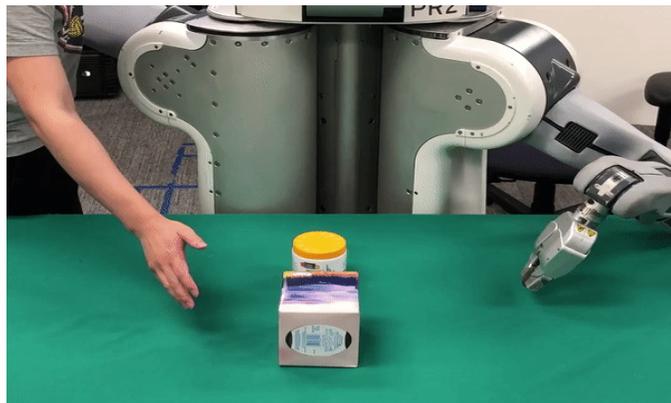# Learning from Visual Observation via Offline Pretrained State-to-Go Transformer

**Bohan Zhou, Ke Li, Jiechuan Jiang, Zongqing Lu**
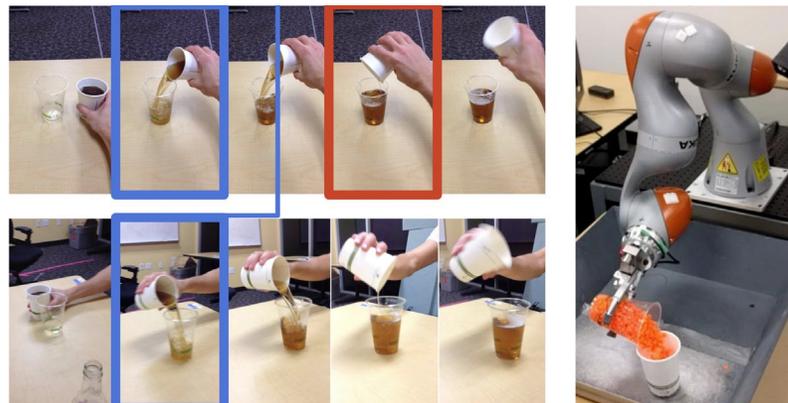
24/10/2023

# Motivation



Learning from Demonstrations
(LfD)

**+** Easy to learn

**–** Hard & expensive annotations



Learning from Visual Ovservations
(LfVO)

**+** No actions or rewards

**+** An ocean of Internet videos

**+** Explore unknown expert policy

**–** Hard to extract useful experience

**From LfD to LfVO**

✓ **Less Supervision**

✓ **Enlarging resource**

✓ **Biologically reasonable**
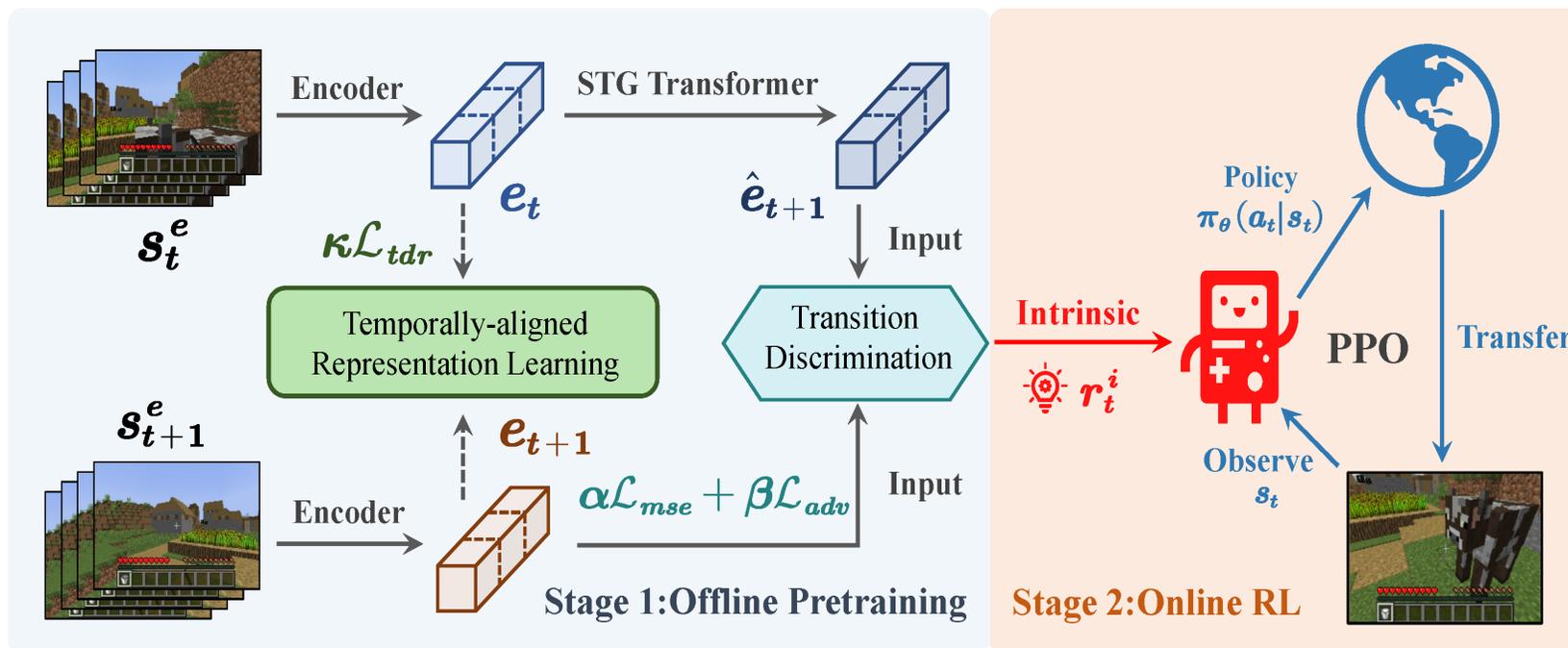
- **IDM-based methods – extra component, compounding error**

- **Adversarial methods – sample-inefficient online learning schemes**

- **Representation-learning-based methods – over-optimistic estimation**

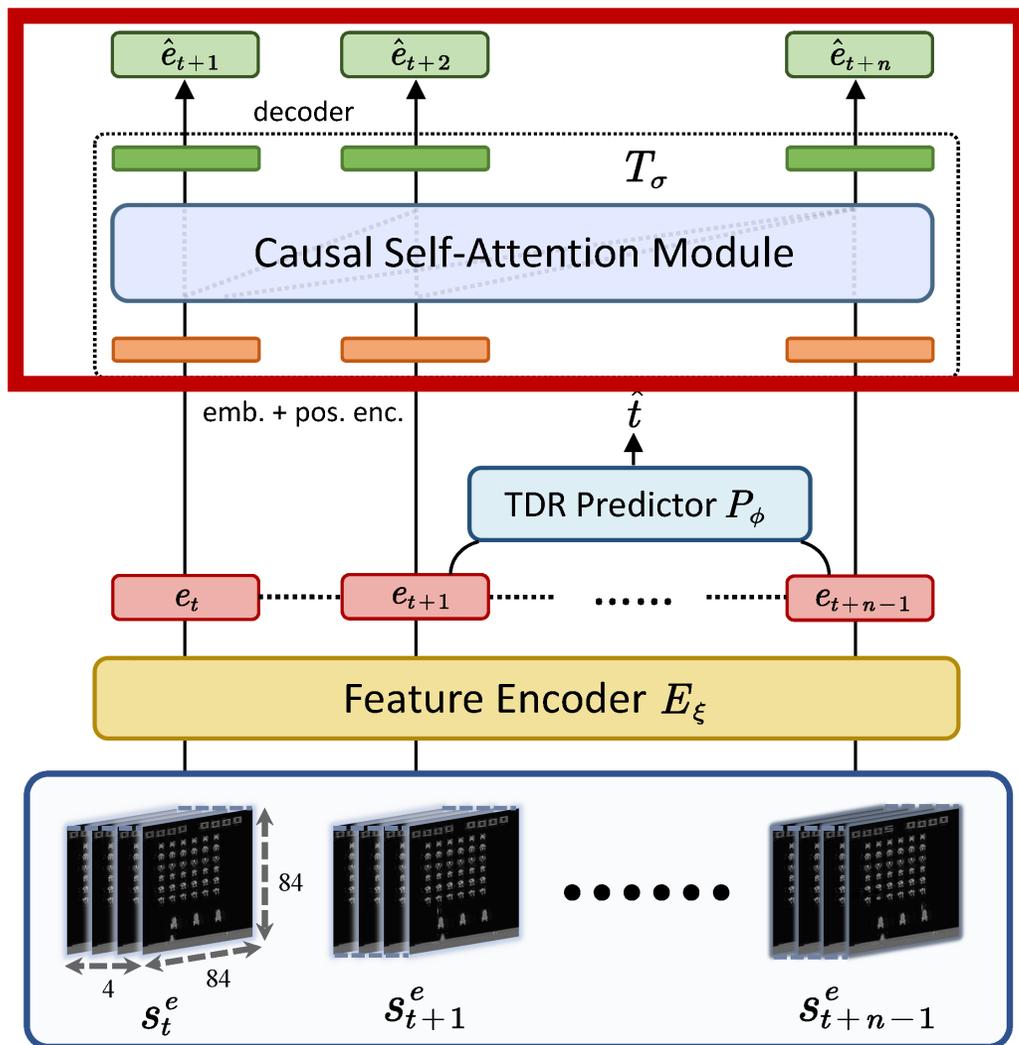- **Goal-oriented methods – extra task-specific information**

Abundant **video-only** data contain useful behavior patterns. How can we effectively leverage them to tackle downstream **reward-free** visual control tasks?

# Two-stage framework



- **Pretraining stage**: we simultaneously learn a **GPT** for latent transition prediction, an expert transition **discriminator** for intrinsic rewards and a temporal distance regressor (**TDR**) for temporally-aligned representations.

- **Reinforcecment learning stage**: agents **merely** learn from generated rewards from discriminator without environmental reward signals.

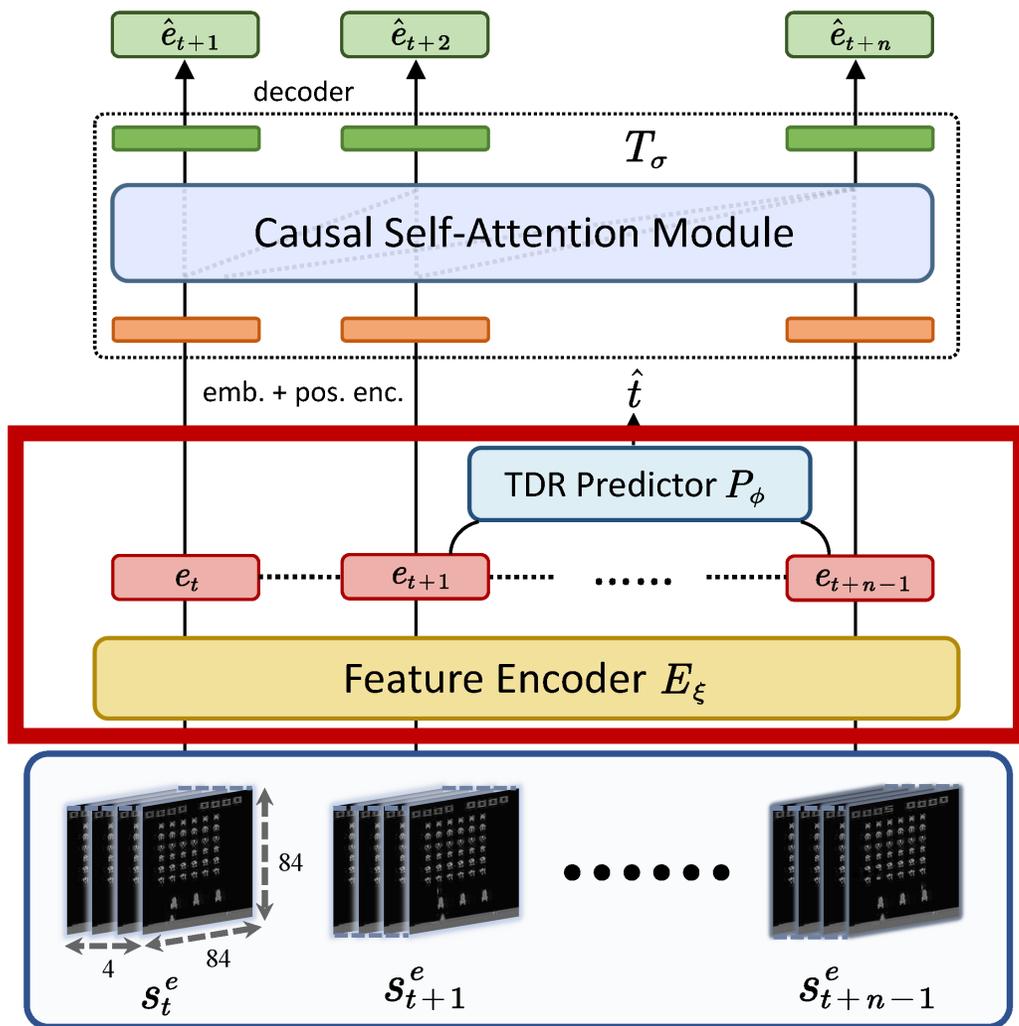## 1. Predicting Latent Transition

Adversarially learn transition module with L2 regularization as well as a WGAN discriminator

for $D_\omega$: $\min\limits_{w \in \mathcal{W}} \mathbb{E}_{\mathcal{D}^e} \left[ D_\omega(e_t, \hat{e}_{t+1}) - D_\omega(e_t, e_{t+1}) \right]$

for $T_\sigma$: $\min\limits_{\xi,\sigma} \mathbb{E}_{\mathcal{D}^e} \left[ -D_\omega(e_t, \hat{e}_{t+1}) + \|\hat{e}_{t+1} - e_{t+1}\|_2^2 \right]$
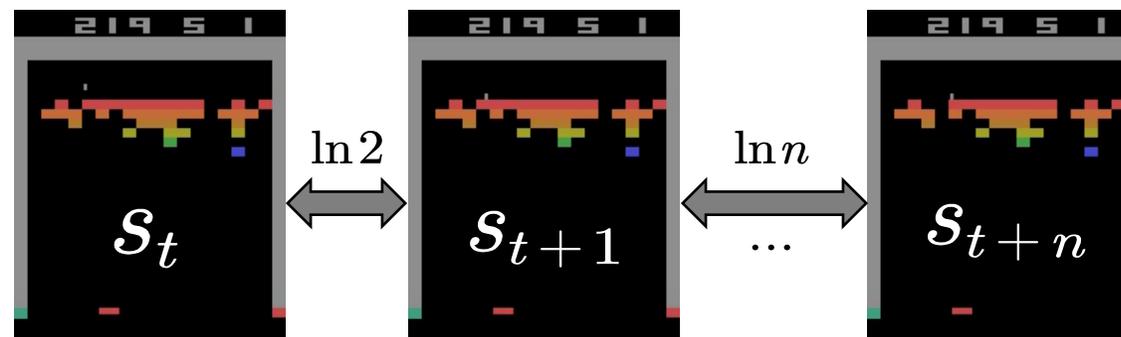
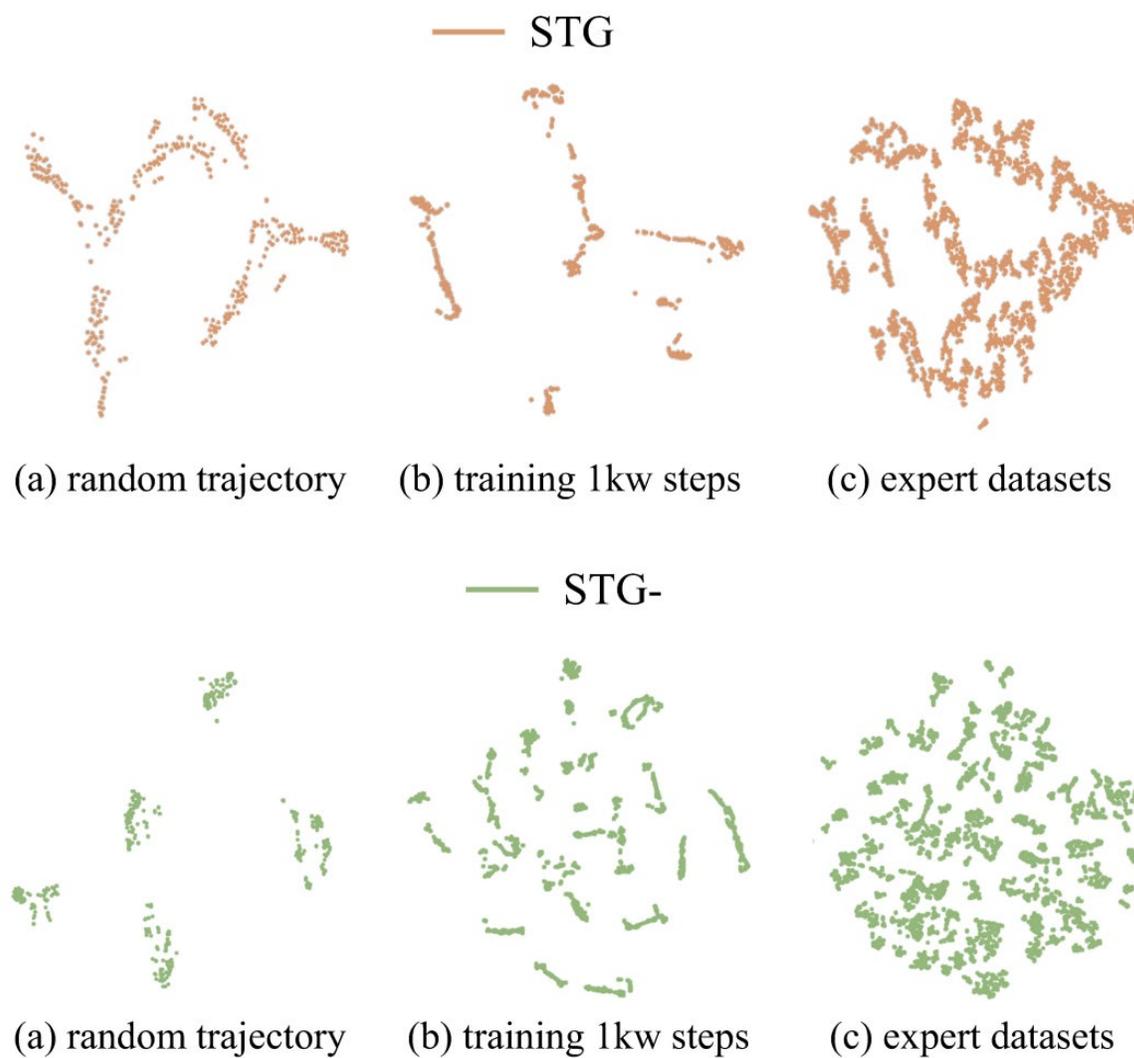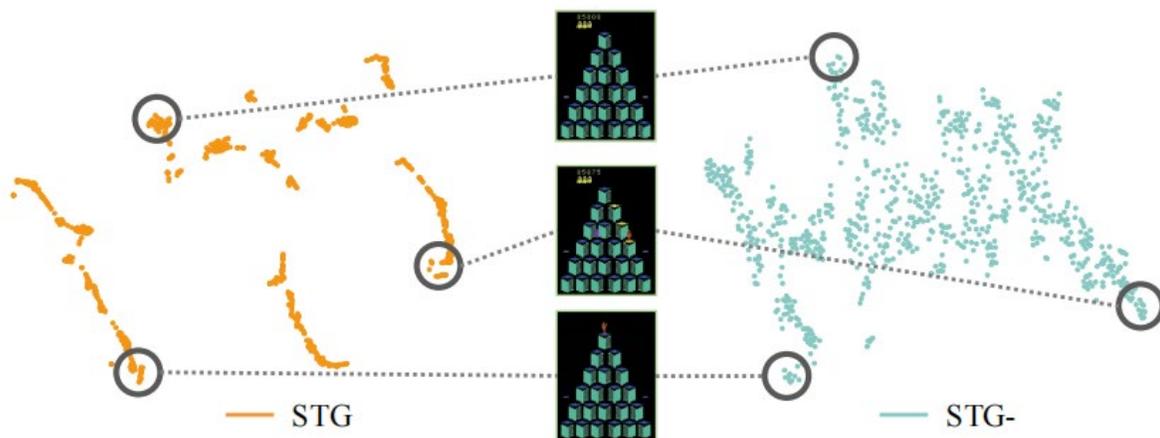$$e_t = E_\xi(s_t), \quad \hat{e}_{t+1} = T_\sigma(e_t)$$

## 2. Learning Temporally-Aligned Representation

Apply symlog temporal distance prior in low-dimensional representation space

$$\min_{\xi,\phi}\mathbb{E}_{\mathcal{D}^e}\left\|P_\phi(e_t,e_{t+j})-\text{sign}(j)\ln(1+|j|)\right\|$$

(a) random trajectory    (b) training 1kw steps    (c) expert datasets

— STG-



(a) random trajectory    (b) training 1kw steps    (c) expert datasets

— STG

— STG-

**Algorithm 1** STG Transformer Offline Pretraining

**Input:** STG Transformer $T_\sigma$, feature encoder $E_\xi$, discriminator $D_\omega$, expert dataset $D^e = \{\tau^1, \tau^2, \ldots, \tau^m\}, \tau^i = \{s_1^i, s_2^i, \ldots\}$, buffer $\mathcal{B}$, loss weights $\alpha, \beta, \kappa$.

1: Initialize parametric network $E_\xi, T_\sigma, D_\omega$ randomly.
2: **for** $e \leftarrow 0, 1, 2 \ldots$ **do**                                            ▷ epoch
3:      Empty buffer $\mathcal{B}$.
4:      **for** $b \leftarrow 0, 1, 2 \ldots |\mathcal{B}|$ **do**                          ▷ batchsize
5:          Stochastically sample state sequence $\tau^i$ from $D^e$.
6:          Stochastically sample timestep $t$ and $n$ adjacent states $\{s_t^i, \ldots, s_{t+n-1}^i\}$ from $\tau^i$.
7:          Store $\{s_t^i, \ldots, s_{t+n-1}^i\}$ in $\mathcal{B}$.
8:      **end for**
9:      Update $D_\omega$: $\omega \leftarrow \text{clip}(\omega - \epsilon \nabla_\omega \mathcal{L}_{dis}, -0.01, 0.01)$.
10:     Update $E_\xi$ and $T_\sigma$ concurrently by minimizing total loss $\alpha \mathcal{L}_{mse} + \beta \mathcal{L}_{adv} + \kappa \mathcal{L}_{tdr}$.
11: **end for**

Pretrained WGAN discriminator works as reward function:

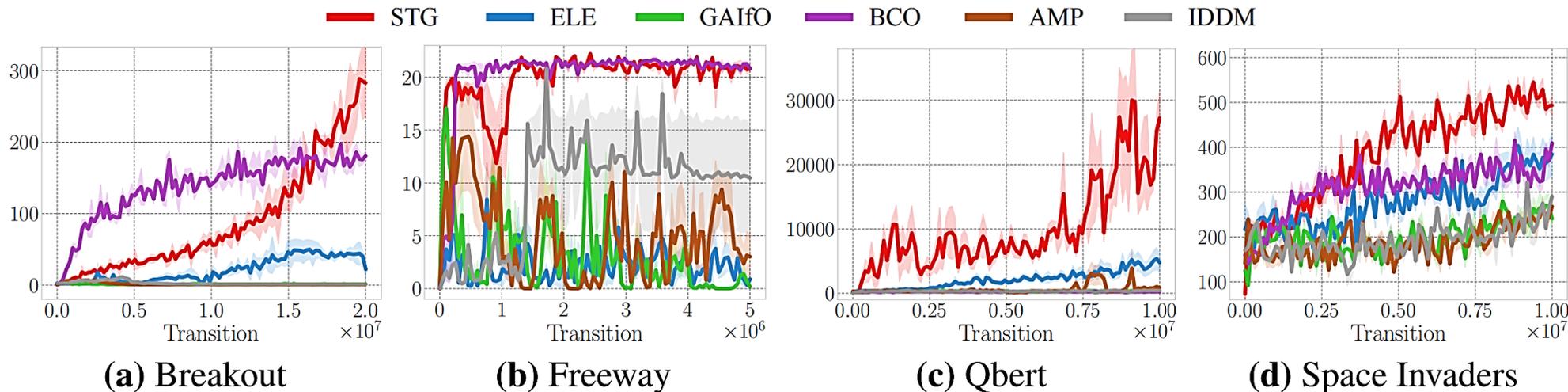$$r_t^i = -\left[ D_\omega\big(E_\xi\left(s_t\right), T_\sigma\left(E_\xi\left(s_t\right)\right)\big) - D_\omega\big(E_\xi\left(s_t\right), E_\xi\left(s_{t+1}\right)\big) \right]$$

---

**Algorithm 2** Online Reinforcement Learning with Intrinsic Rewards

---

**Input:** pretrained $E_\xi, T_\sigma, D_\omega$, policy $\pi_\theta$, MDP $\mathcal{M}$, intrinsic coefficient $\eta$.
1: Initialize parametric policy $\pi_\theta$ with random $\theta$ randomly and reset $\mathcal{M}$.
2: **while** updating $\pi_\theta$ **do**                                    ▷ policy improvement
3:      Execute $\pi_\theta$ and store the resulting $n$ state transitions $\{(s, s')\}_t^{t+n}$.
4:      Use $E_\xi$ to obtain $n$ real latent transitions $\{(e, e')\}_t^{t+n}$.
5:      Use $T_\sigma$ to obtain $n$ predicted latent transitions $\{(e, \hat{e}')\}_t^{t+n}$.
6:      Use $D_\omega$ to calculate intrinsic rewards: $\Delta_t^{t+n} = \{D_\omega(e, \hat{e}')\}_t^{t+n} - \{D_\omega(e, e')\}_t^{t+n}$.
7:      Perform PPO update to improve $\pi_\theta$ with respect to $r^i = -\eta\Delta$.
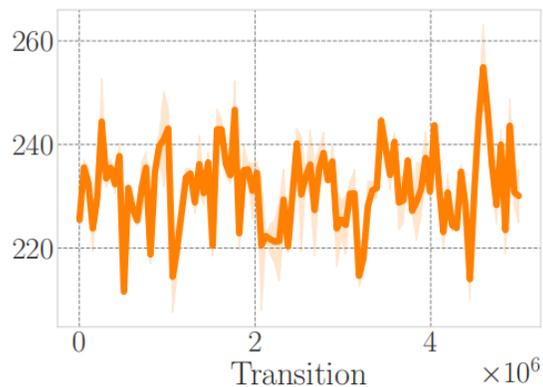8: **end while**

---

**(a)** Breakout  **(b)** Freeway  **(c)** Qbert  **(d)** Space Invaders

| Environment | GAIfO | AMP | IDDM | ELE | BCO | STG | Expert | PPO |
|---|---|---|---|---|---|---|---|---|
| Breakout | 1.5 | 0.6 | 1.2 | 22.0 | 180.4 | **288.8** | 212.5 | 274.8 |
| Freeway | 0.6 | 3.0 | 10.5 | 2.7 | 21.6 | 21.8 | 31.9 | 32.5 |
| Qbert | 394.4 | 874.9 | 423.3 | 4698.6 | 234.1 | **27234.1** | 15620.7 | 14293.3 |
| Space Invaders | 260.2 | 268.1 | 290.4 | 384.6 | 402.2 | 502.1 | 1093.9 | 942.5 |

Learning from **50** trajectories for each task, STG demonstrates **superiority among baselines** and even **surpass expert level**.

(a) Breakout

(b) Freeway

(c) Qbert

(d) Space Invaders

The rising trend of **intrinsic return** proves that online collected observation distribution is getting **closer** to expert observation distribution during training.

# Minecraft Experiments



**(a)** Pick a flower  **(b)** Milk a cow  **(c)** Harvest tallgrass  **(d)** Gather wool

In challenging **open-ended** Minecraft tasks, shows superiority over baselines!
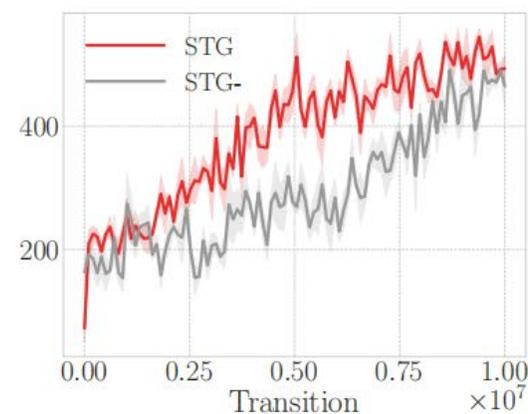
**(a)** Breakout      **(b)** Freeway      **(c)** Qbert      **(d)** Space Invaders

Pretrained on whole Atari datasets, STG-Multi shows comparable performance.

**(a)** Breakout  **(b)** Freeway  **(c)** Qbert  **(d)** Space Invaders
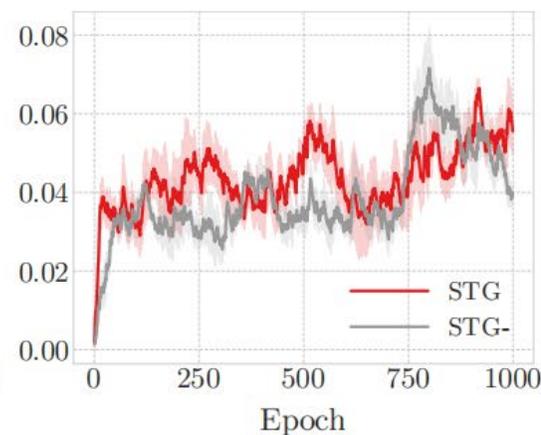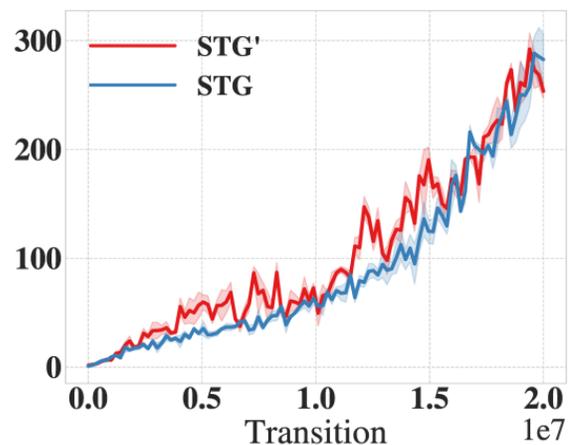
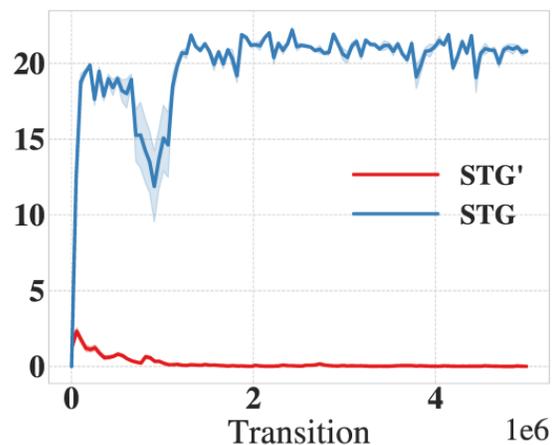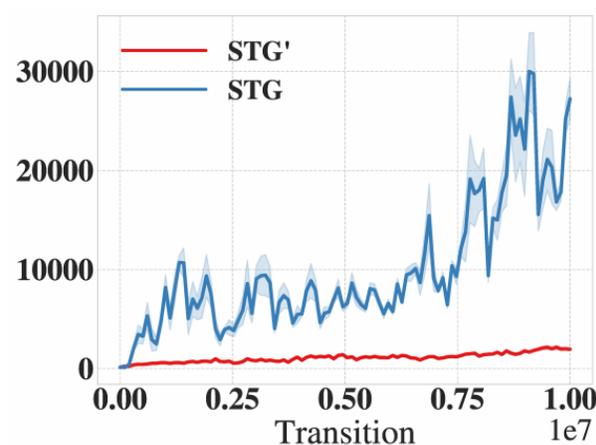**(e)** Pick a flower  **(f)** Milk a cow  **(g)** Harvest tallgrass  **(h)** Gather wool

$$r_t^i = D_\omega\big(E_\xi(s_t), E_\xi(s_{t+1})\big) - D_\omega\big(E_\xi(s_t), T_\sigma(E_\xi(s_t))\big) = r_t^{guide} - r_t^{base}$$
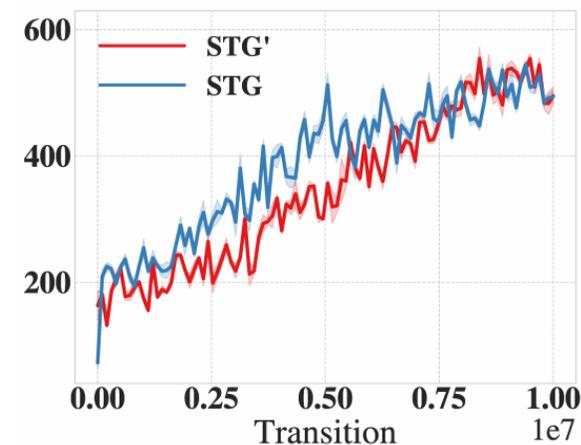


**(a)** Breakout  **(b)** Freeway  **(c)** Qbert  **(d)** Space Invaders
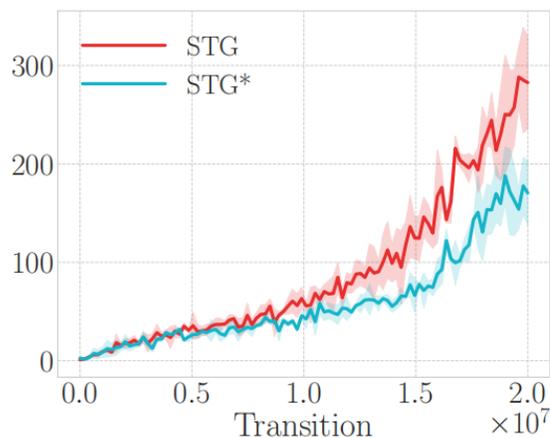
**Figure 4:** Atari experiments comparing using $r^{guide}$ (STG') and $r^i$ (STG) as intrinsic reward.
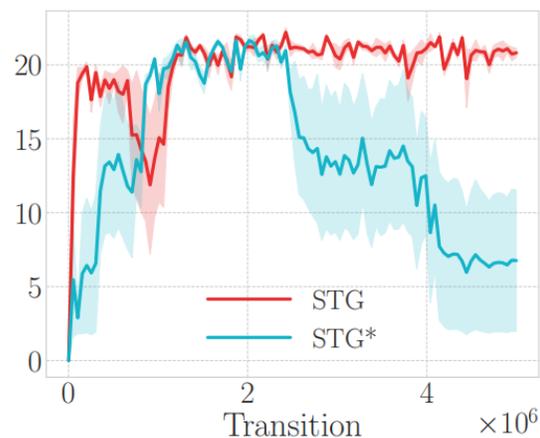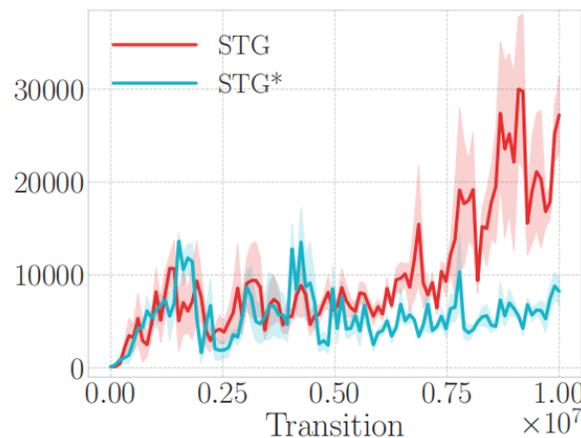
$$\begin{cases} \text{discrimination reward: } D_\omega(e_t, e_{t+1}) - D_\omega(e_t, \hat{e}_{t+1}) \\ \text{progression reward: } P_\phi(e_t, e_{t+k}), \ k=1 \end{cases}$$



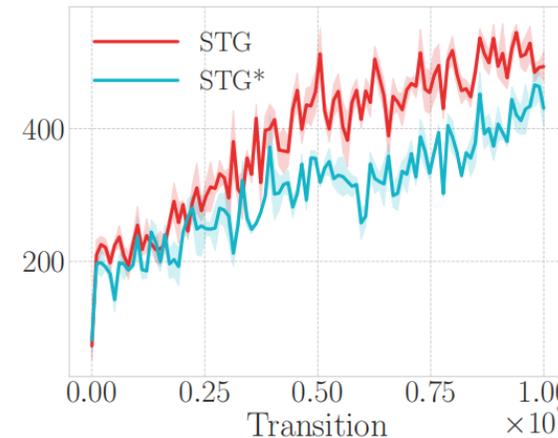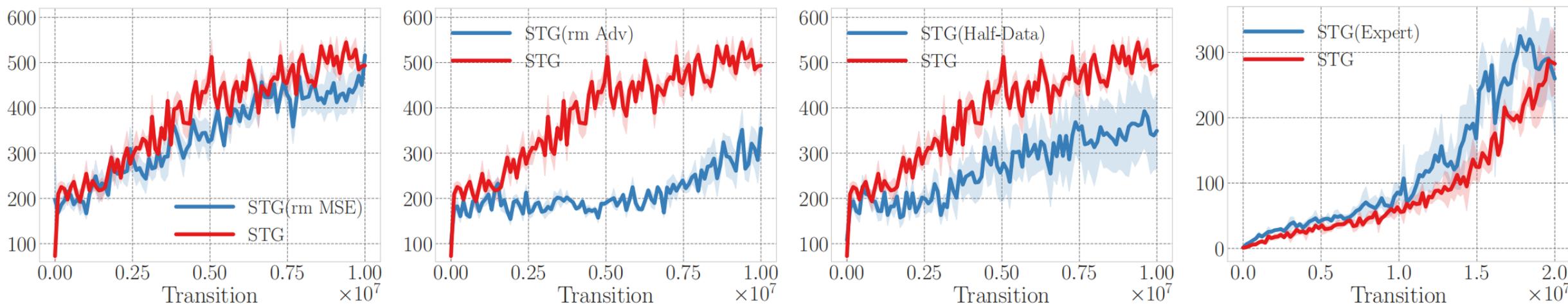**(a)** Breakout  **(b)** Freeway  **(c)** Qbert  **(d)** Space Invaders

**Figure 9:** Atari experiments comparing using discriminative rewards (STG) and using both discriminative rewards and progression rewards (STG*).

**(a)** Ablate removing $\mathcal{L}_{mse}$ **(b)** Ablate removing $\mathcal{L}_{adv}$ **(c)** Ablate dataset size **(d)** Ablate dataset quality

**Figure 8:** Learning curves of four pre-training ablations: (a) removing $\mathcal{L}_{mse}$ in SpaceInvaders; (b) removing $\mathcal{L}_{adv}$ in SpaceInvaders; (c) using half dataset to train STG in SpaceInvaders; (d) using expert dataset to train STG in Breakout.

# Extensions

STG offers an effective solution in situations with plentiful video demonstrations, limited environment interactions, and inaccessible labeled action or rewards.

In future work, STG is likely to benefit from:

- more powerful large-scale vision foundation models to facilitate generalization across a broader range of related tasks, domains or embodiments.
- hierarchical framework where one-step predicted rewards can be employed for training low-level policies and multi-step rewards for a high-level policy to tackle long-horizon tasks.