# PlaNet

## Learning Latent Dynamics
## for Planning from Pixels

# PlaNet Deep Planning Network

- Scalable Model-based RL

- Efficint planning in latent space with large batch size

- Reaches top performance using 200X fewer episodes

# PlaNet Deep Planning Network

- **Recurrent State Space** Model (RSSM)

  Deterministic & Stochastic Components

- **Latent Overshooting**

  Latent Space Result로 Latent Sequence Model을 굴리는게 장기예측에 적합

# Recurrent State Space Model

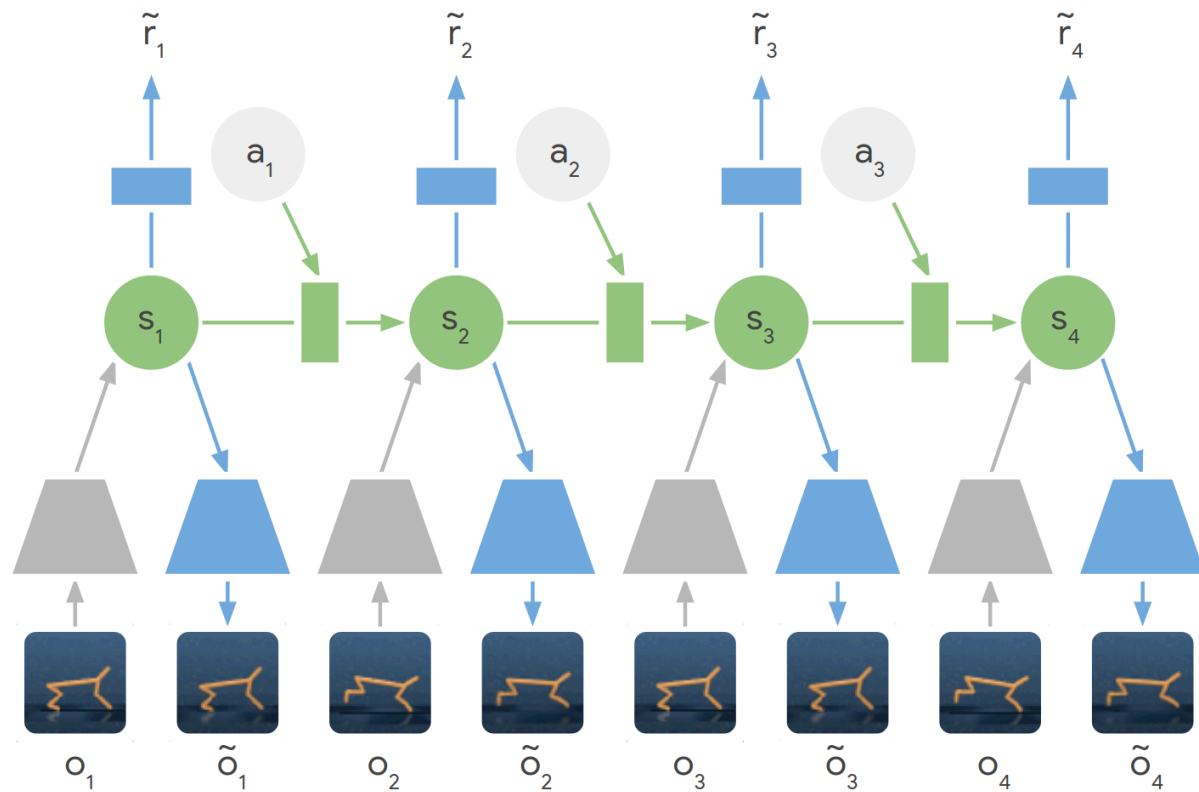# Recurrent State Space Model (RSSM)

보이는게 State의 전부가 아닐 수 있다. (Non-Marcov)

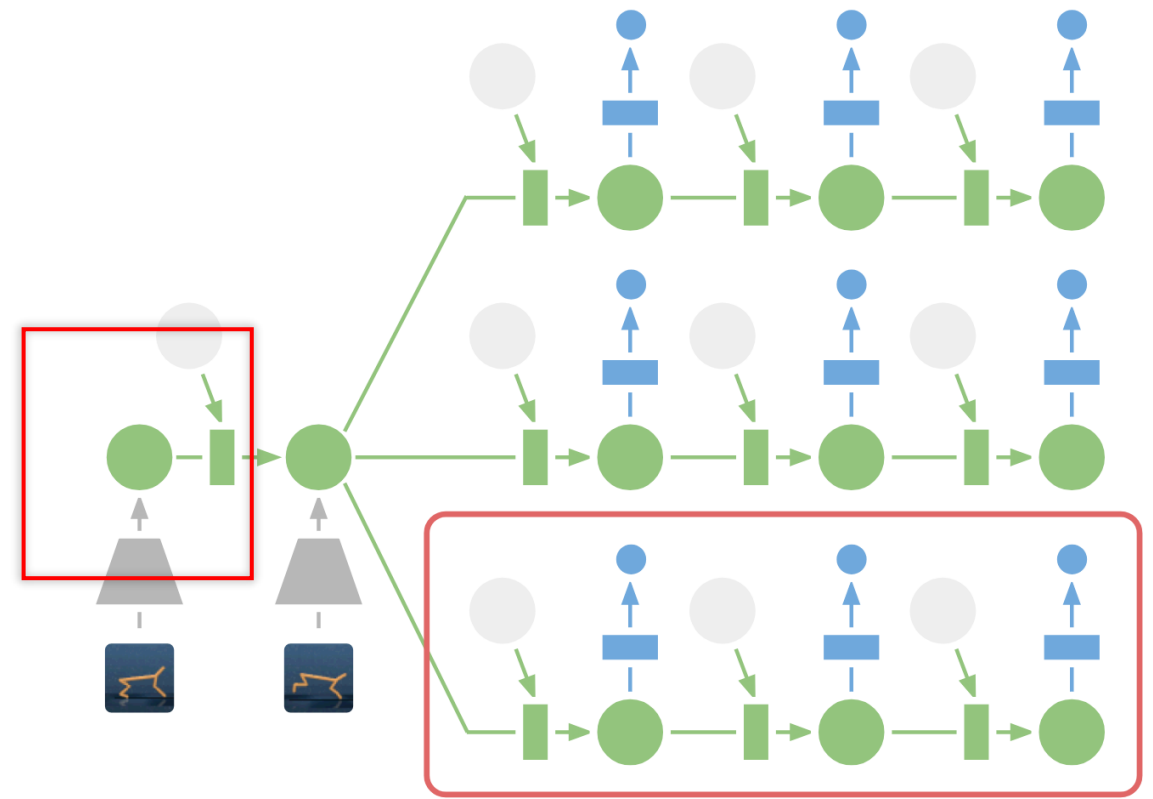Partially Observable Markov Decision Process를 Marcov Decision Process로 바꾸는 기법



DQN에서는 비슷하게 **이미지 네 장을 묶어 하나의 State**로 취급하는 것과 유사하게 (Non-Marcov를 Marcov하게 바꾸는 기법으로 제시하는 것이다.

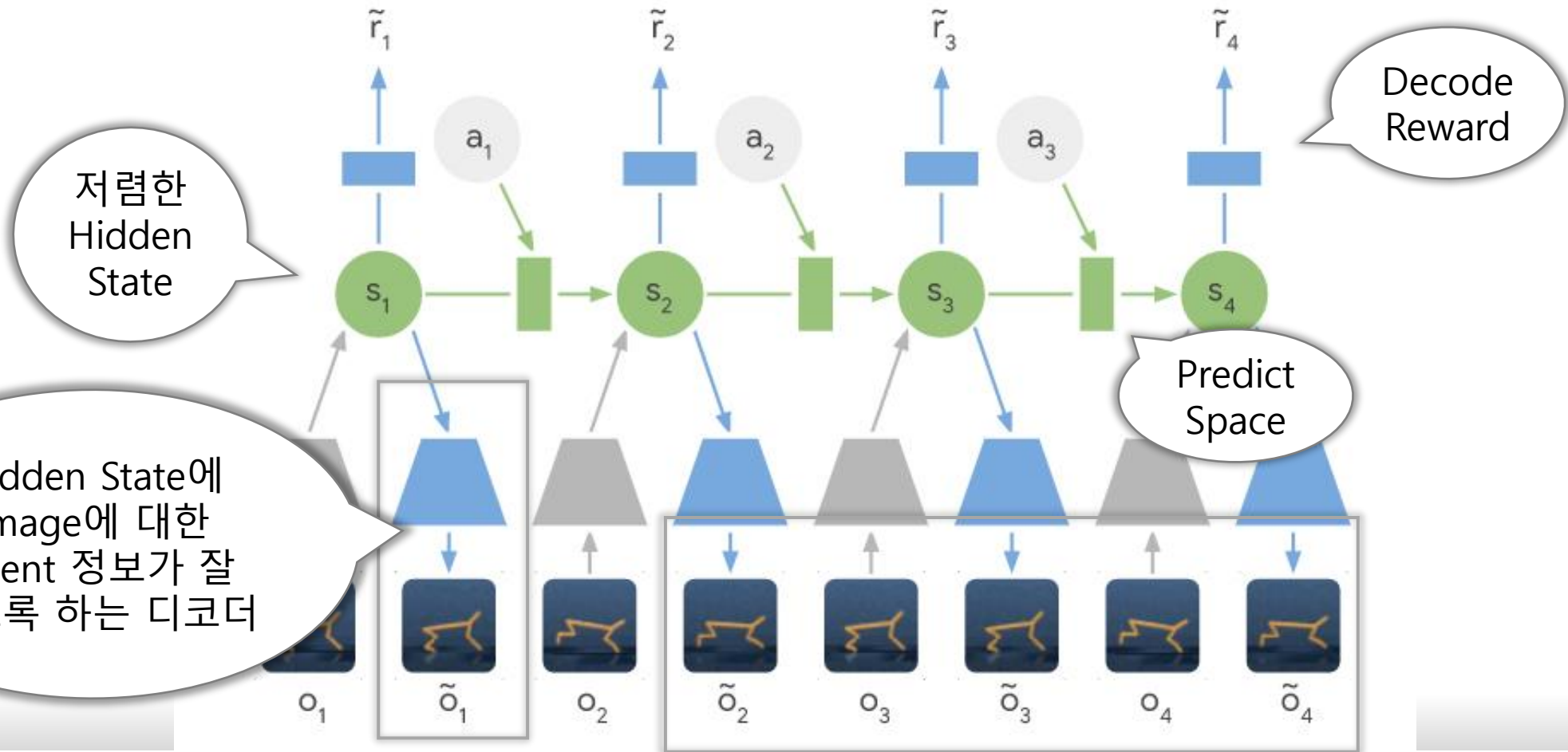# Recurrent State Space Model (RSSM)
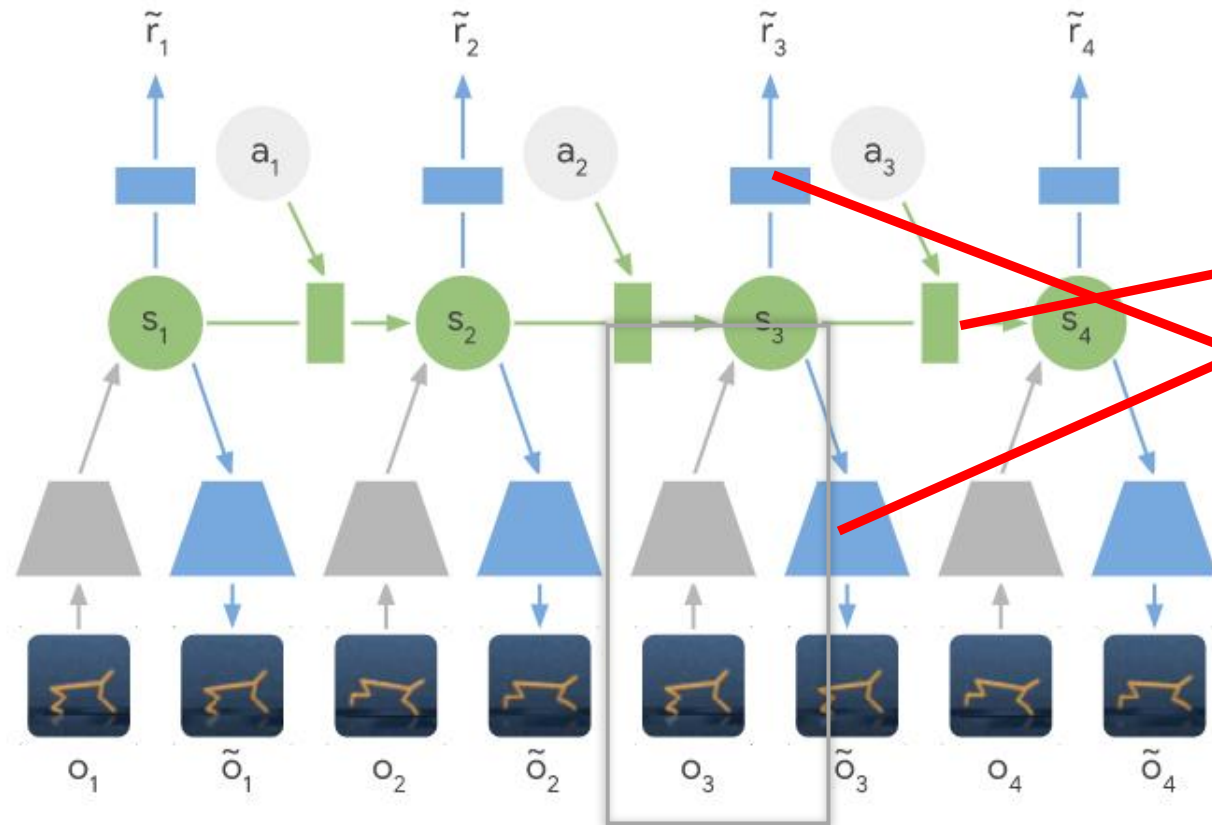


Learned Latent Dynamics Model

Planning in Latent Space

# Recurrent State Space Model (RSSM)

## ■ Learned Latent Dynamics Model

# Recurrent State Space Model

## ■ Learned Latent Dynamics Model



Transition function: $\quad s_t \sim \mathrm{p}(s_t \mid s_{t-1}, a_{t-1})$

Observation function: $\quad o_t \sim \mathrm{p}(o_t \mid s_t)$

Reward function: $\quad r_t \sim \mathrm{p}(r_t \mid s_t)$ $\qquad (1)$

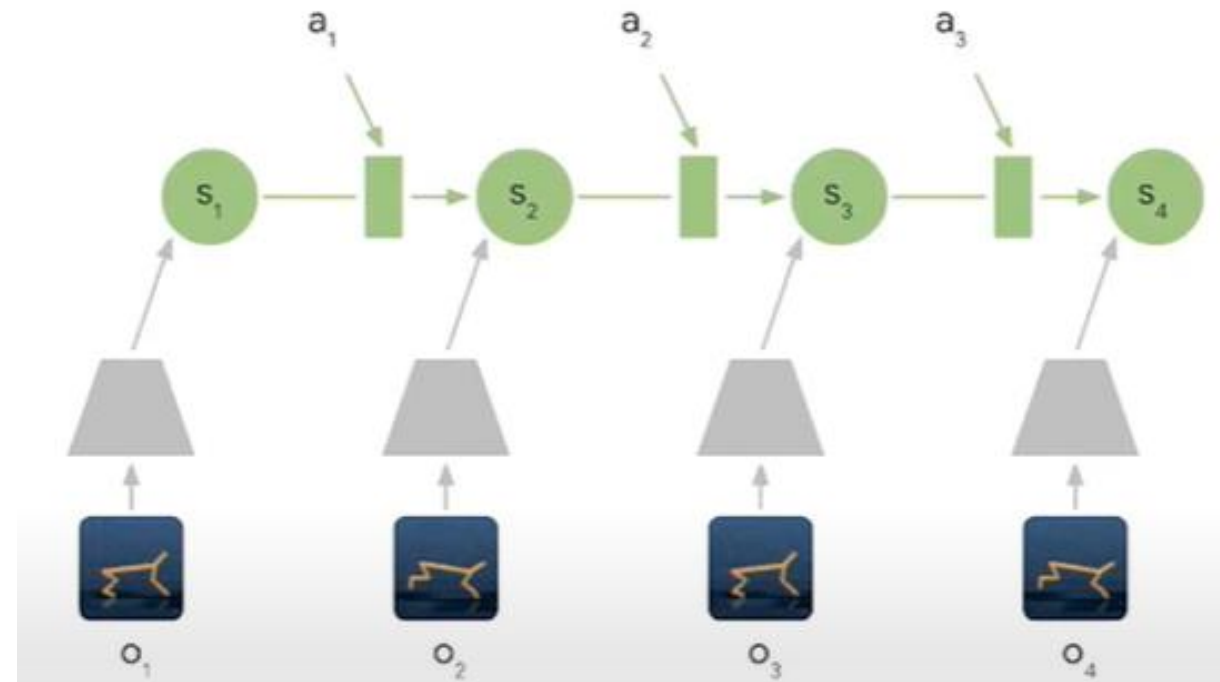Policy: $\quad a_t \sim \mathrm{p}(a_t \mid o_{\leq t}, a_{<t}),$
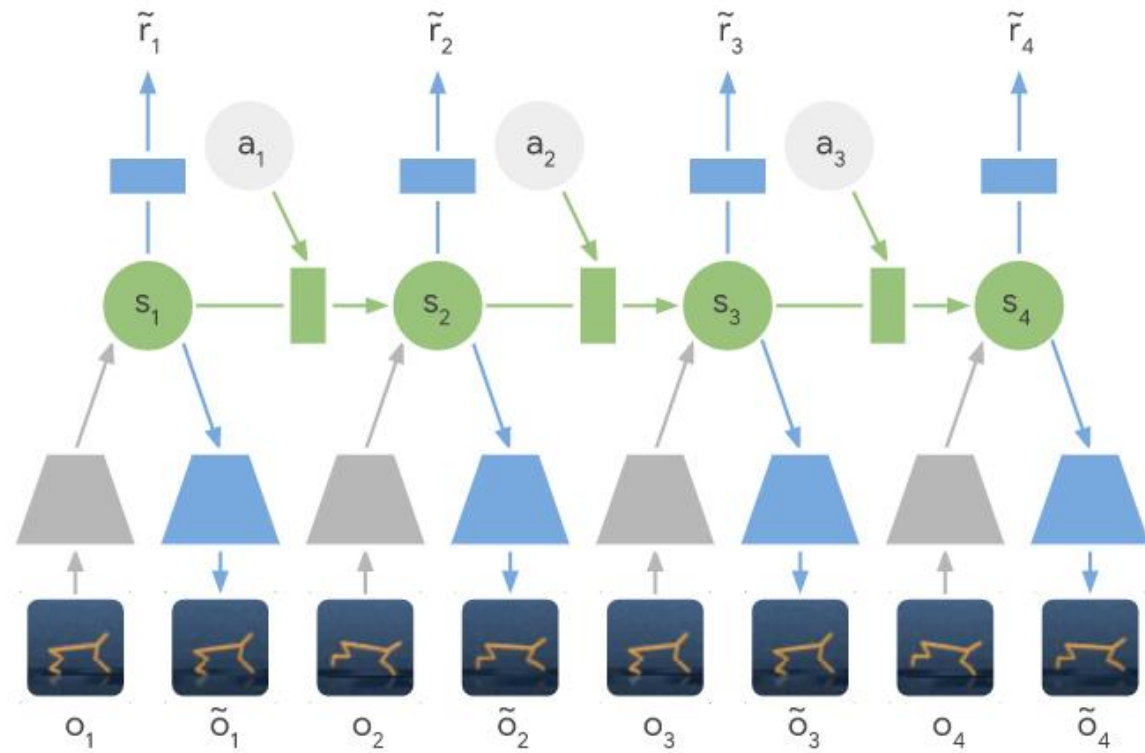
Policy는?
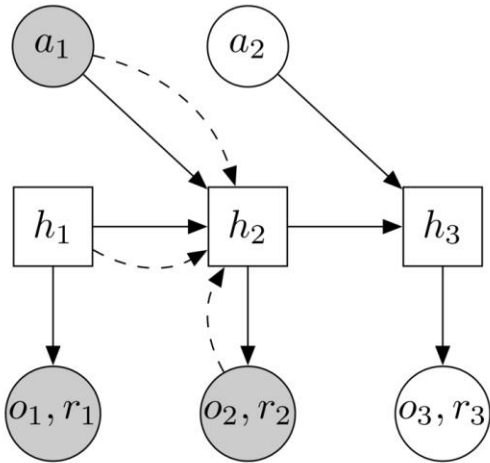Model-free처럼 Policy Net을 사용하는 게 아니고
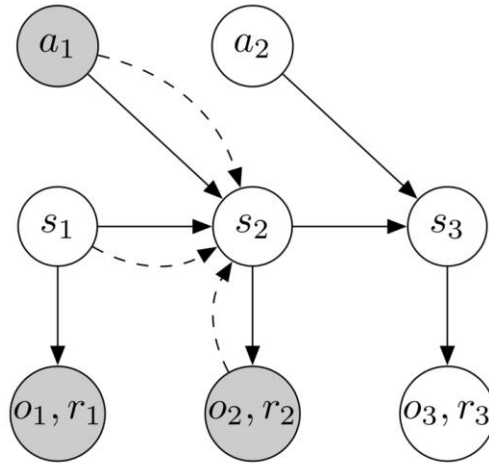굴려온 RNN이용해서 그때그때 Planning한다.

# Recurrent State Space Model (RSSM)
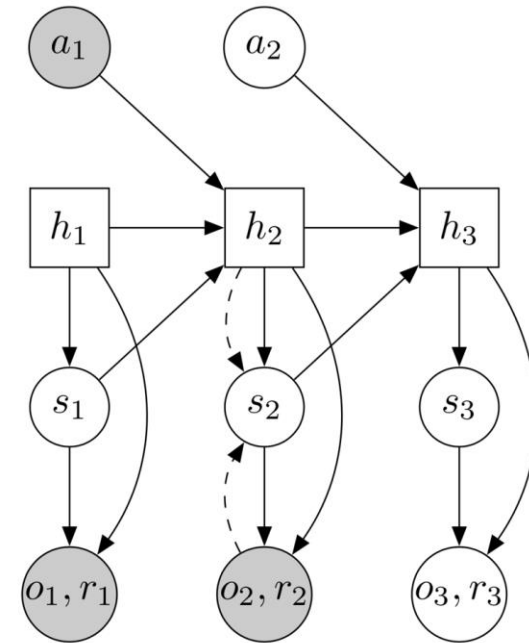
■ Learned Latent Dynamics Model

# Recurrent State Space Model (RSSM)



(a) Deterministic model (RNN)

(b) Stochastic model (SSM)

(c) Recurrent state-space model (RSSM)

# Recurrent State Space Model (RSSM)

RSSM에서 State는 Deterministic & Stochastic Components 으로 이루어져 있다.

## Deterministic Components

Discrete Time step        $t$
Hidden states        $s_t$
Image observations        $o_t$
Continuous action vectors        $a_t$
Scalar rewards        $r_t$

## Stochastic Components

$$
\begin{aligned}
\text{Transition function:} \quad & s_t \sim \mathrm{p}(s_t \mid s_{t-1}, a_{t-1}) \\
\text{Observation function:} \quad & o_t \sim \mathrm{p}(o_t \mid s_t) \\
\text{Reward function:} \quad & r_t \sim \mathrm{p}(r_t \mid s_t) \\
\text{Policy:} \quad & a_t \sim \mathrm{p}(a_t \mid o_{\leq t}, a_{<t}),
\end{aligned} \quad (1)
$$

# Recurrent State Space Model (RSSM)

## Motivation

### Deterministic

Stochastic transitions make it difficult for the transition model to reliably remember information for multiple time steps.
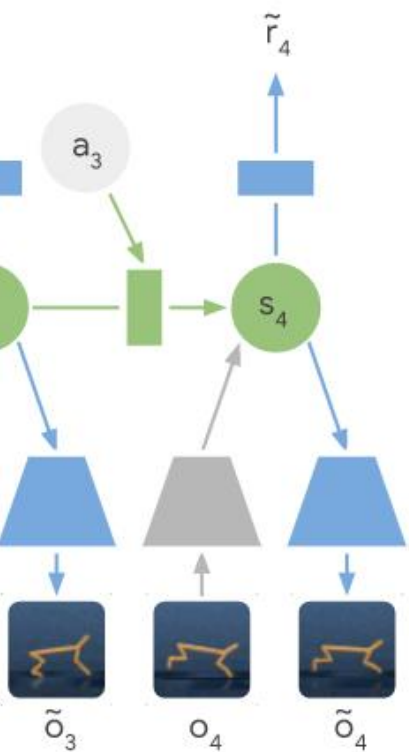
### Stochastic

Split the state into stochastic and deterministic parts, allowing the model to robustly learn to predict multiple futures.

## Detail

All information about the observations must pass through the sampling step of the encoder

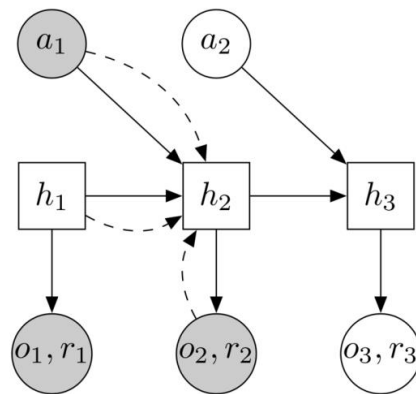to avoid a deterministic shortcut from inputs to reconstructions.

## ■ Recurrent State Space Model (RSSM)



(a) Deterministic model (RNN)   (b) Stochastic model (SSM)   (c) Recurrent state-space model (RSSM)

Transition function: $\quad s_t \sim \mathrm{p}(s_t \mid s_{t-1}, a_{t-1})$

Observation function: $\quad o_t \sim \mathrm{p}(o_t \mid s_t)$

Reward function: $\quad r_t \sim \mathrm{p}(r_t \mid s_t)$

Policy: $\quad a_t \sim \mathrm{p}(a_t \mid o_{\leqslant t}, a_{<t})$,

$(1)$

Deterministic state model: $\quad h_t = f(h_{t-1}, s_{t-1}, a_{t-1})$

Stochastic state model: $\quad s_t \sim p(s_t \mid h_t)$

Observation model: $\quad o_t \sim p(o_t \mid h_t, s_t)$

Reward model: $\quad r_t \sim p(r_t \mid h_t, s_t)$,

$(4)$

# Latent Dynamics (PlaNet, Deep Planning Network)

■ Planning in Latent Space



Past Experience

searches for the best sequence of future actions.

# Latent Dynamics (PlaNet, Deep Planning Network)

**Algorithm 1:** Deep Planning Network (PlaNet)

**Input :**

| | | | |
|---|---|---|---|
| $R$ | Action repeat | $p(s_t \mid s_{t-1}, a_{t-1})$ | Transition model |
| $S$ | Seed episodes | $p(o_t \mid s_t)$ | Observation model |
| $C$ | Collect interval | $p(r_t \mid s_t)$ | Reward model |
| $B$ | Batch size | $q(s_t \mid o_{\leq t}, a_{<t})$ | Encoder |
| $L$ | Chunk length | $p(\epsilon)$ | Exploration noise |
| $\alpha$ | Learning rate | | |

1  Initialize dataset $\mathcal{D}$ with $S$ random seed episodes.
2  Initialize model parameters $\theta$ randomly.
3  **while** *not converged* **do**

    `// Model fitting`
4      **for** *update step* $s = 1..C$ **do**
5          Draw sequence chunks $\{(o_t, a_t, r_t)_{t=k}^{L+k}\}_{i=1}^{B} \sim \mathcal{D}$ uniformly at random from the dataset.
6          Compute loss $\mathcal{L}(\theta)$ from Equation 8.
7          Update model parameters $\theta \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}(\theta)$.

    `// Data collection`
8      $o_1 \leftarrow$ `env.reset()`
9      **for** *time step* $t = 1..\lceil \frac{T}{R} \rceil$ **do**
10         Infer belief over current state $q(s_t \mid o_{\leq t}, a_{<t})$ from the history.
11         $a_t \leftarrow$ `planner`$(q(s_t \mid o_{\leq t}, a_{<t}), p)$, see Algorithm 2 in the appendix for details.
12         Add exploration noise $\epsilon \sim p(\epsilon)$ to the action.
13         **for** *action repeat* $k = 1..R$ **do**
14             $r_t^k, o_{t+1}^k \leftarrow$ `env.step`$(a_t)$
15         $r_t, o_{t+1} \leftarrow \sum_{k=1}^{R} r_t^k, o_{t+1}^R$
16     $\mathcal{D} \leftarrow \mathcal{D} \cup \{(o_t, a_t, r_t)_{t=1}^{T}\}$

PlaNet은 Planning을 위해
Cross Entropy Method (CEM) 을 사용했다.

**Algorithm 2:** Latent planning with CEM

**Input :**

| | | | |
|---|---|---|---|
| $H$ | Planning horizon distance | $q(s_t \mid o_{\leq t}, a_{<t})$ | Current state belief |
| $I$ | Optimization iterations | $p(s_t \mid s_{t-1}, a_{t-1})$ | Transition model |
| $J$ | Candidates per iteration | $p(r_t \mid s_t)$ | Reward model |
| $K$ | Number of top candidates to fit | | |

1  Initialize factorized belief over action sequences $q(a_{t:t+H}) \leftarrow \text{Normal}(0, \mathbb{I})$.
2  **for** *optimization iteration* $i = 1..I$ **do**

    `// Evaluate J action sequences from the current belief.`
3      **for** *candidate action sequence* $j = 1..J$ **do**
4          $a_{t:t+H}^{(j)} \sim q(a_{t:t+H})$
5          $s_{t:t+H+1}^{(j)} \sim q(s_t \mid o_{1:t}, a_{1:t-1}) \prod_{\tau=t+1}^{t+H+1} p(s_\tau \mid s_{\tau-1}, a_{\tau-1}^{(j)})$
6          $R^{(j)} = \sum_{\tau=t+1}^{t+H+1} \mathbb{E}[p(r_\tau \mid s_\tau^{(j)})]$

    `// Re-fit belief to the K best action sequences.`
7      $\mathcal{K} \leftarrow \text{argsort}(\{R^{(j)}\}_{j=1}^{J})_{1:K}$
8      $\mu_{t:t+H} = \frac{1}{K} \sum_{k \in \mathcal{K}} a_{t:t+H}^{(k)}, \quad \sigma_{t:t+H} = \frac{1}{K-1} \sum_{k \in \mathcal{K}} |a_{t:t+H}^{(k)} - \mu_{t:t+H}|.$
9      $q(a_{t:t+H}) \leftarrow \text{Normal}(\mu_{t:t+H}, \sigma_{t:t+H}^2 \mathbb{I})$
10 **return** *first action mean* $\mu_t$.

## ■ Objective Function

$$\ln p(o_{1:T} \mid a_{1:T}) \triangleq \ln \int \prod_t p(s_t \mid s_{t-1}, a_{t-1}) p(o_t \mid s_t) \, ds_{1:T}$$

$$\geq \sum_{t=1}^{T} \Big( \underbrace{\mathrm{E}_{q(s_t \mid o_{\leq t}, a_{<t})} \big[ \ln p(o_t \mid s_t) \big]}_{\text{reconstruction}} \ \hookleftarrow$$

$$- \underbrace{\mathrm{E}_{q(s_{t-1} \mid o_{\leq t-1}, a_{<t-1})} \big[ \mathrm{KL} \big[ q(s_t \mid o_{\leq t}, a_{<t}) \parallel p(s_t \mid s_{t-1}, a_{t-1}) \big] \big]}_{\text{complexity}} \Big). \qquad (3)$$

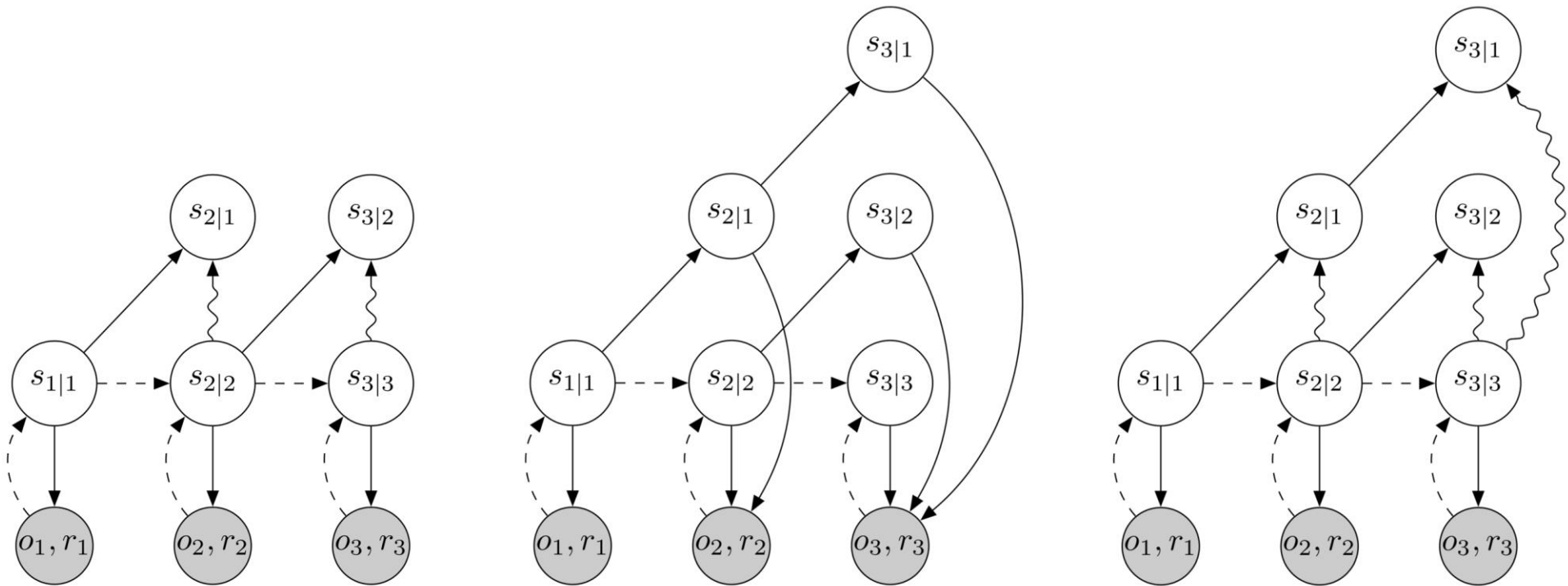Limitation    Gradient가 한 step만 흐르기 때문에, **1-step prediction**만 보장한다

# Latent Overshooting

# Latent OverShooting

## Motivation

multi-step prediction은 추가 image를 가하지 않았을 때 **latent space에서의 손실**로 개선할 수 있다.



(a) Standard variational bound  (b) Observation overshooting  (c) Latent overshooting

# Latent OverShooting

■ Define **Multi Step Prediction**

$$p(s_t \mid s_{t-d}) \triangleq \int \prod_{\tau=t-d+1}^{t} p(s_\tau \mid s_{\tau-1}) \, ds_{t-d+1:t-1} \qquad (5)$$

$$= \mathrm{E}_{p(s_{t-1} \mid s_{t-d})}[p(s_t \mid s_{t-1})].$$

*d = 1 이면  1 - s t e p   p r e d i c t i o n*

■ **Objective Function** (given distance d)

$$\ln p_d(o_{1:T}) \triangleq \ln \int \prod_{t=1}^{T} p(s_t \mid s_{t-d})p(o_t \mid s_t) \, ds_{1:T}$$

$$\geq \sum_{t=1}^{T} \left( \underbrace{\mathrm{E}_{q(s_t \mid o_{\leq t})}[\ln p(o_t \mid s_t)]}_{\text{reconstruction}} \hookleftarrow \right. \qquad (6)$$

$$\left. - \underbrace{\mathrm{E}_{p(s_{t-1} \mid s_{t-d})q(s_{t-d} \mid o_{\leq t-d})}\big[\mathrm{KL}[q(s_t \mid o_{\leq t}) \parallel p(s_t \mid s_{t-1})]\big]}_{\text{multi-step prediction}} \right).$$

와 ~ Marcov하다 !

# Latent OverShooting

■ **Objective Function** (all distances up to the planning horizon.)

$$\frac{1}{D}\sum_{d=1}^{D}\ln p_d(o_{1:T}) \geqslant \sum_{t=1}^{T}\Bigg(\underbrace{\mathrm{E}_{q(s_t|o_{\leqslant t})}[\ln p(o_t \mid s_t)]}_{\text{reconstruction}} \hookleftarrow$$

$$-\frac{1}{D}\sum_{d=1}^{D}\beta_d\,\mathrm{E}\underbrace{\big[\mathrm{KL}[q(s_t \mid o_{\leqslant t}) \parallel p(s_t \mid s_{t-1})]\big]}_{p(s_{t-1}|s_{t-d})q(s_{t-d}|o_{\leqslant t-d})}\Bigg). \quad (7)$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{\text{latent overshooting}}$$

# Latent OverShooting

■ **Objective Function** (all distances up to the planning horizon.)

$$\frac{1}{D}\sum_{d=1}^{D}\ln p_d(o_{1:T}) \geqslant \sum_{t=1}^{T}\left(\underbrace{\mathrm{E}_{q(s_t|o_{\leqslant t})}[\ln p(o_t \mid s_t)]}_{\text{reconstruction}} \;\hookleftarrow\right.$$

$$\left.-\frac{1}{D}\sum_{d=1}^{D}\beta_d\,\underbrace{\mathrm{E}\big[\mathrm{KL}[q(s_t \mid o_{\leqslant t}) \parallel p(s_t \mid s_{t-1})]\big]}_{p(s_{t-1}|s_{t-d})q(s_{t-d}|o_{\leqslant t-d})}\right). \quad (7)$$

$$\underbrace{\phantom{-\frac{1}{D}\sum_{d=1}^{D}\beta_d\,\mathrm{E}\big[\mathrm{KL}[q(s_t \mid o_{\leqslant t}) \parallel p(s_t \mid s_{t-1})]\big]}}_{\text{latent overshooting}}$$
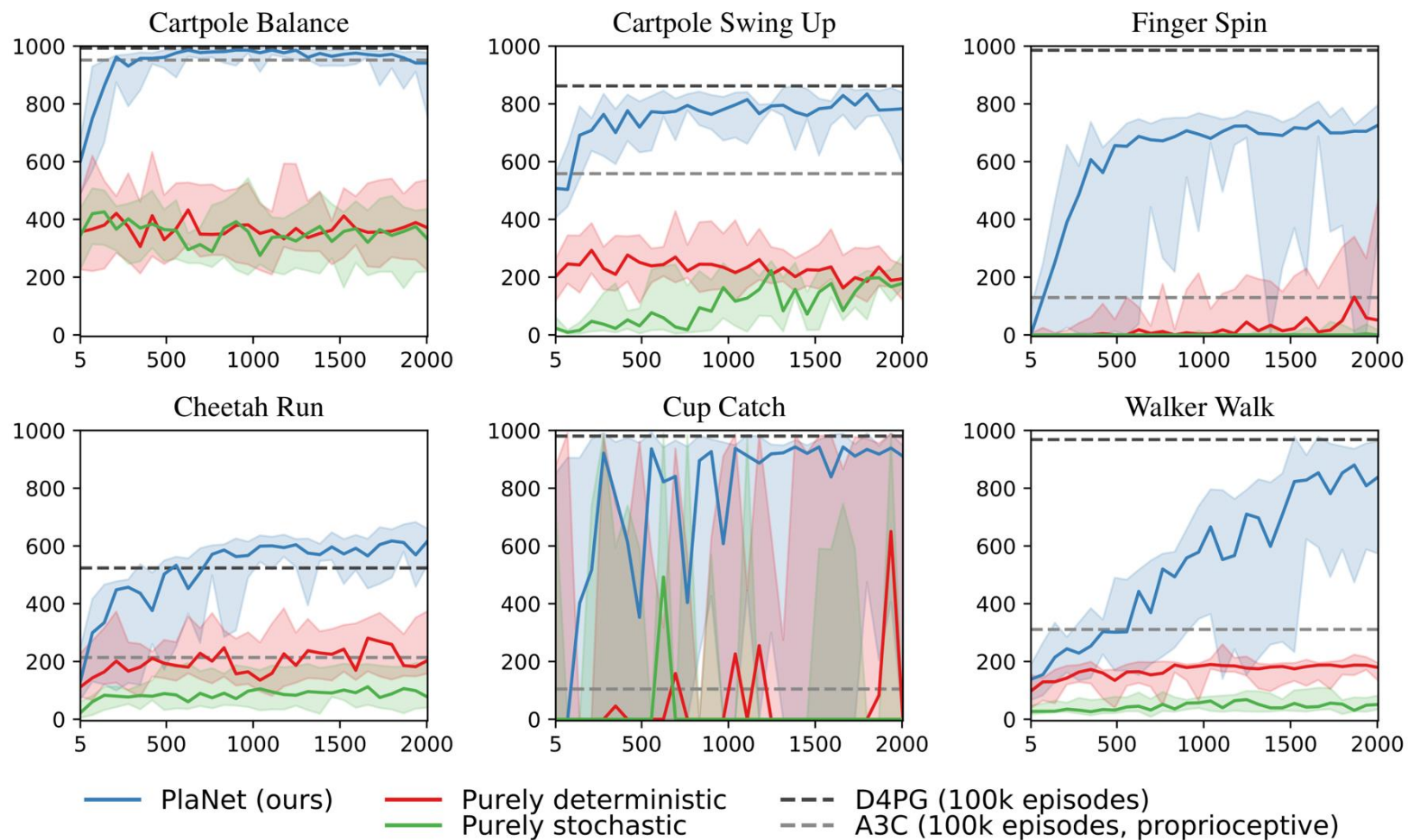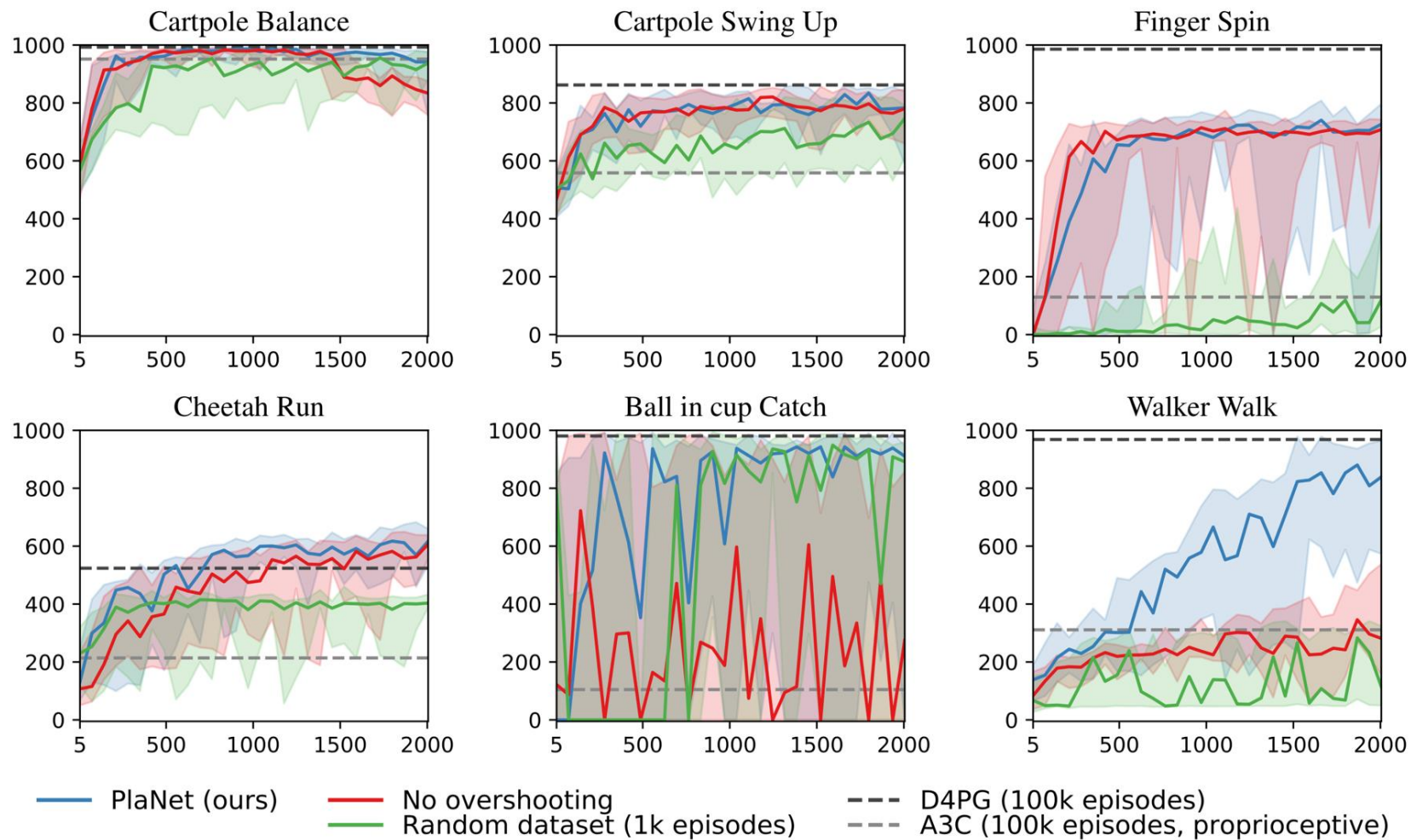
# Experiments

# Comparison (Model-Free Method)

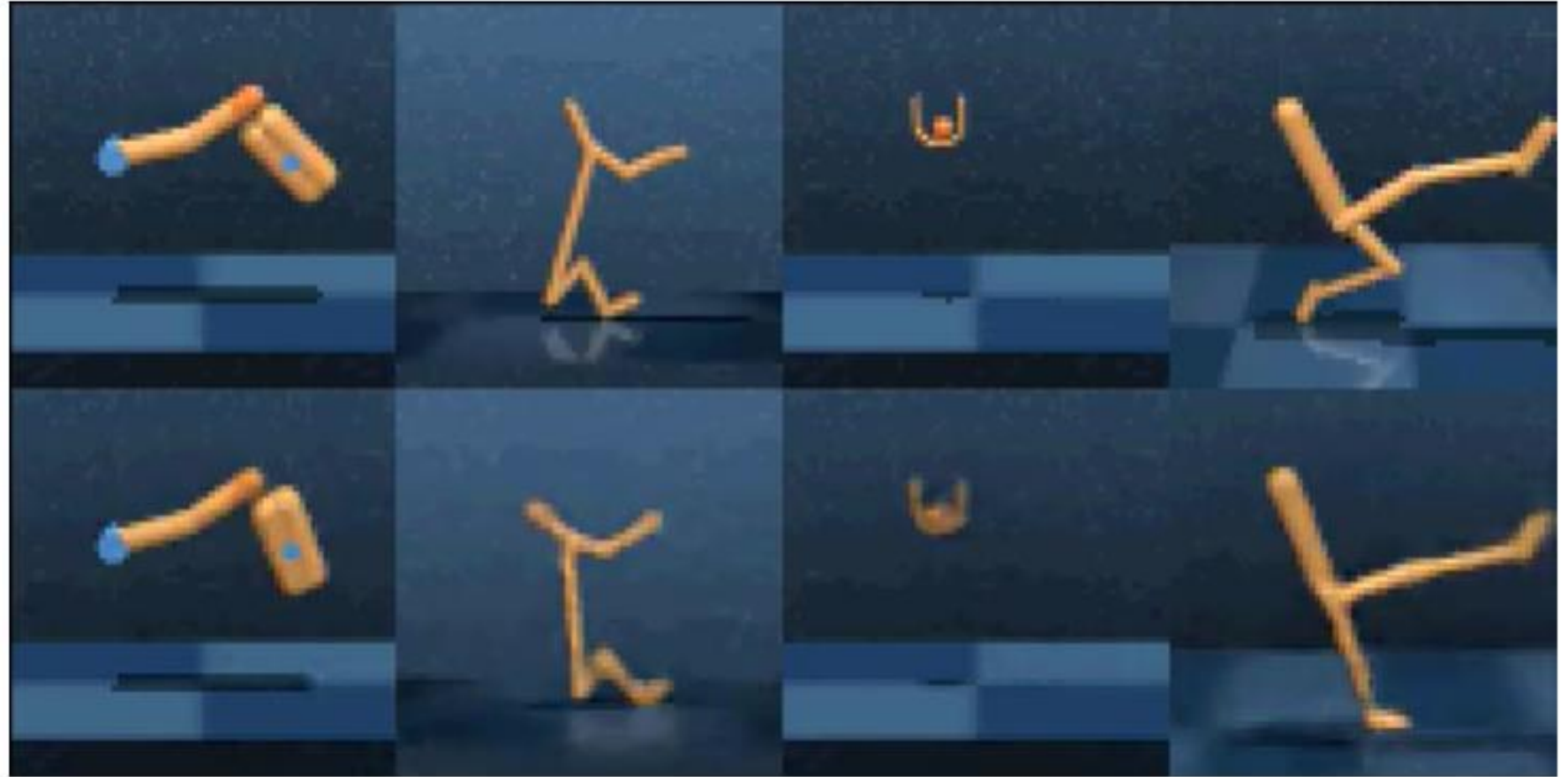| Method | Modality | Episodes | Cartpole Balance | Cartpole Swingup | Finger Spin | Cheetah Run | Ball in cup Catch | Walker Walk |
|---|---|---|---|---|---|---|---|---|
| A3C | proprioceptive | 100,000 | 952 | 558 | 129 | 214 | 105 | 311 |
| D4PG | pixels | 100,000 | 993 | 862 | 985 | 524 | 980 | 968 |
| PlaNet (ours) | pixels | 2,000 | 986 | 831 | 744 | 650 | 914 | 890 |
| CEM + true simulator | simulator state | 0 | 998 | 850 | 825 | 656 | 993 | 994 |
| Data efficiency gain PlaNet over D4PG (factor) | | | 100 | 180 | 16 | 50+ | 20 | 11 |

# Model



Cartpole Balance · Cartpole Swing Up · Finger Spin · Cheetah Run · Cup Catch · Walker Walk

PlaNet (ours) · Purely deterministic · Purely stochastic · D4PG (100k episodes) · A3C (100k episodes, proprioceptive)

# Agent



Cartpole Balance · Cartpole Swing Up · Finger Spin · Cheetah Run · Ball in cup Catch · Walker Walk

PlaNet (ours) — No overshooting — D4PG (100k episodes) — Random dataset (1k episodes) — A3C (100k episodes, proprioceptive)

# One agent all tasks

Episode

Prediction

# Relative Works

# Relative Works

1. **Planning in state space**

2. **Hybrid agents**

3. **Multi-step predictions**

4. **Latent sequence models**

5. **Video prediction**

# End !