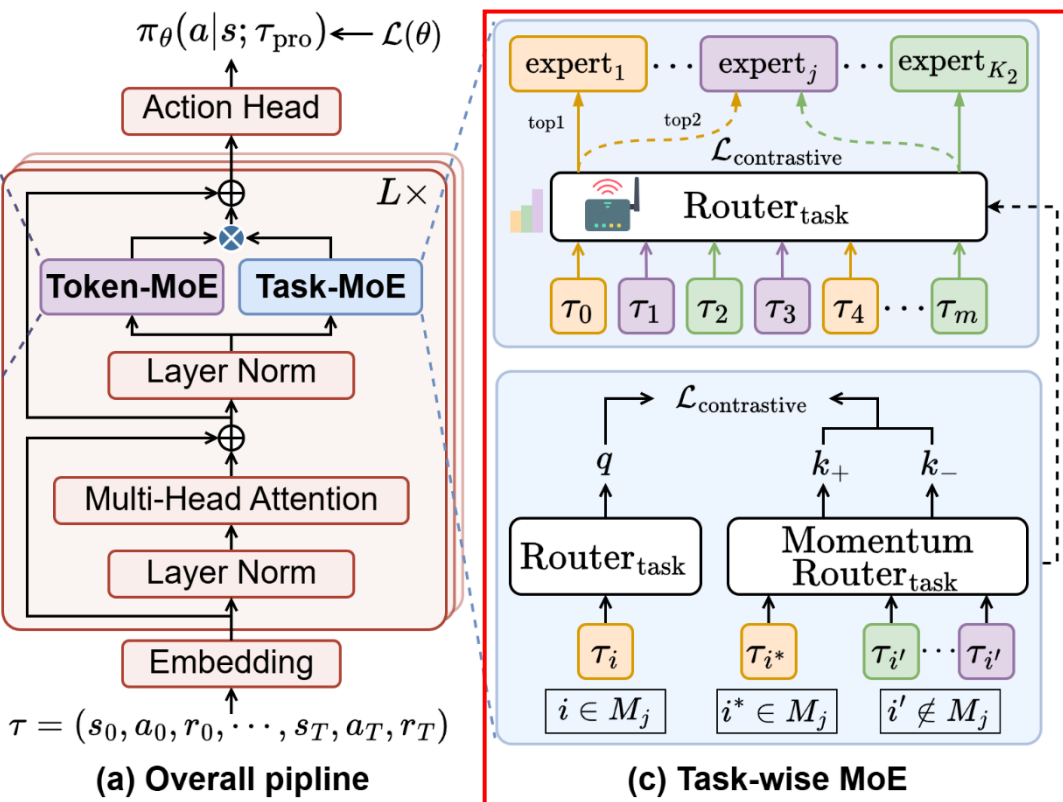
$i \in M_j$   $i^* \in M_j$   $i' \notin M_j$

**(c) Task-wise MoE**

- ➤ **motivation: task diversity and heterogeneity**

- some tasks are similar and others differ significantly

- learning efficiency can be impeded by intrinsic gradient conflicts in scenarios with significant task variation

- router $G_{\text{task}}$ learns to assign tokens to specialized experts at the task level

$$y_{\text{task}} = \sum_{i=1}^{K_2} w_{\text{task}}(i; \bar{h}) \cdot E_{\text{task}}(i|\bar{h})$$

$$w_{\text{task}}(i; \bar{h}) = \text{softmax}\left(\text{topk}\left(G_{\text{task}}(i|\bar{h})\right)\right)[i]$$

- view $G_{\text{task}}$ as a task encoder, balance expert utilization with InfoNCE loss

$$\mathcal{L}_{\text{contrastive}} = \frac{\sum_{i^* \in M_j} \exp(z_i^\top W z_{i^*})}{\sum_{i^* \in M_j} \exp(z_i^\top W z_{i^*}) + \sum_{i' \notin M_j} \exp(z_i^\top W z_{i'})}$$
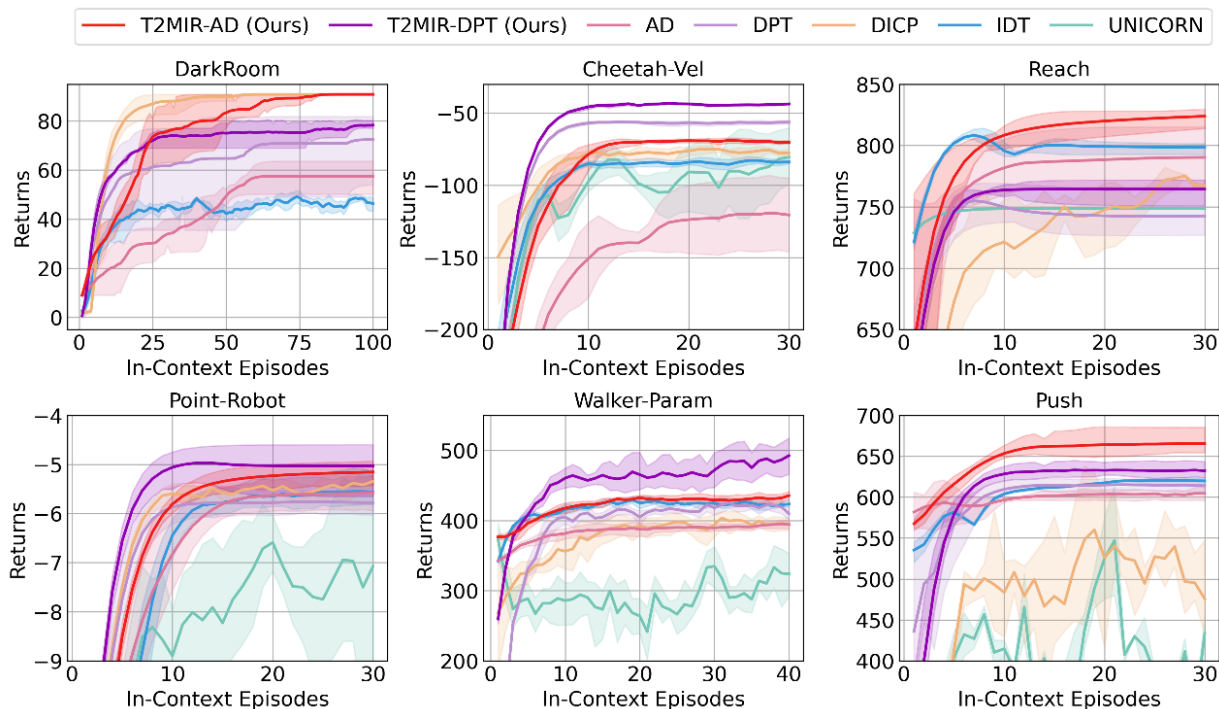
Figure 3: Test return curves of two T2MIR implementations against baselines using Mixed datasets
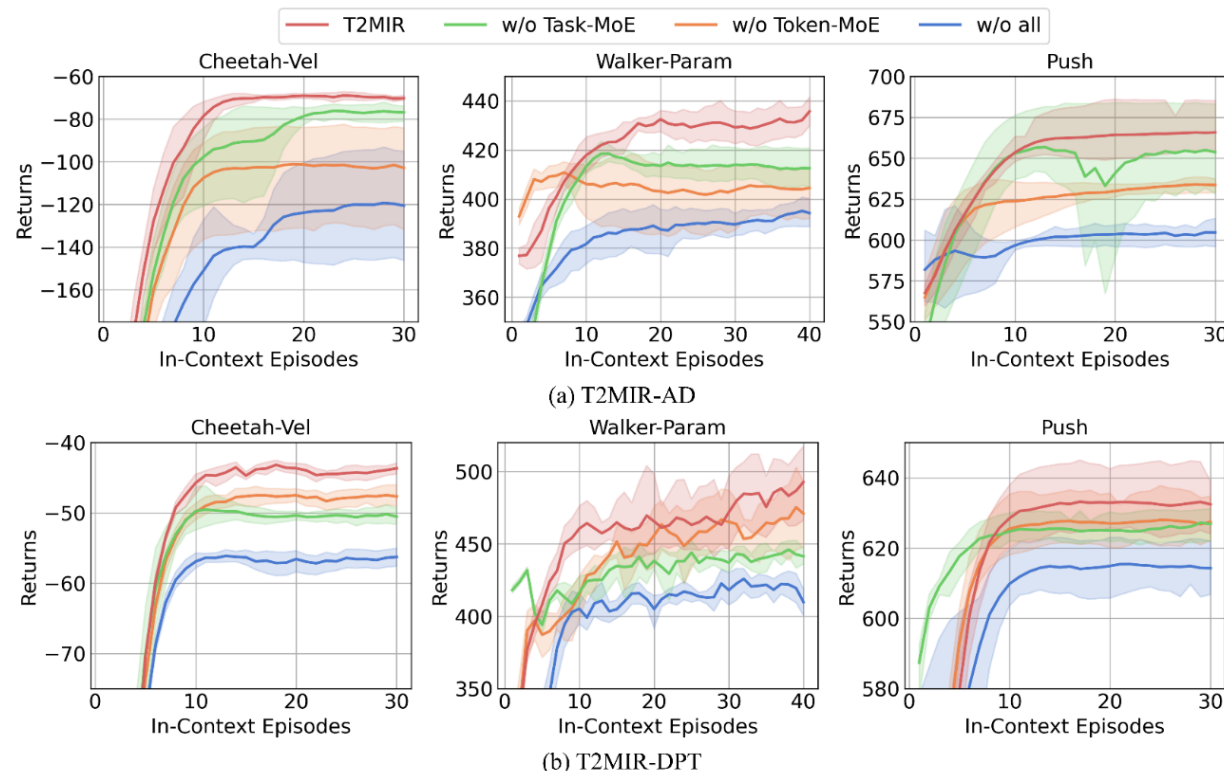
Figure 4: Ablation results of both T2MIR-AD and T2MIR-DPT architectures using Mixed datasets
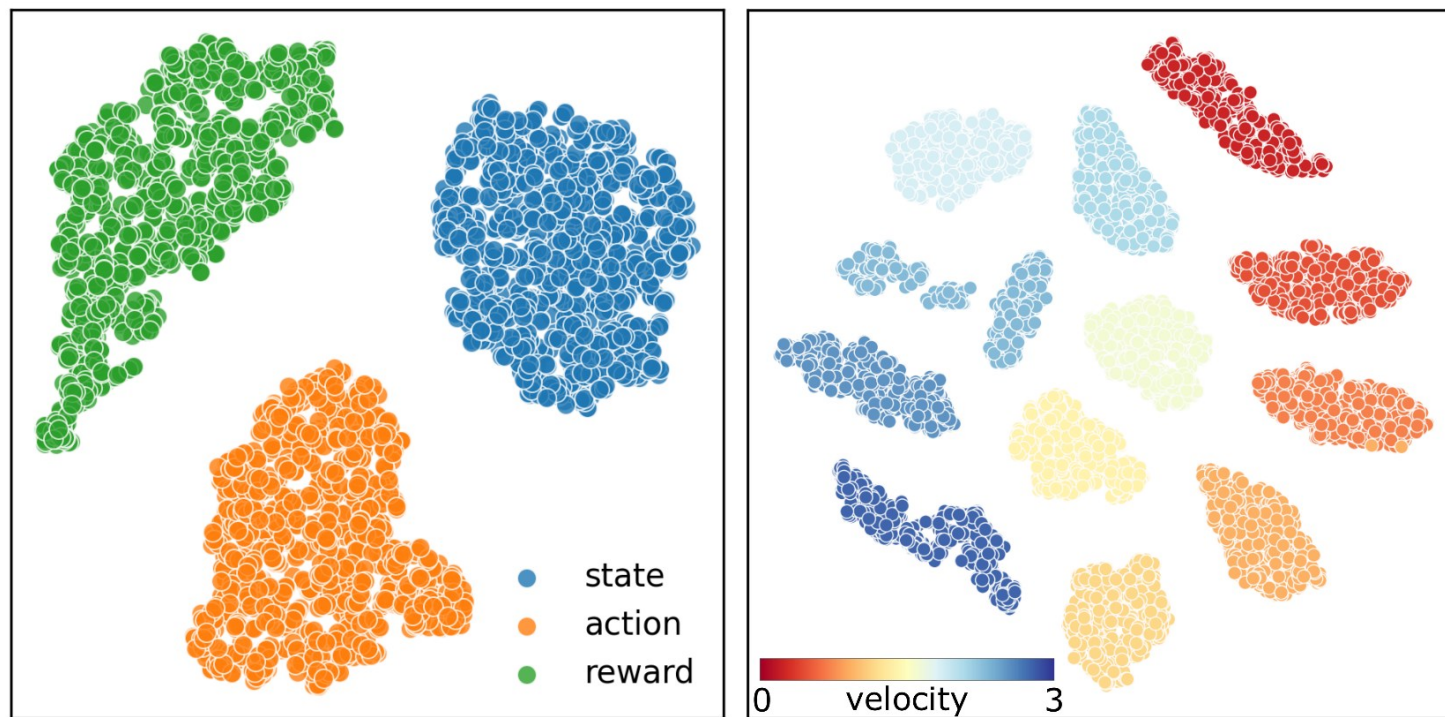
Figure 1: t-SNE visualization of expert assignments on Cheetah-Vel where tasks differ in target velocities. Left: token-wise MoE enables different experts to process tokens with distinct semantics. Right: task-wise MoE effectively manages a broad task distribution, where the difference between expert assignments is positively related to the difference between tasks.
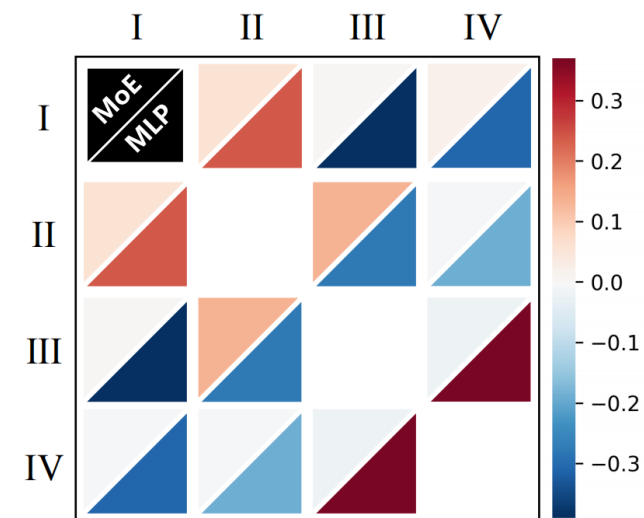


Figure 8: Cosine similarity of gradients between Point-Robot tasks in four quadrants (I-IV), comparing T2MIR-AD (MoE) with AD (MLP).

- **Innovative Framework**: This study introduces T2MIR, a novel framework that integrates the mixture-of-experts (MoE) architecture into transformer-based decision models for in-context reinforcement learning (ICRL).

- **Key Contributions**: The proposed Token-wise MoE effectively handles multi-modal inputs by capturing distinct semantics of input tokens. The Task-wise MoE manages a broad task distribution and reduces gradient conflicts through specialized experts and contrastive learning-enhanced task routing.

- **Significant Advantages**: T2MIR significantly boosts in-context learning capacity. It demonstrates superior performance over various baselines across multiple benchmarks, proving its effectiveness in advancing ICRL.

- **Future Prospects**: An urgent improvement is to evaluate on more complex environments such as XLand-MiniGrid with huge datasets, unlocking the scaling properties of MoE in ICRL domains. Another step is to deploy our method to vision-language-action (VLA) tasks that naturally involve more complex input multi-modality and task diversity.

# Thank you!

Wenhao Wu (吴文浩)

Department of Control Science and Intelligent Engineering
Nanjing University

Email: wenhaowu@smail.nju.edu.cn

Code: https://github.com/NJU-RL/T2MIR

Paper: https://arxiv.org/abs/2506.05426