# Main focus of Meta-Learning

- Excels in adaptation to unseen but similar tasks



**Env Distribution A**
⋮
**Meta Learner A**
**→ Learning Algorithm A**
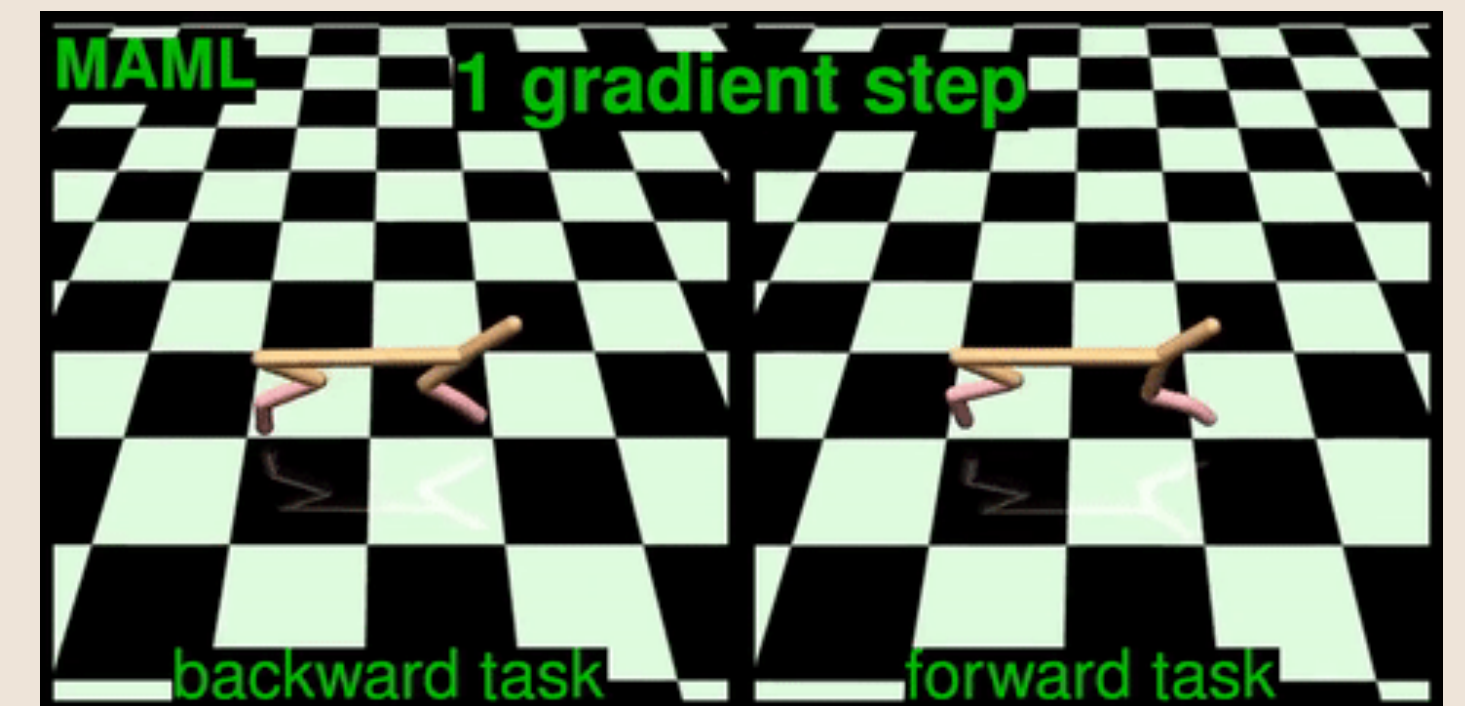
**Env Distribution B**
⋮
**Meta Learner B**
**→ Learning Algorithm B**
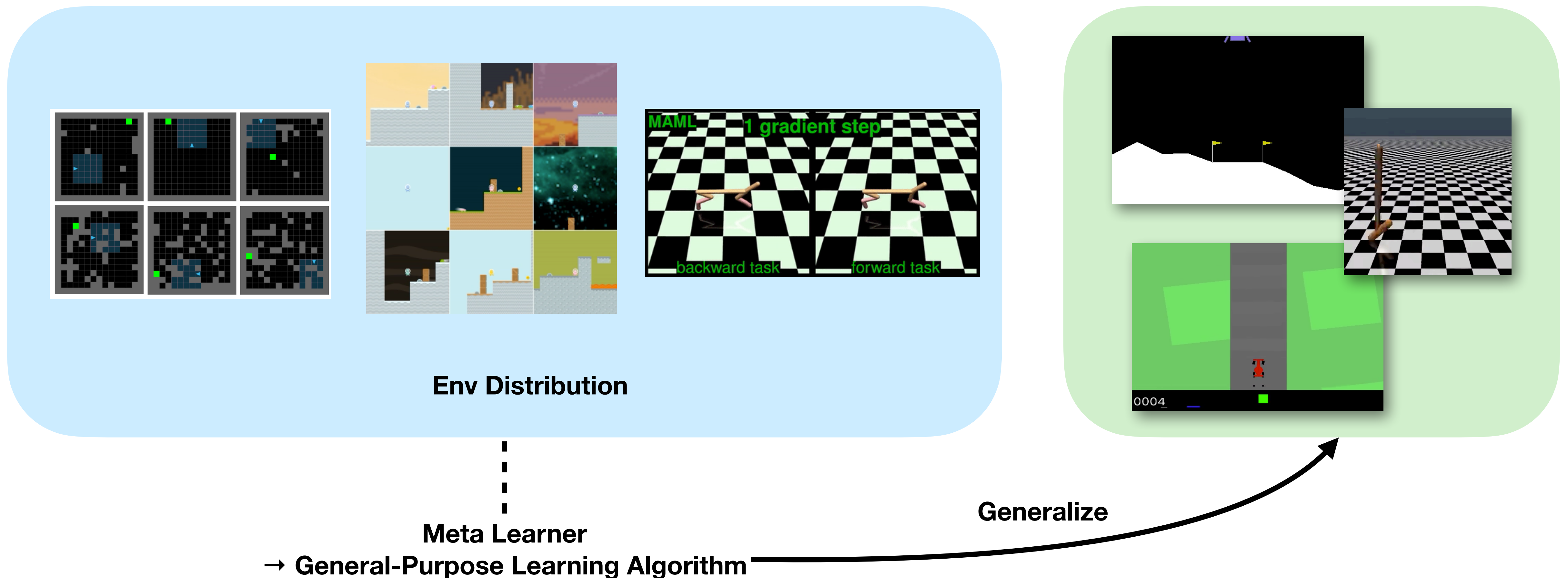
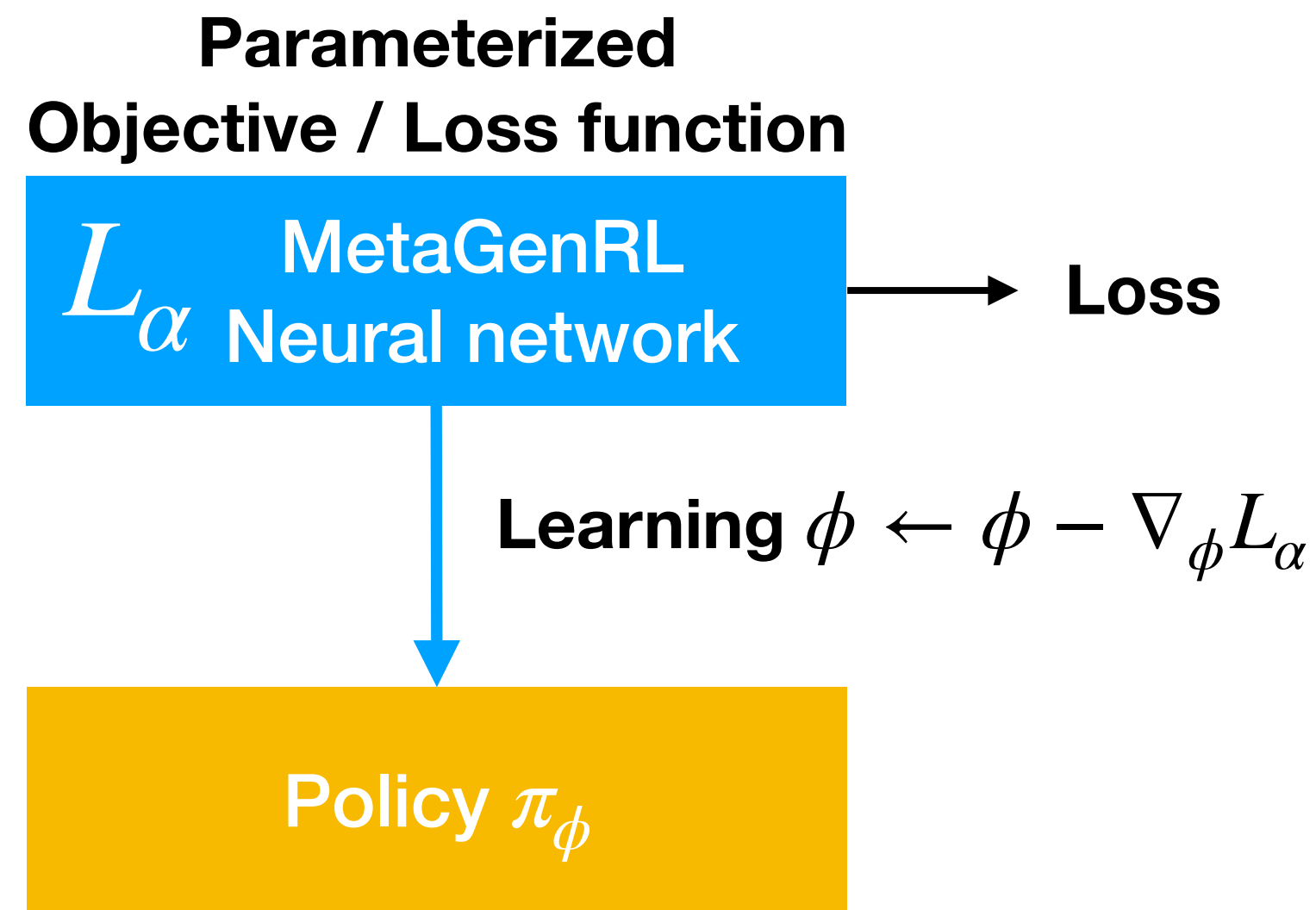**Env Distribution C**
⋮
**Meta Learner C**
**→ Learning Algorithm C**

# General-Purpose Meta-Learning

Can we train an agent that can efficiently in-context **learn and act in any environment**?



**Env Distribution**

**Meta Learner**
→ **General-Purpose Learning Algorithm**

**Generalize**

# MetaGenRL

**Parameterized Objective / Loss function**

$L_\alpha$ MetaGenRL Neural network → **Loss**

**Learning** $\phi \leftarrow \phi - \nabla_\phi L_\alpha$

Policy $\pi_\phi$
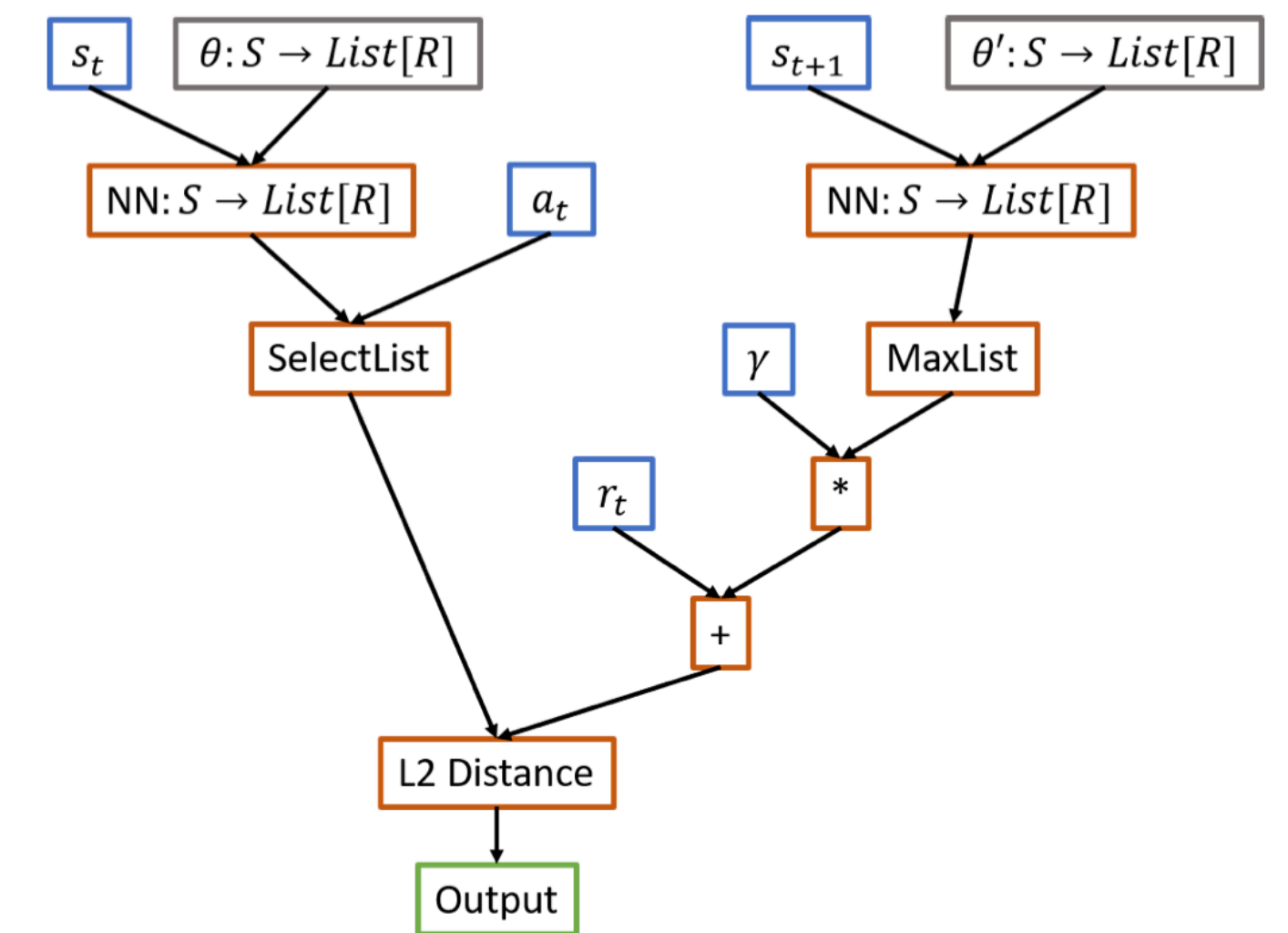
**Kirsch et al [ICLR 2020]**

# Learned Policy Gradient

**Oh et al [NeurIPS 2020]**

# Evolving RL Algorithms

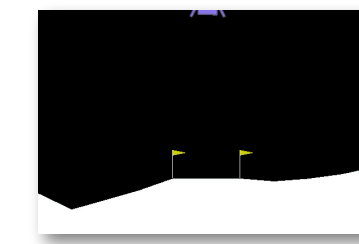**Co-Reyes et al [ICLR 2021]**

**MetaGenRL, LPG, etc**　　　　**In-context meta-RL**

$\{(o_i,)a_i, r_i\}_{i=1}^{t}$

**Gray-box**
**Learning** $\phi \leftarrow \phi - \nabla_\phi L_\alpha$

Better $\pi_\phi$

$\{o_i, a_i, r_i\}_{i=1}^{t}$

**Black-box function**
**approximator**
**e.g. LSTM, Transformer, FWP**
$\pi(a_{t+1} \,|\, o_{t+1}, \{o_i, a_i, r_i\}_{i=1}^{t})$

**= in-context learning**

Better $\pi$

**[Schmidhuber since 1992,**
**RL² Duan et al 2016,**
**Wang et al 2016,**
**SymLA Kirsch et al 2022,**
**Adaptive Agents 2023, etc]**

**Good generalization**　　**Difficult generalization**

**This paper:**
**How to fix this?**

5

**Env Distribution**

TOUCH

LIFT

(x,y) location

REPEAT

MAML    1 gradient step

backward task    forward task

Tuesday, Mar 16

0004

# GPICL: From memorization to general learn-to-learn

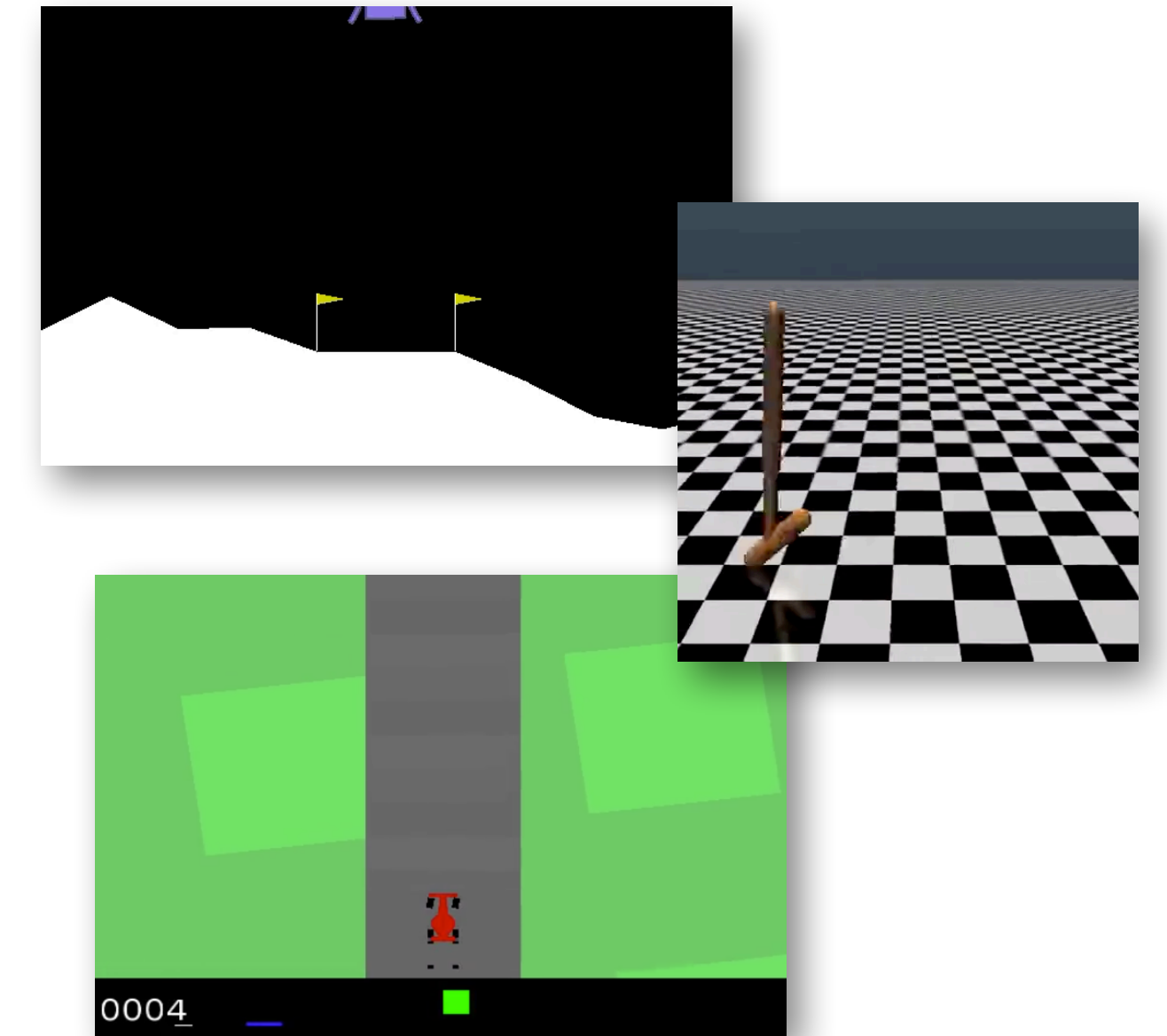**Transformers exhibit three different phases in terms of meta-learned behavior.**

| Learning | Generalization | Algorithm Description |
|---|---|---|
| × No | × No | Task memorization |
| ✓ Yes | × No | Task identification |
| ✓ Yes | ✓ Yes | General-purpose learning algorithm |

**Kirsch et al [2022]**

**Adaptation to RL**

Instead of training on a supervised dataset $\left( \{x_i, y_i\}_{i=1}^{N_D}, x' \right) \mapsto y'$

Train on RL data $\left( \{s_i, a_i, r_i, d_i\}_{i=1}^{N_D}, s \right) \mapsto a$

Our recipe: **Offline meta-training + augmented RL data $\rightarrow$ Generalization**

**Training is offline, but agent can learn online!**

# How it works – PPO Data collection

**Data Collection**



Environment Distribution

PPO Training → Training Dataset $\bar{D} = \{\tau\}$

PPO update    PPO      PPO      PPO      PPO      PPO      PPO

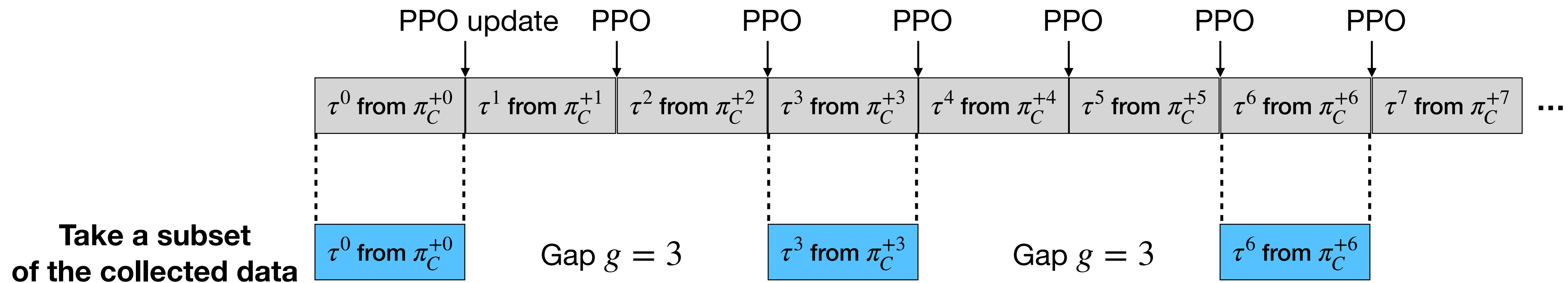| $\tau^0$ from $\pi_C^{+0}$ | $\tau^1$ from $\pi_C^{+1}$ | $\tau^2$ from $\pi_C^{+2}$ | $\tau^3$ from $\pi_C^{+3}$ | $\tau^4$ from $\pi_C^{+4}$ | $\tau^5$ from $\pi_C^{+5}$ | $\tau^6$ from $\pi_C^{+6}$ | $\tau^7$ from $\pi_C^{+7}$ | $\cdots$ |

**Take a subset of the collected data**

| $\tau^0$ from $\pi_C^{+0}$ | Gap $g = 3$ | $\tau^3$ from $\pi_C^{+3}$ | Gap $g = 3$ | $\tau^6$ from $\pi_C^{+6}$ |

# How it works - Model learning algorithm

# How it works - Augmentations

**Data Collection**

Environment Distribution

PPO Training →

Training Dataset $\bar{D} = \{\tau\}$

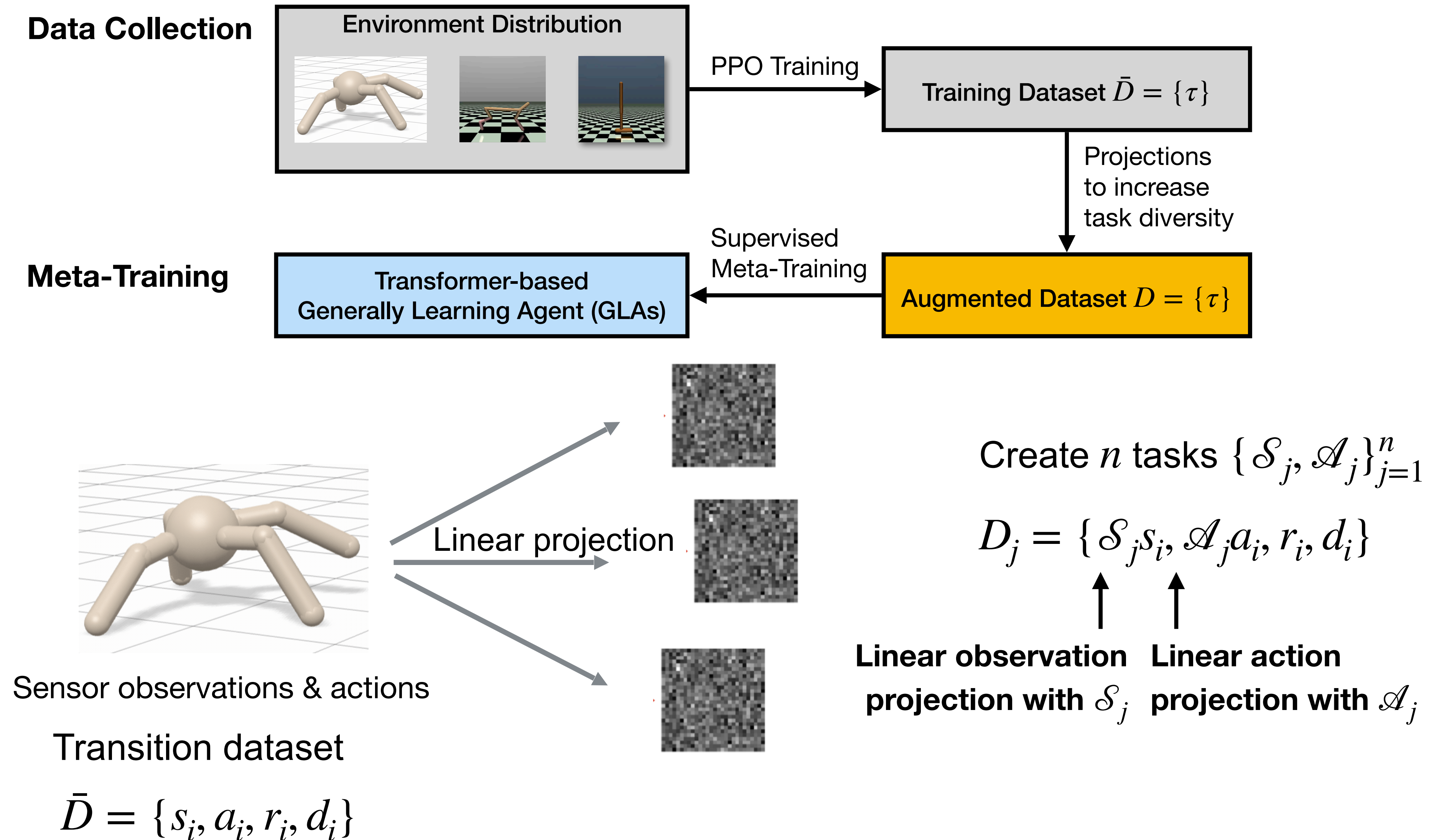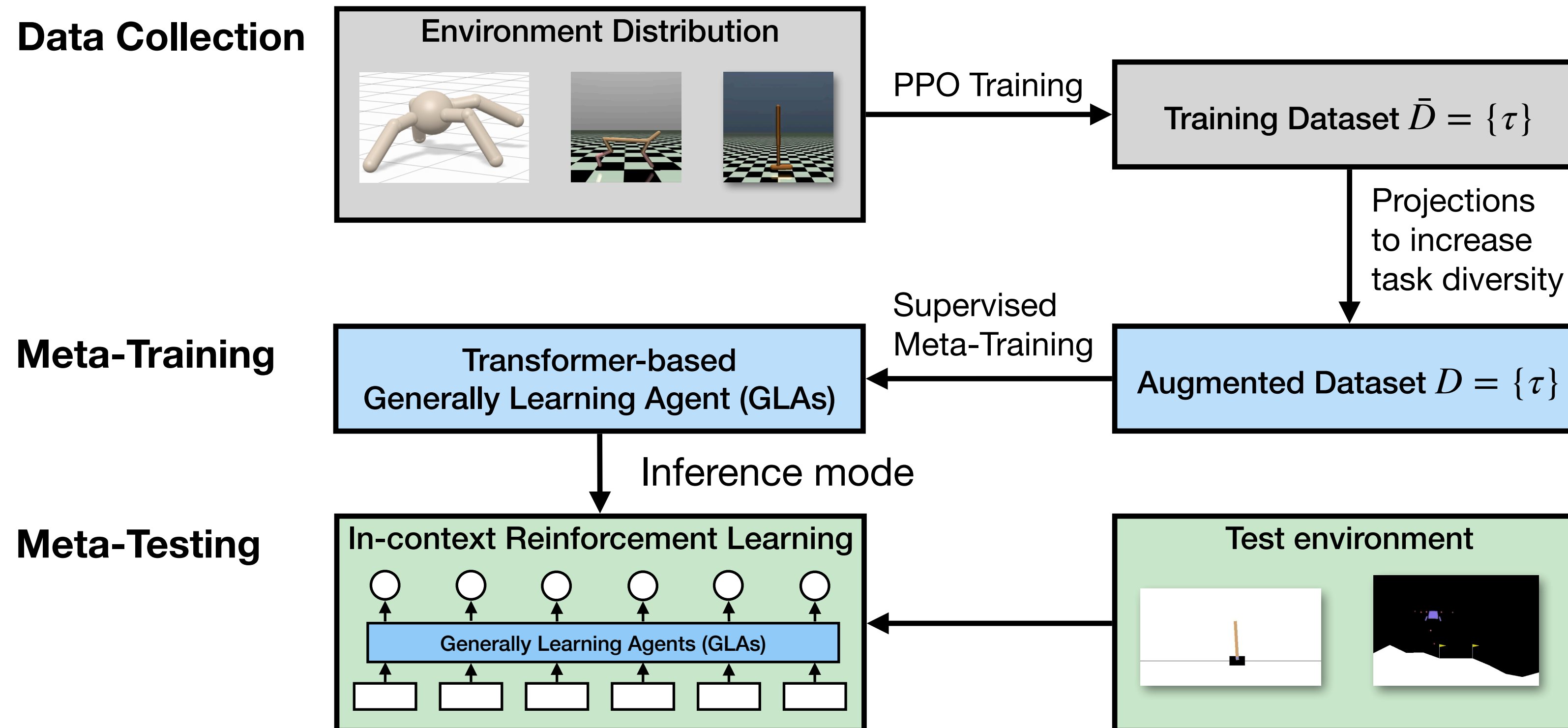Projections to increase task diversity

**Meta-Training**

Transformer-based Generally Learning Agent (GLAs)

← Supervised Meta-Training

Augmented Dataset $D = \{\tau\}$

Linear projection

Sensor observations & actions

Transition dataset

$$\bar{D} = \{s_i, a_i, r_i, d_i\}$$

Create $n$ tasks $\{\mathscr{S}_j, \mathscr{A}_j\}_{j=1}^{n}$

$$D_j = \{\mathscr{S}_j s_i, \mathscr{A}_j a_i, r_i, d_i\}$$

**Linear observation projection with** $\mathscr{S}_j$    **Linear action projection with** $\mathscr{A}_j$

# Evaluation: Single Environment

In-context learning policies that encode a learning algorithm on a specific task

Training Dataset $\bar{D} = \{\tau\}$
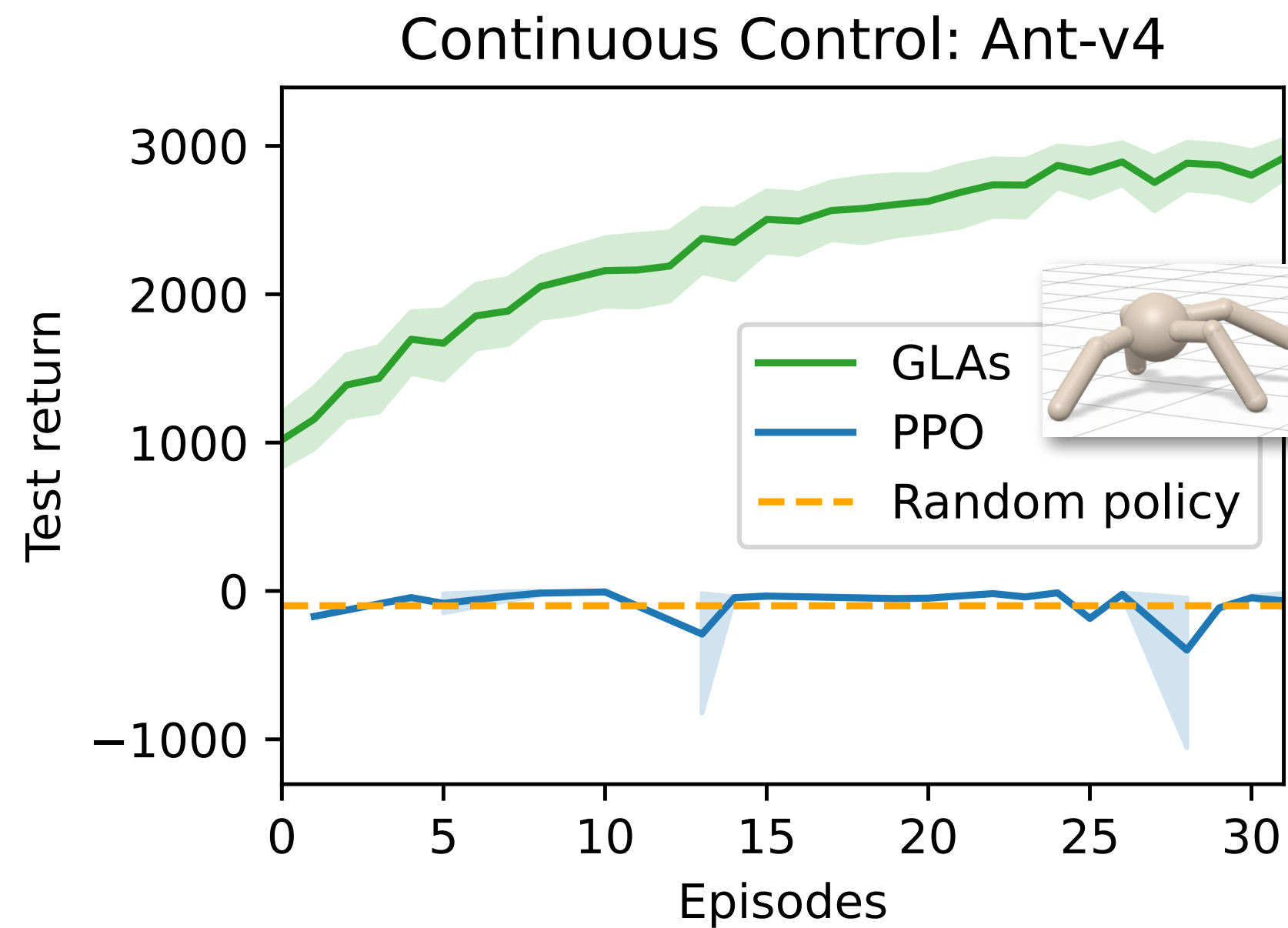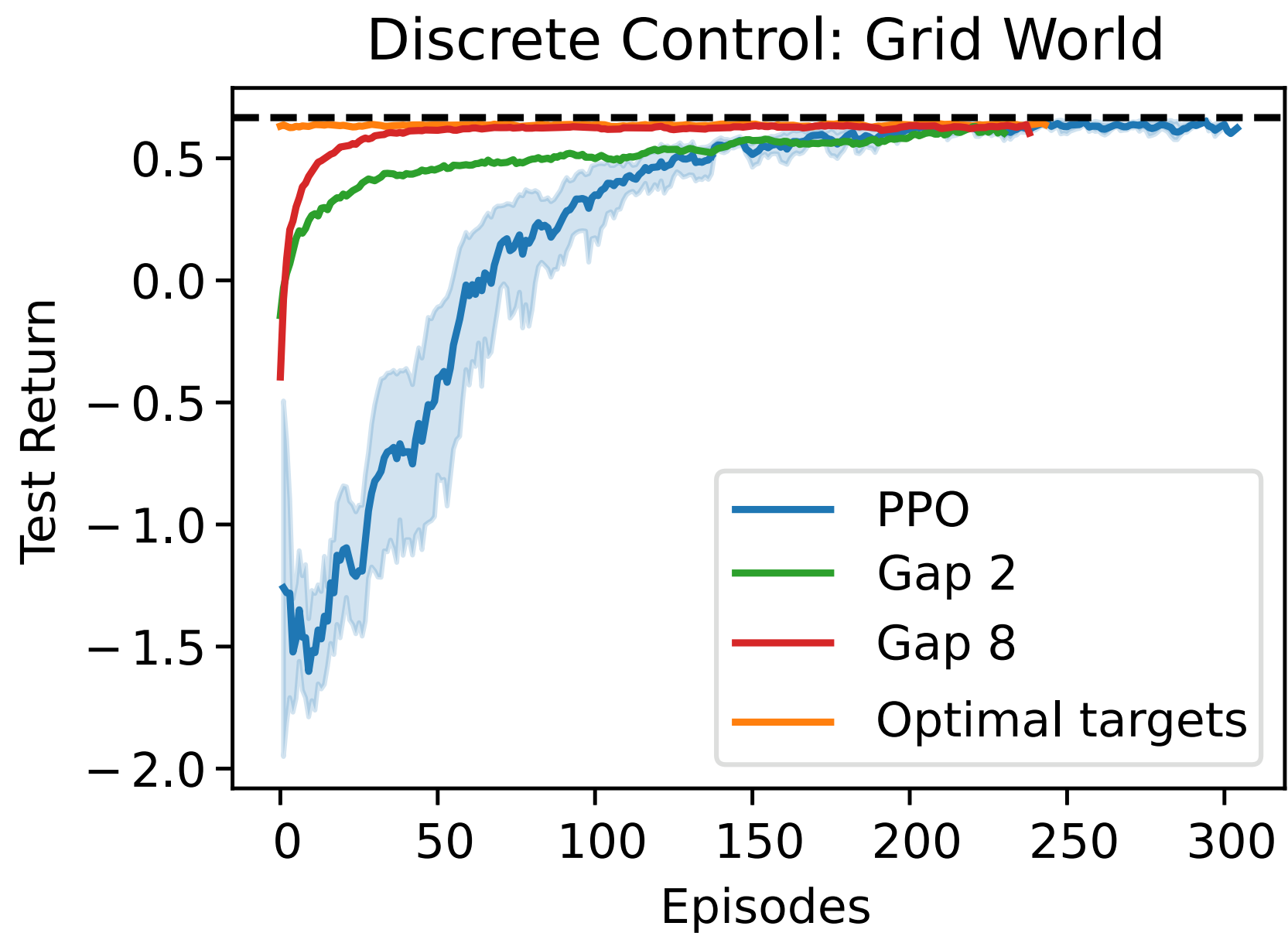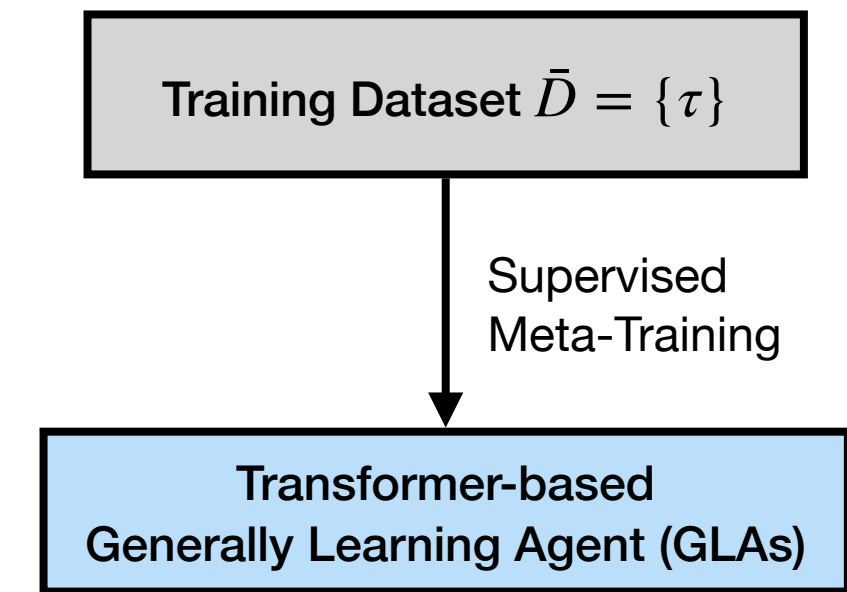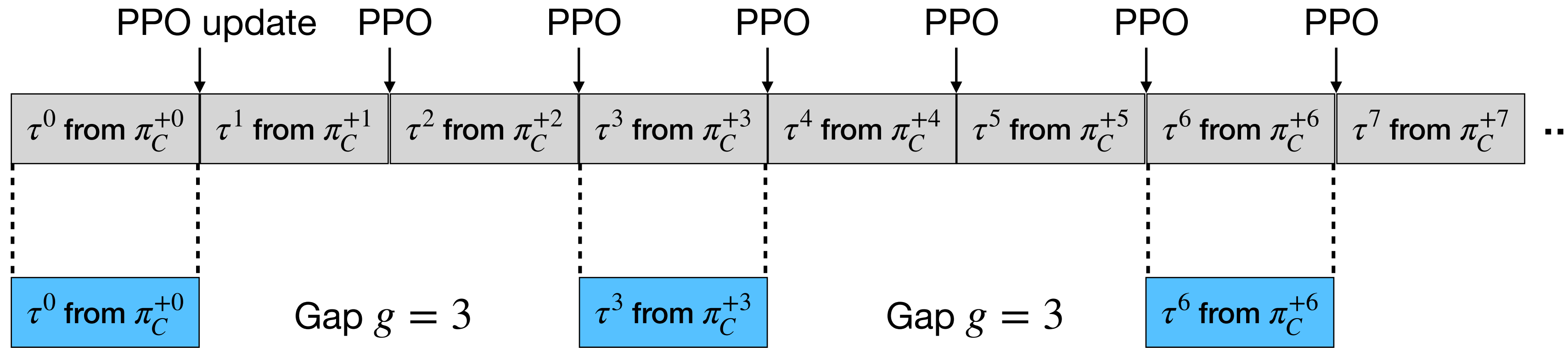
Supervised Meta-Training

Transformer-based Generally Learning Agent (GLAs)



Discrete Control: Grid World

Continuous Control: Ant-v4

# Effect of the gap



PPO update    PPO    PPO    PPO    PPO    PPO    PPO

| $\tau^0$ from $\pi_C^{+0}$ | $\tau^1$ from $\pi_C^{+1}$ | $\tau^2$ from $\pi_C^{+2}$ | $\tau^3$ from $\pi_C^{+3}$ | $\tau^4$ from $\pi_C^{+4}$ | $\tau^5$ from $\pi_C^{+5}$ | $\tau^6$ from $\pi_C^{+6}$ | $\tau^7$ from $\pi_C^{+7}$ | $\cdots$ |

$\tau^0$ from $\pi_C^{+0}$    Gap $g = 3$    $\tau^3$ from $\pi_C^{+3}$    Gap $g = 3$    $\tau^6$ from $\pi_C^{+6}$

Training Dataset $\bar{D} = \{\tau\}$

Supervised
Meta-Training

Transformer-based
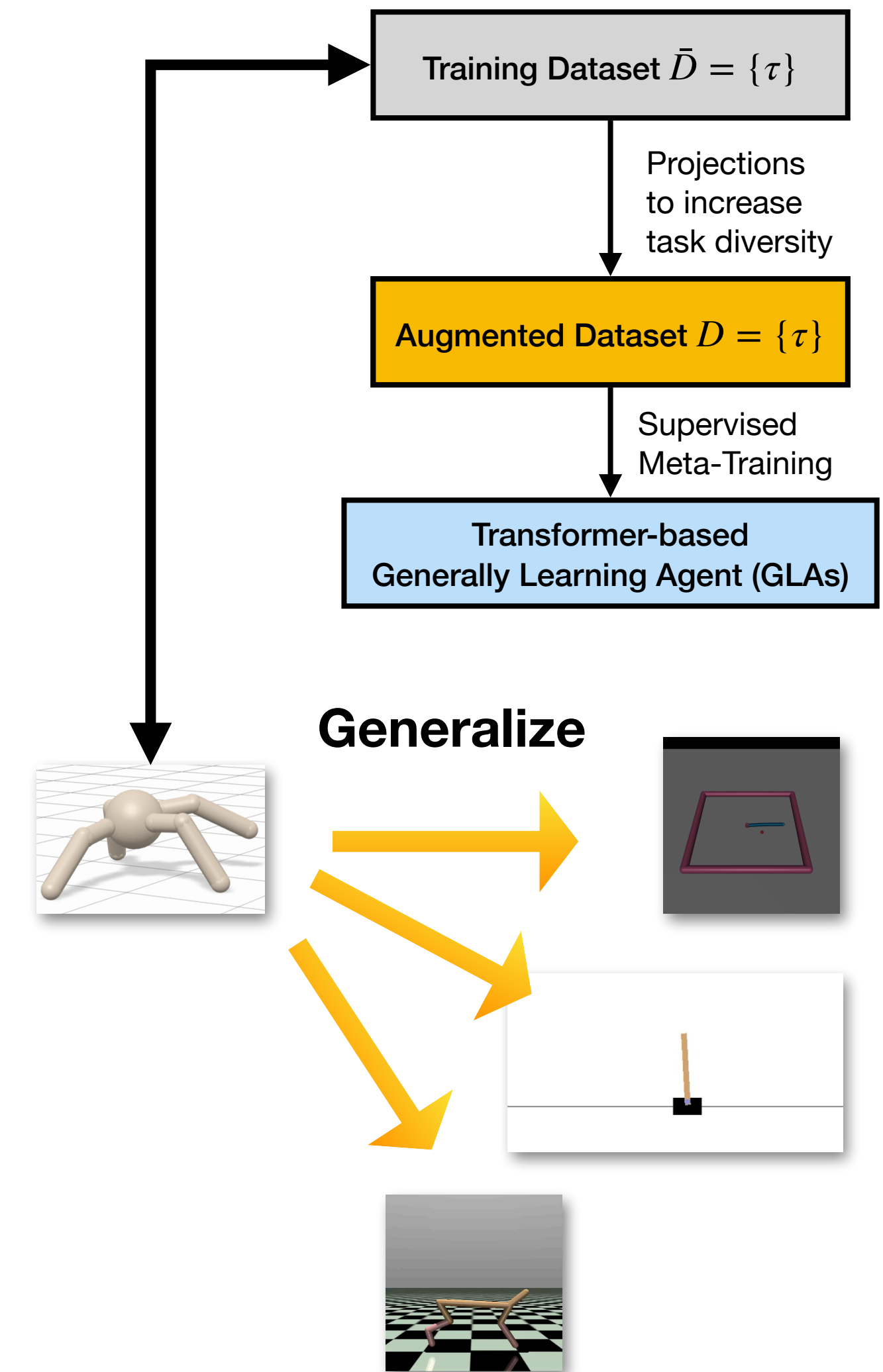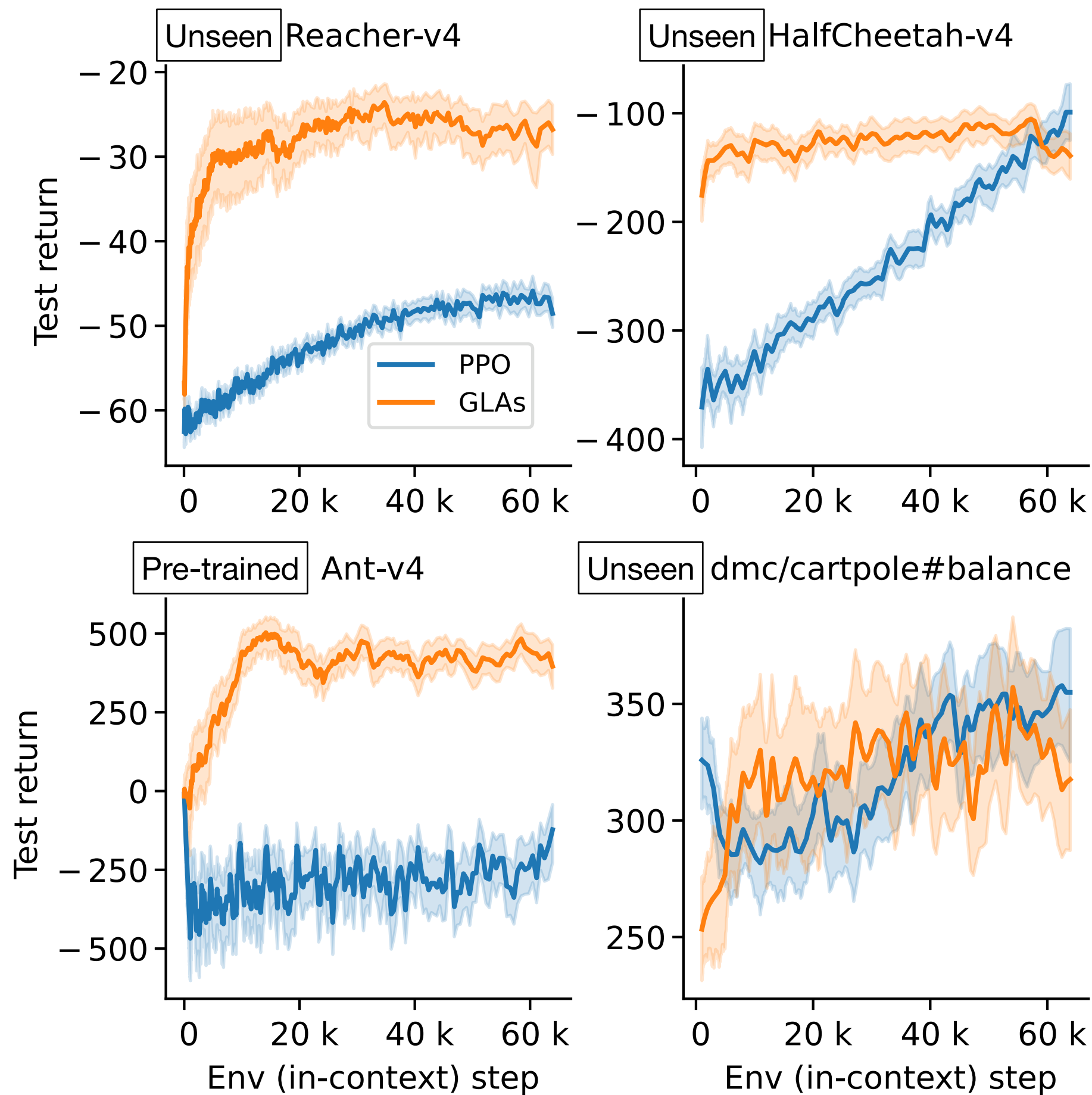Generally Learning Agent (GLAs)

**Discrete Control: Grid World**

**Continuous Control: Ant-v4**

**Generalize across environments?**

# Introducing random projections

**Generalize across environments?**

**Unseen** Reacher-v4

**Unseen** HalfCheetah-v4

**Pre-trained** Ant-v4

**Unseen** dmc/cartpole#balance

PPO
GLAs

Test return

Env (in-context) step

Test return

Env (in-context) step

**Early results!**
* **Longer contexts**
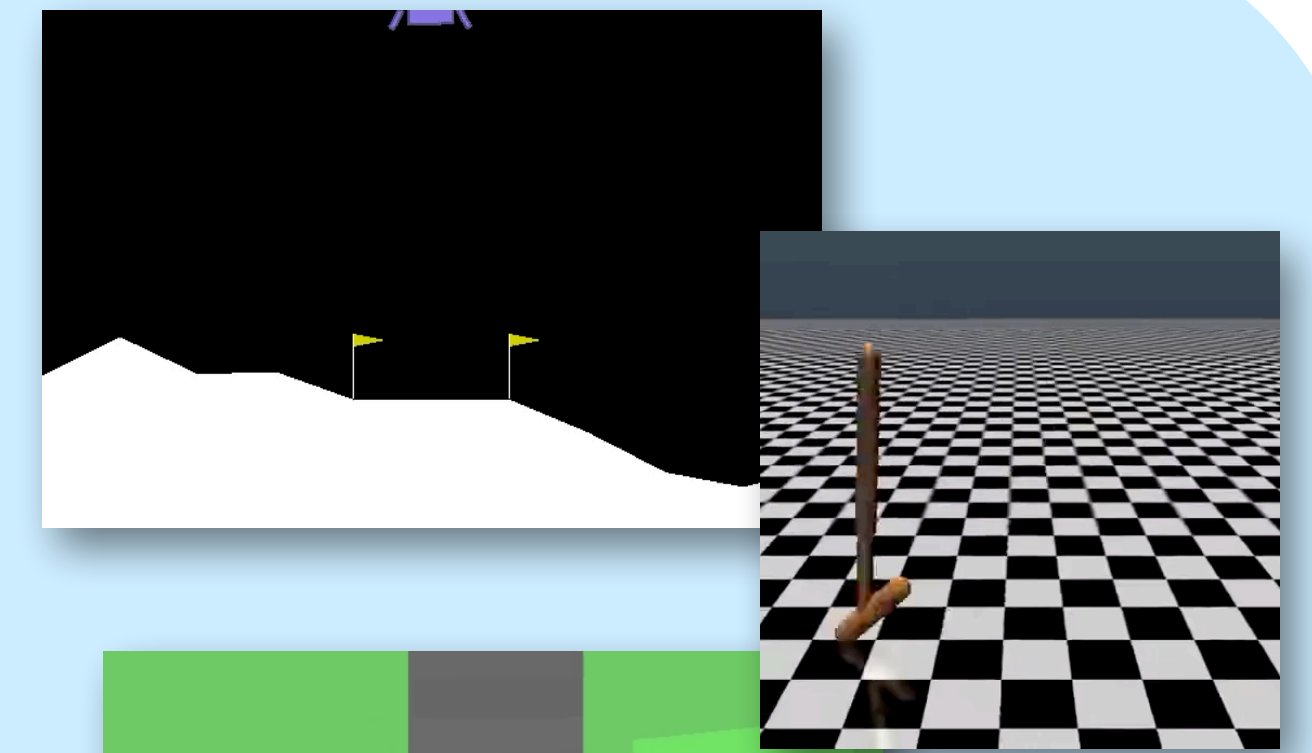* **Training dynamics**

Training Dataset $\bar{D} = \{\tau\}$

Projections to increase task diversity

Augmented Dataset $D = \{\tau\}$

Supervised Meta-Training

Transformer-based Generally Learning Agent (GLAs)

**Generalize**

# Automated task / environment generation



**Env Distribution**

→ **General-Purpose Learning Algorithm**

# Where to find me

louiskirsch.com

@LouisKirschAI

## Summary:

- **We distill an accelerated PPO into a Transformer**
- **Strong data augmentation helps the Transformer implement a general-purpose online RL algorithm generalizing across domains**



**Data Collection**

Environment Distribution

PPO Training → Training Dataset $\bar{D} = \{\tau\}$

Projections to increase task diversity

**Meta-Training**

Transformer-based Generally Learning Agent (GLAs) ← Supervised Meta-Training ← Augmented Dataset $D = \{\tau\}$

**Meta-Testing** — Inference mode

In-context Reinforcement Learning

Generally Learning Agents (GLAs)

Test environment