

# Plan2Explore

---

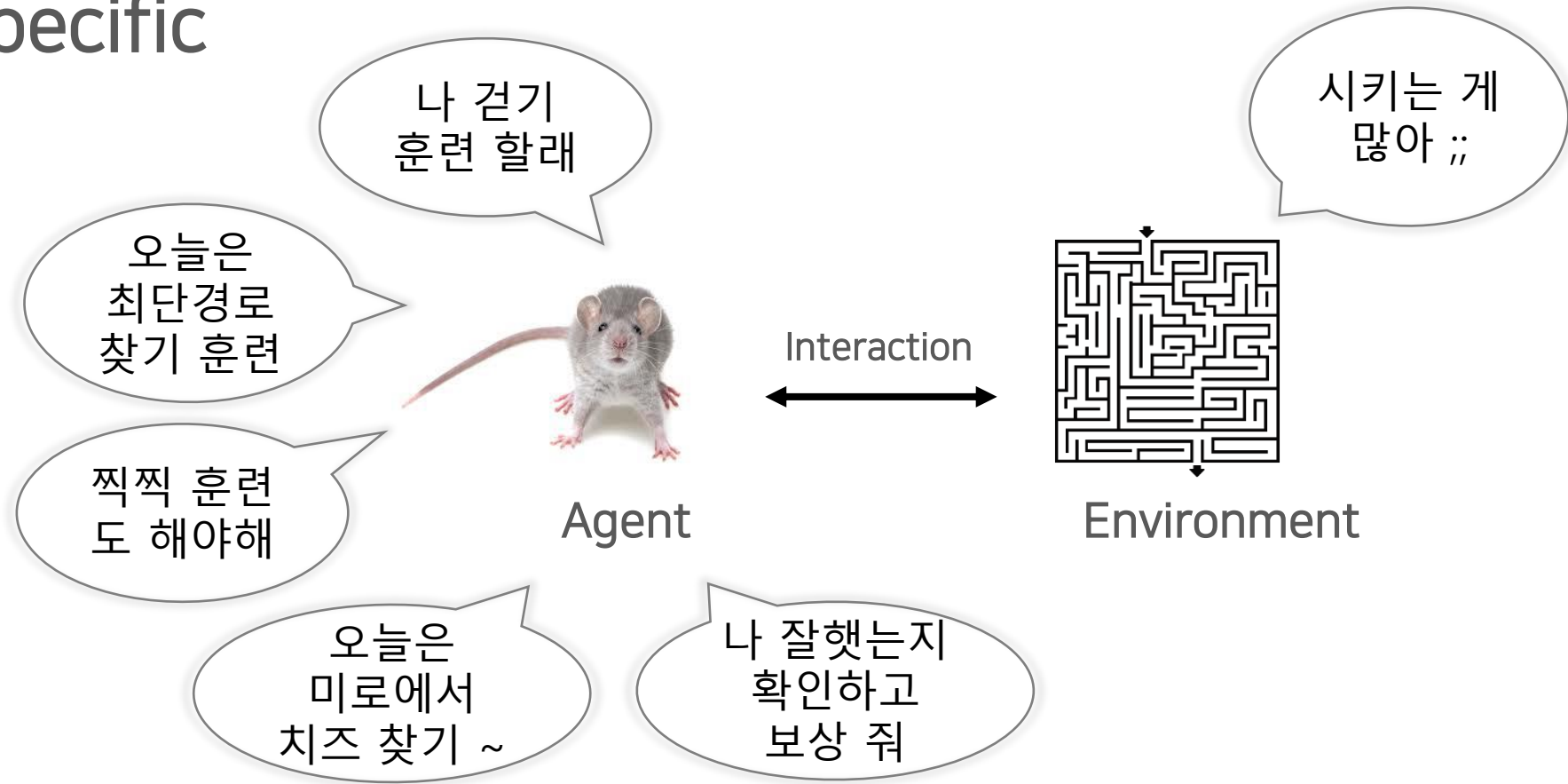
Planning to Explore via Self-Supervised World Models

# Plan2Explore

- Self Supervised Model Based Learning
- Work directly from images
- Interact with environment without agent, to collect new data
- Train the World Model with Intrinsic Motivation

# Supervised Learning

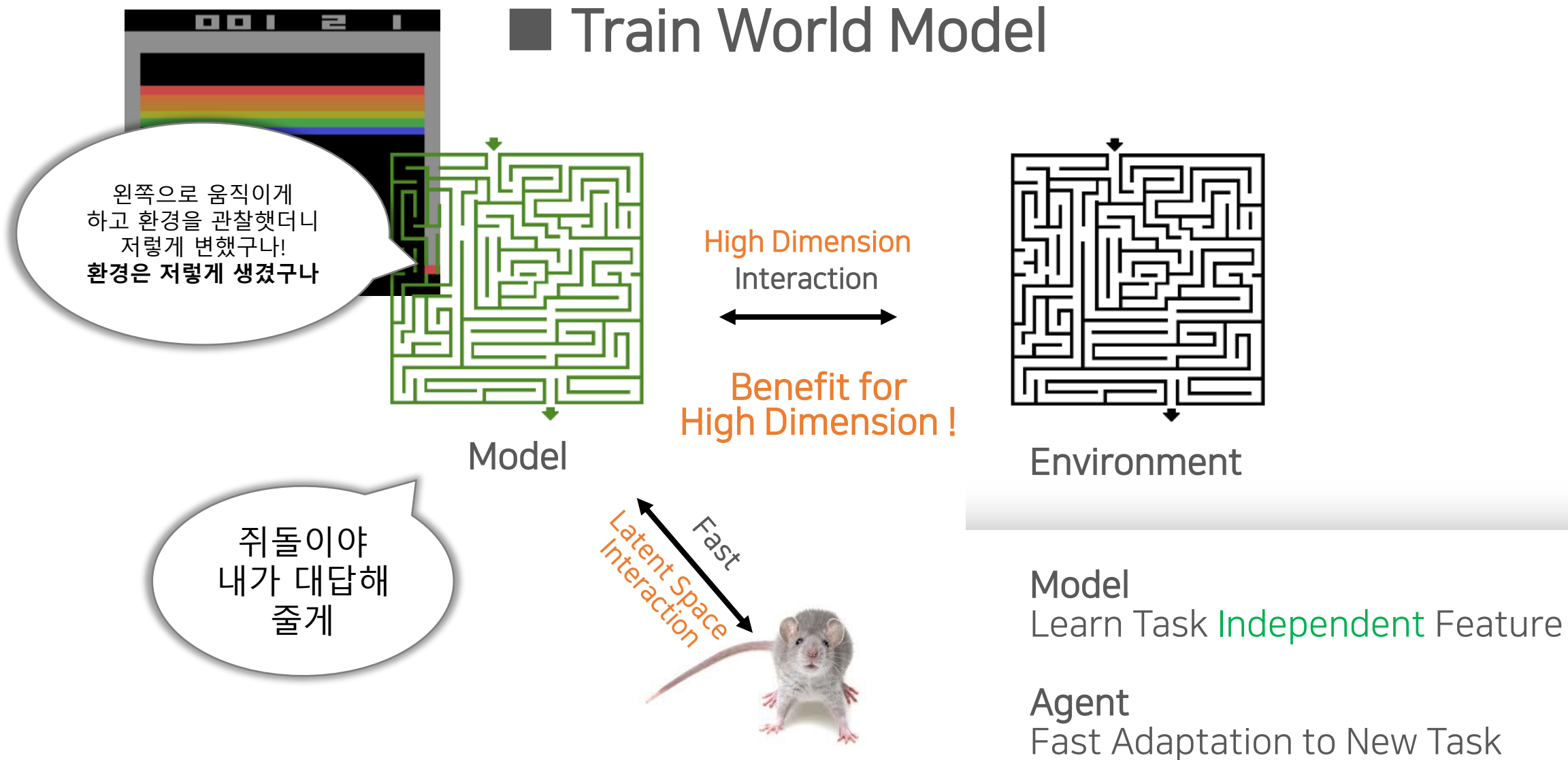
## ■ Task Specific



Number of Task ↑ → Require Large Amount of Experience

# Self Supervised Model Based Learning

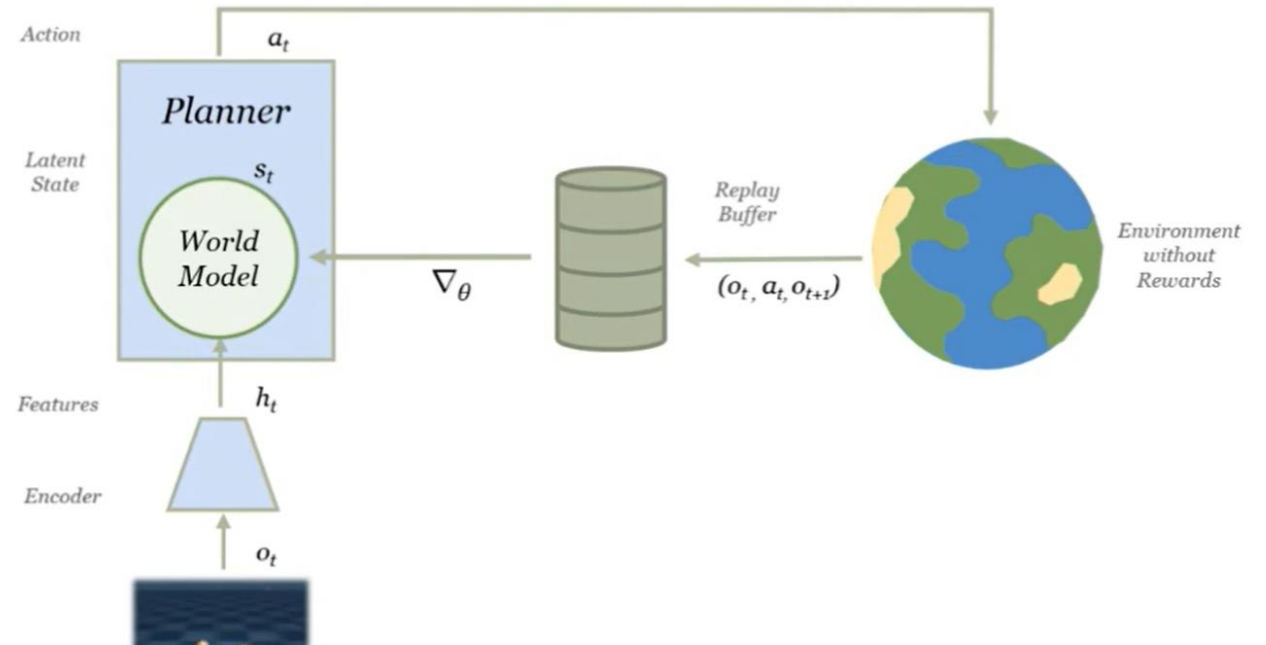
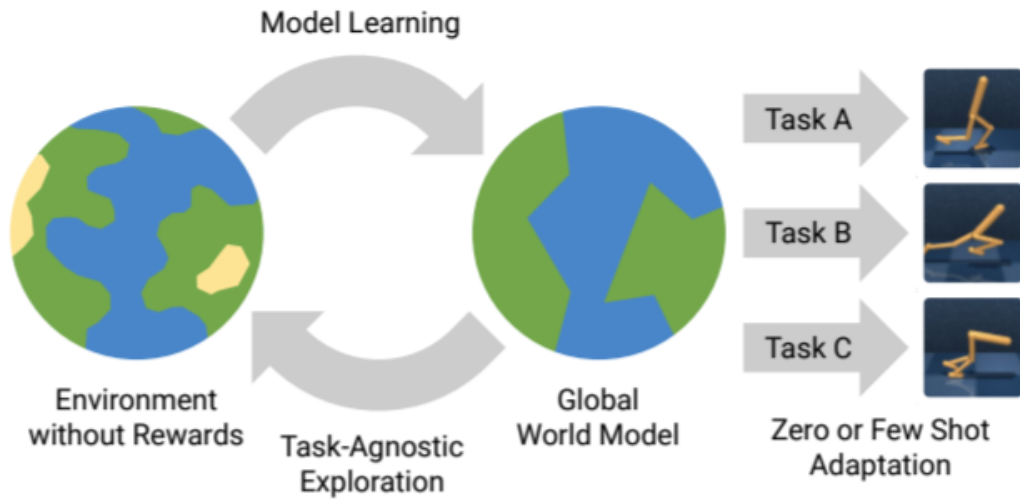
## ■ Train World Model



# Exploration

---

# Exploration



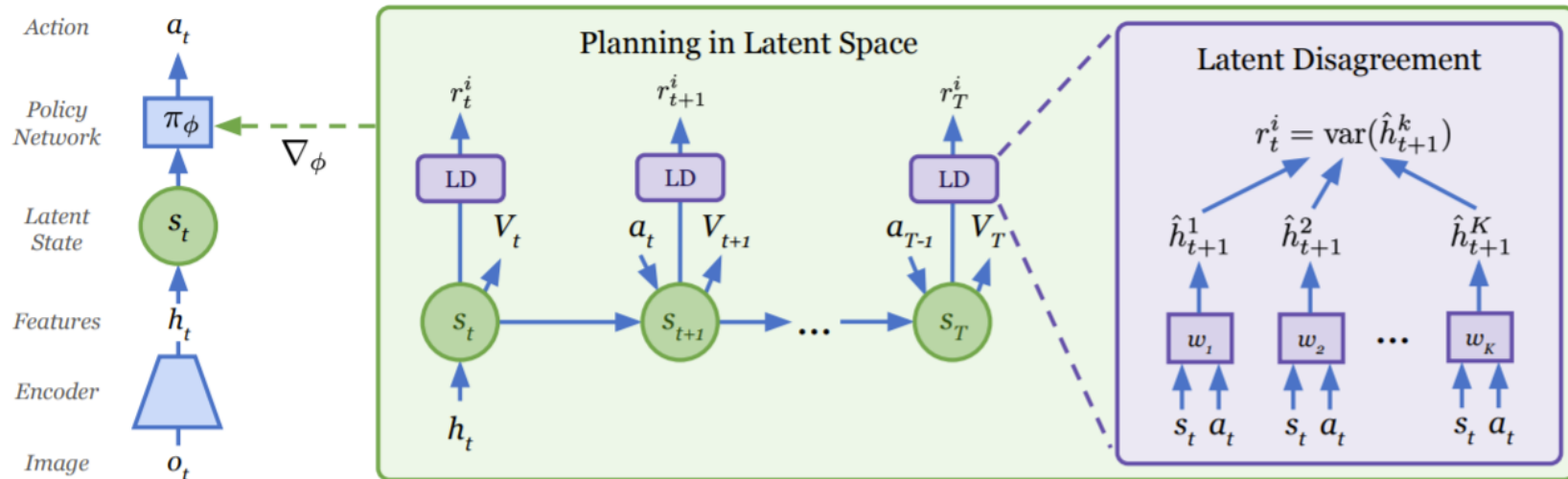
> Environment과의 Interaction으로 World Model 학습

> Planner로 다음 Interaction 결정

Key Idea      Use long-term planning to collecting novel data

# Exploration

## ■ Planning in Latent Space



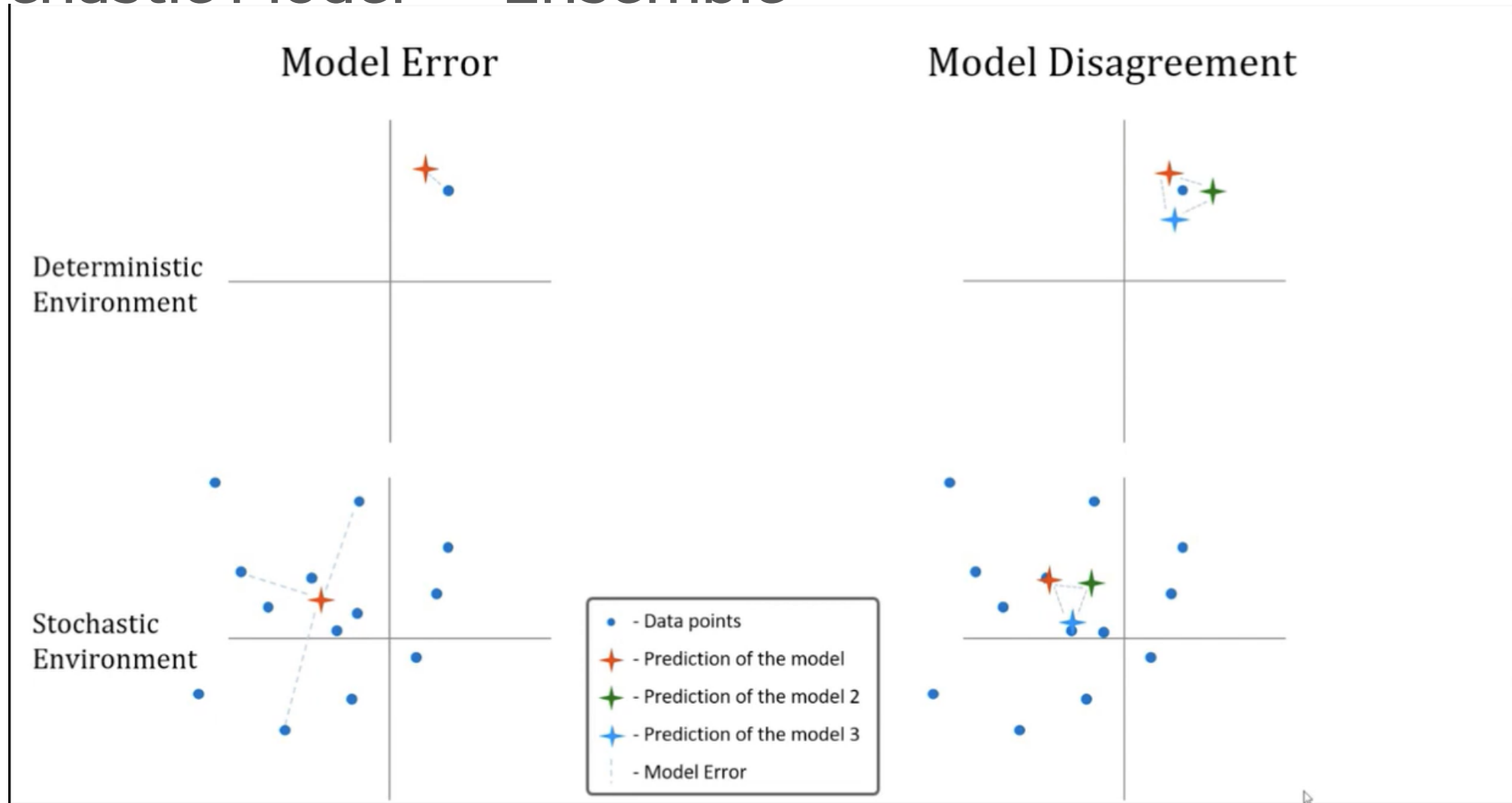
Latent State를 이용하는  
가벼운 1-step model들을 Ensemble

### Motivation

직관적으로, 초기화도 다르고, 관찰 순서도 다른 앙상블은 처음에는 예측이 다르나, 데이터가 늘어날수록 모델이 동일한 예측으로 수렴하며 불일치가 감소한다.

# Exploration

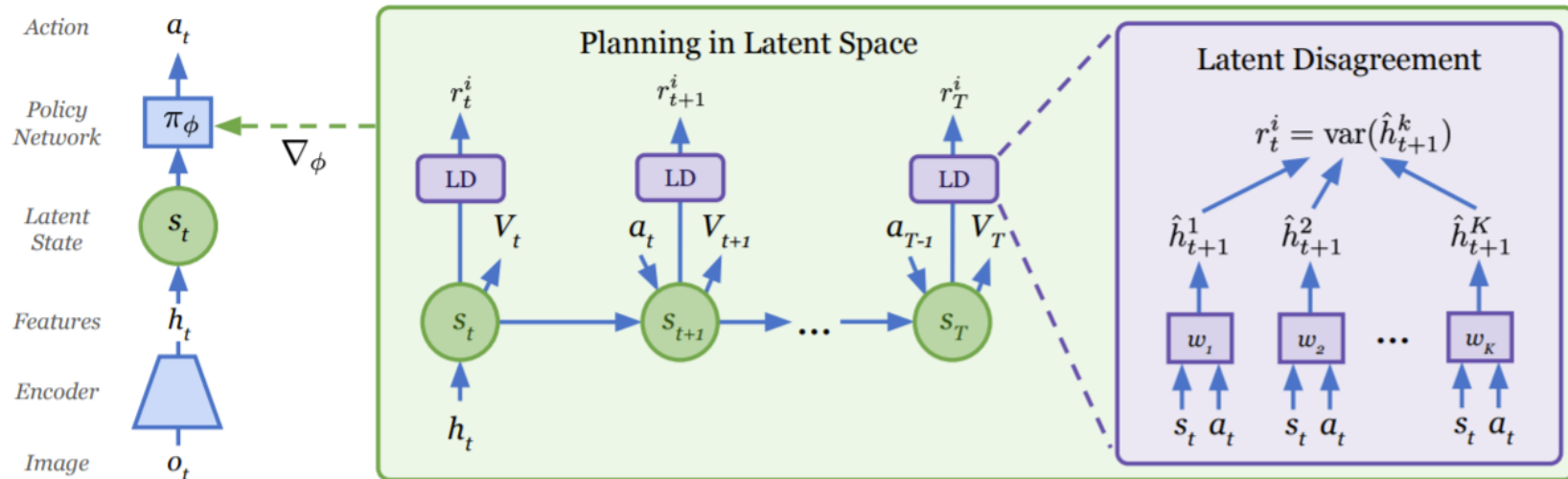
## ■ Stochastic Model -> Ensemble





# Exploration

## ■ Planning in Latent Space



Latent State를 이용하는  
가벼운 1-step model들을 Ensemble

Ensemble predictors:  $q(h_{t+1} \mid w_k, s_t, a_t)$

$q(h_{t+1} \mid w_k, s_t, a_t) \triangleq \mathcal{N}(\mu(w_k, s_t, a_t), 1)$ .

### Motivation

직관적으로, 초기화도 다르고, 관찰 순서도 다른 앙상블은 처음에는 예측이 다르나, 데이터가 늘어날수록 모델이 동일한 예측으로 수렴하며 불일치가 감소한다.

# Exploration

## ■ Planning in Latent Space

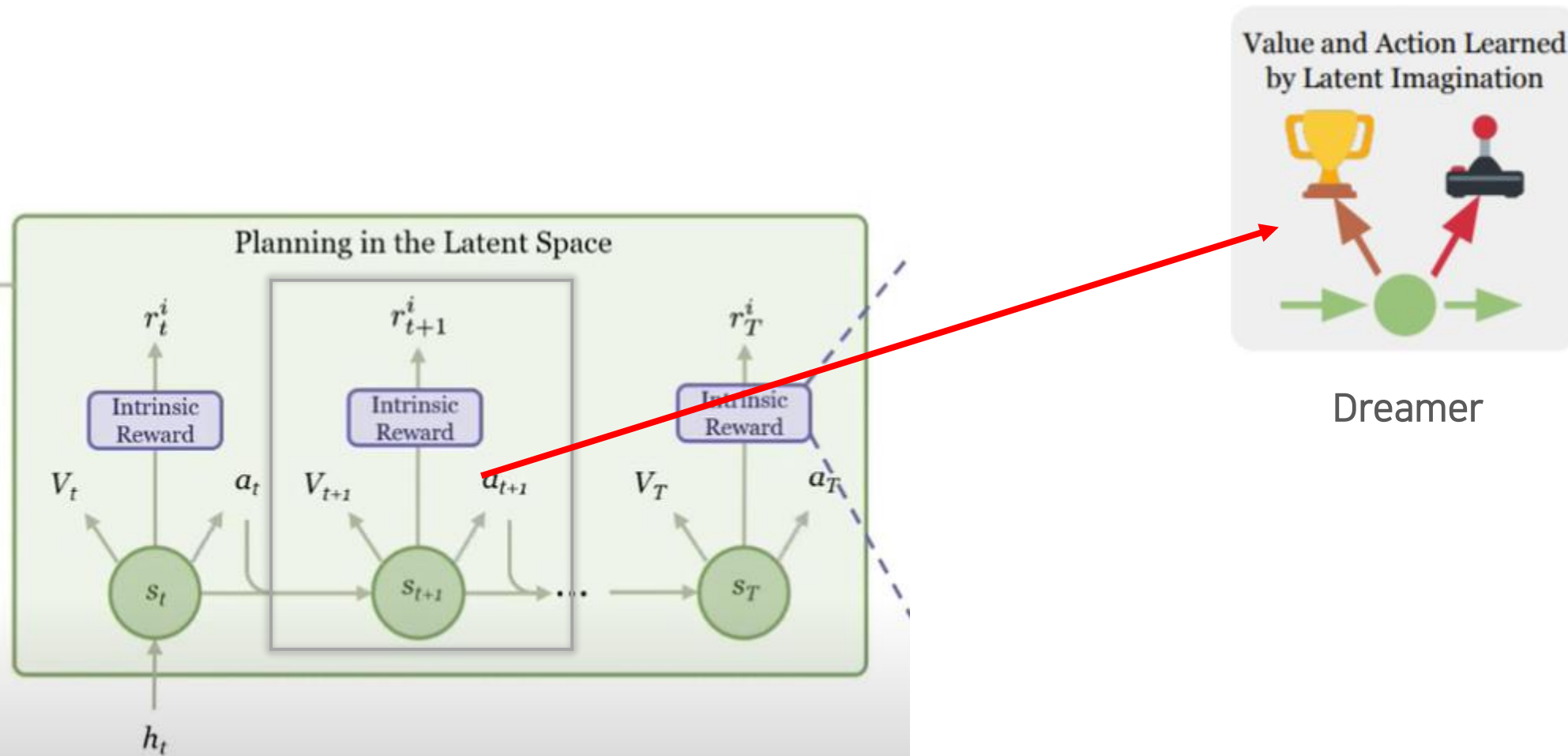
Latent Disagreement를 예측한 평균에 대한 분산으로 정의

$$\begin{aligned} D(s_t, a_t) &\triangleq \text{Var}(\{\mu(w_k, s_t, a_t) \mid k \in [1; K]\}) \\ &= \frac{1}{K-1} \sum_k (\mu(w_k, s_t, a_t) - \mu')^2, \\ \mu' &\triangleq \frac{1}{K} \sum_k \mu(w_k, s_t, a_t). \end{aligned} \quad (4)$$

Latent Disagreement는  
Exploration policy를 훈련시키기 위한 Intrinsic Reward로 사용됨

# Exploration

## ■ Exploration Policy : Train by Dreamer



# Exploration

---

**Algorithm 1** Planning to Explore via Latent Disagreement

---

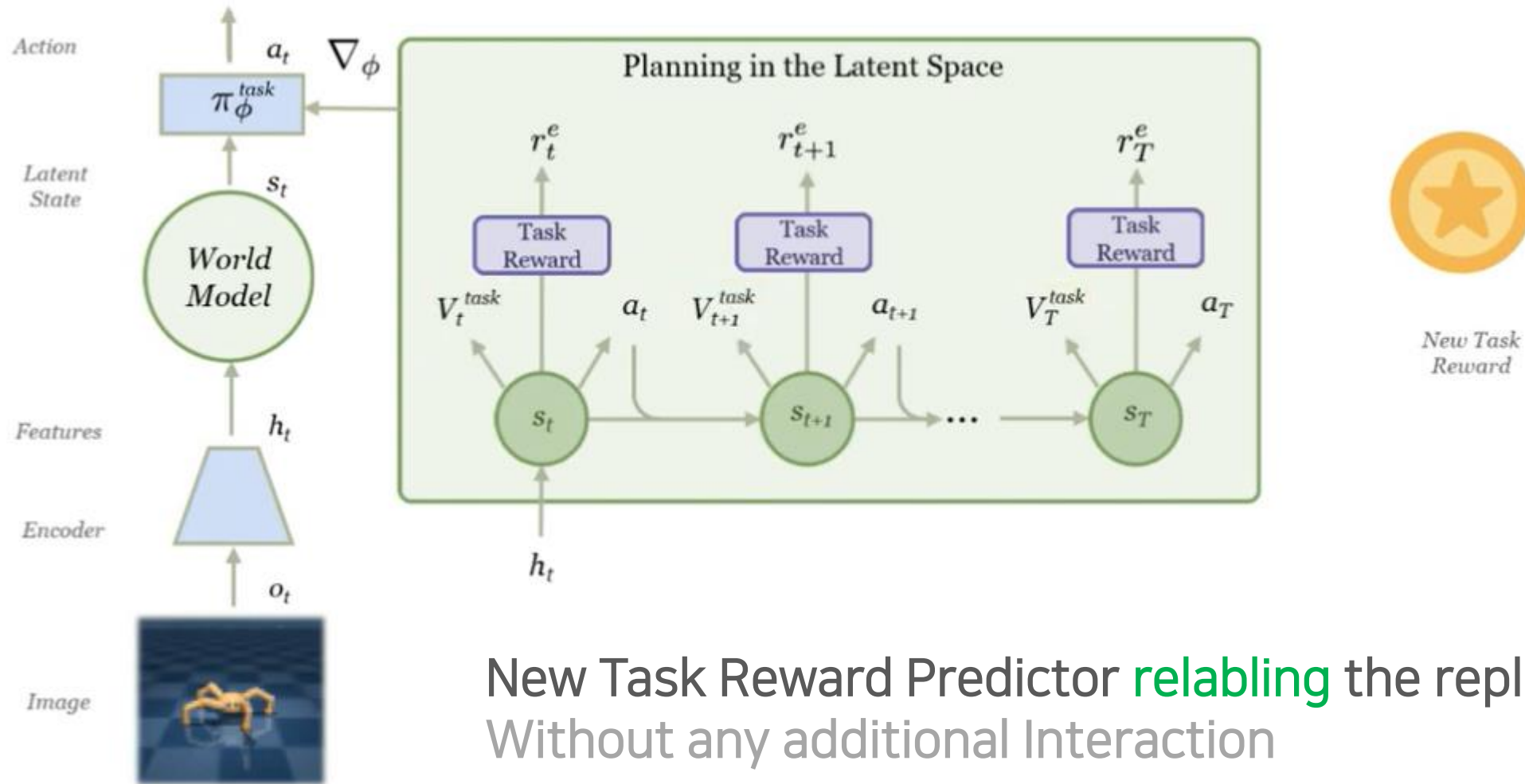
```
1: initialize: Dataset  $D$  from a few random episodes.  
2:               World model  $M$ .  
3:               Latent disagreement ensemble  $E$ .  
4:               Exploration actor-critic  $\pi_{LD}$ .  
5: while exploring do  
6:   Train  $M$  on  $D$ .  
7:   Train  $E$  on  $D$ .  
8:   Train  $\pi_{LD}$  on LD reward in imagination of  $M$ .  
9:   Execute  $\pi_{LD}$  in the environment to expand  $D$ .  
10: end while  
11: return Task-agnostic  $D$  and  $M$ .
```

---

# Solving Task

---

# Solving Task



New Task Reward Predictor **relabelling** the replay buffer  
Without any additional Interaction

# Solving Task

---

**Algorithm 2** Zero and Few-Shot Task Adaptation

---

```
1: input:      World model  $M$ .
2:             Dataset  $D$  without rewards.
3:             Reward function  $R$ .
4: initialize: Latent-space reward predictor  $\hat{R}$ .
5:             Task actor-critic  $\pi_R$ .
6: while adapting do
7:   Distill  $R$  into  $\hat{R}$  for sequences in  $D$ .
8:   Train  $\pi_R$  on  $\hat{R}$  in imagination of  $M$ .
9:   Execute  $\pi_R$  for the task and report performance.
10:  Optionally, add task-specific episode to  $D$  and repeat.
11: end while
12: return Task actor-critic  $\pi_R$ .
```

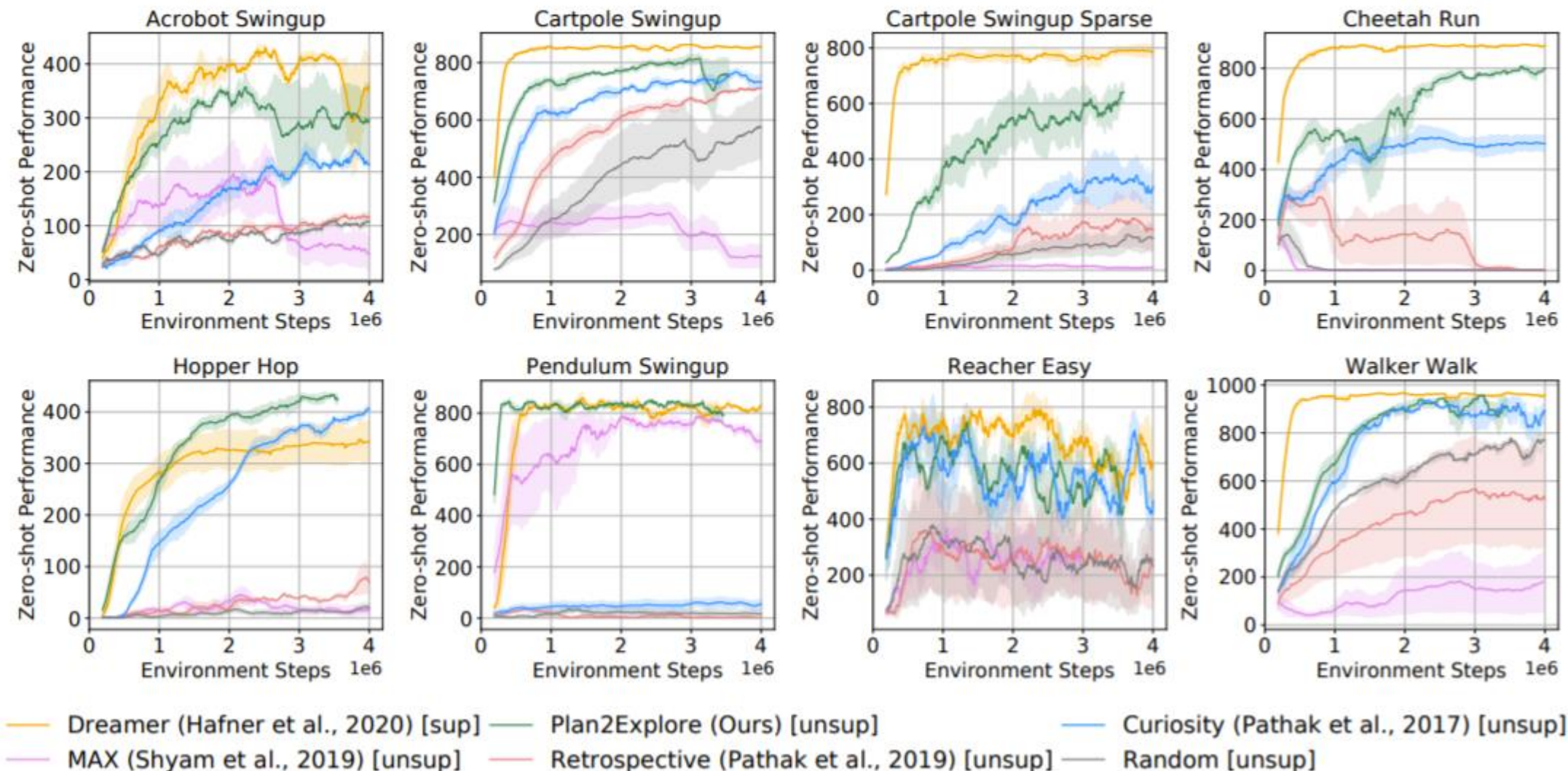
---

# Experiment

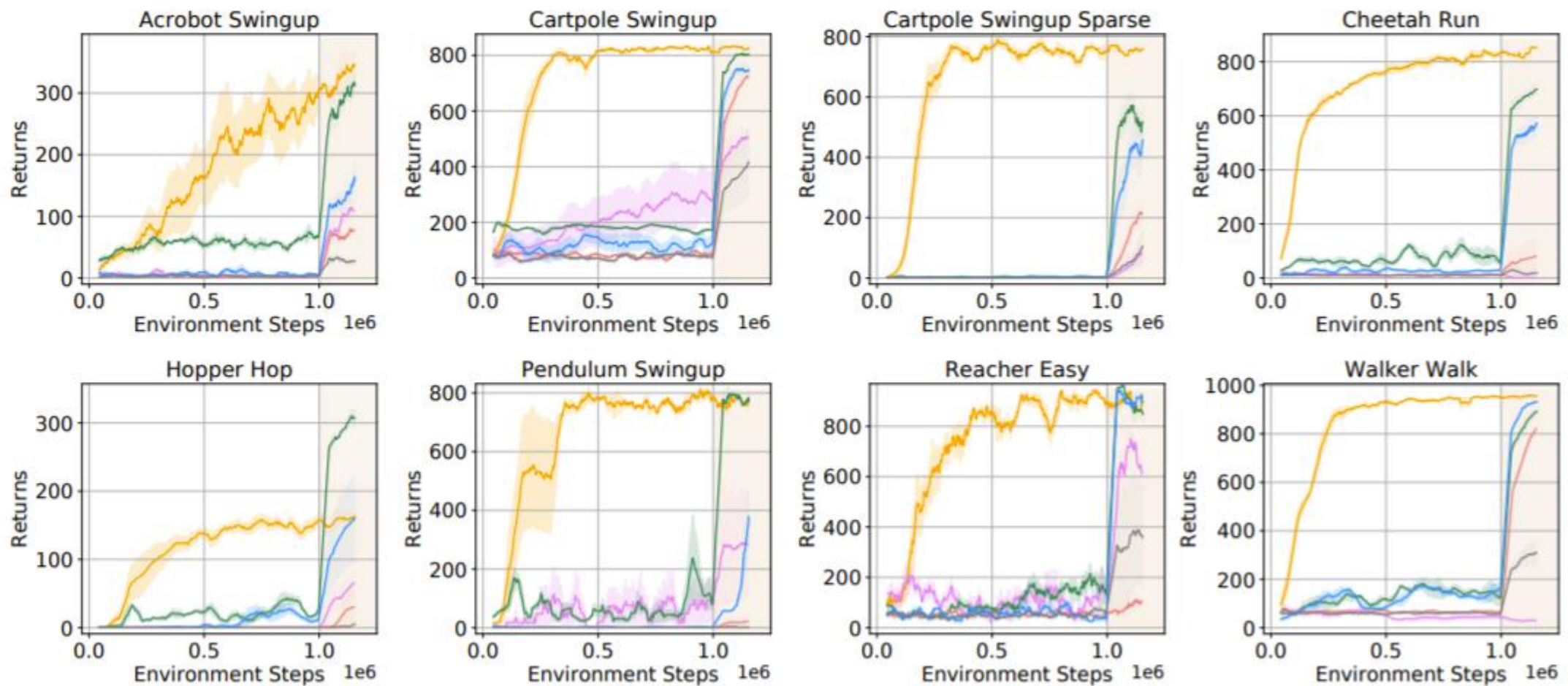
---



# Solving a new-task in zero-shot



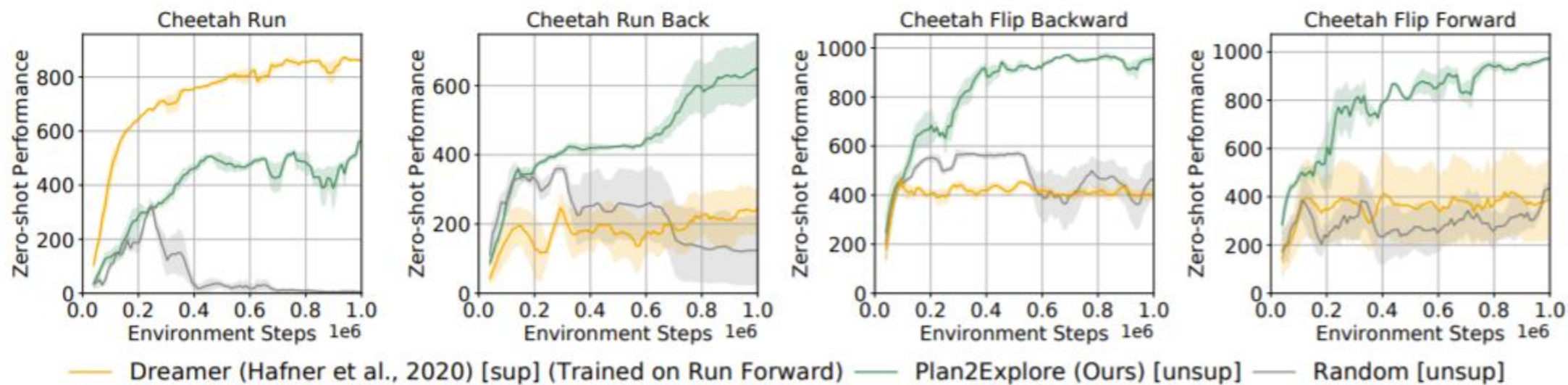
# Few shot adaptation



— Dreamer (Hafner et al., 2020) [sup] — Plan2Explore (Ours) [unsup] — Curiosity (Pathak et al., 2017) [unsup]  
— MAX (Shyam et al., 2019) [unsup] — Retrospective (Pathak et al., 2019) [unsup] — Random [unsup]



# Multi Task Performance



왜 Latent Disagreement ?

# **Expected Information Gain**

---

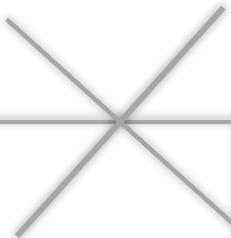
# Latent Disagreement

Maximize Information Gain

$$a_t^* \triangleq \arg \max_{a_t} I(h_{t+1}; w \mid s_t, a_t).$$

Information Gain

$$\begin{aligned} I(h_{t+1}; w \mid s_t, a_t) \\ = H(h_{t+1} \mid s_t, a_t) - H(h_{t+1} \mid w, s_t, a_t). \end{aligned}$$


$$p(w) \triangleq \frac{1}{K} \sum_k \delta(w - w_k)$$

$$p(h_{t+1} \mid w_k, s_t, a_t) \triangleq \mathcal{N}(h_{t+1} \mid \mu(w_k, s_t, a_t), \sigma^2).$$

$$D(s_t, a_t) \triangleq \frac{1}{K-1} \sum_k (\mu(w_k, s_t, a_t) - \mu')^2,$$

$$\mu' \triangleq \frac{1}{K} \sum_k \mu(w_k, s_t, a_t).$$

**End**

---