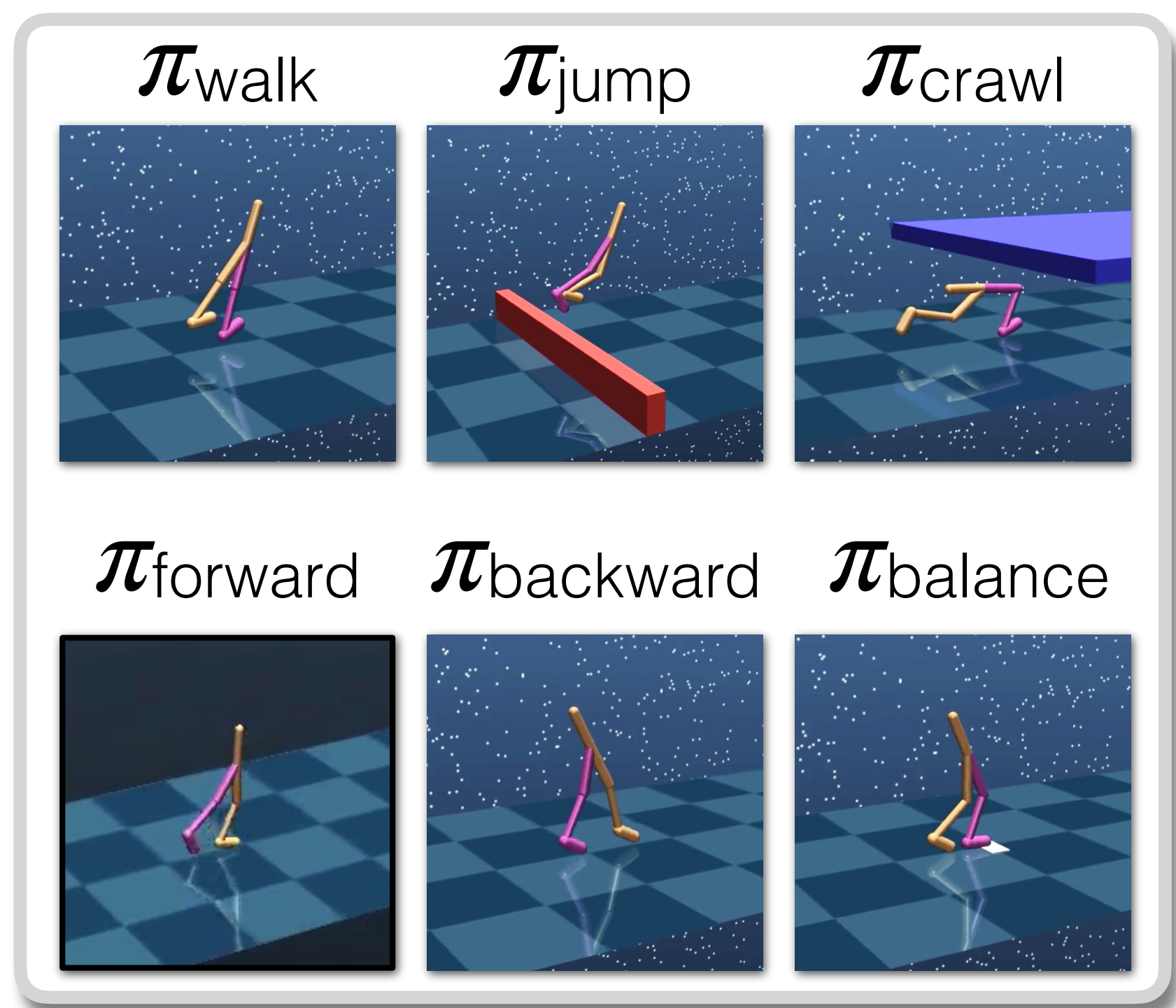


# Composing Complex Skills by Learning Transition Policies

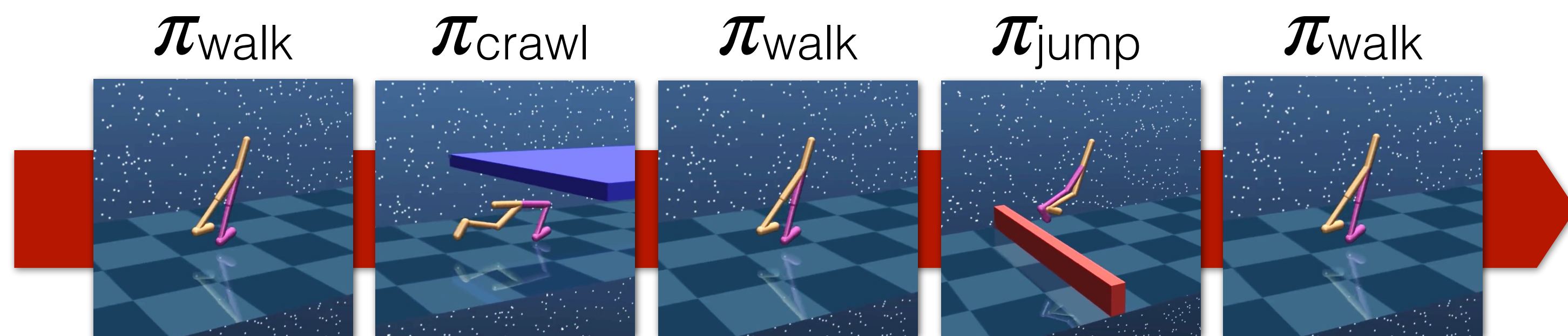
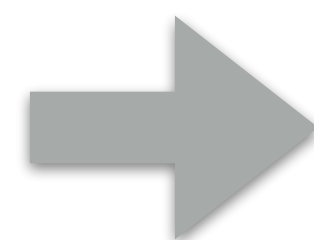
Youngwoon Lee\*, Shao-Hua Sun\*, Sriram Somasundaram, Edward S. Hu, Joseph J. Lim

Presented in ICLR 2019



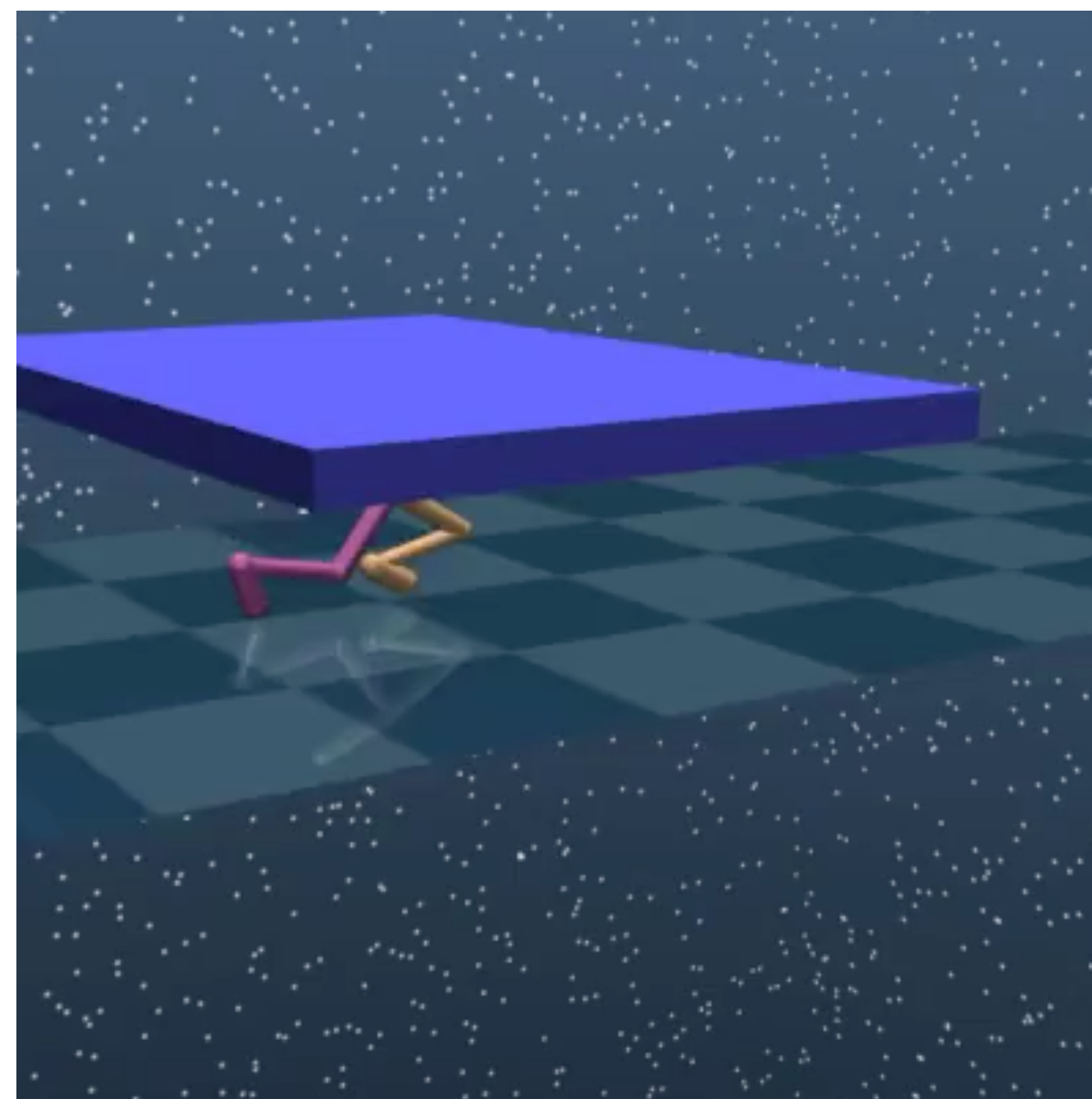
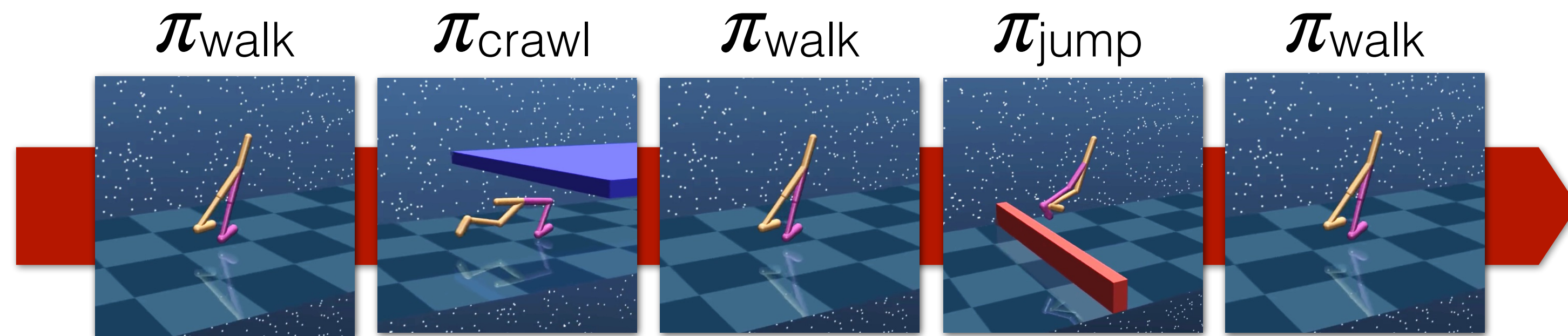


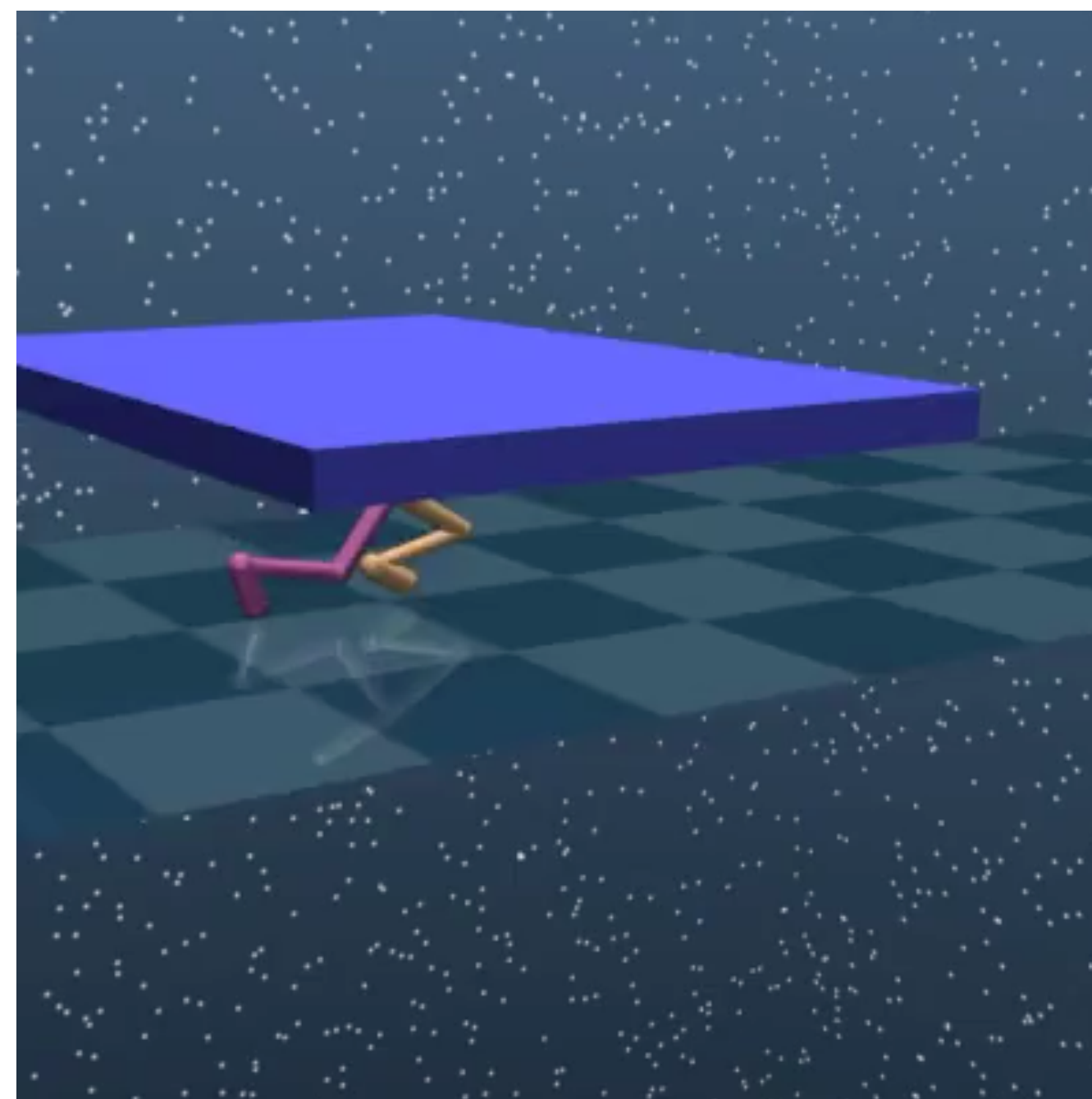
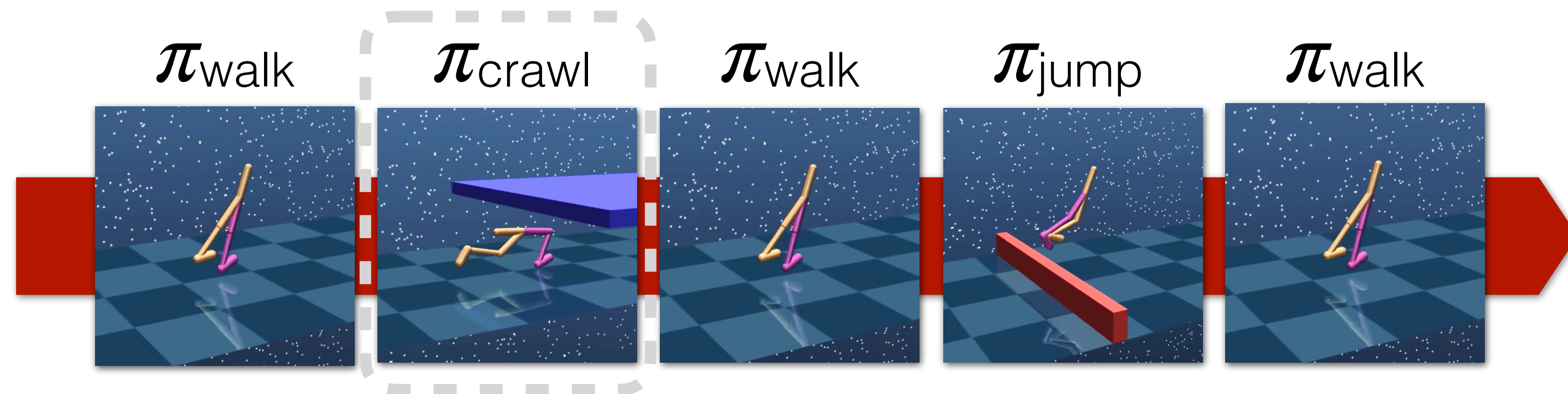
Reusable Skills



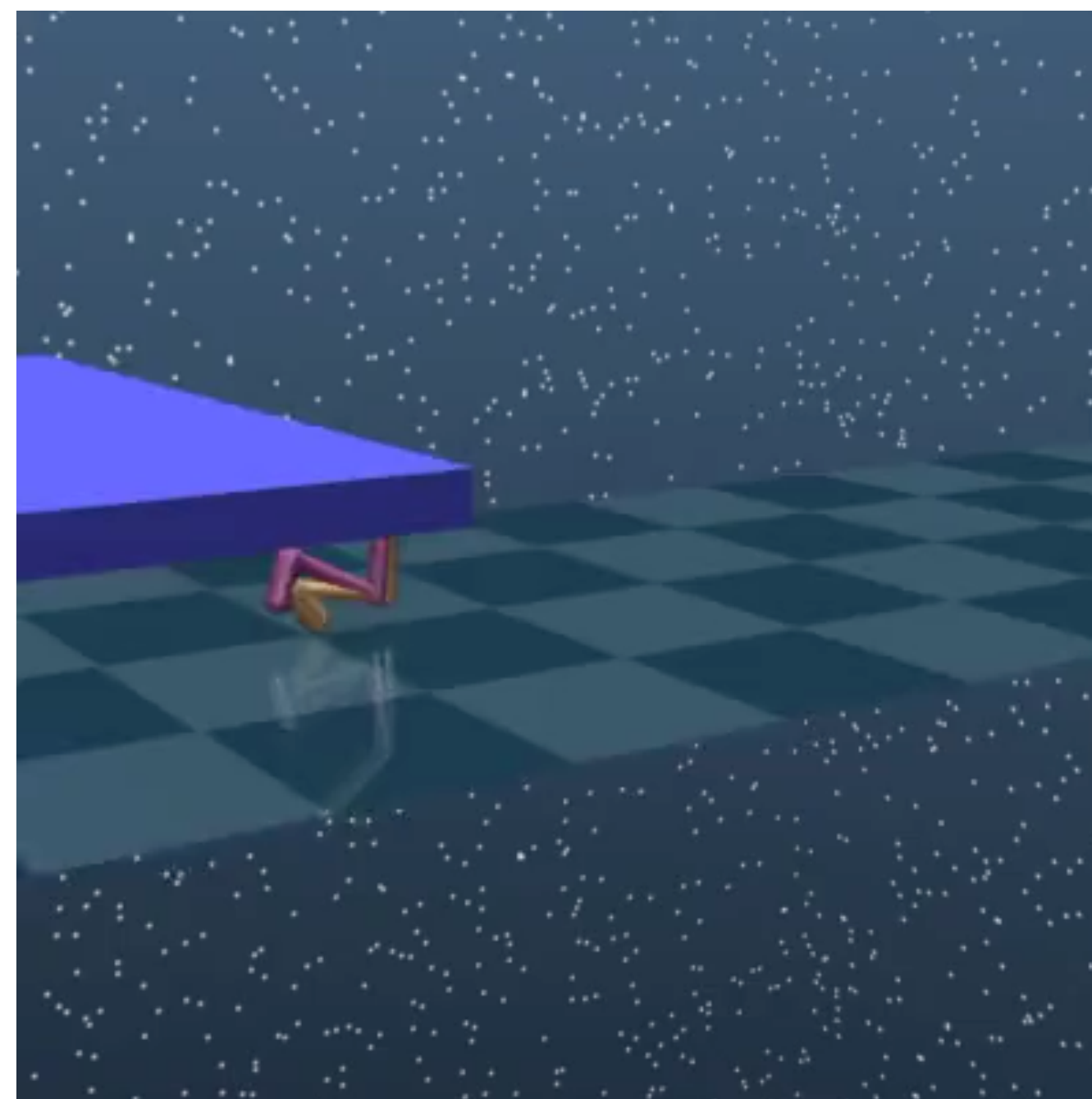
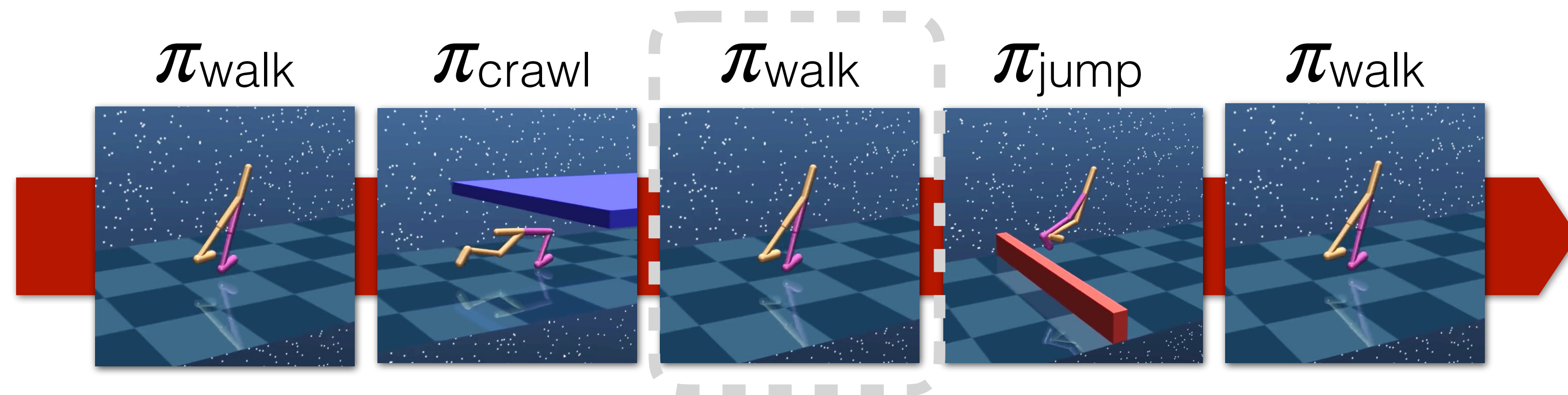
Compositional Task

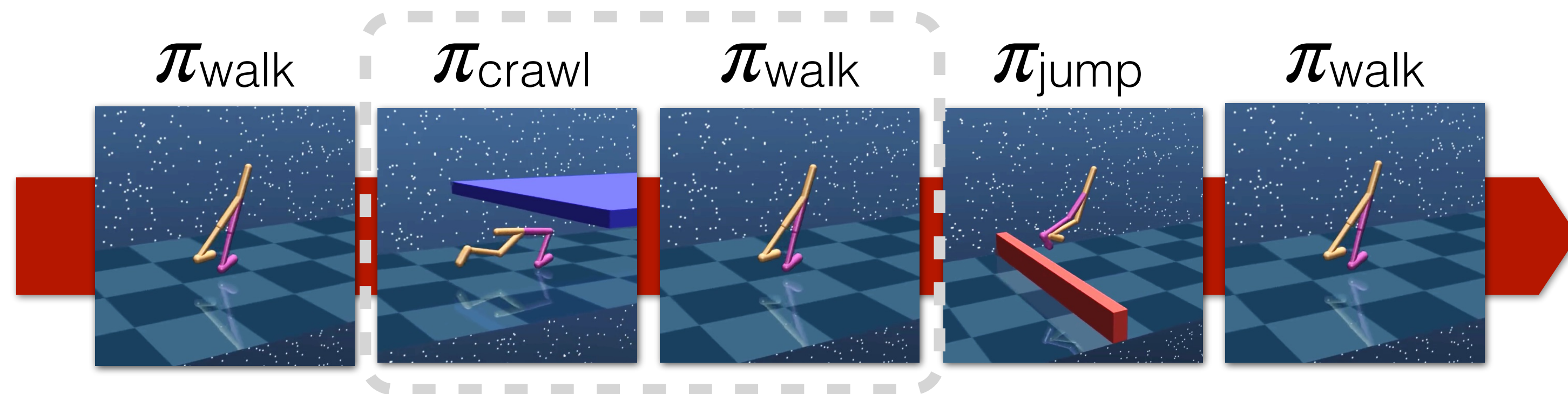




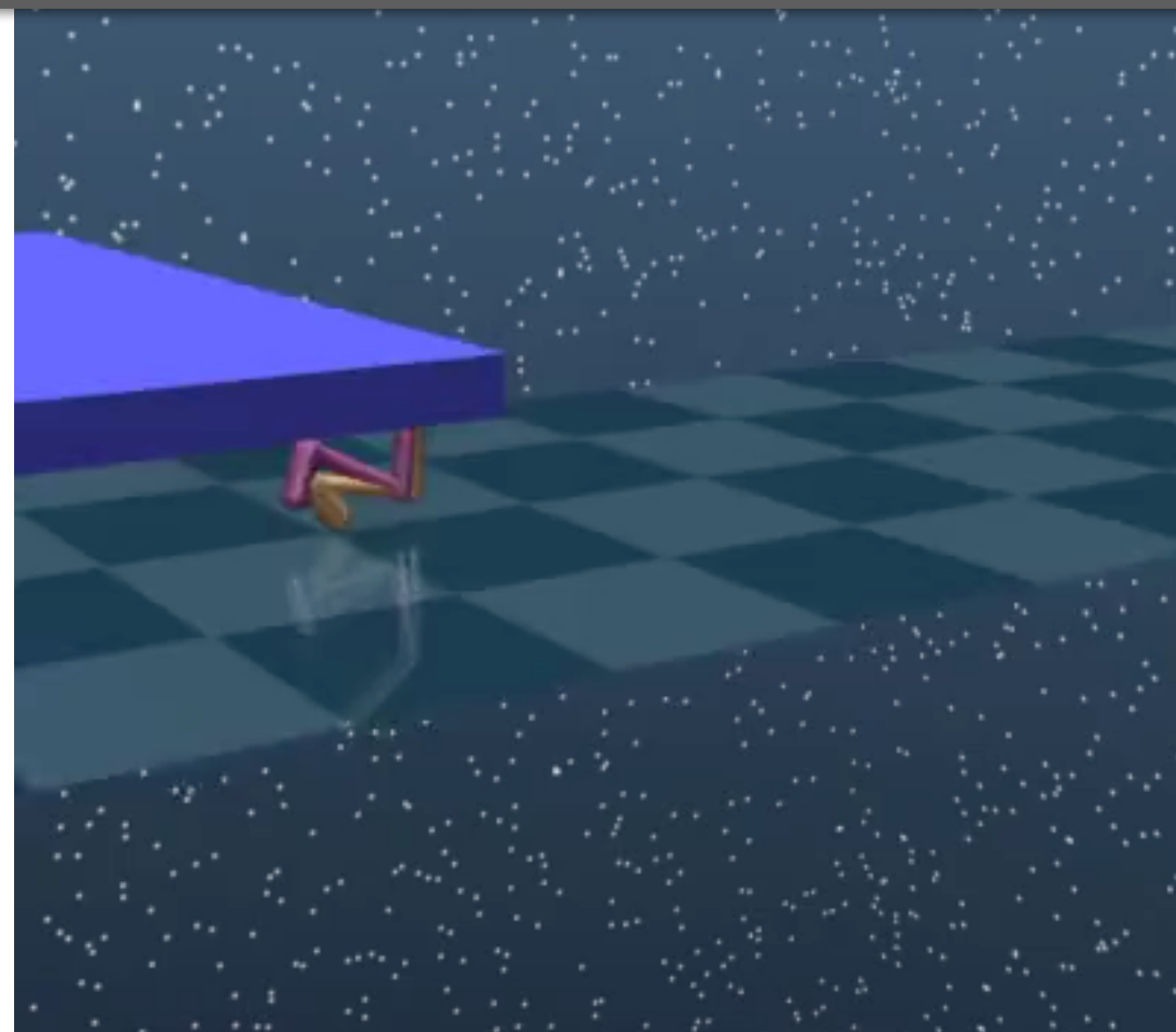




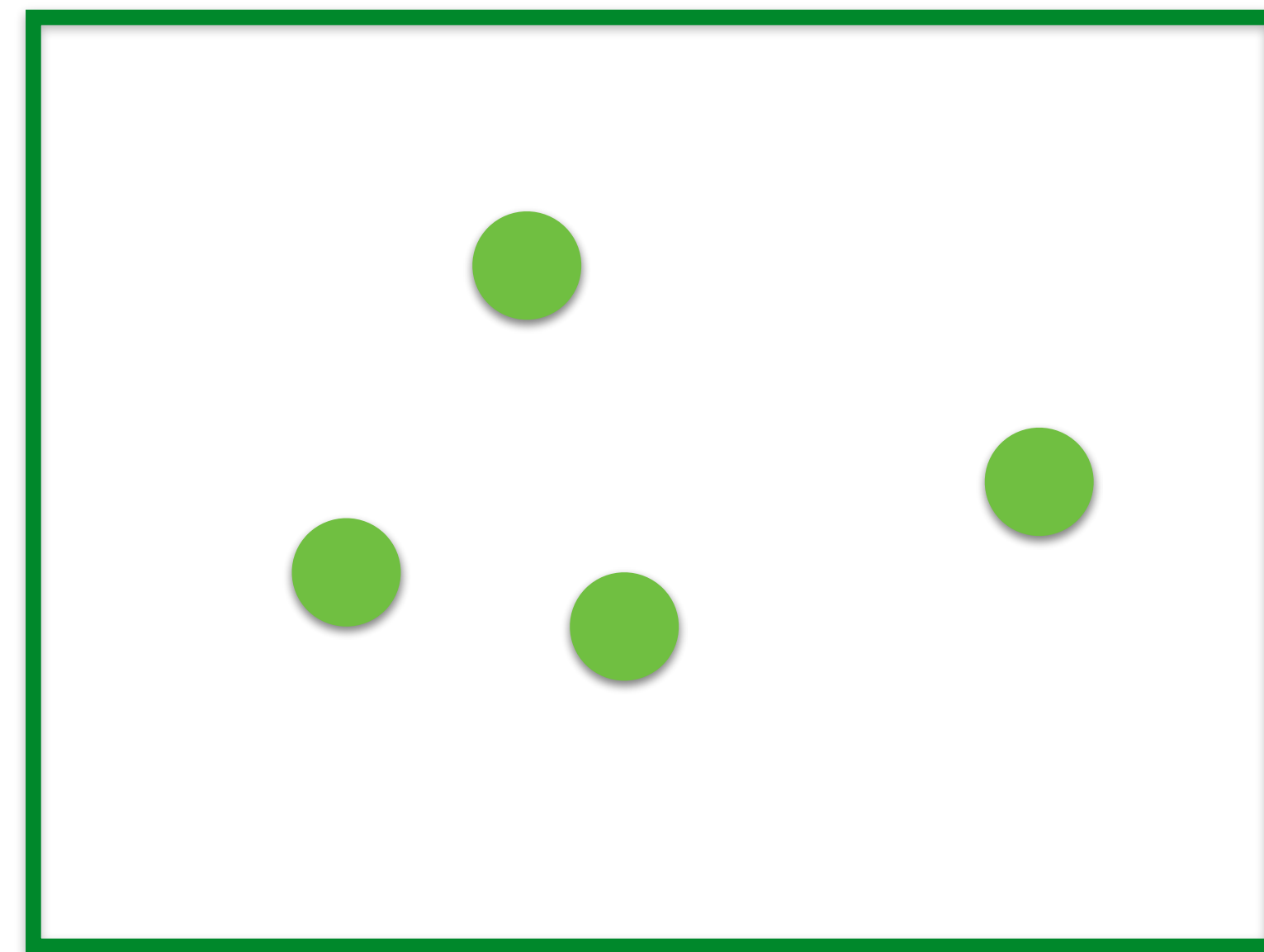
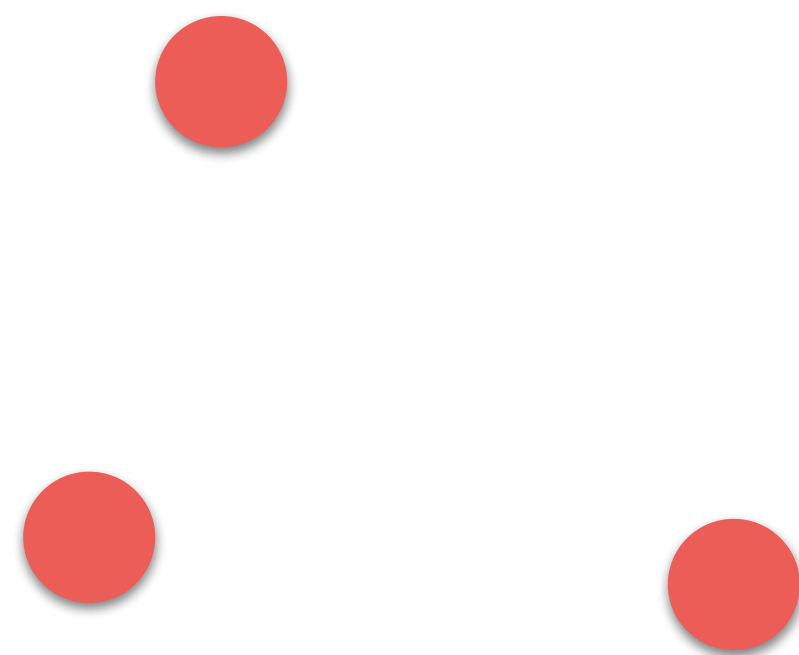




**Fail** since these skills never learned to connect

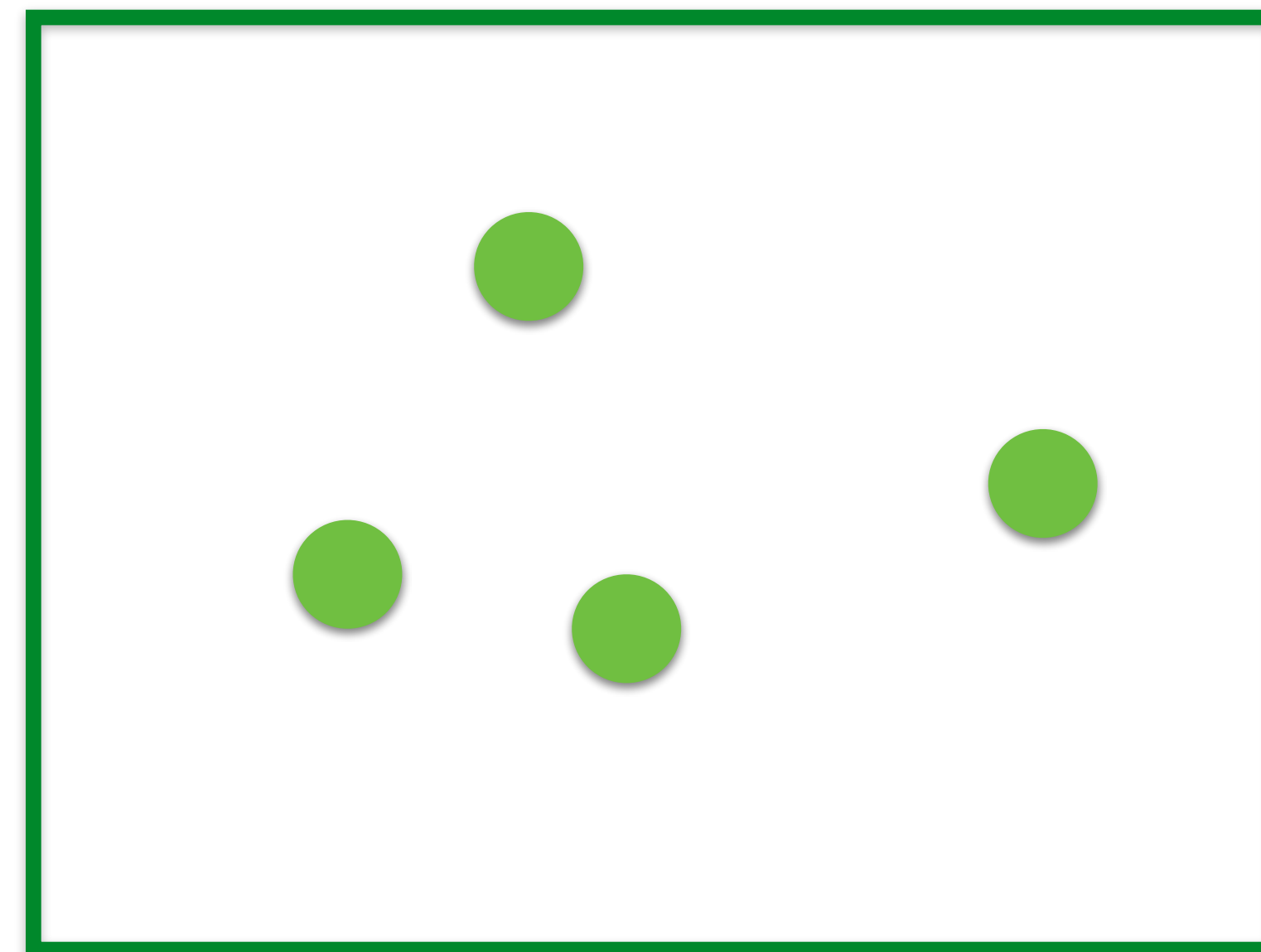
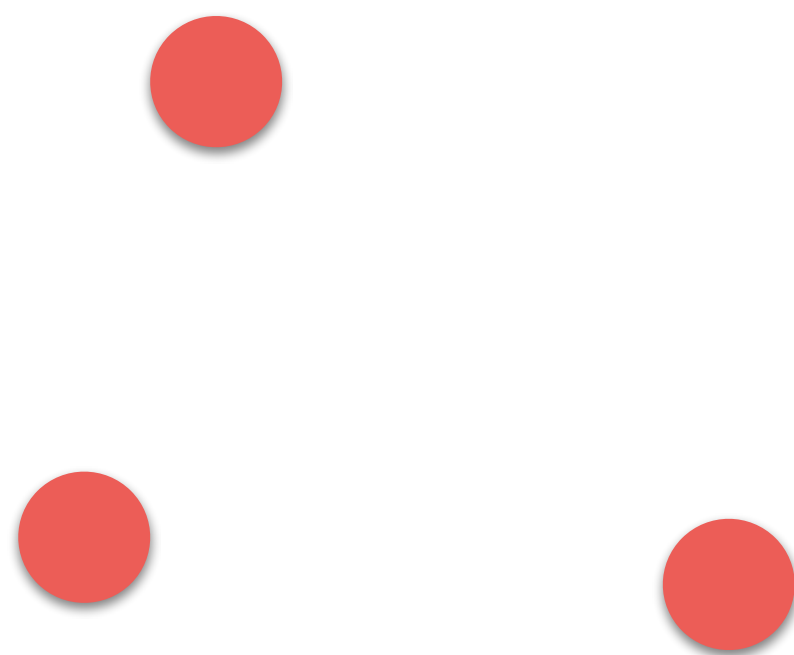
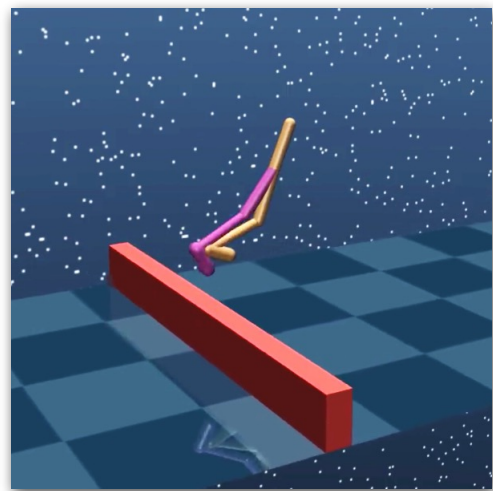






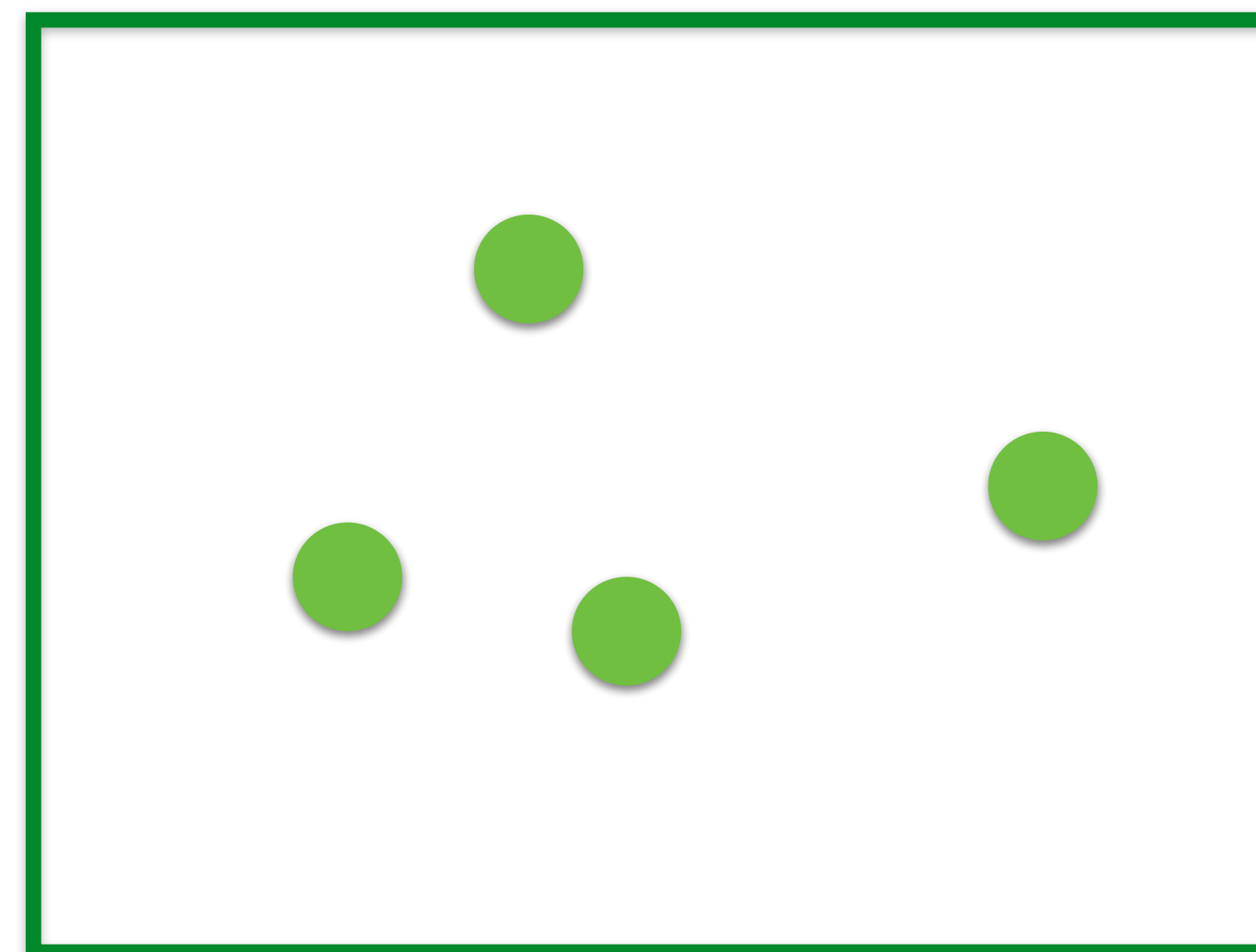
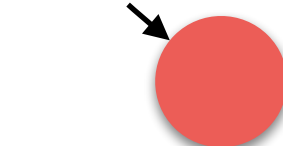
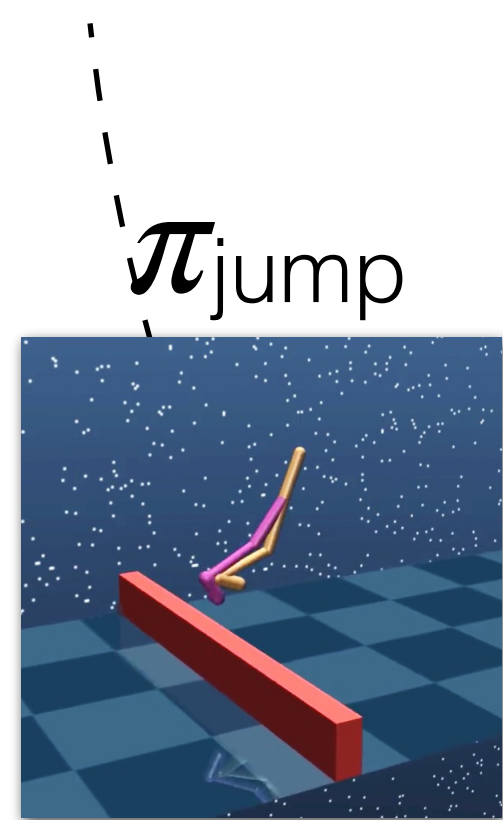
Good initial states for  $\pi_{\text{walk}}$

$\pi_{\text{jump}}$

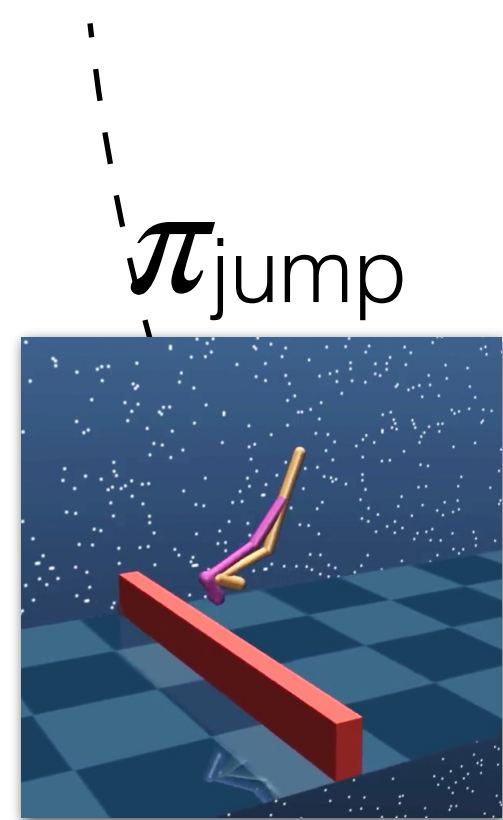


Good initial states for  $\pi_{\text{walk}}$

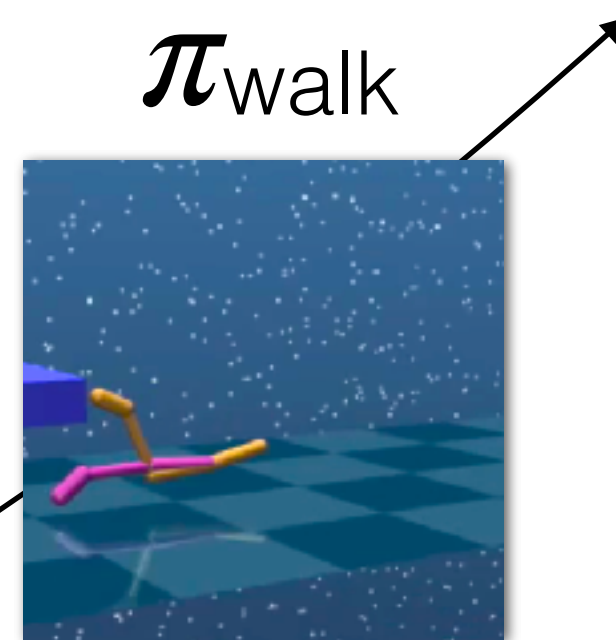




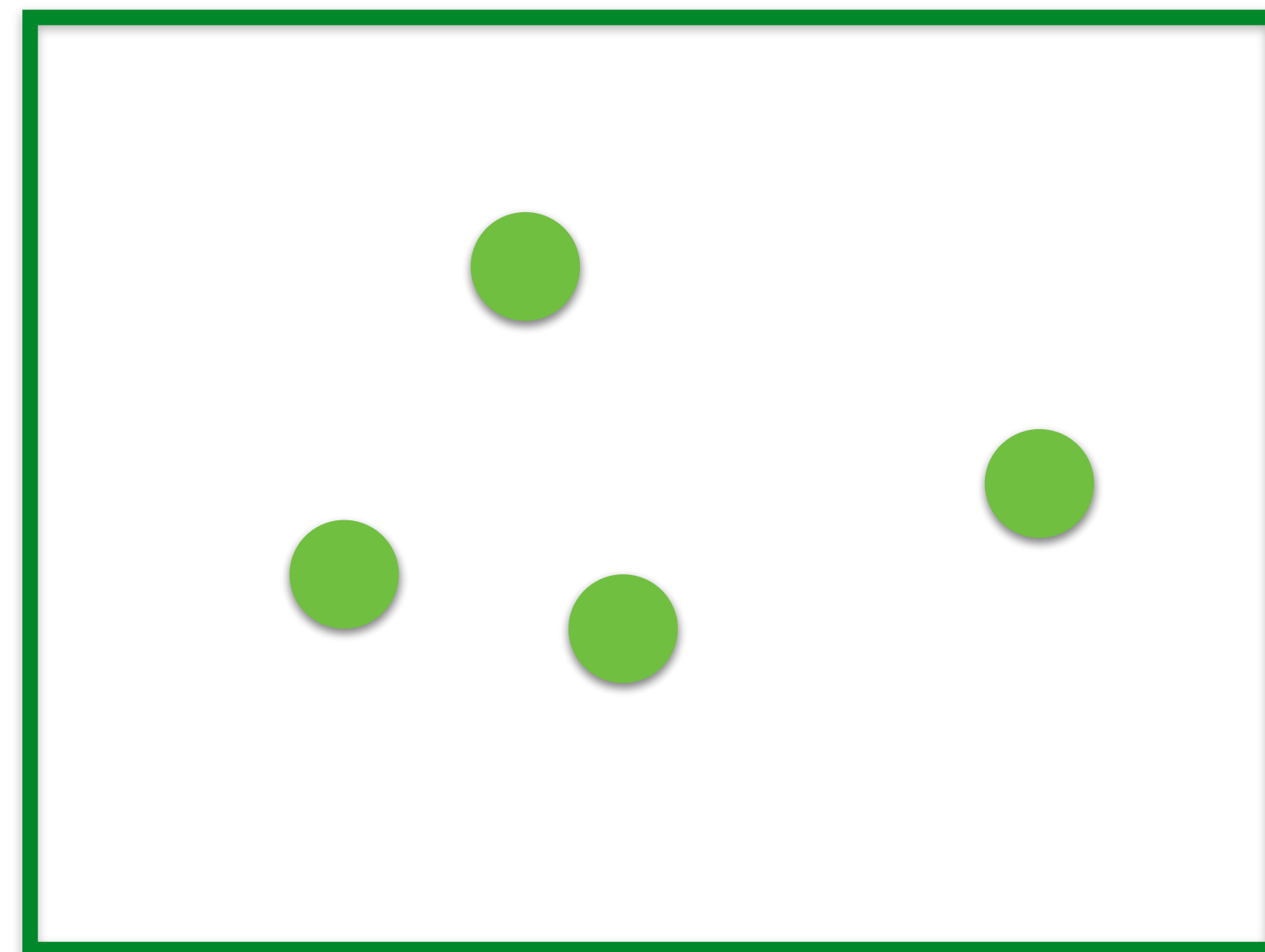
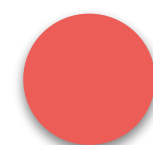
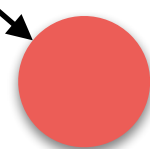
Good initial states for  $\pi_{\text{walk}}$



$\pi_{\text{jump}}$

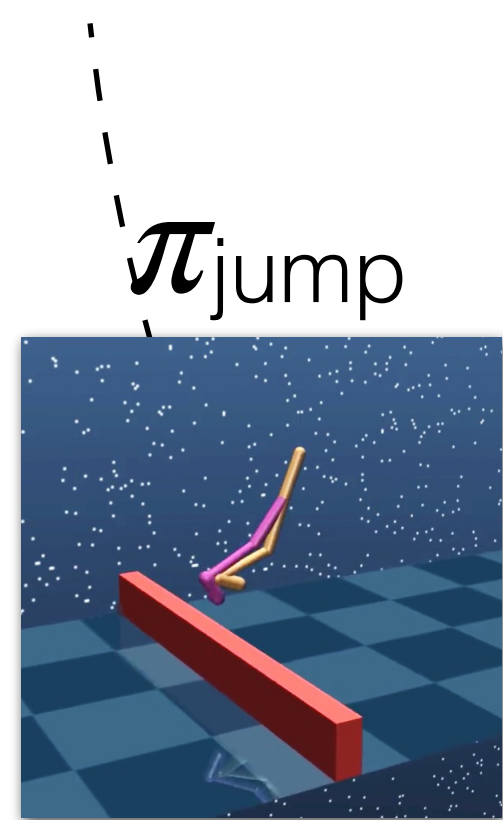


$\pi_{\text{walk}}$

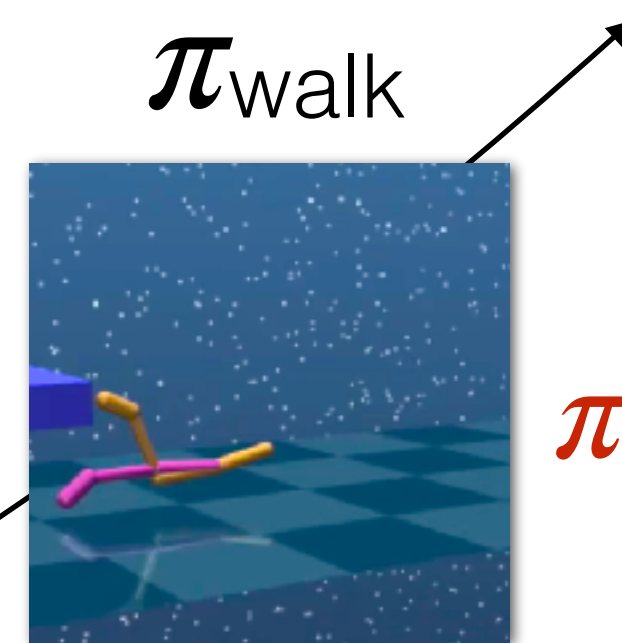


Good initial states for  $\pi_{\text{walk}}$



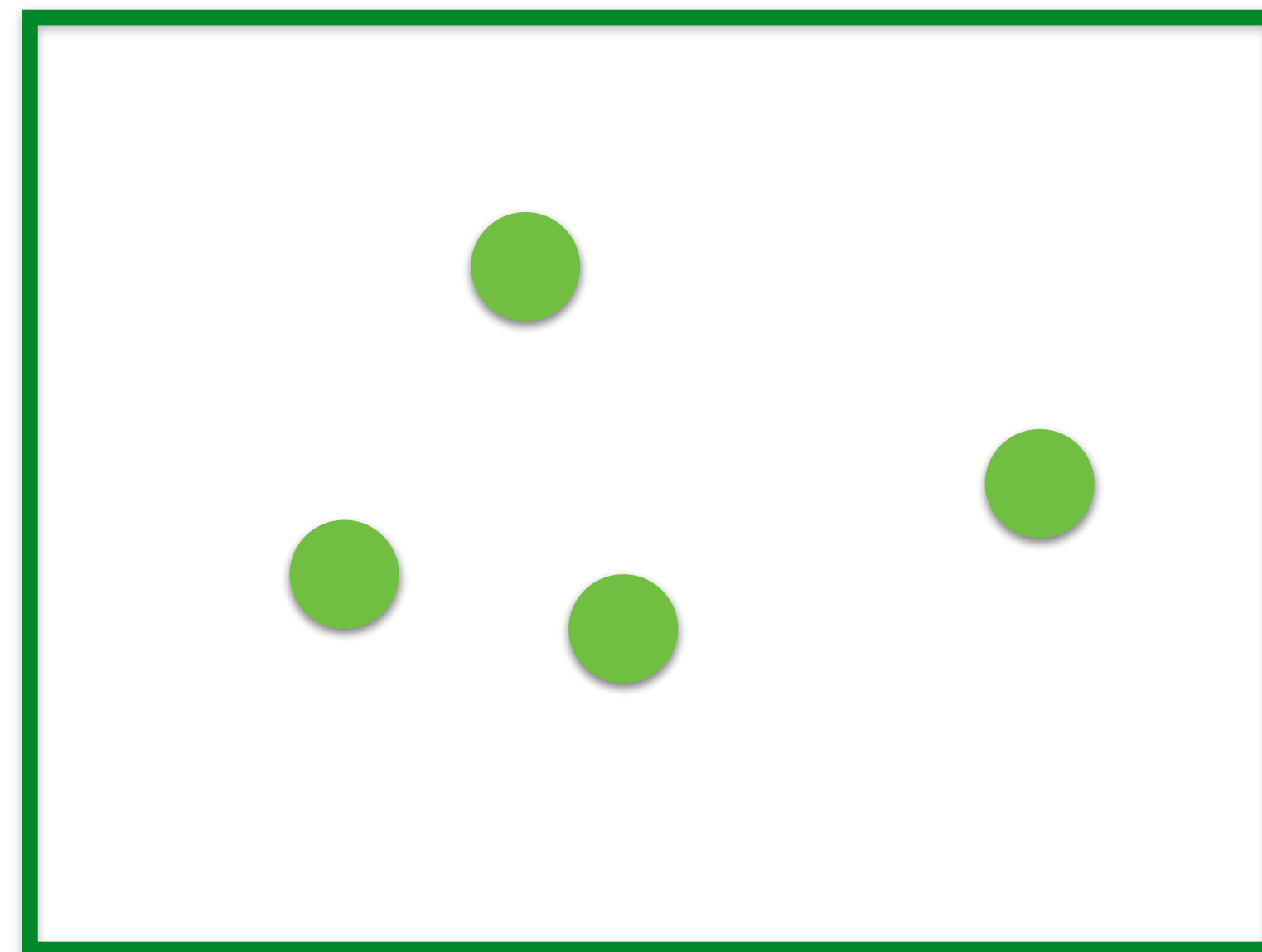
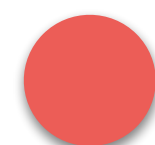
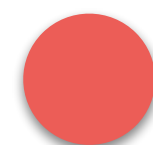
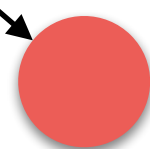


$\pi_{\text{jump}}$

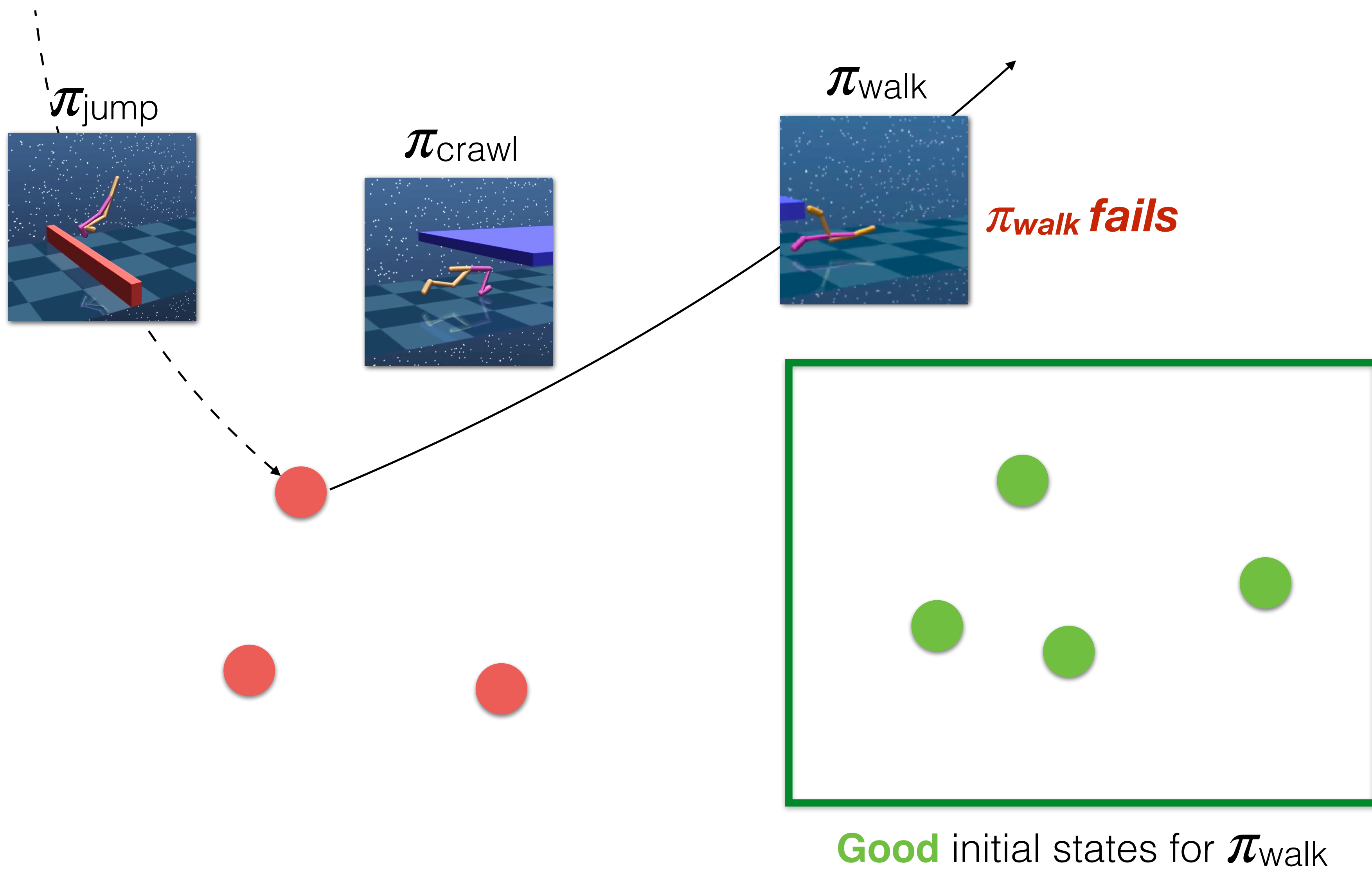


$\pi_{\text{walk}}$

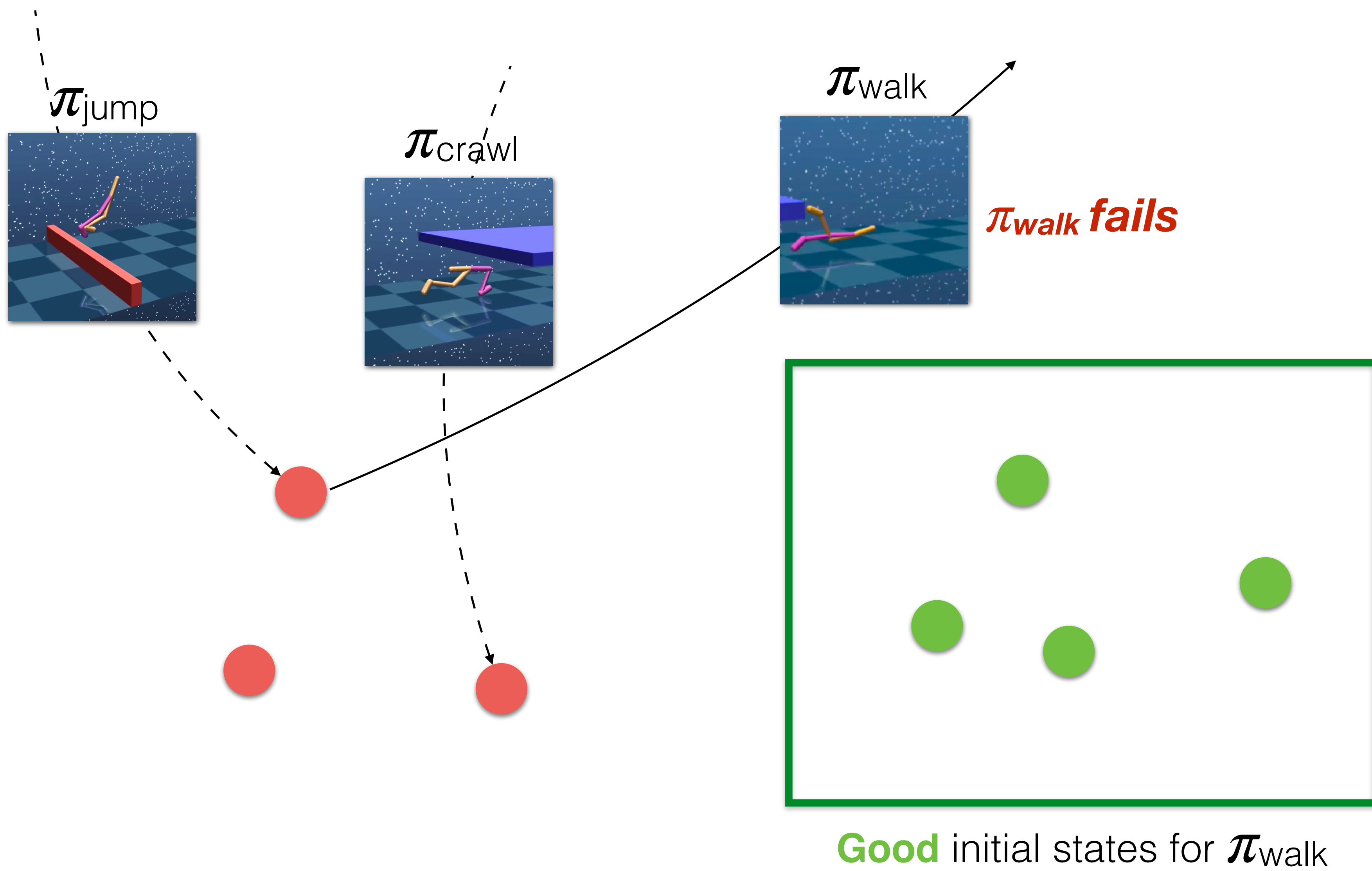
$\pi_{\text{walk fails}}$

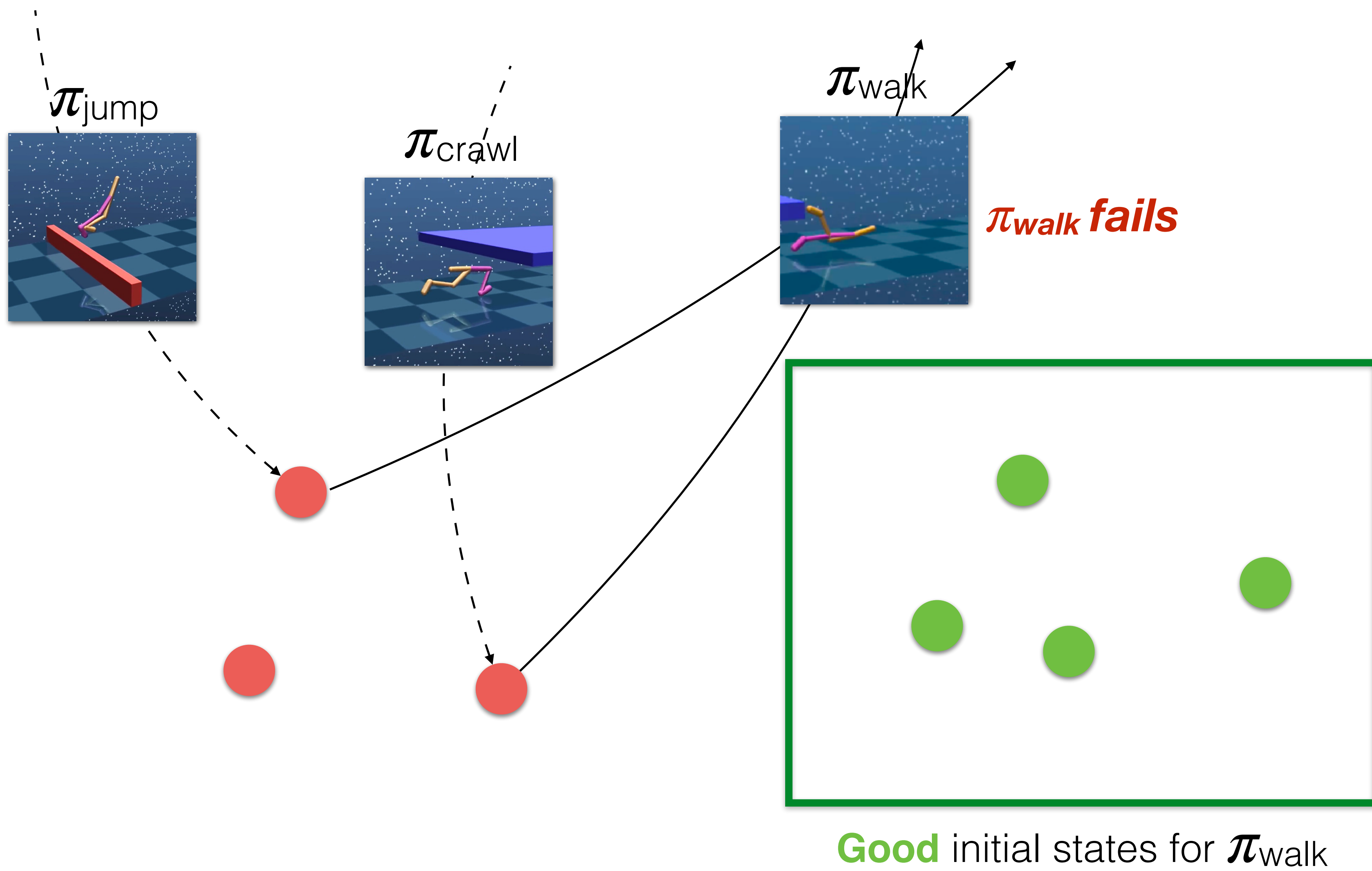


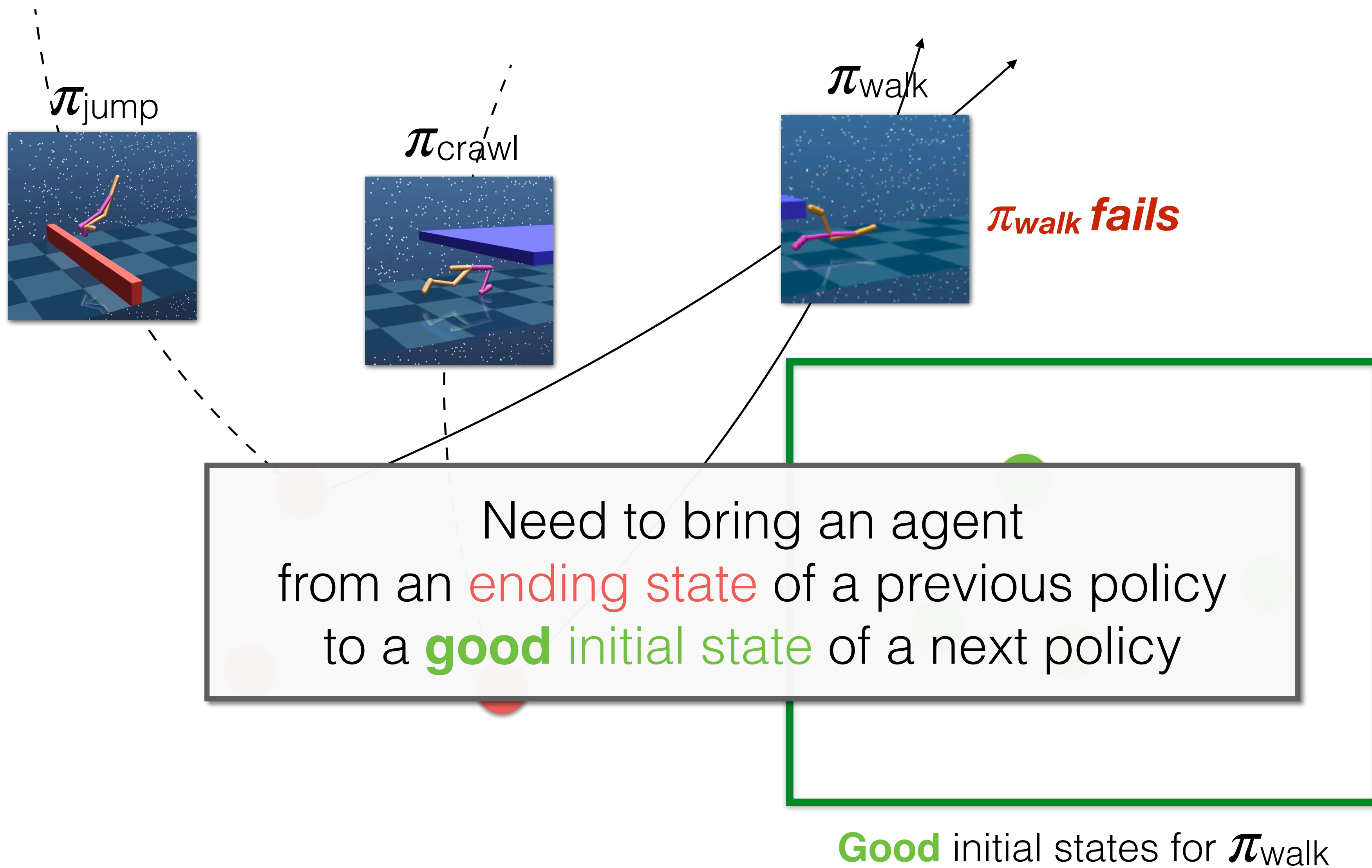
Good initial states for  $\pi_{\text{walk}}$

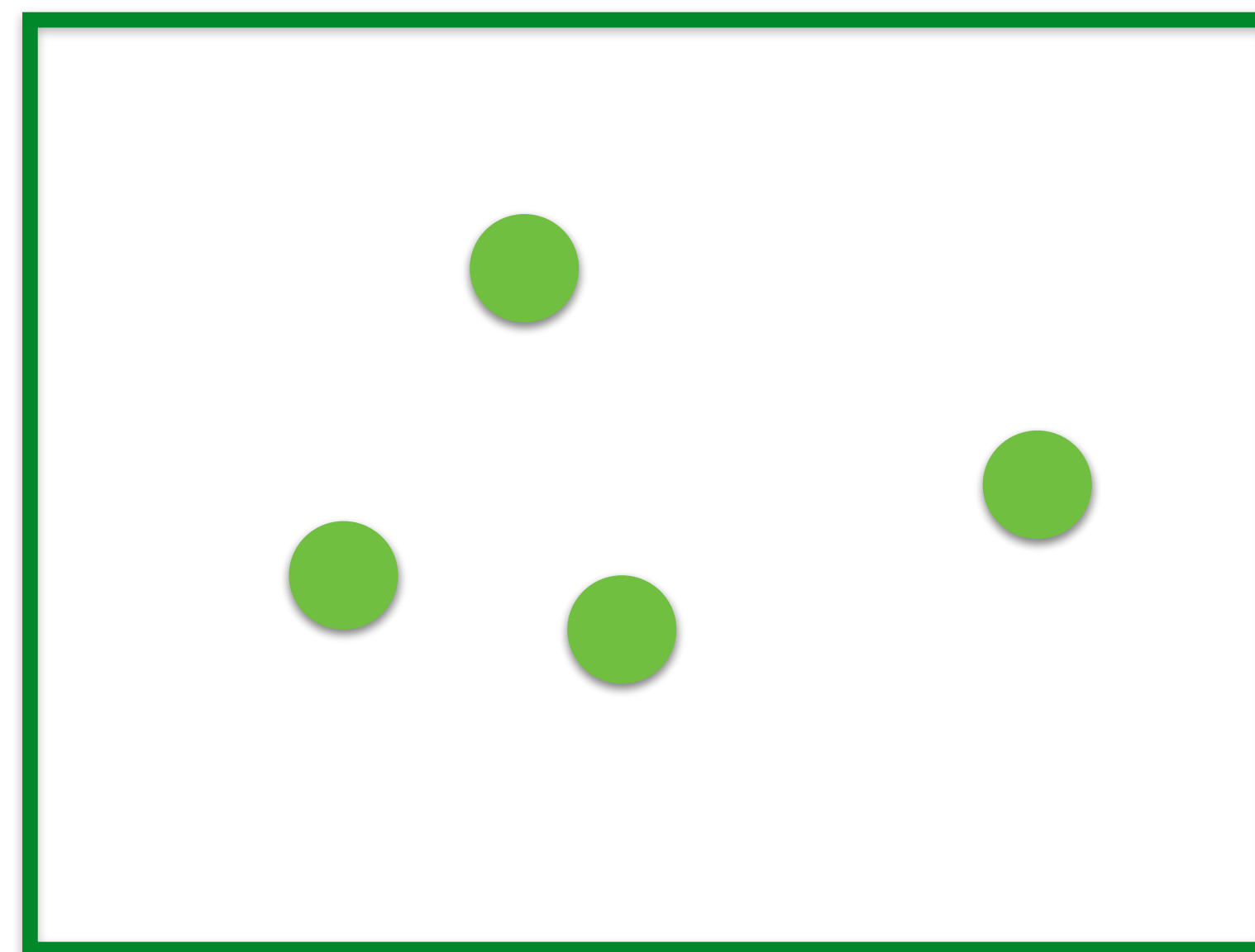
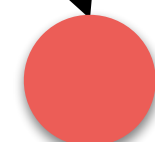
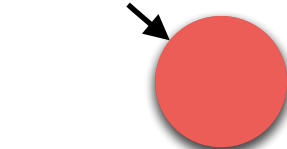
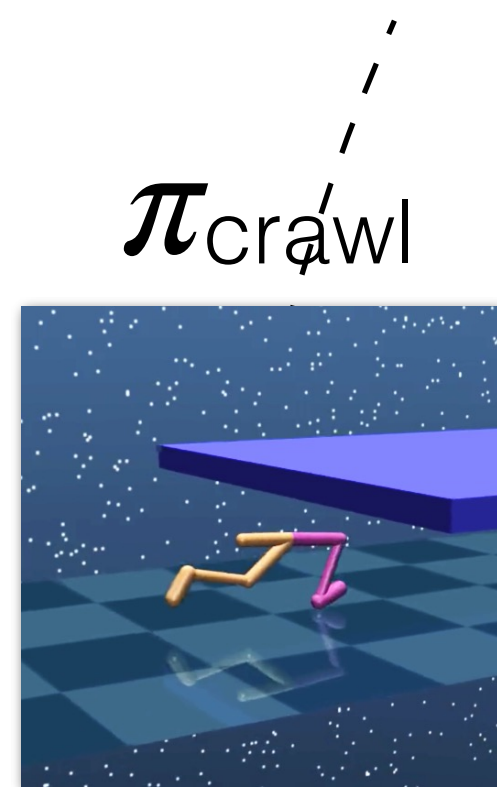
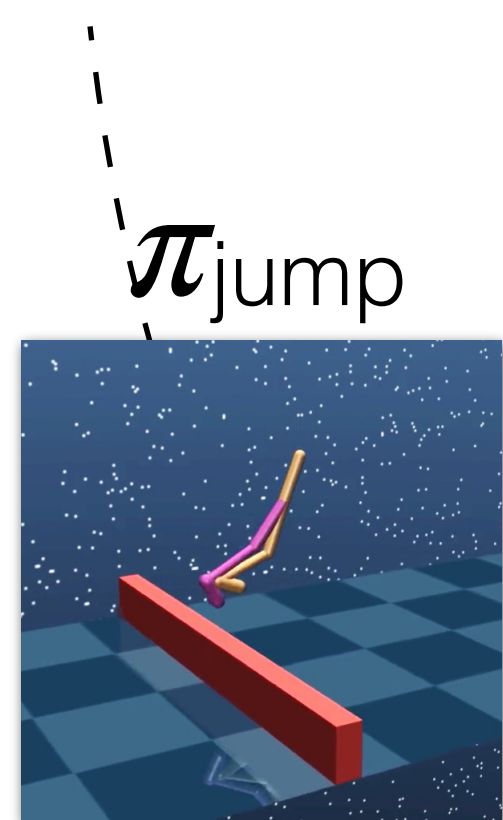






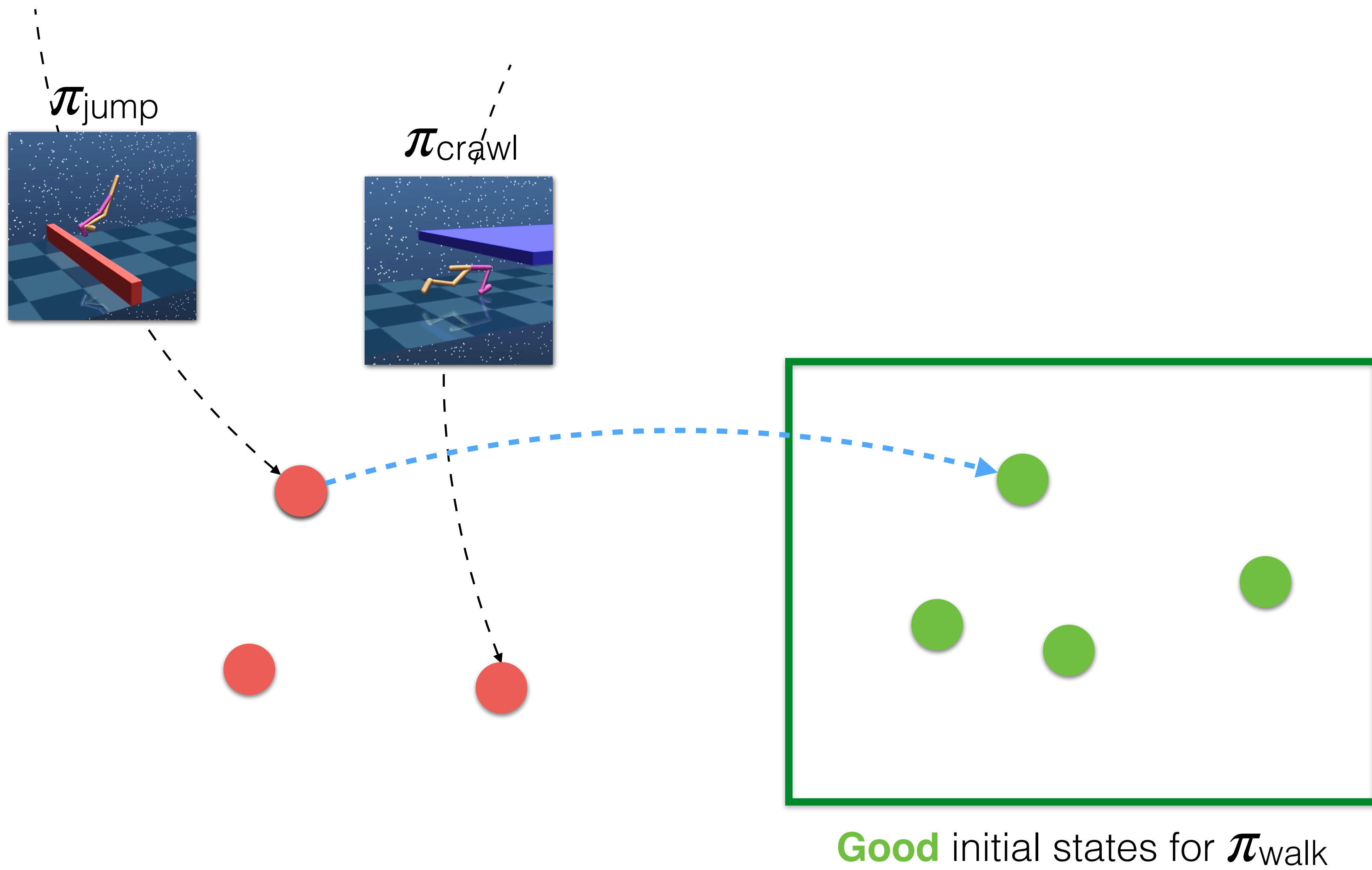


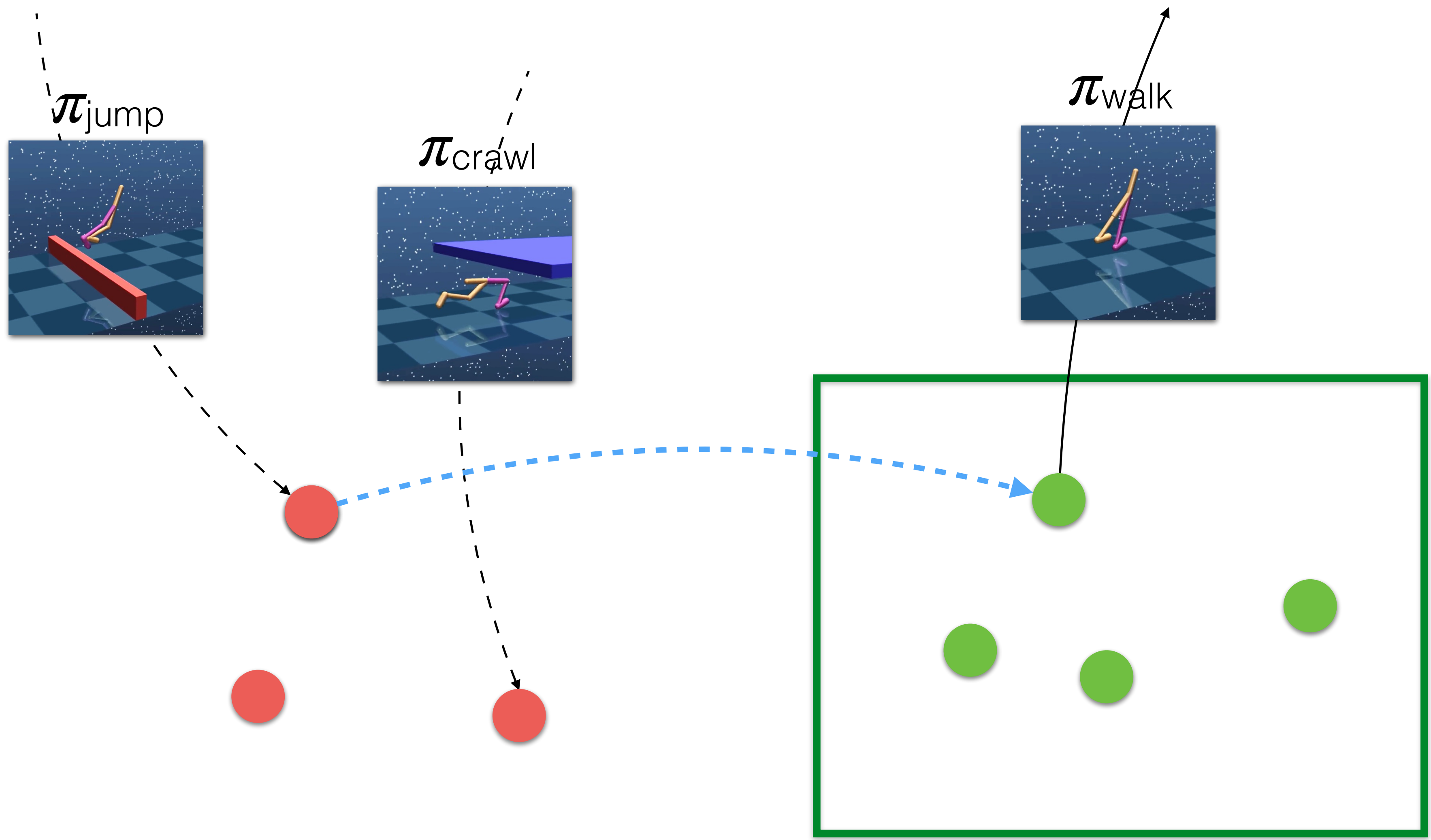




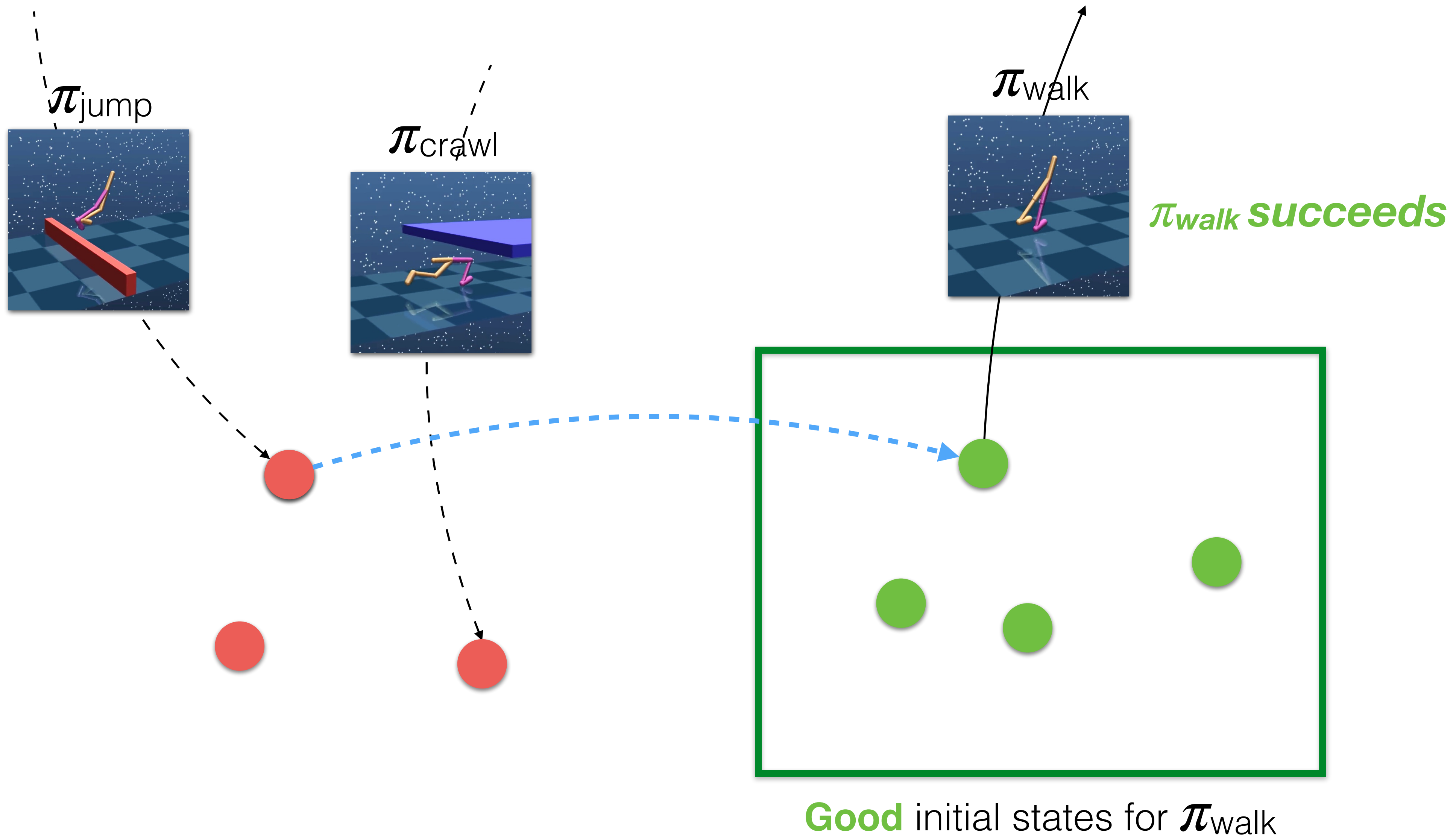
Good initial states for  $\pi_{\text{walk}}$

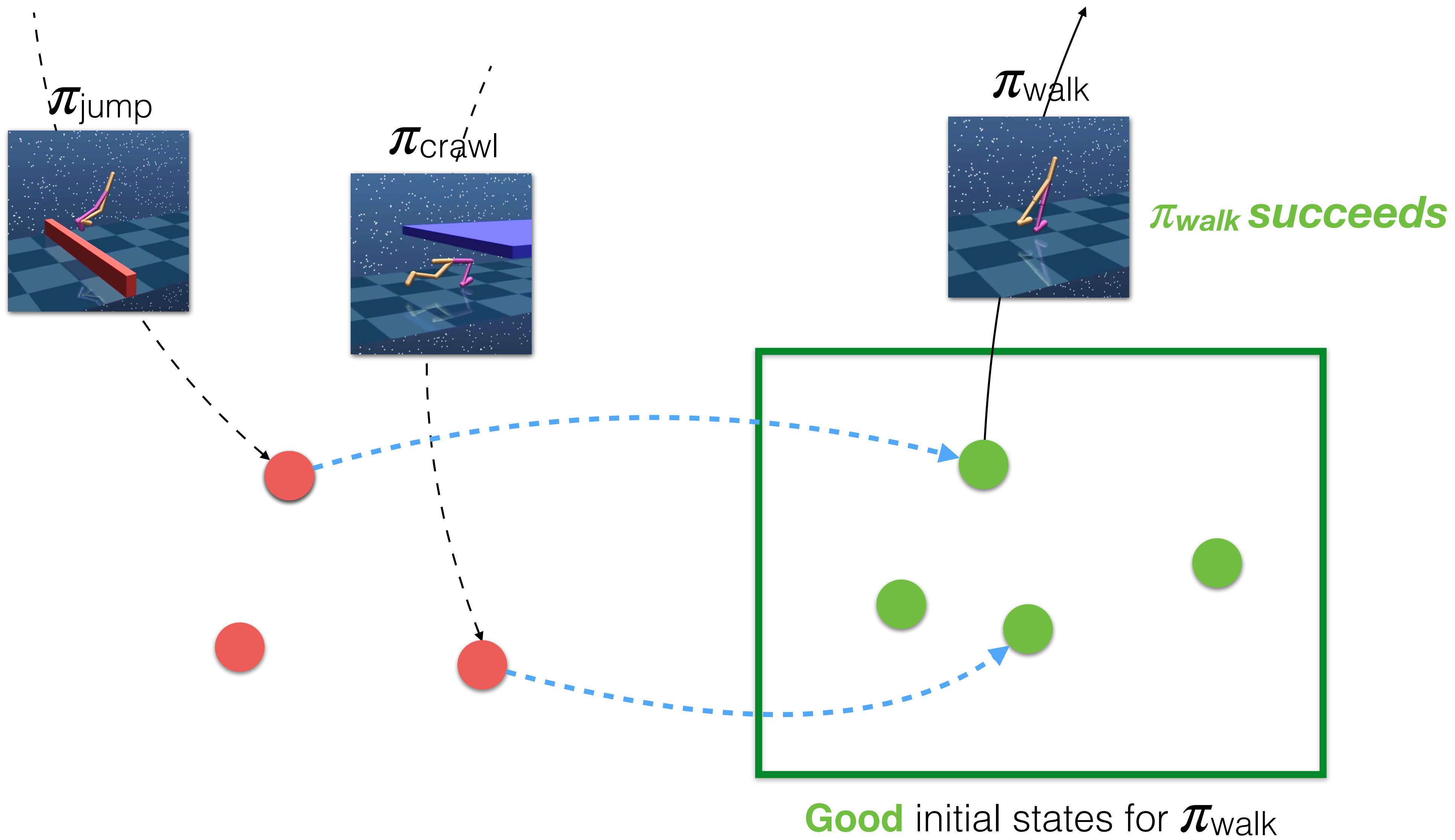




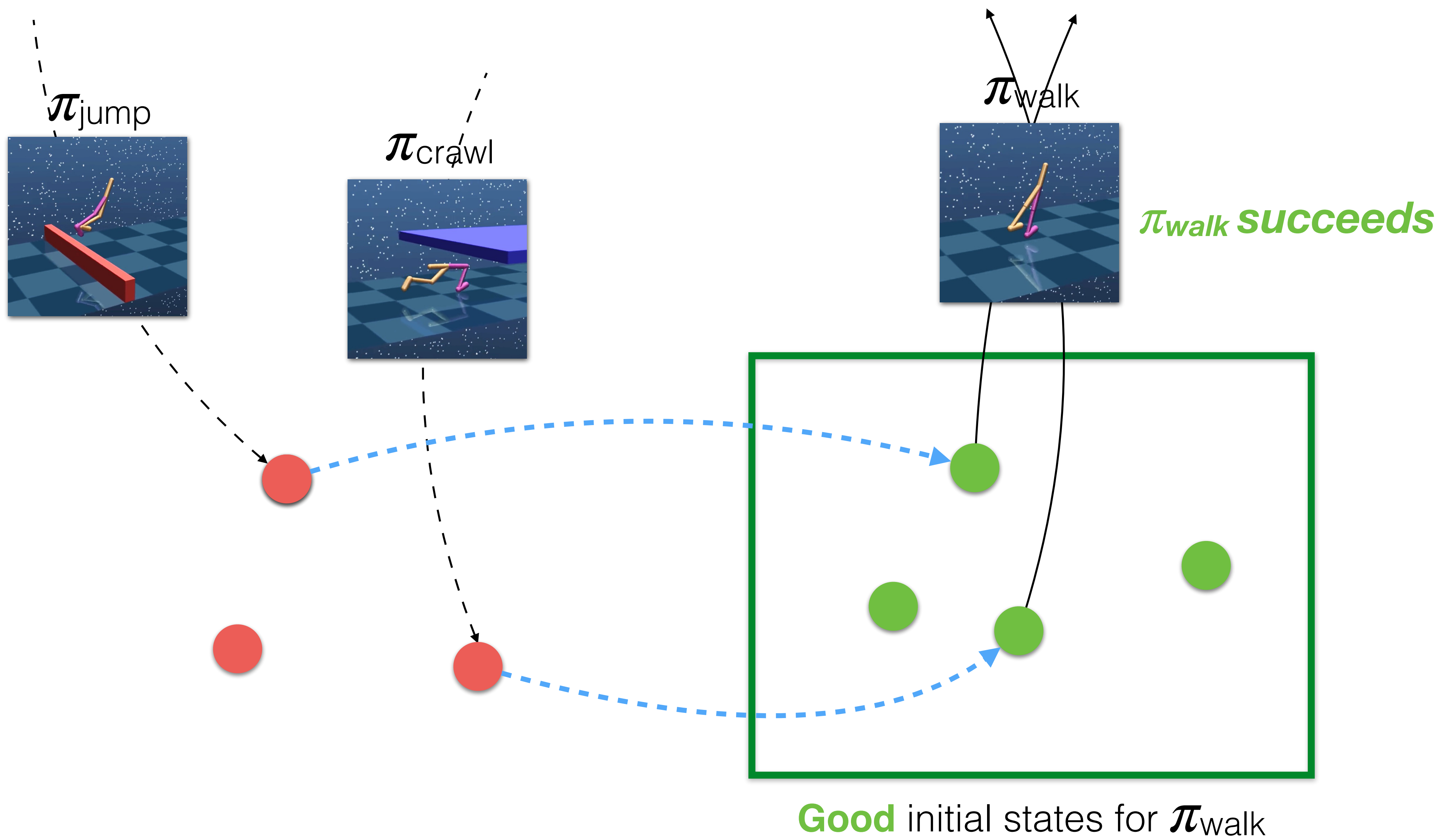


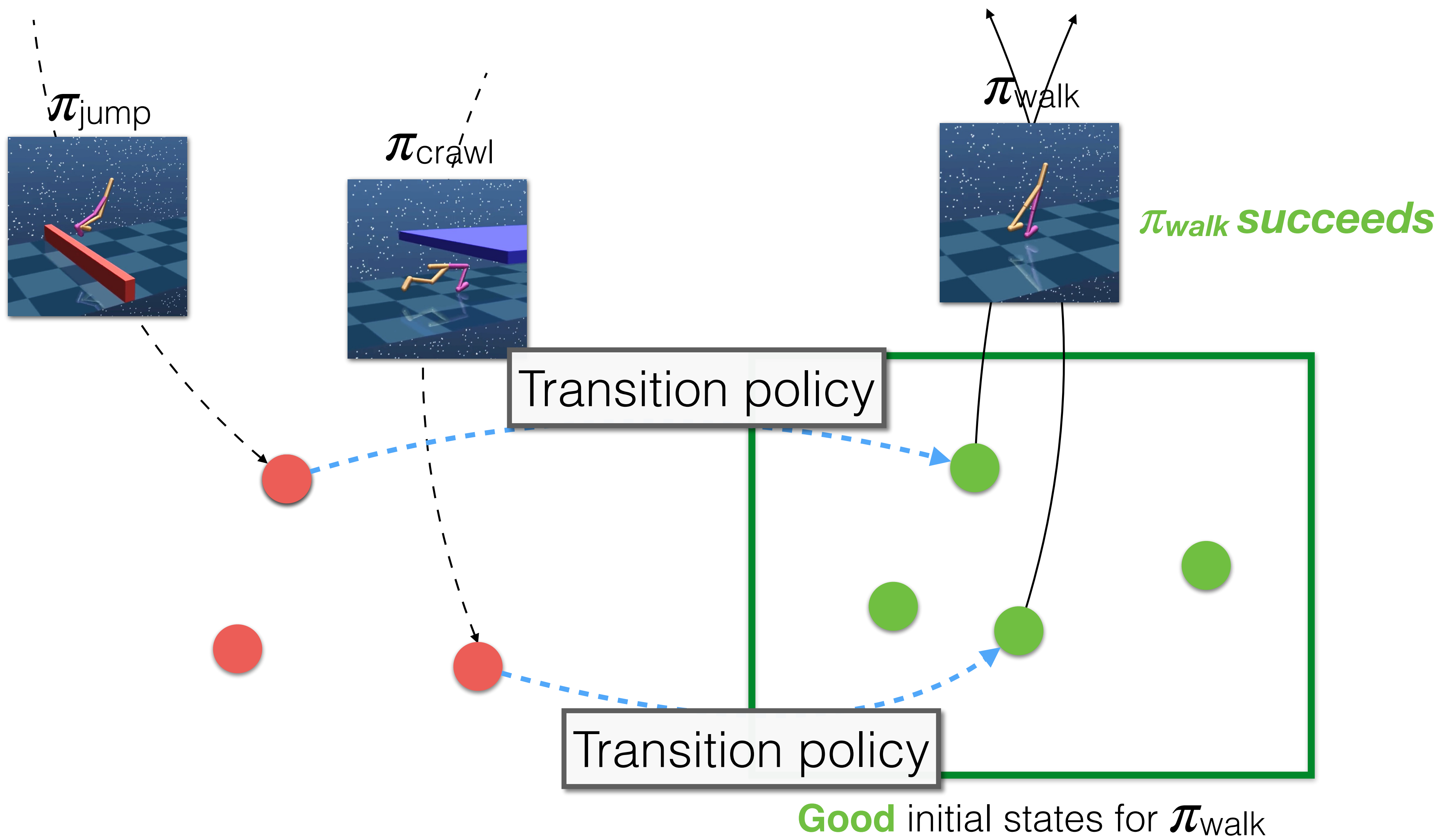
Good initial states for  $\pi_{\text{walk}}$

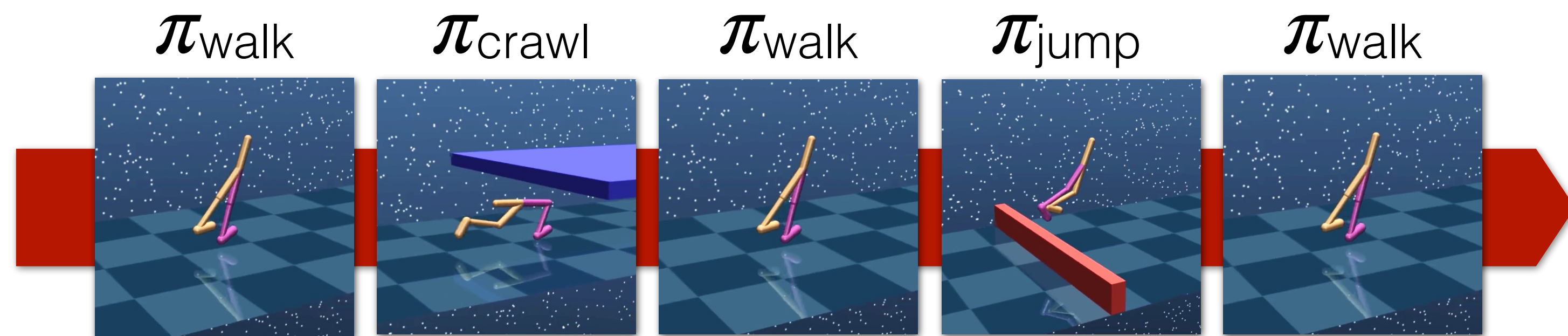




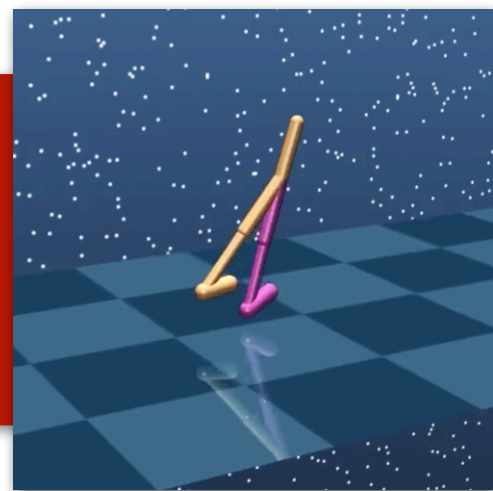




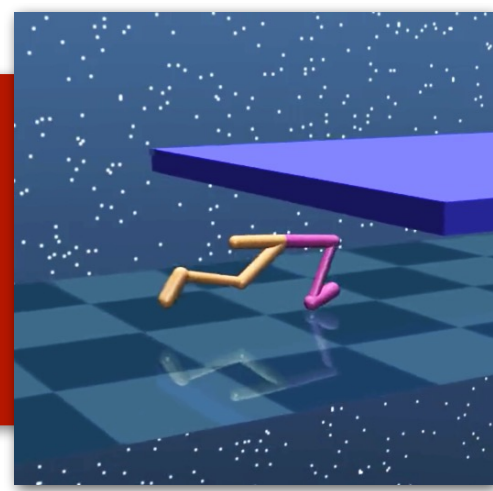




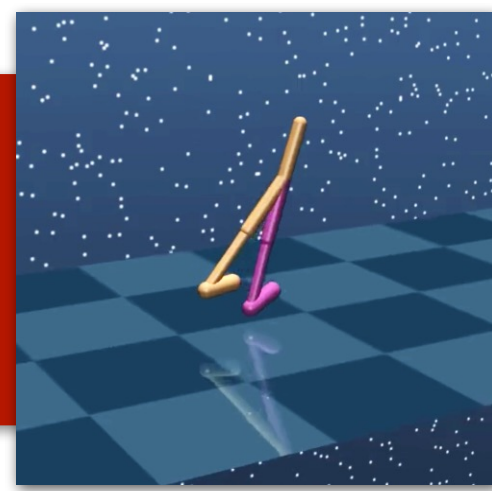
$\pi_{\text{walk}}$



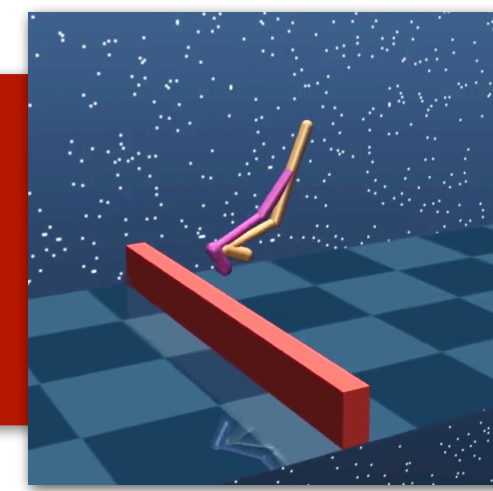
$\pi_{\text{crawl}}$



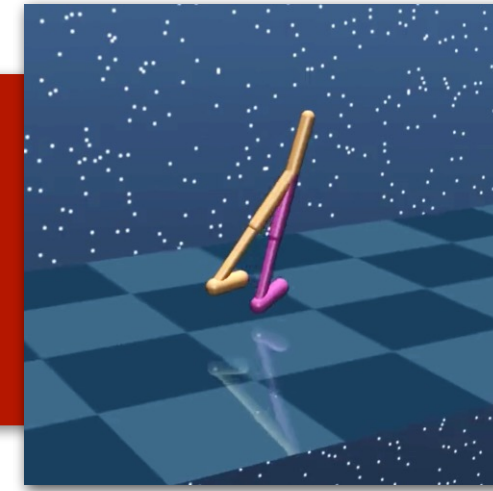
$\pi_{\text{walk}}$



$\pi_{\text{jump}}$

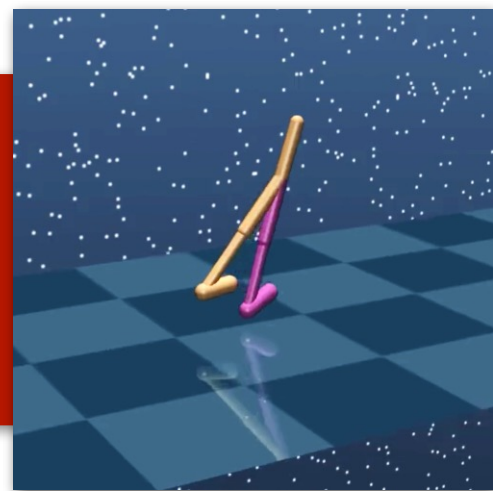


$\pi_{\text{walk}}$



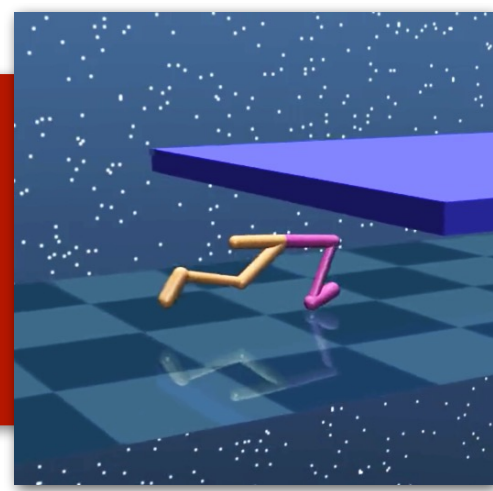


$\pi_{\text{walk}}$



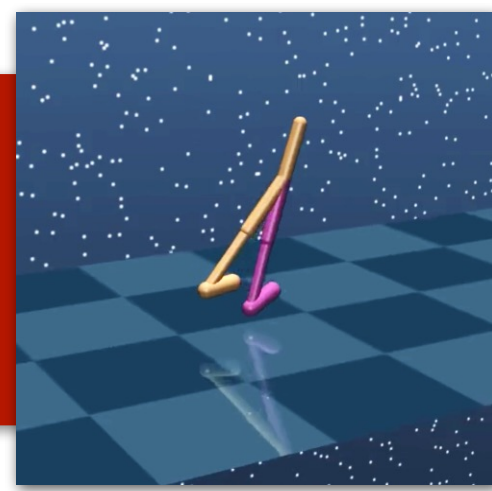
$\pi_{\text{tr}}$

$\pi_{\text{crawl}}$



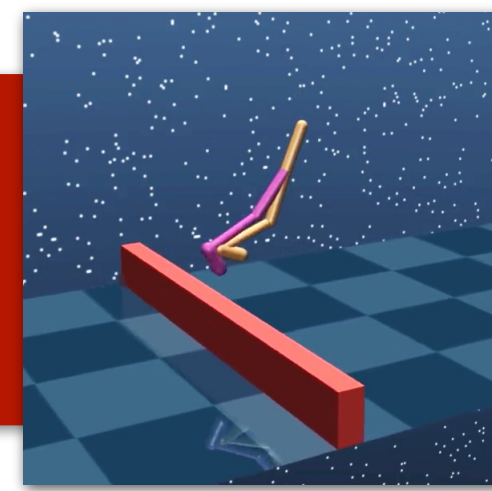
$\pi_{\text{tr}}$

$\pi_{\text{walk}}$



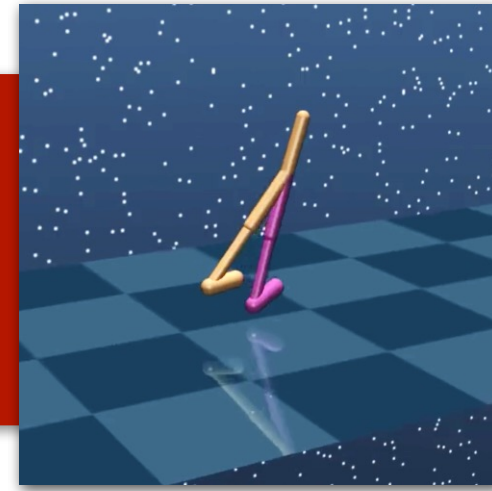
$\pi_{\text{tr}}$

$\pi_{\text{jump}}$



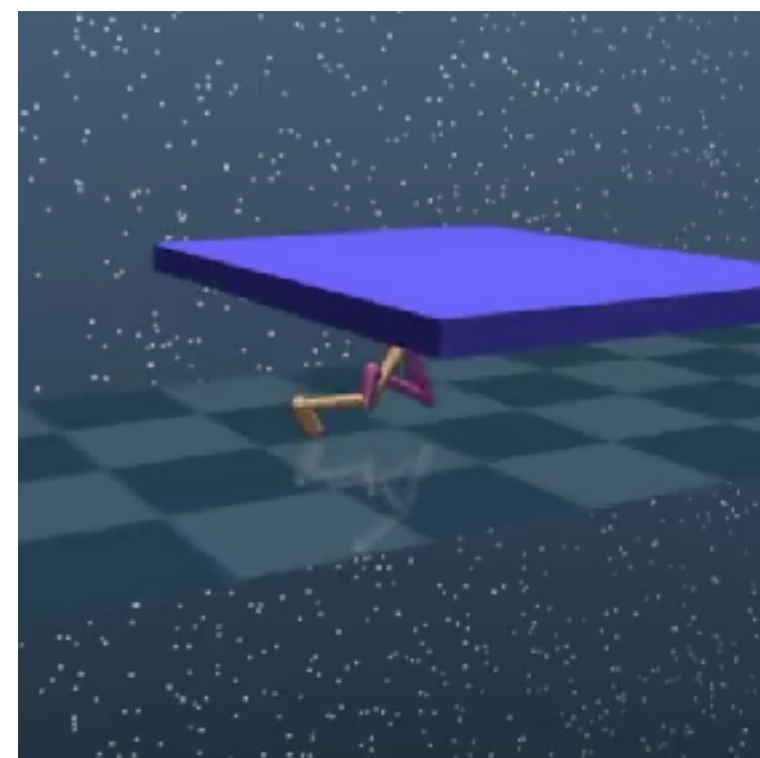
$\pi_{\text{tr}}$

$\pi_{\text{walk}}$

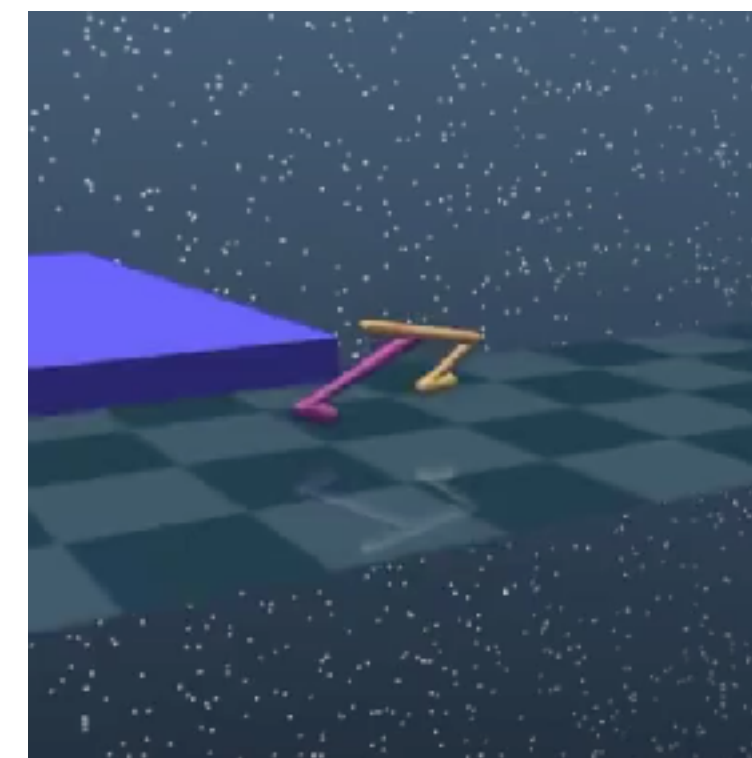


## Obstacle course

$\pi_{\text{crawl}}$

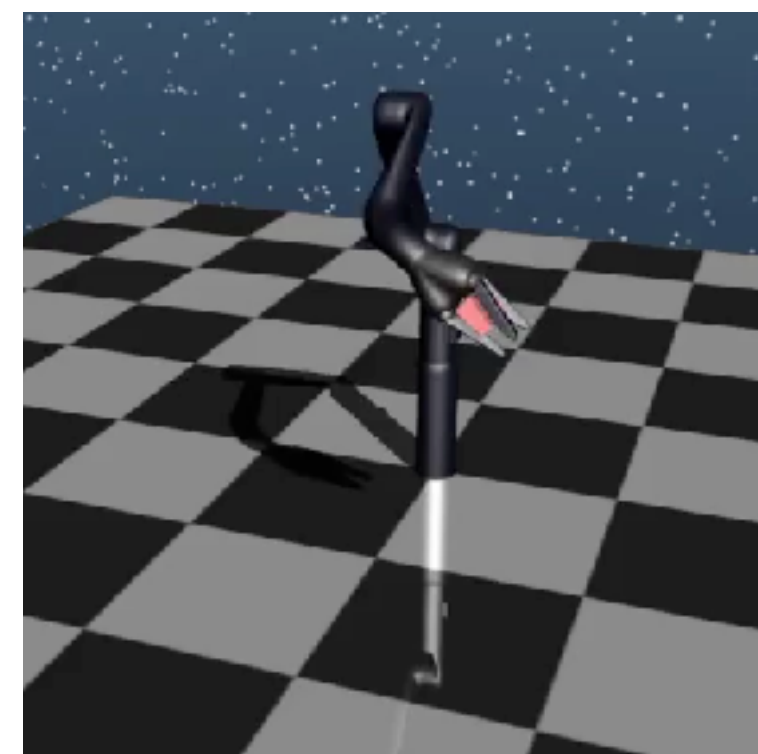
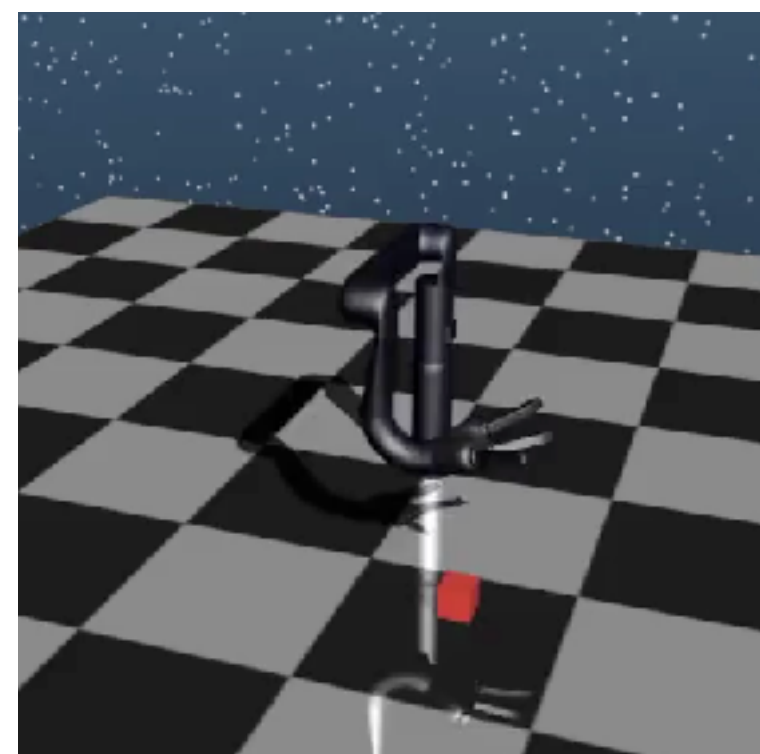


$\pi_{\text{walk}}$

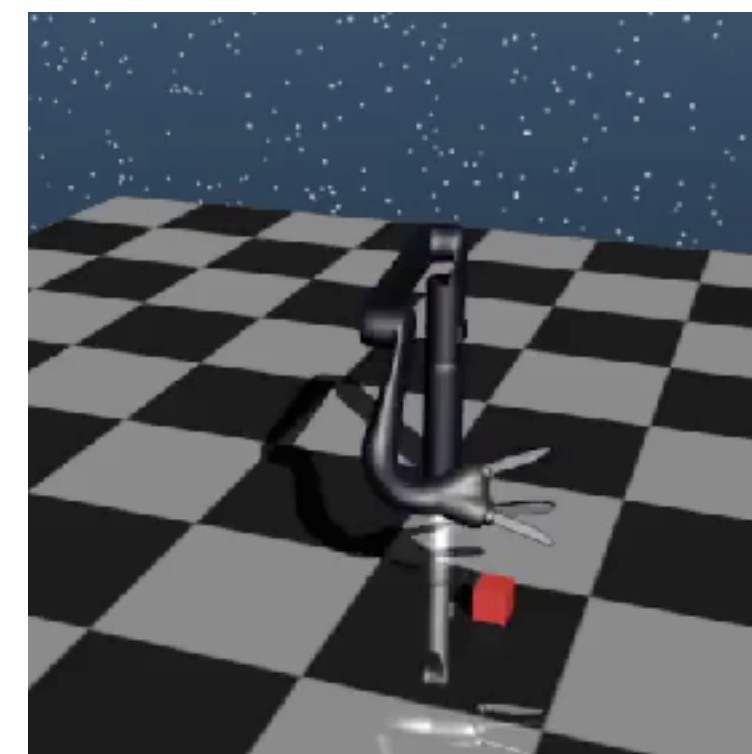


## Repetitive pick

$\pi_{\text{pick}}$

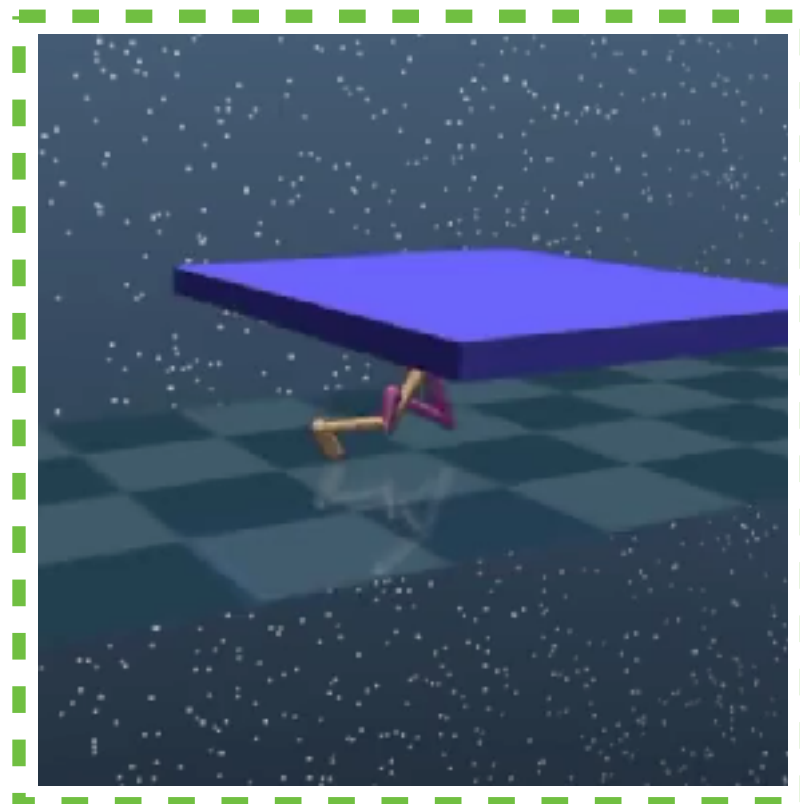


$\pi_{\text{pick}}$

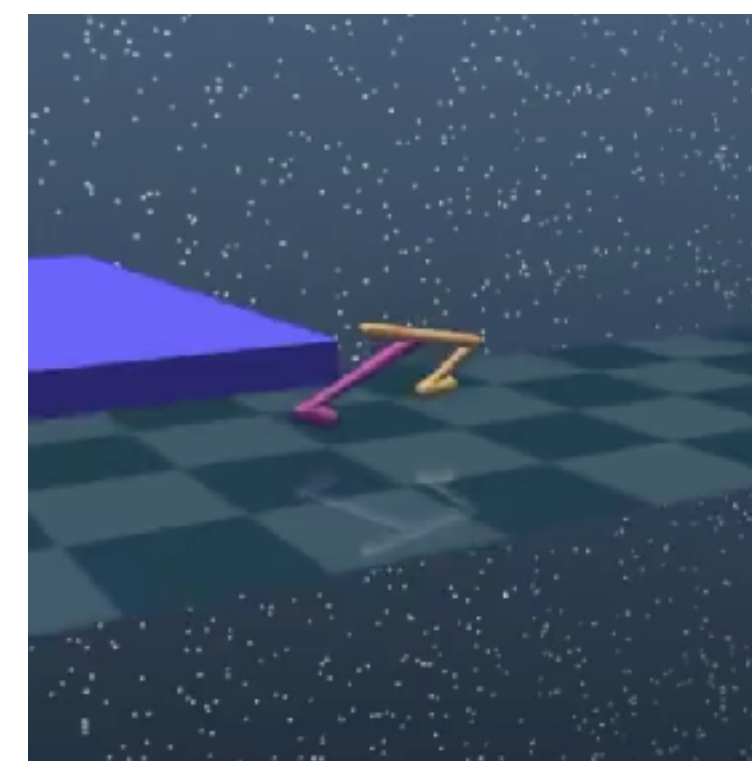


## Obstacle course

$\pi_{\text{crawl}}$

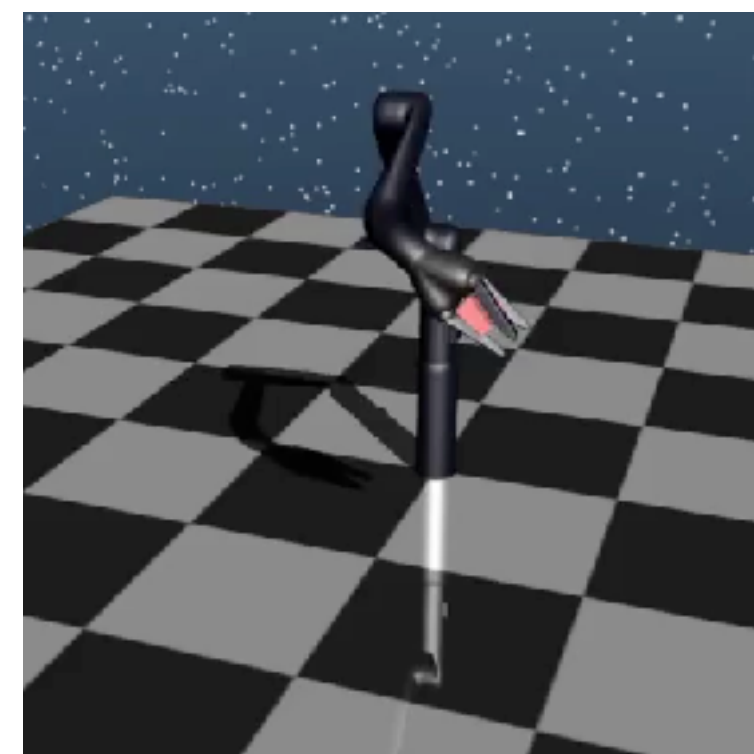
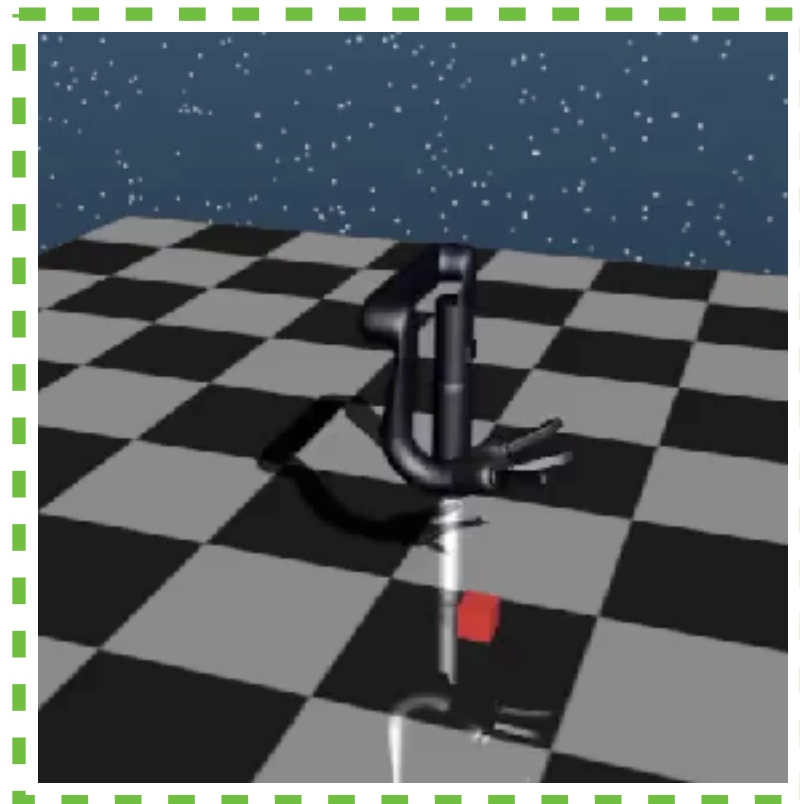


$\pi_{\text{walk}}$

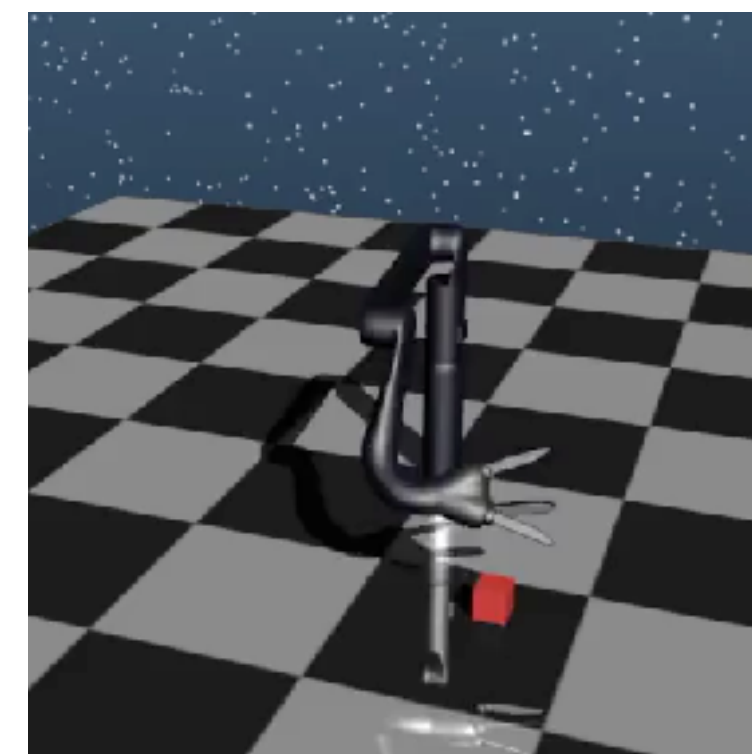


## Repetitive pick

$\pi_{\text{pick}}$



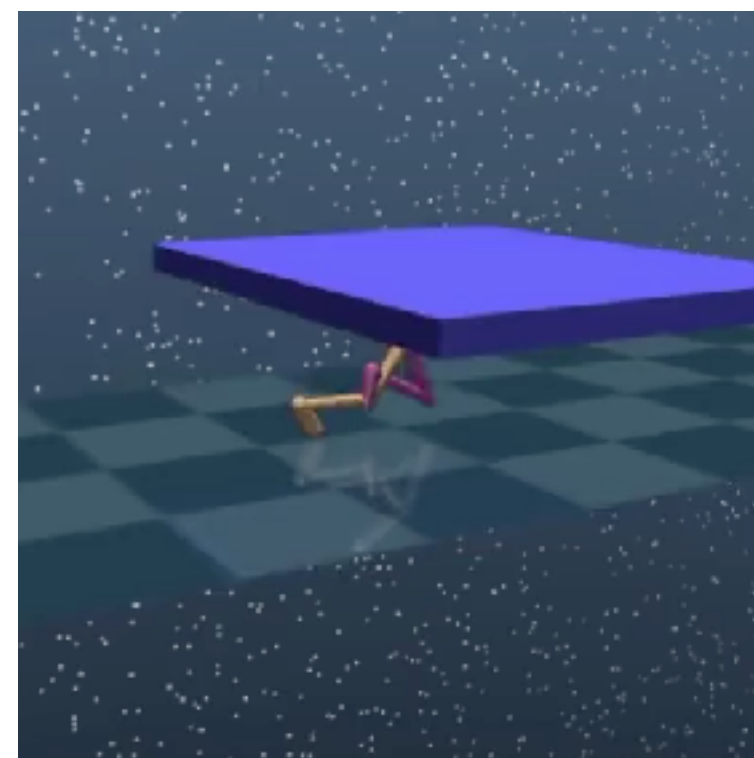
$\pi_{\text{pick}}$



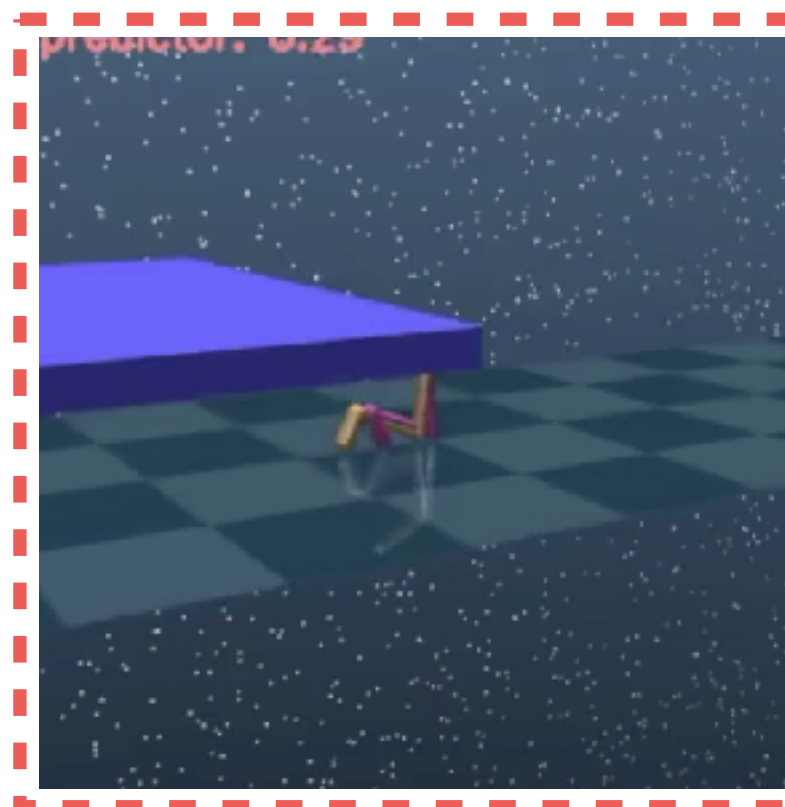


## Obstacle course

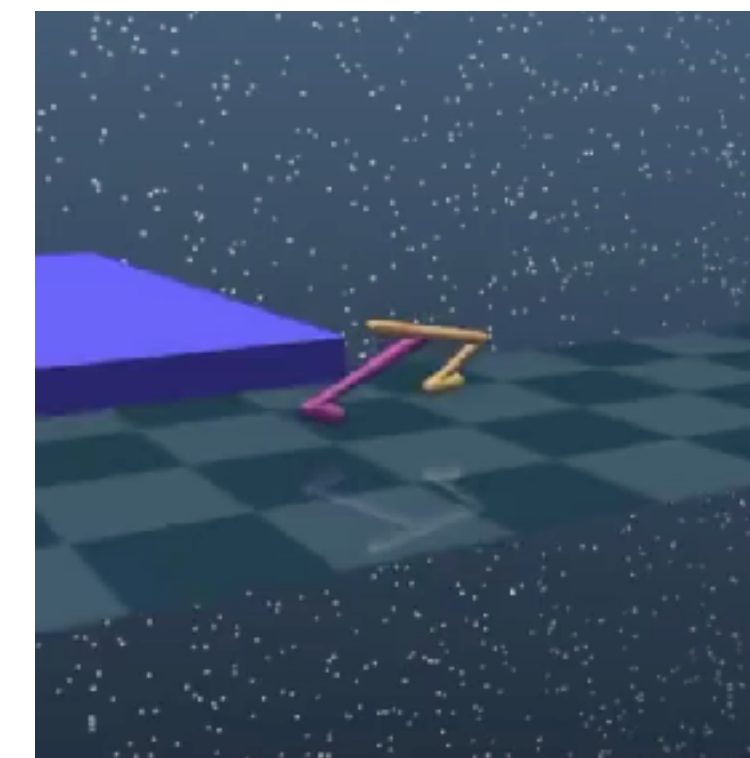
$\pi_{\text{crawl}}$



***Transition policy***

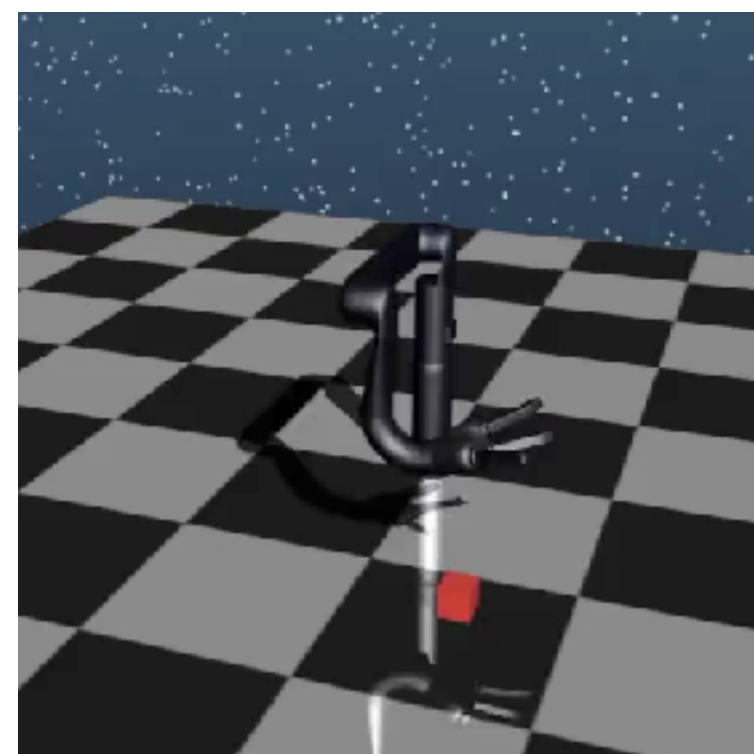


$\pi_{\text{walk}}$

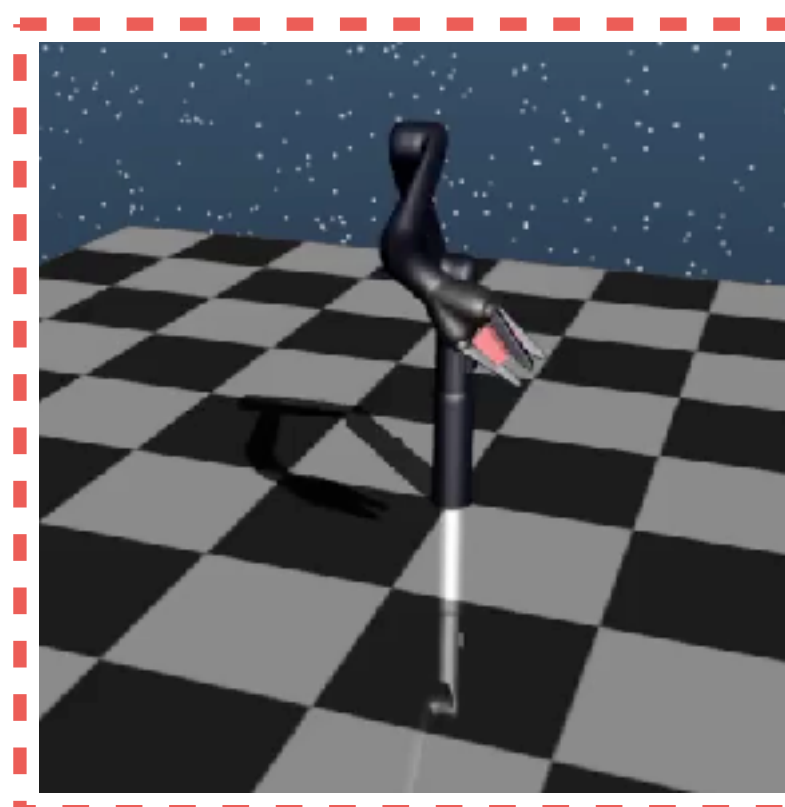


## Repetitive pick

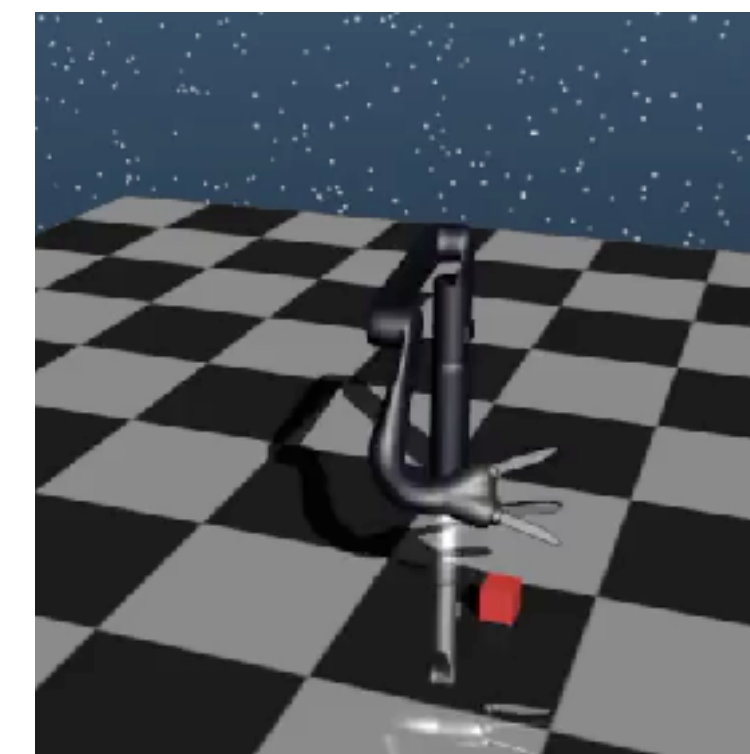
$\pi_{\text{pick}}$



***Transition policy***



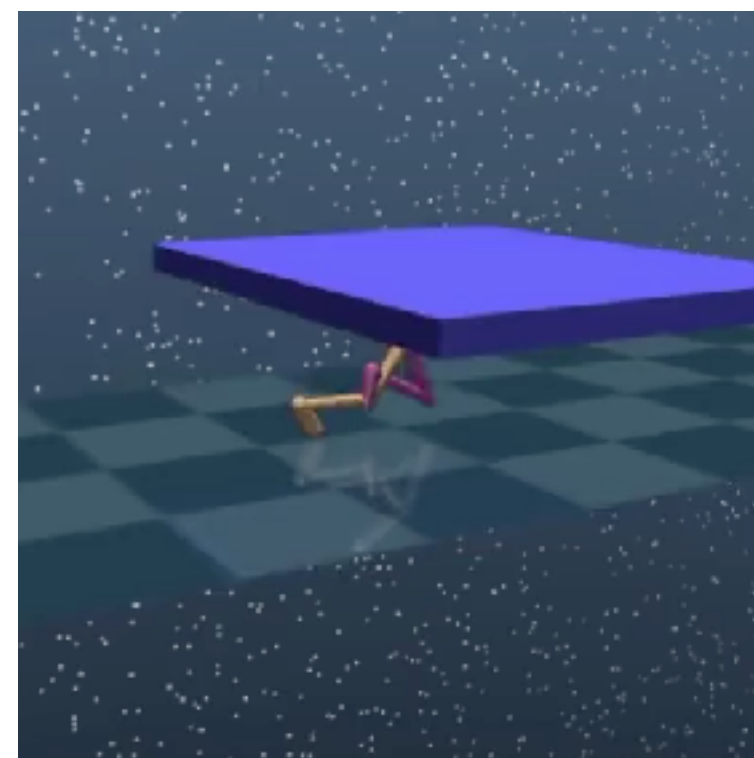
$\pi_{\text{pick}}$





## Obstacle course

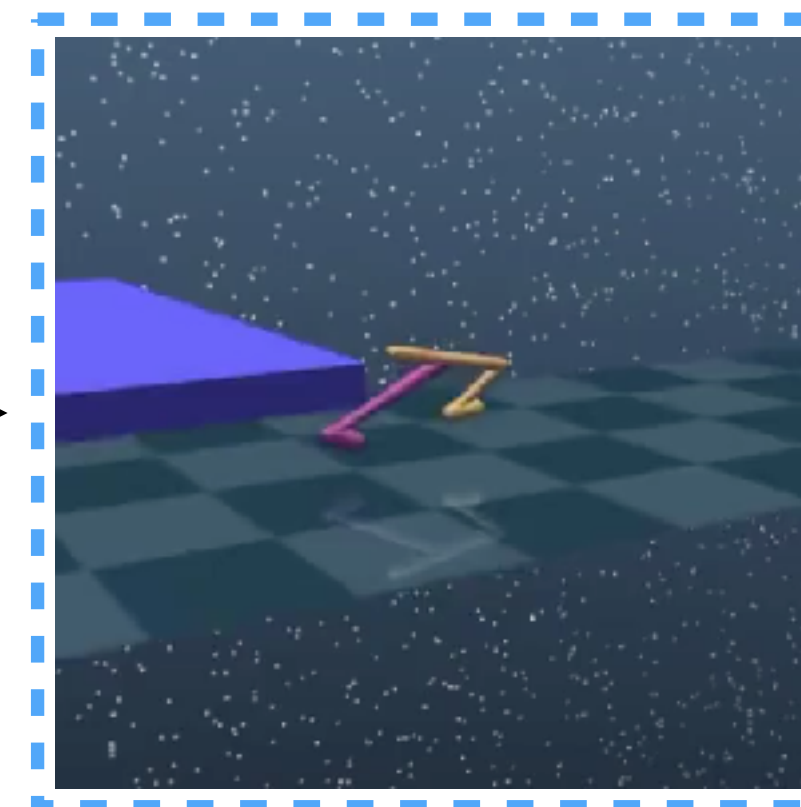
$\pi_{\text{crawl}}$



***Transition policy***

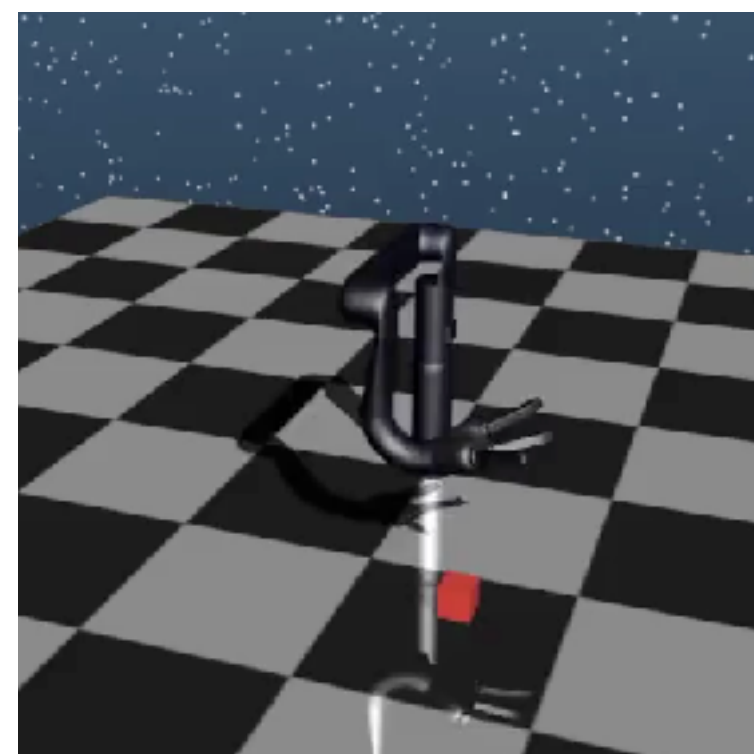


$\pi_{\text{walk}}$

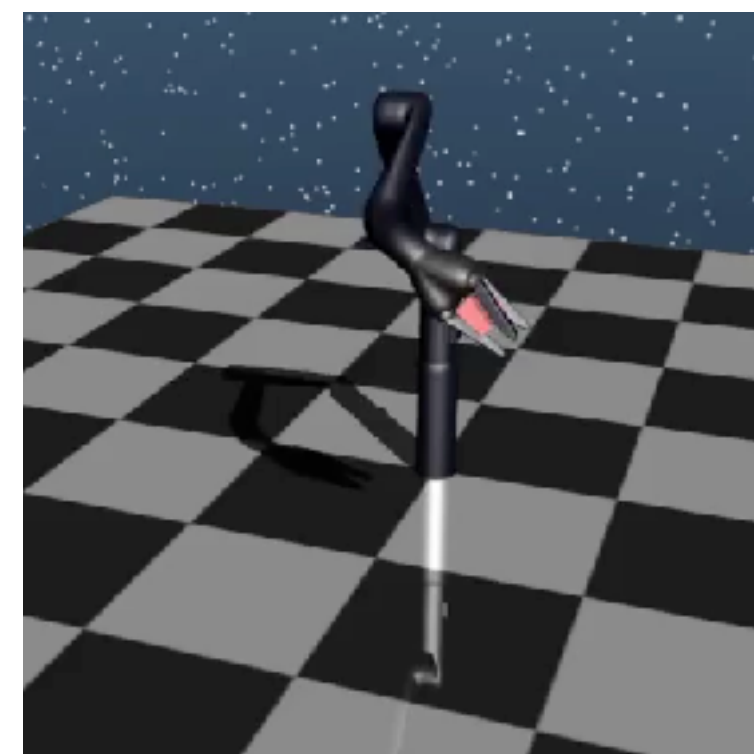


## Repetitive pick

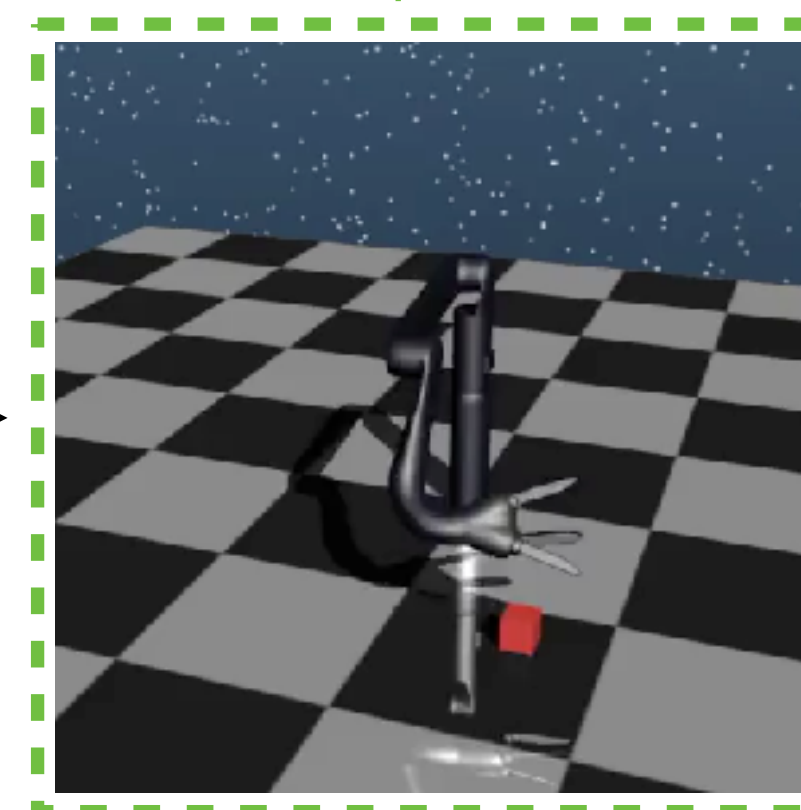
$\pi_{\text{pick}}$



***Transition policy***



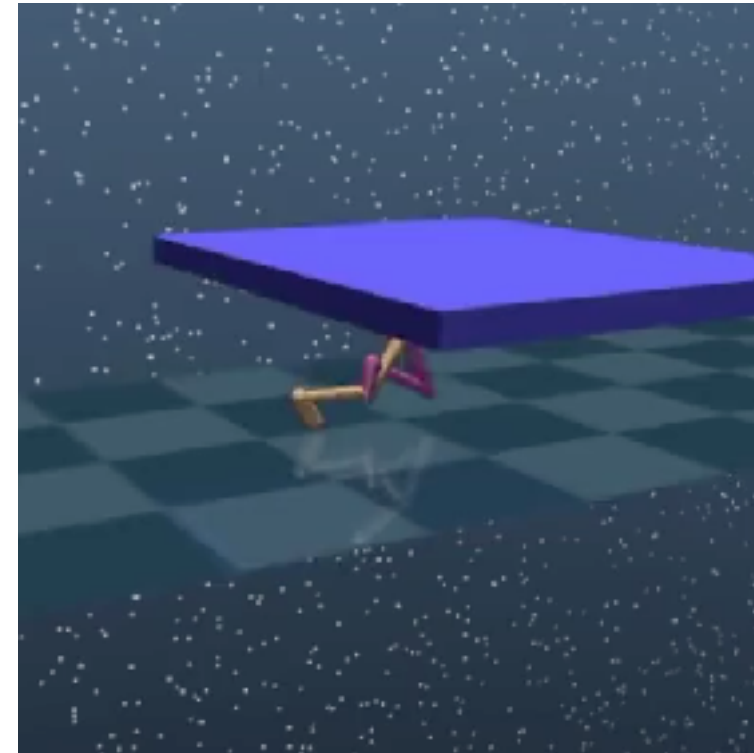
$\pi_{\text{pick}}$



## *Smoothly connect skills*

Obstacle course

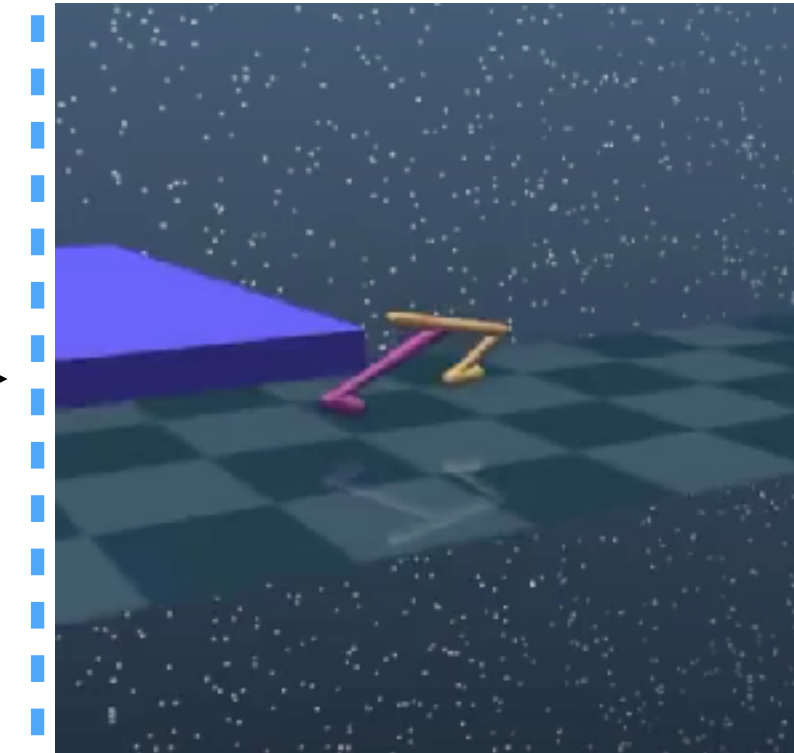
$\pi_{\text{crawl}}$



***Transition policy***

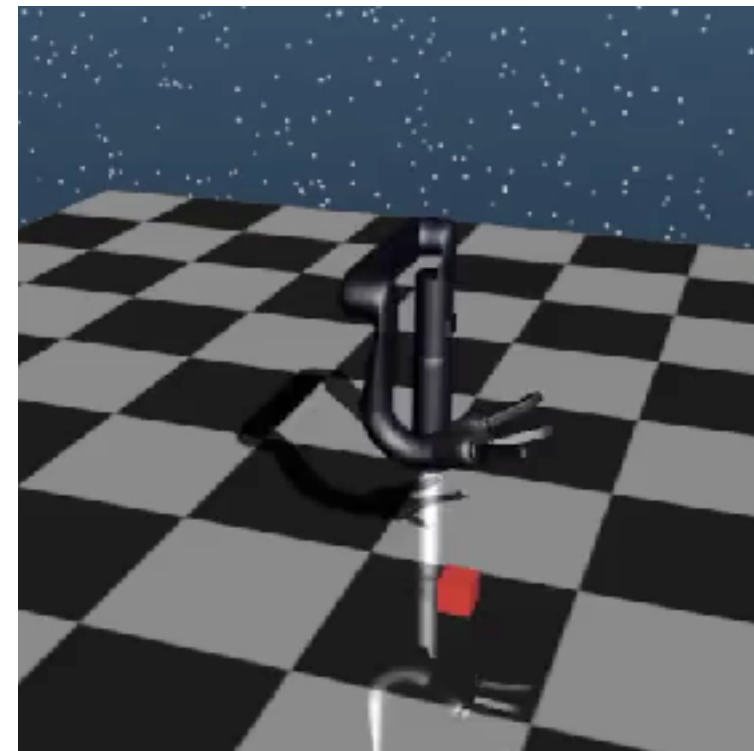


$\pi_{\text{walk}}$

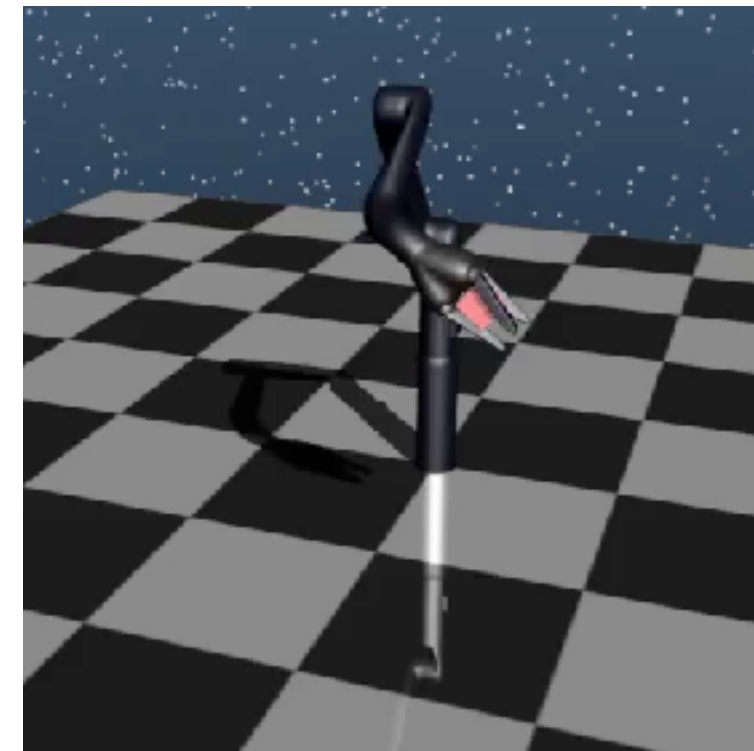


Repetitive pick

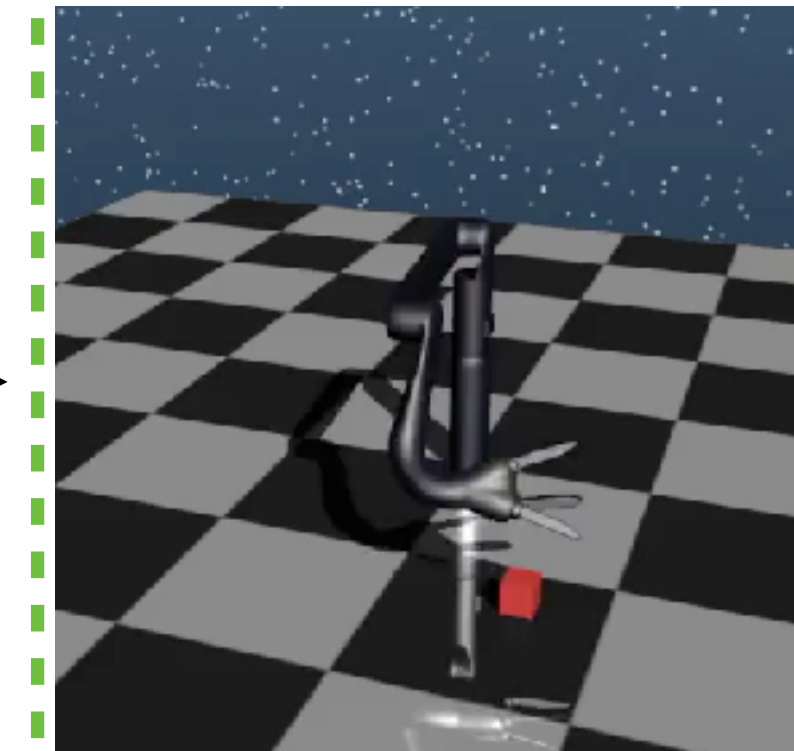
$\pi_{\text{pick}}$



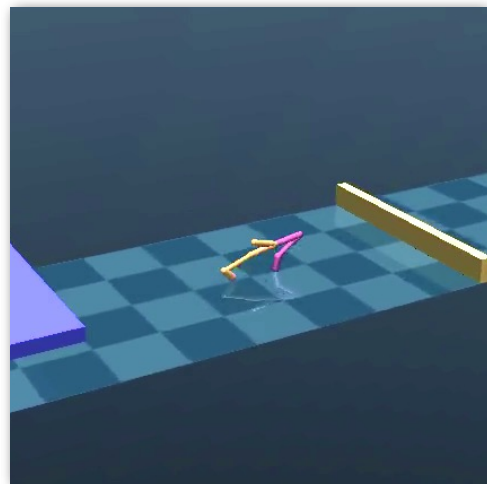
***Transition policy***



$\pi_{\text{pick}}$



# Model



Observation



Meta policy

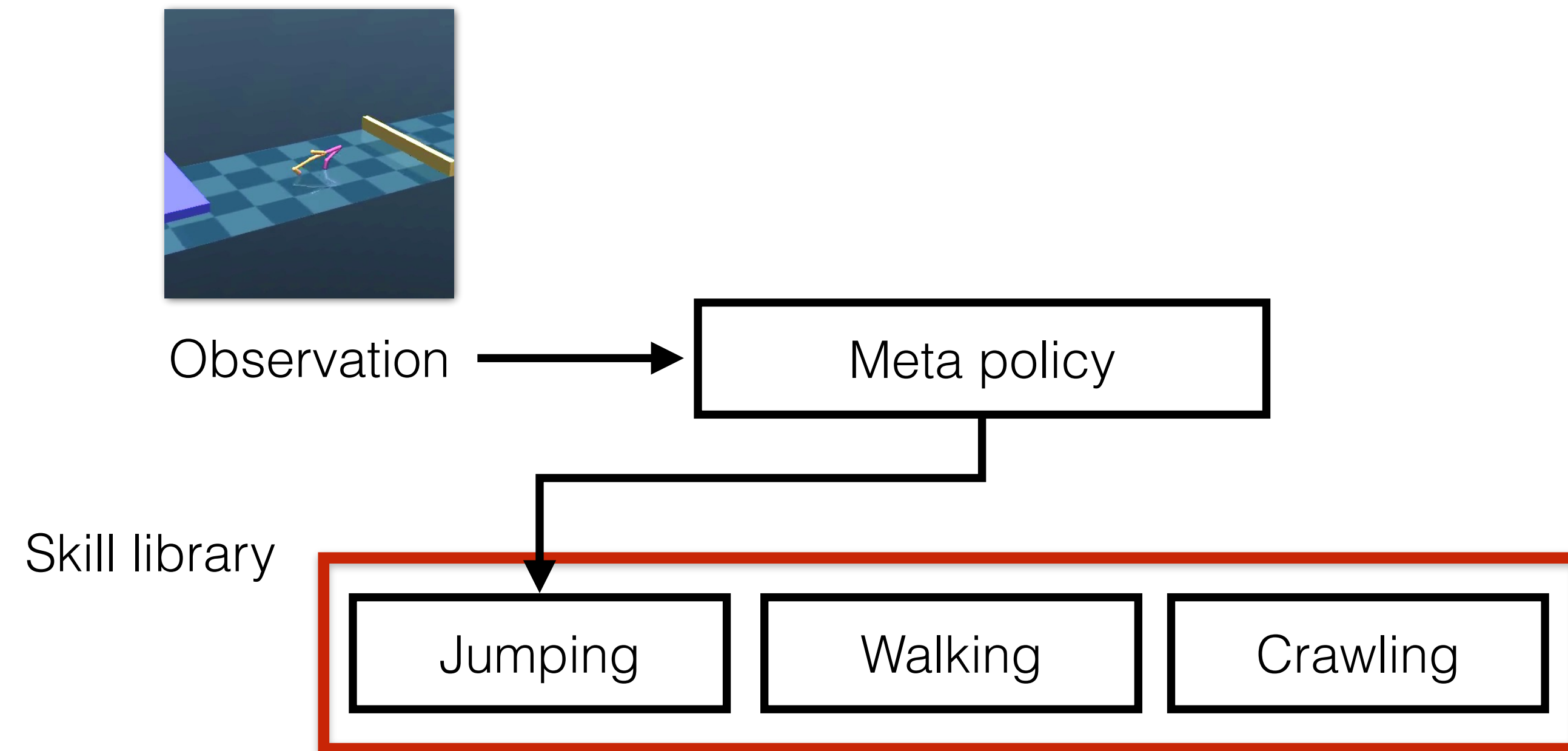
Skill library

Jumping

Walking

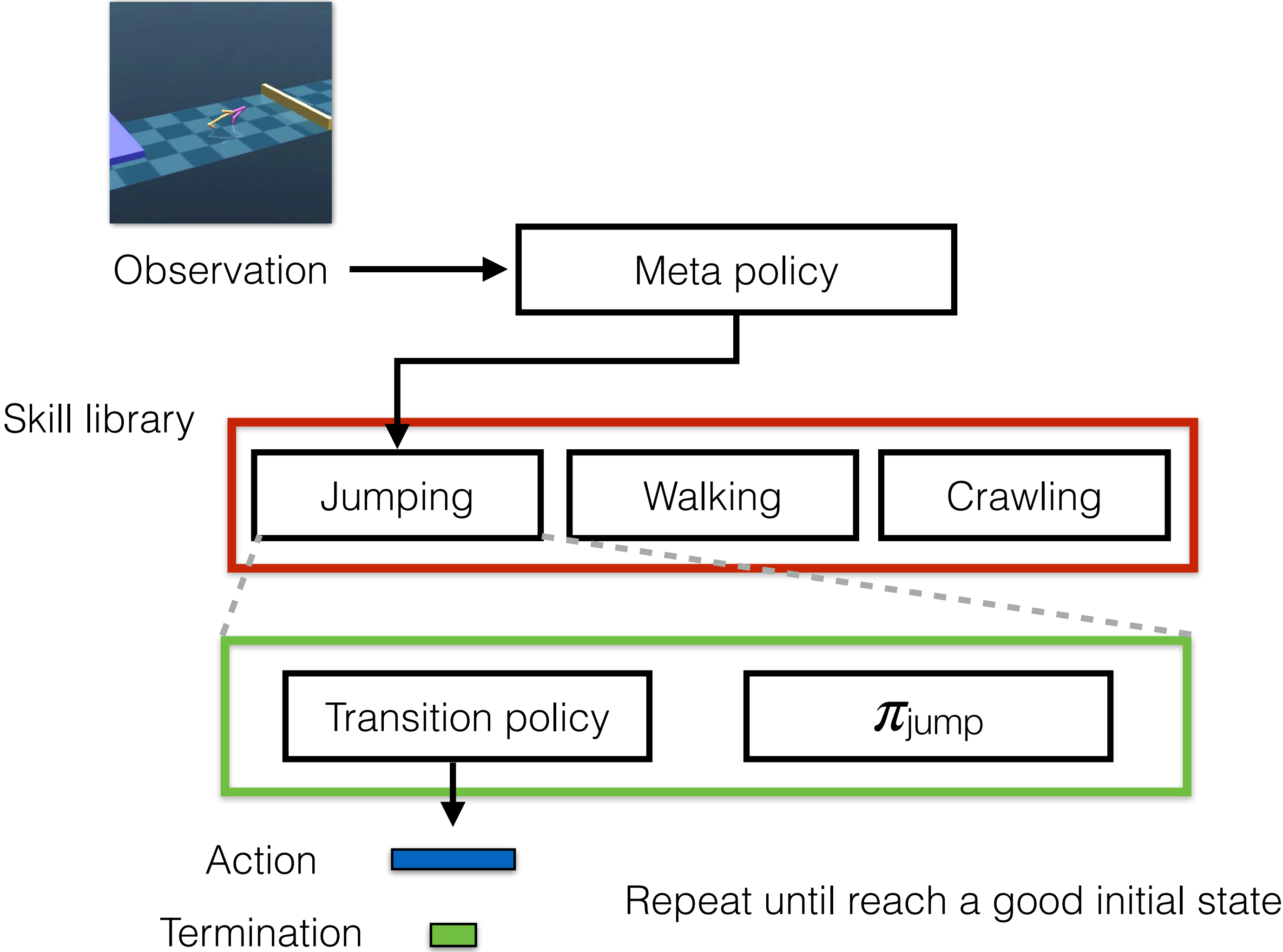
Crawling

# Model

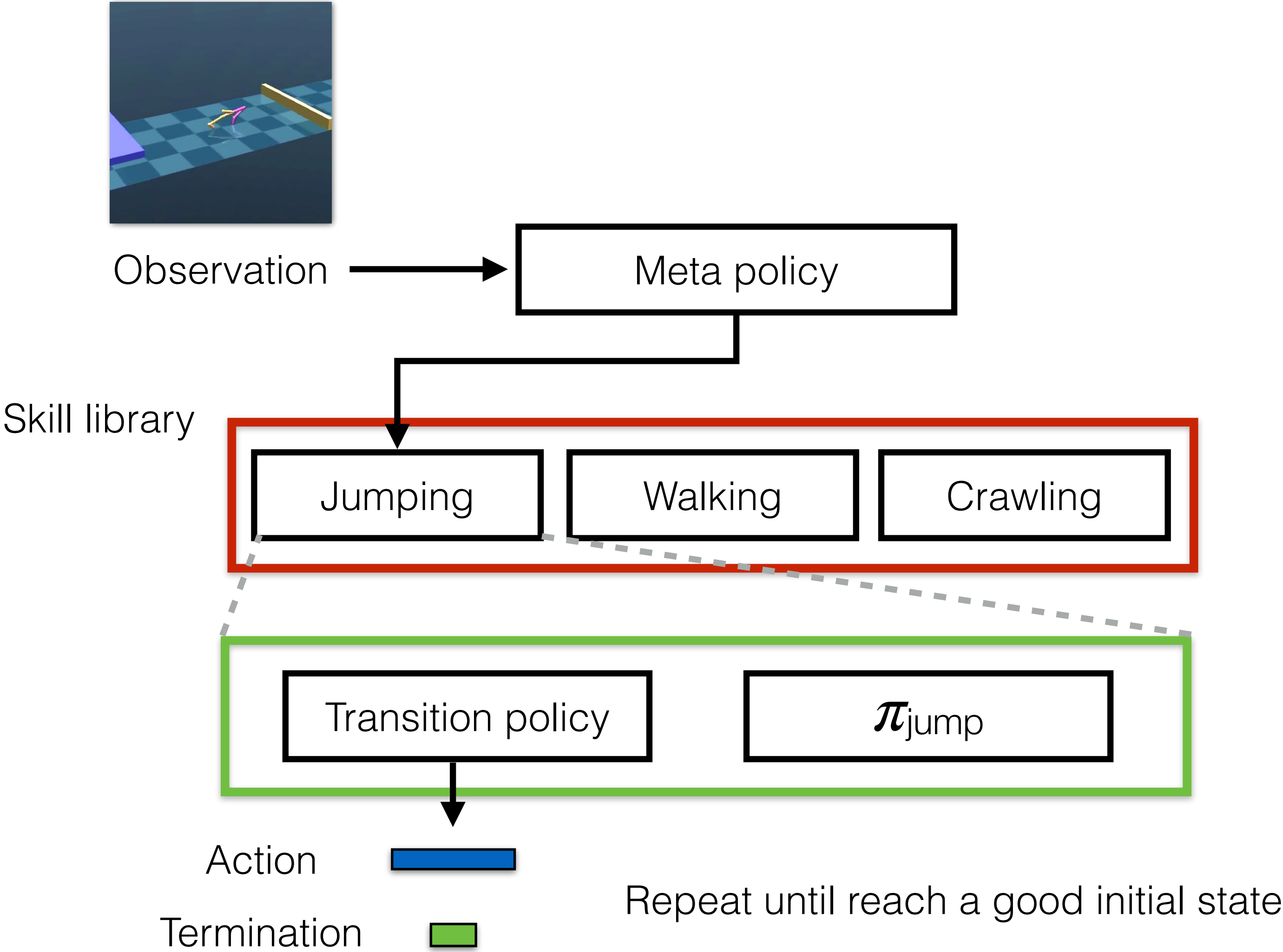




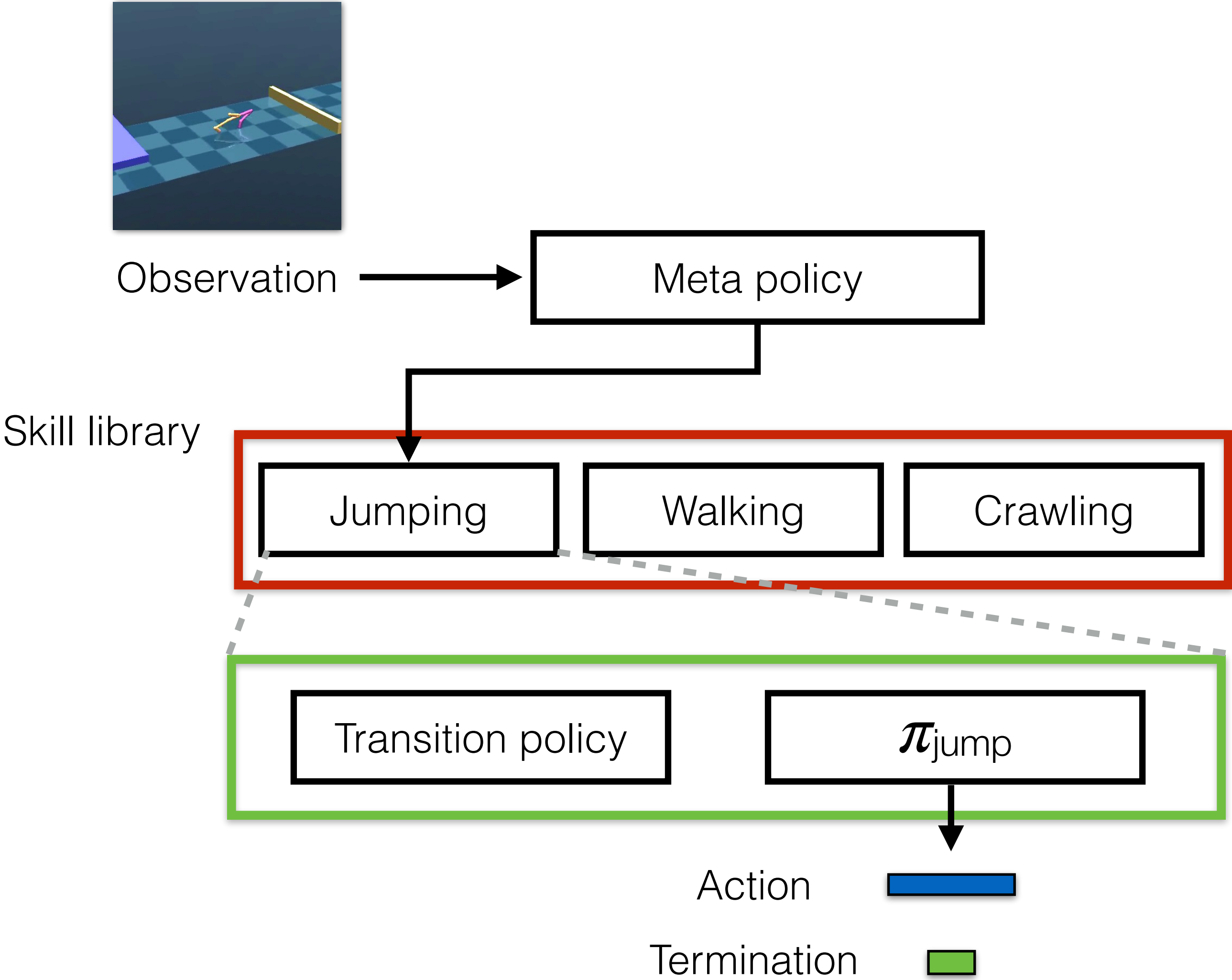
# Model



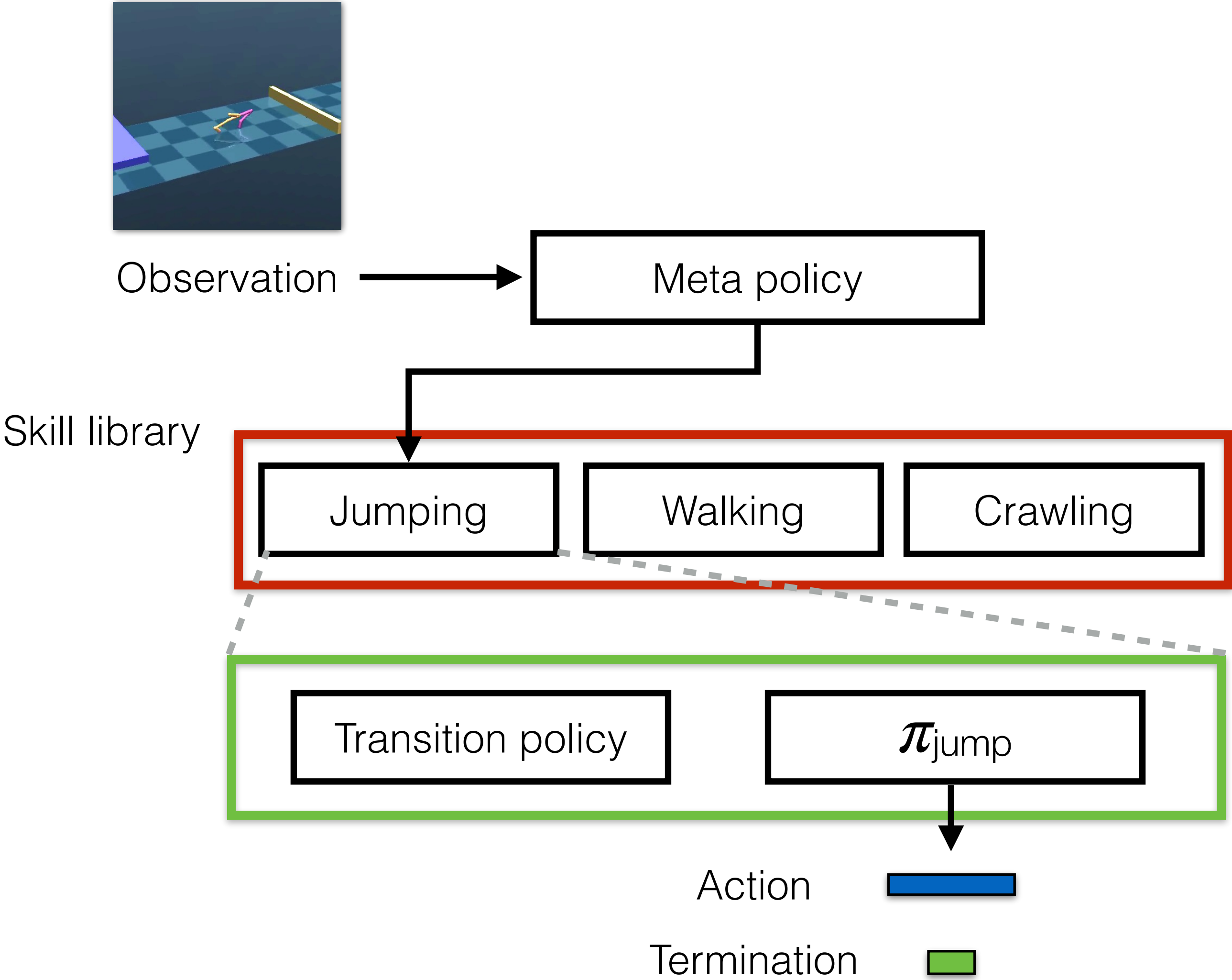
# Model



# Model

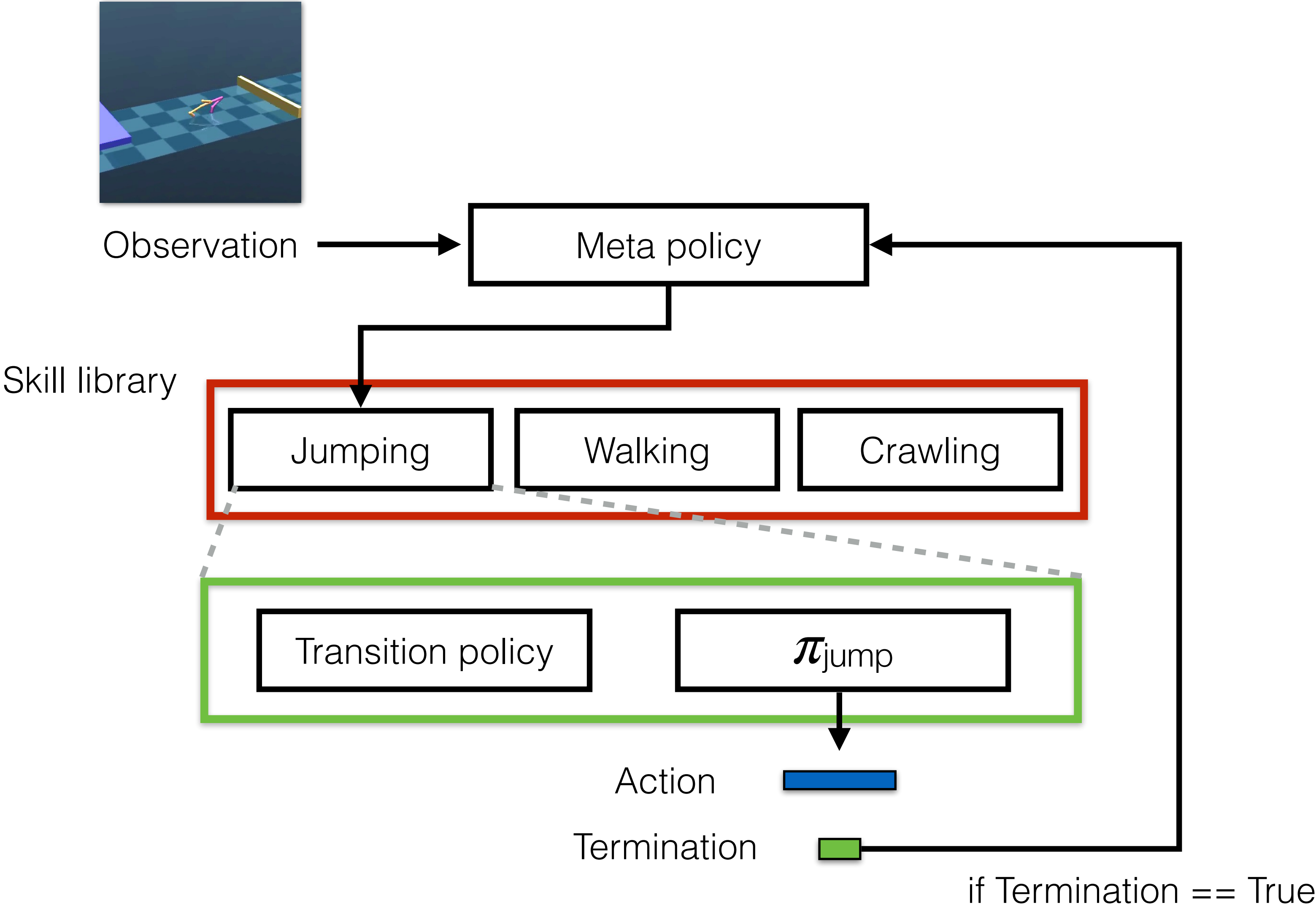


# Model

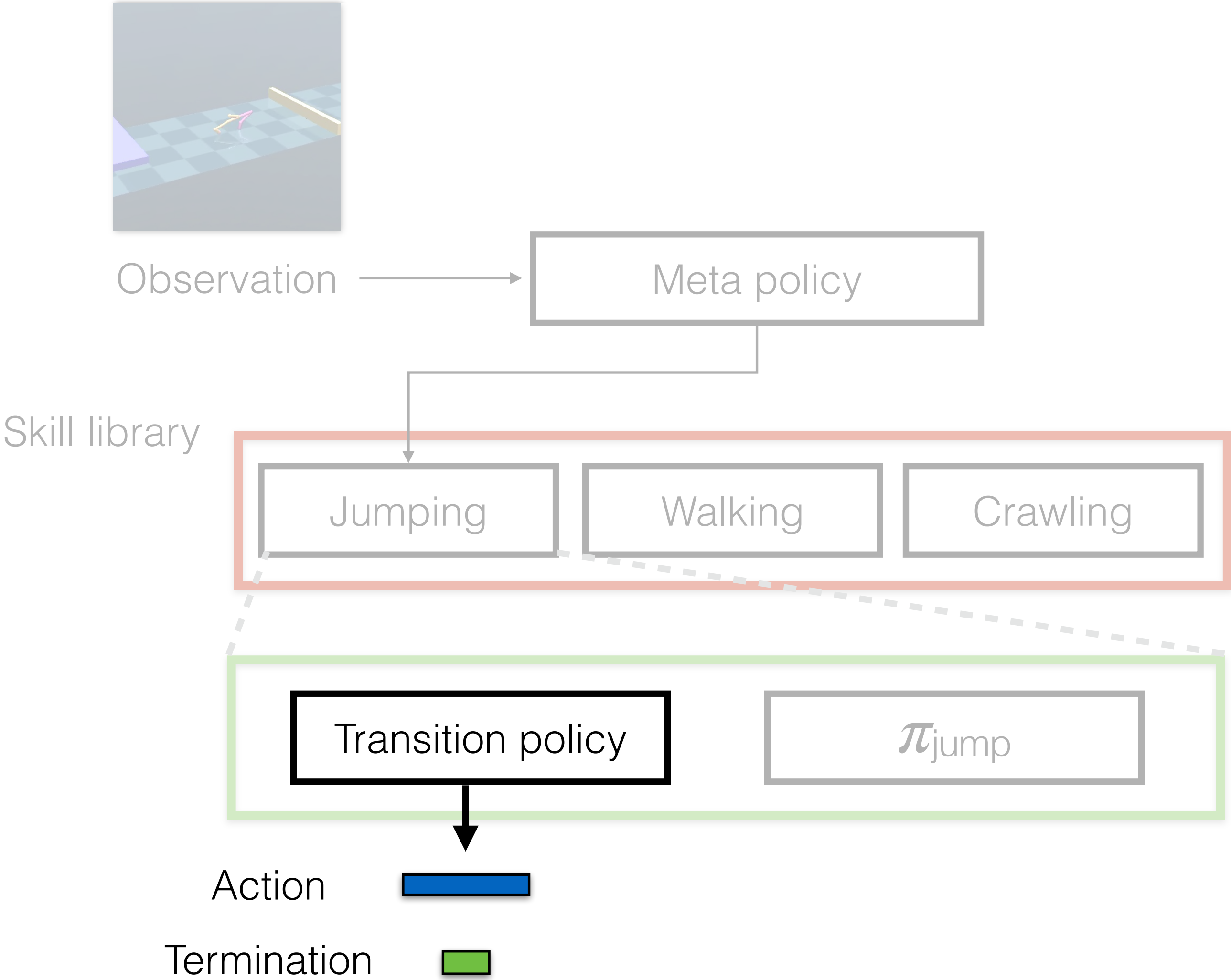




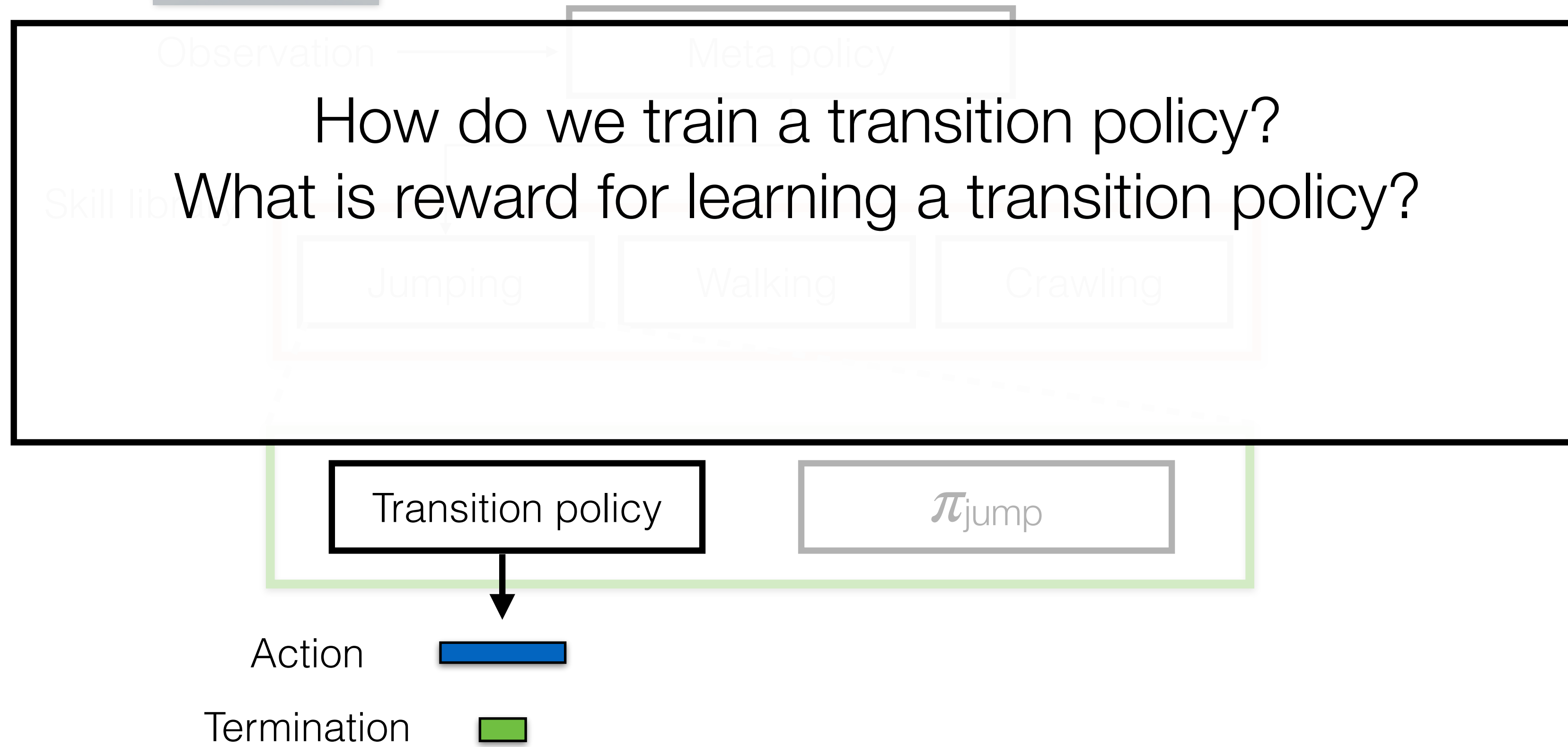
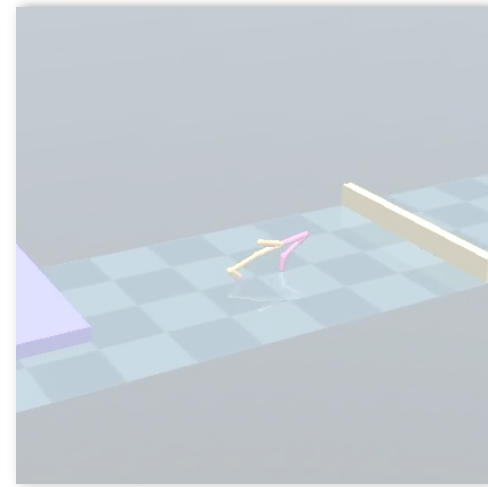
# Model



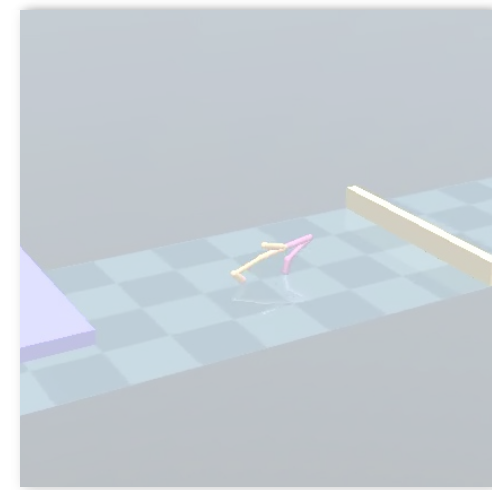
# Model



# Model



# Model



Observation

Meta policy

How do we train a transition policy?  
What is reward for learning a transition policy?

Skill lib

Jumping

Walking

Crawling

**- Success of the following skill**

Transition policy

$\pi_{\text{jump}}$

Action

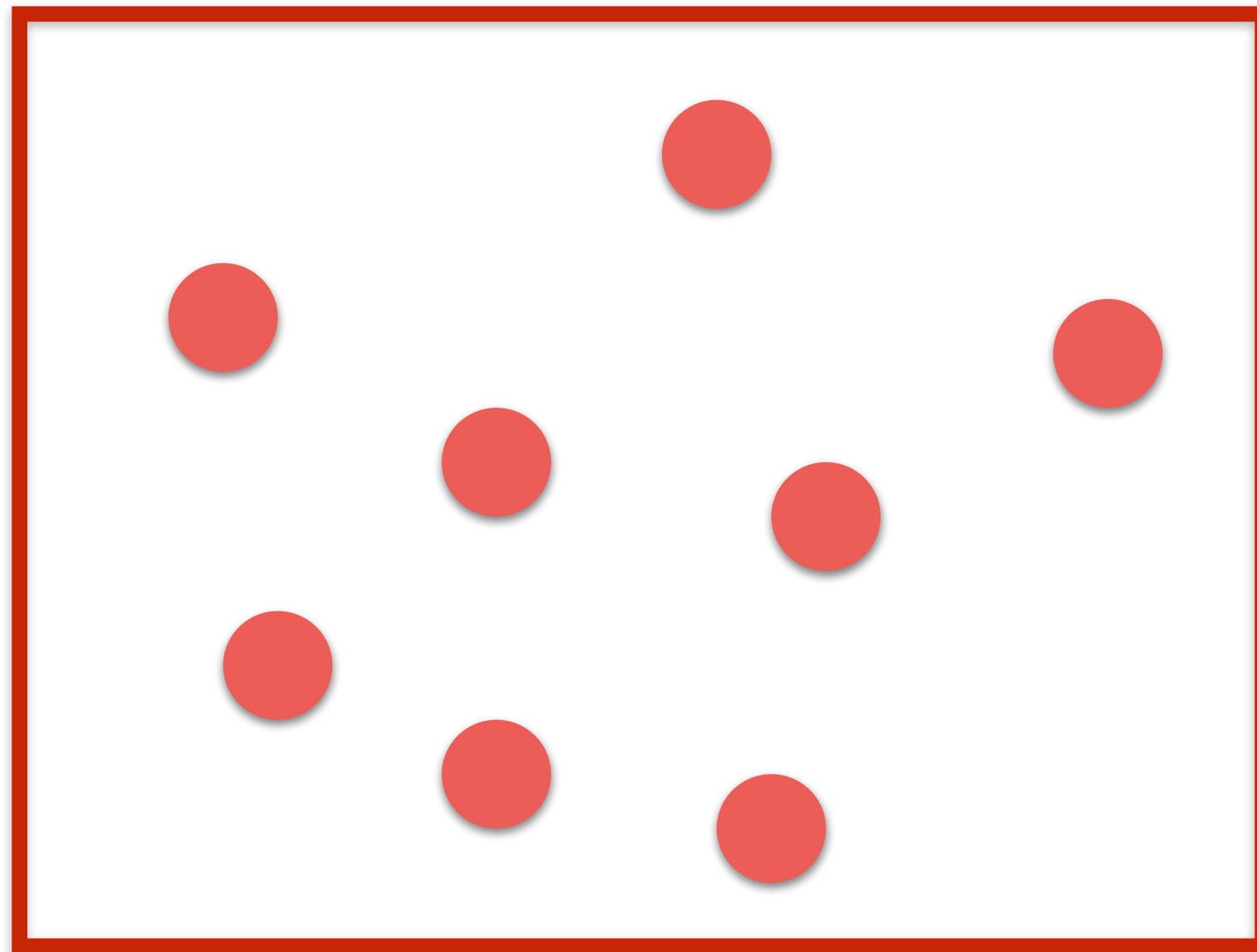


Termination

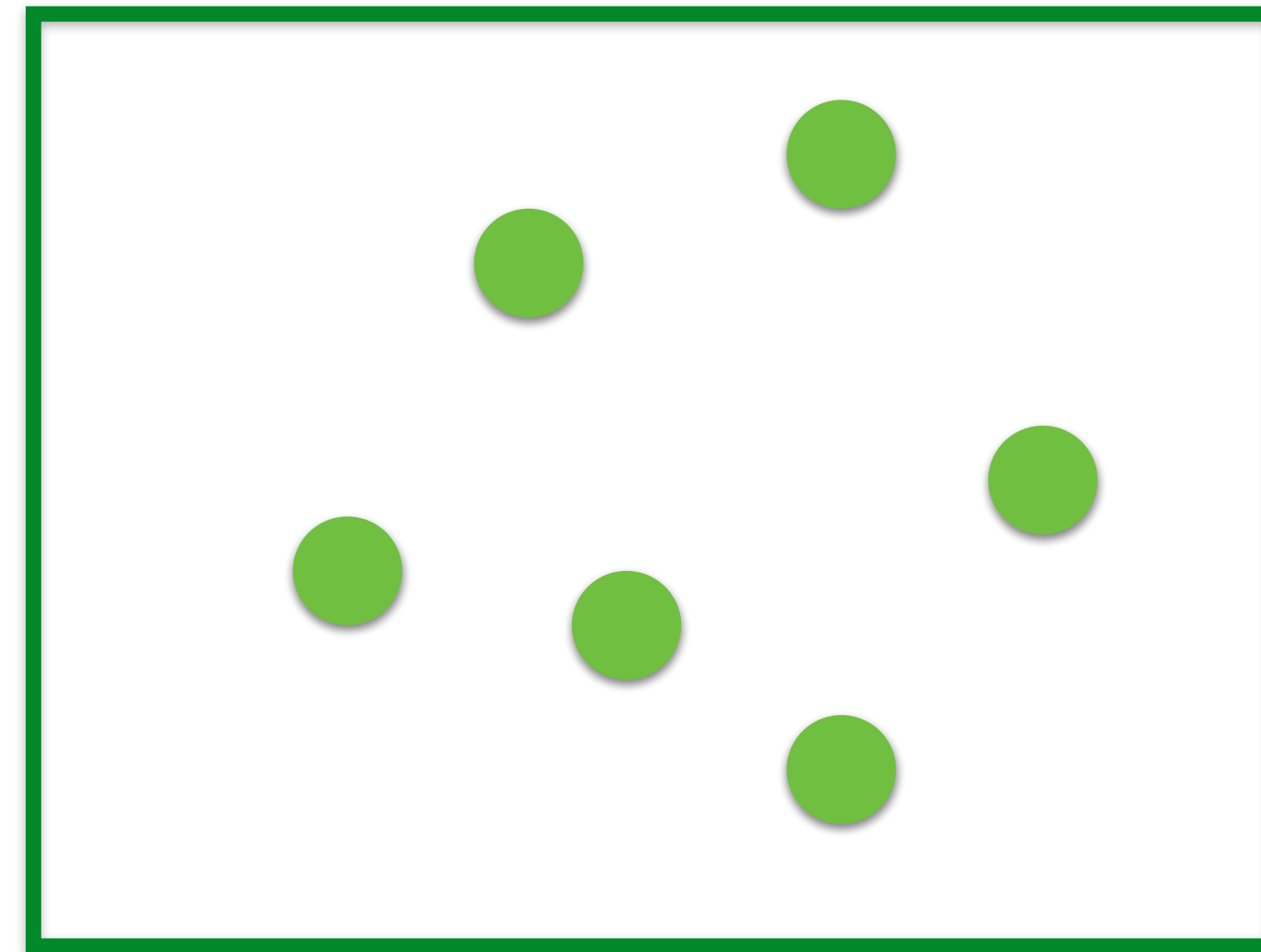




# Learning Transition Policy

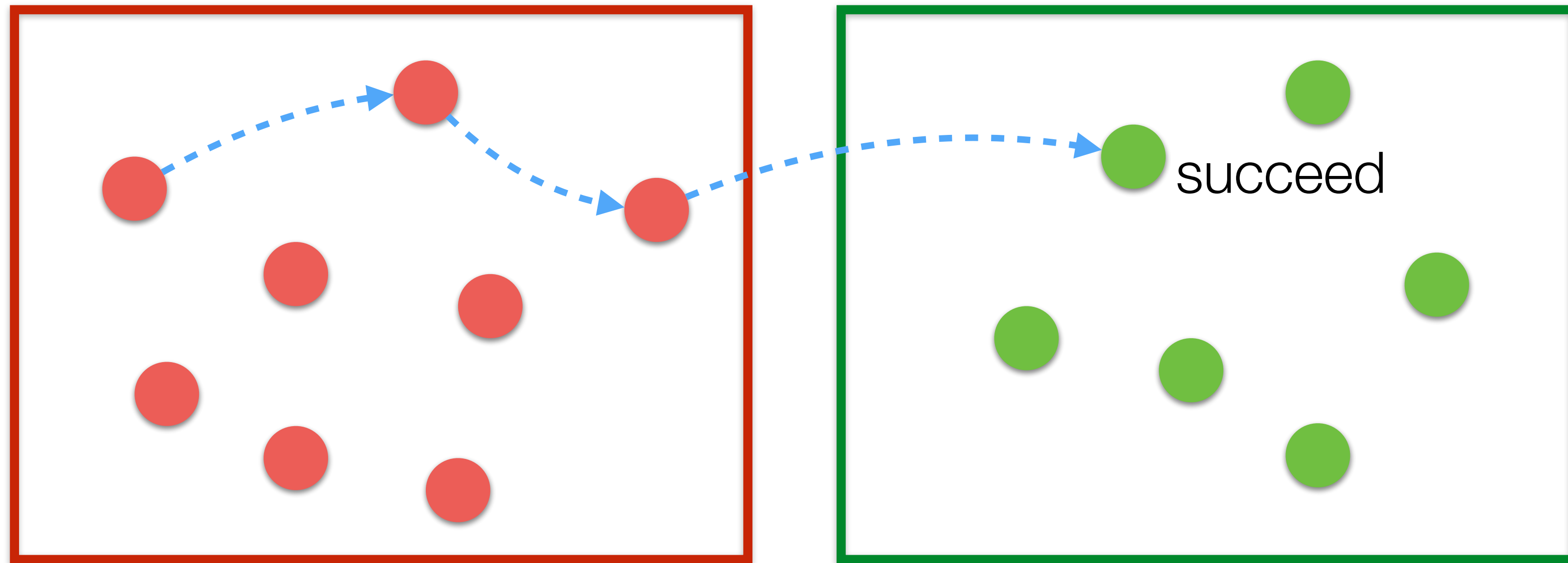


**Bad** initial states for  $\pi_{\text{walk}}$



**Good** initial states for  $\pi_{\text{walk}}$

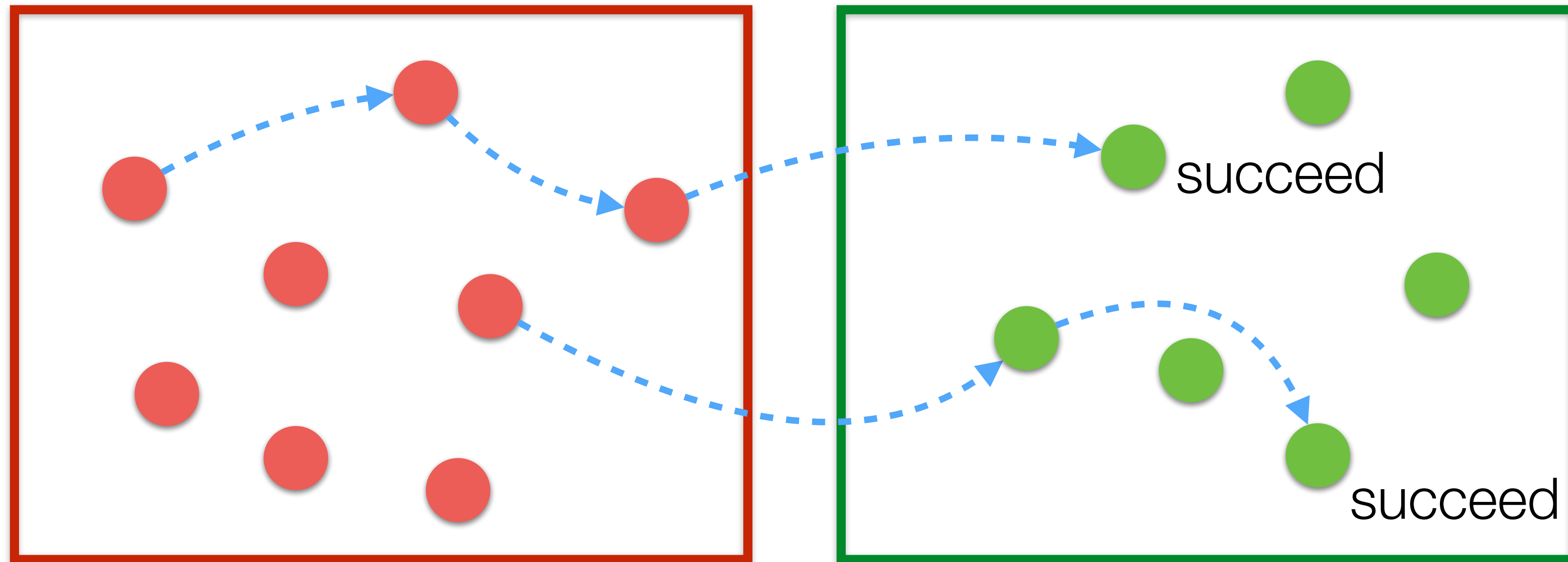
# Learning Transition Policy



**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

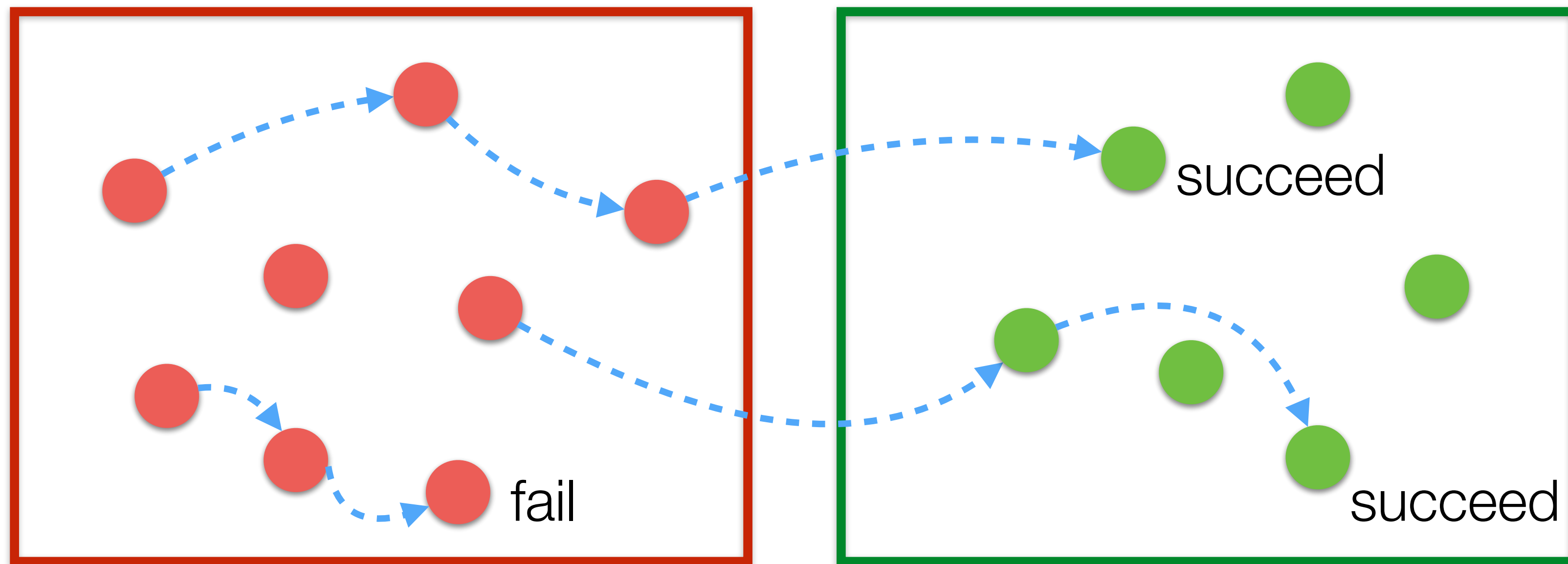
# Learning Transition Policy



**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

# Learning Transition Policy



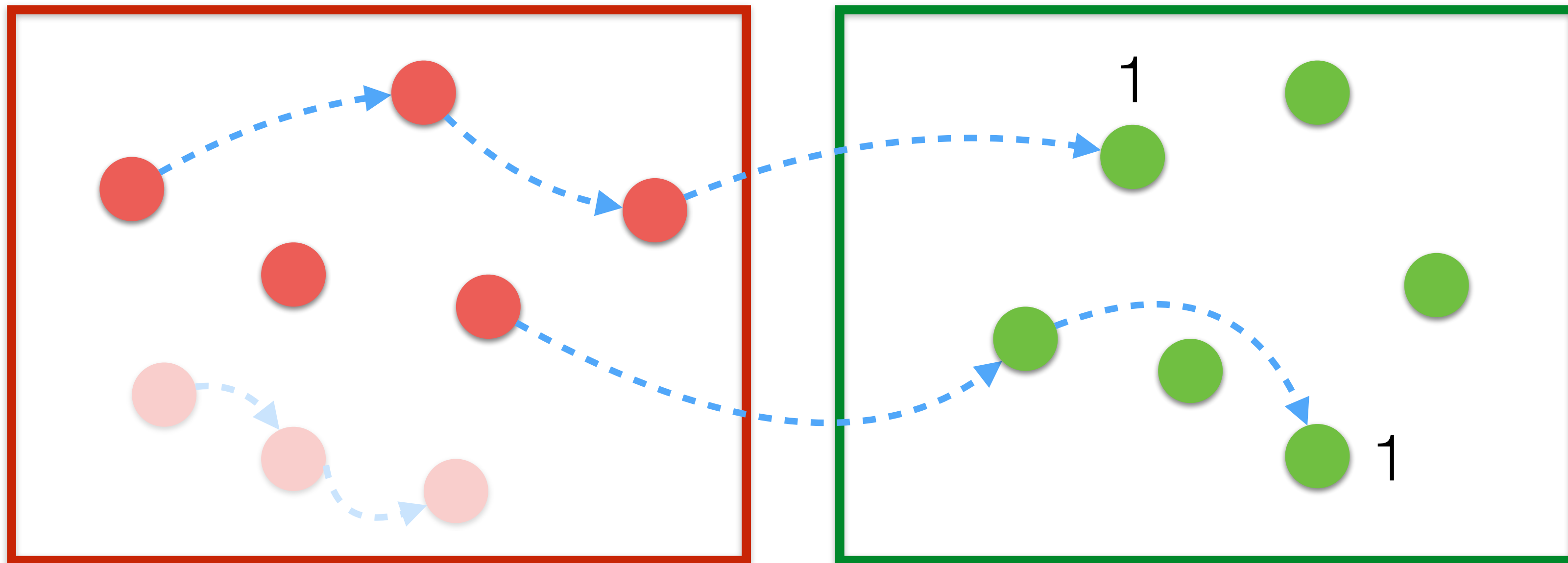
**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$



# Learning Transition Policy

Successful execution of the following skill: +1



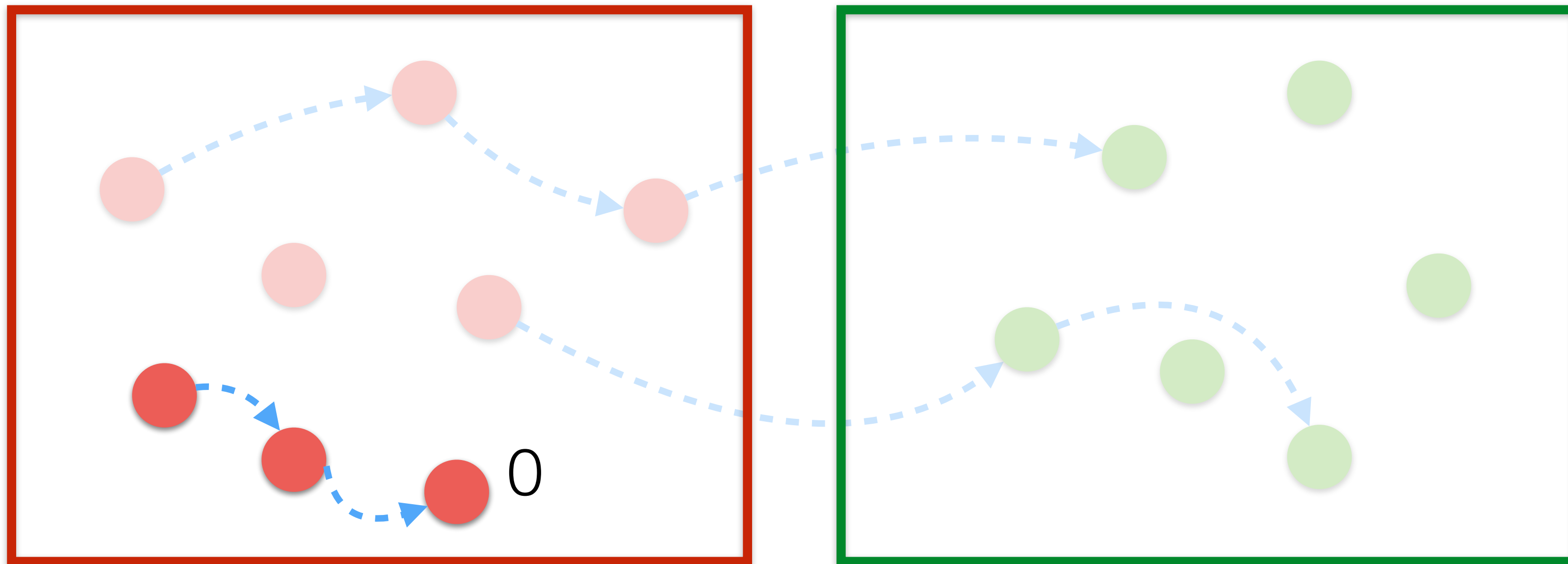
**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

# Learning Transition Policy

Successful execution of the following skill: +1

Failing execution of the following skill: 0



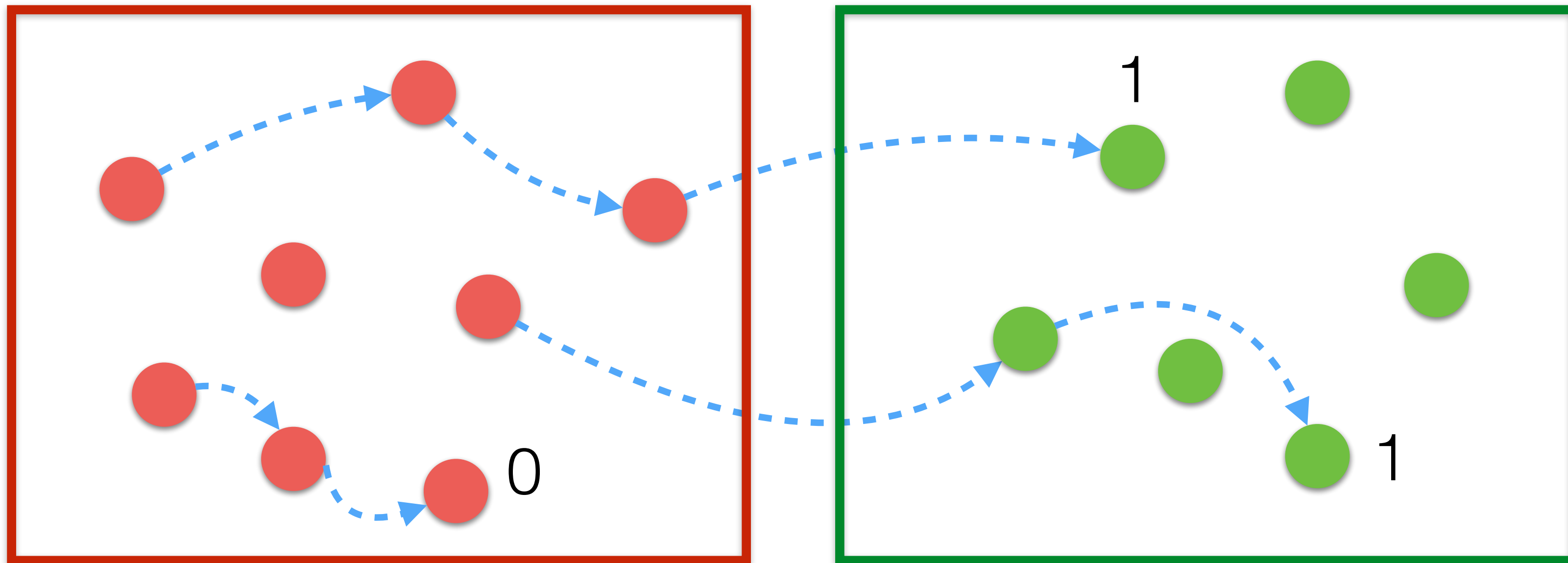
**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

# Learning Transition Policy

Successful execution of the following skill: +1

Failing execution of the following skill: 0



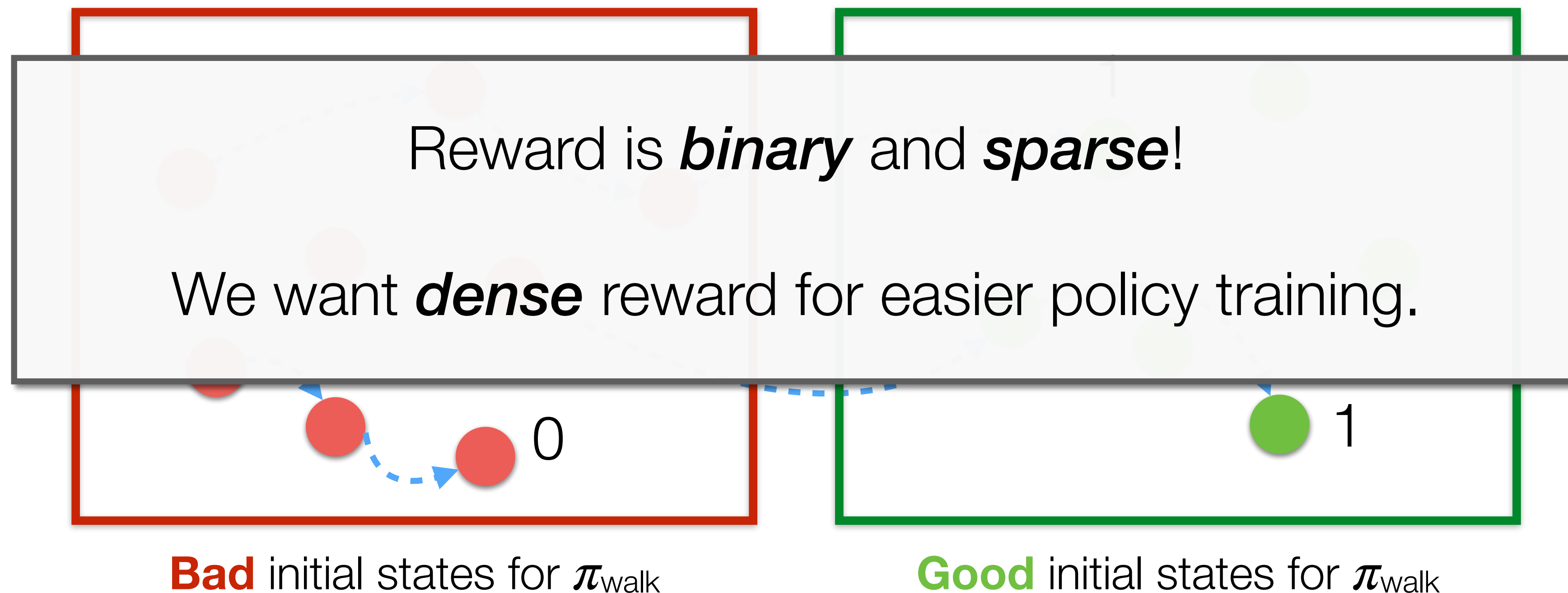
**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

# Learning Transition Policy

Successful execution of the following skill: +1

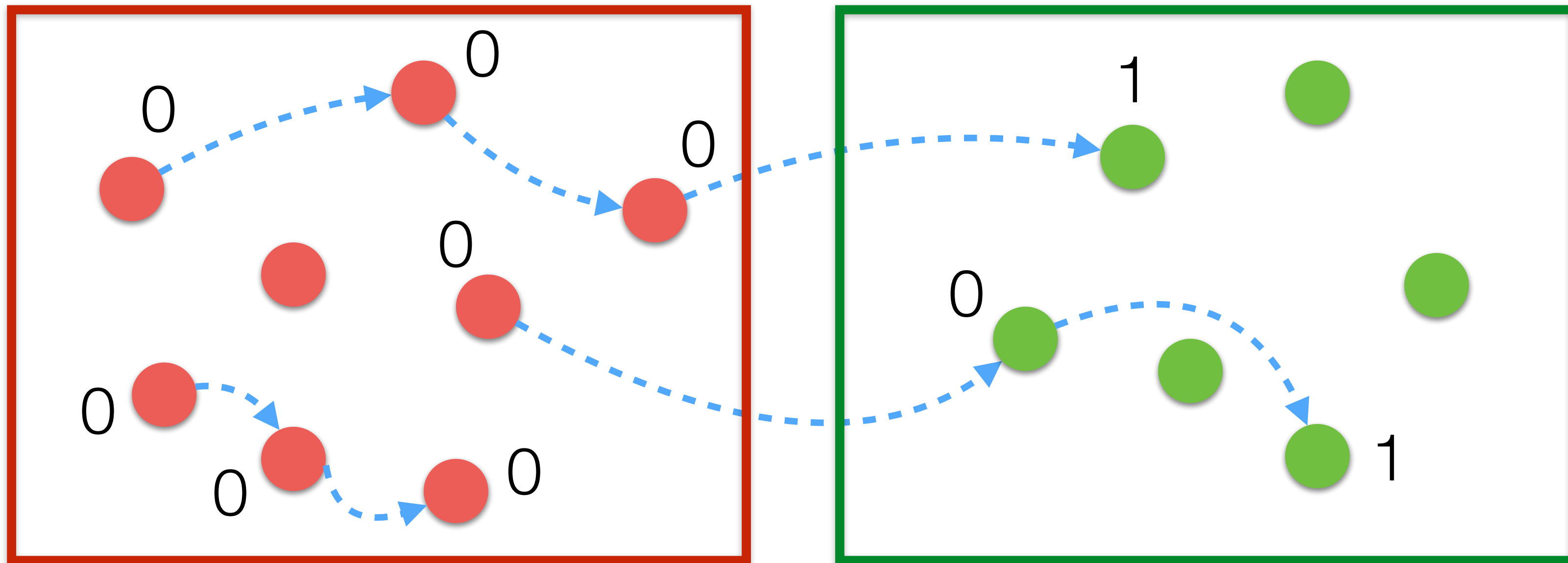
Failing execution of the following skill: 0





# Proximity Reward

Instead of binary reward

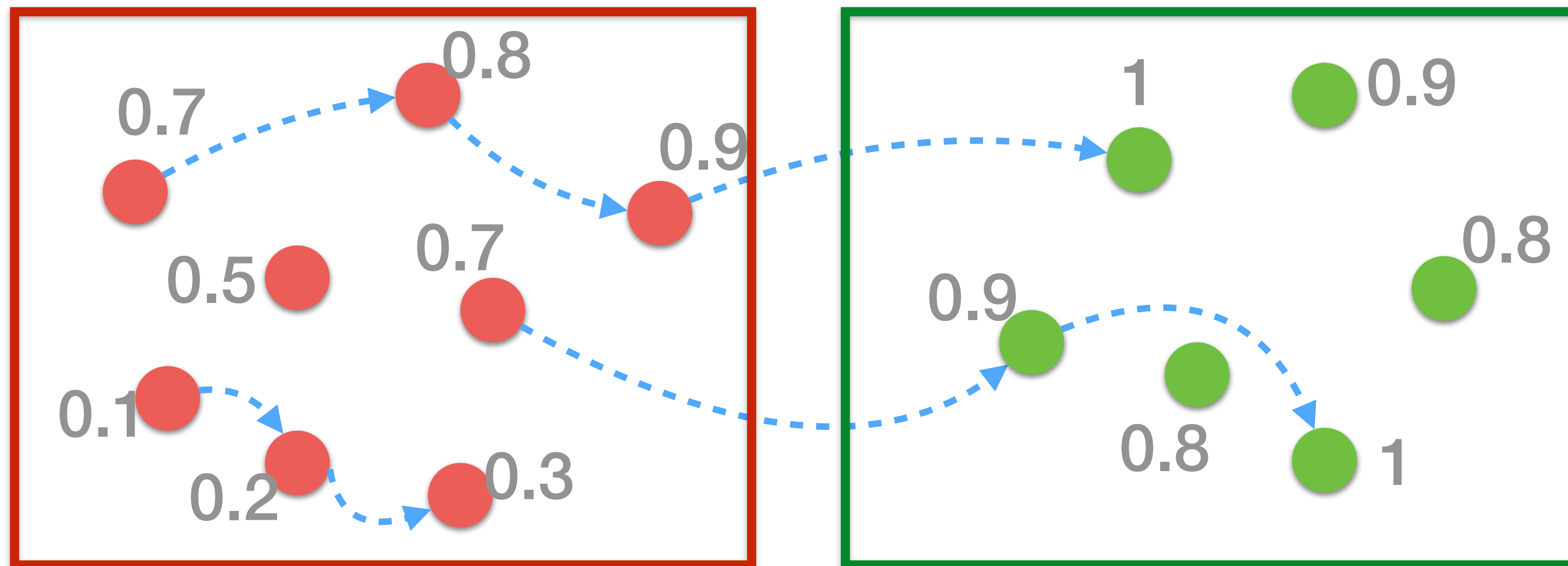


**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

# Proximity Reward

Instead of binary reward, use *“proximity prediction”*, which estimates how close to good initial states

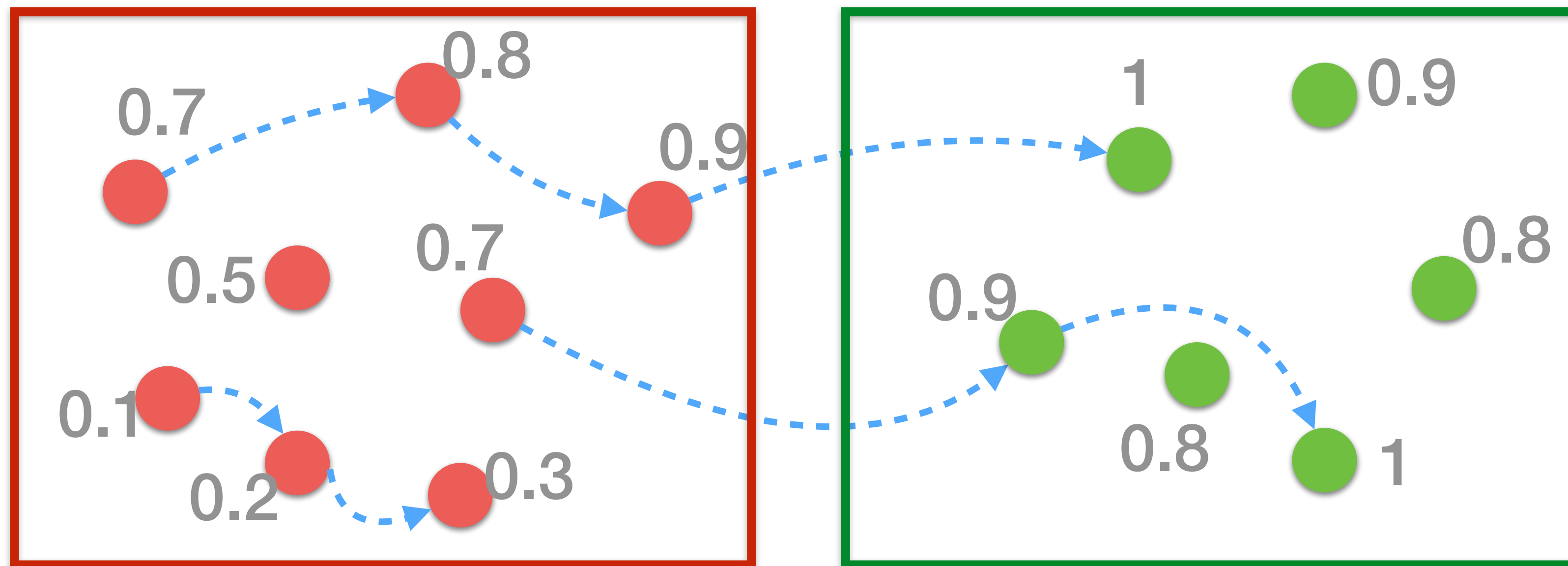


**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

# Proximity Reward

Instead of binary reward, use “*proximity prediction*”, which estimates how close to good initial states



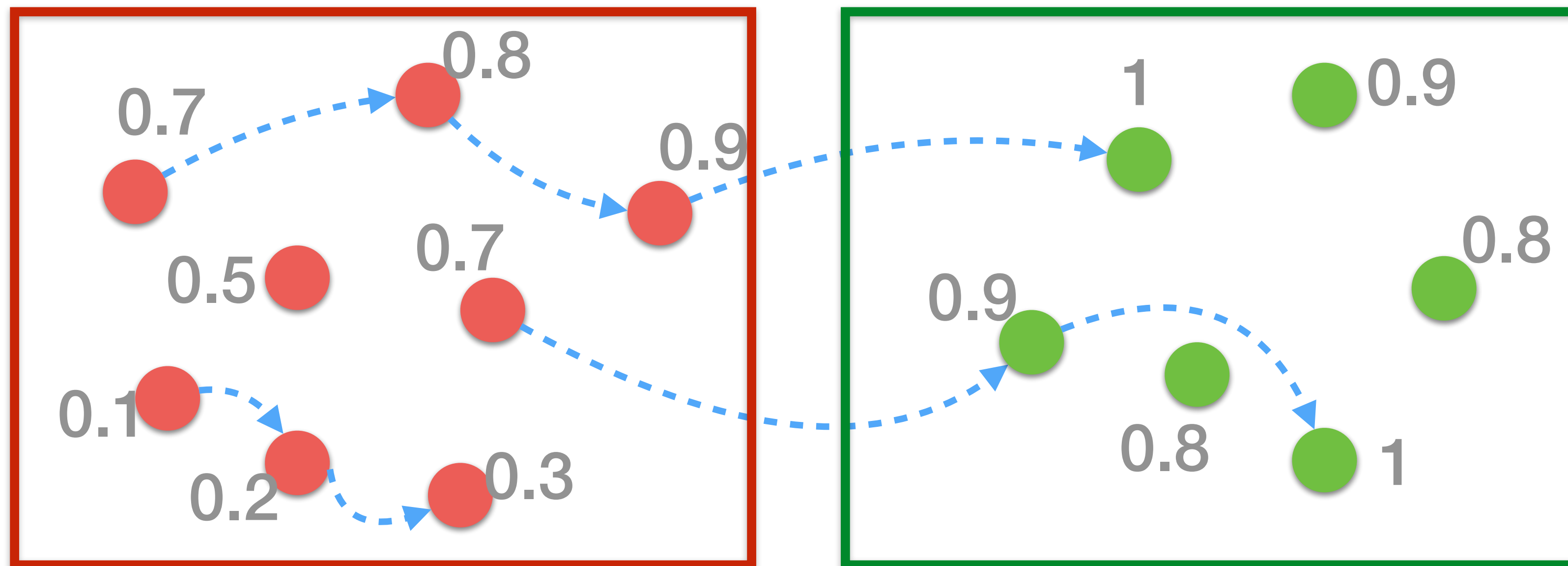
**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

We define *proximity* as:  $P(s) = \delta^{\text{step}}$

# Proximity Reward

Instead of binary reward, use “*proximity prediction*”, which estimates how close to good initial states



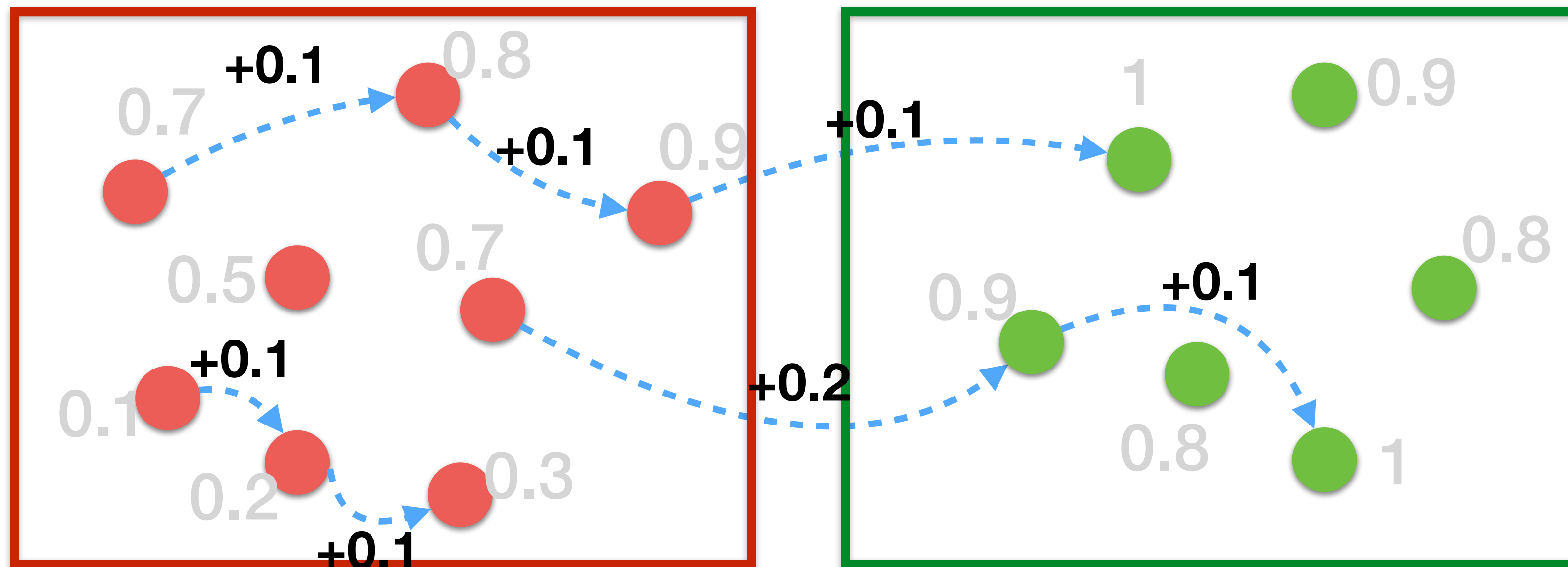
**Bad** initial states for  $\pi_{\text{walk}}$

**Good** initial states for  $\pi_{\text{walk}}$

We define *proximity* as:  $P(s) = \delta^{\text{step}}$

# Proximity Reward

Instead of binary reward, use “*proximity prediction*”, which estimates how close to good initial states



**Bad** initial states for  $\pi_{\text{walk}}$

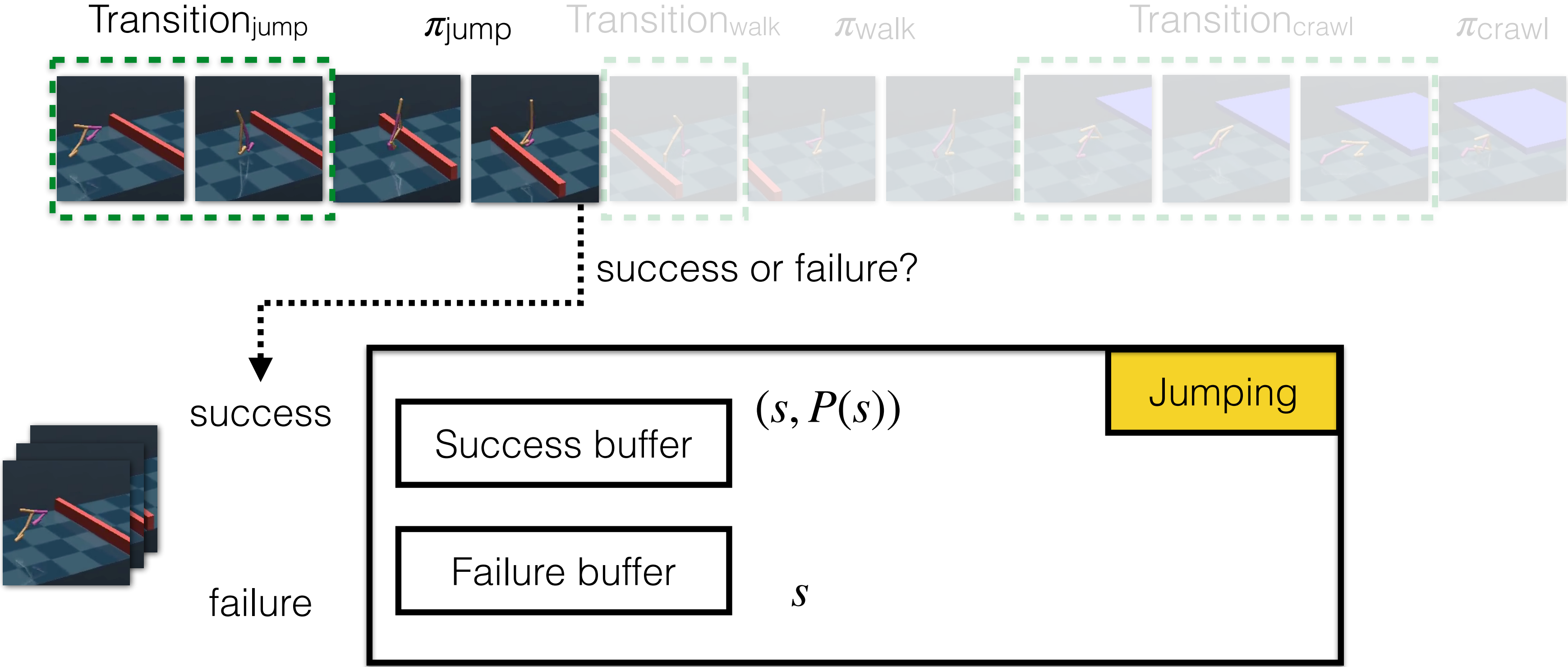
**Good** initial states for  $\pi_{\text{walk}}$

We define *proximity* as:  $P(s) = \delta^{\text{step}}$

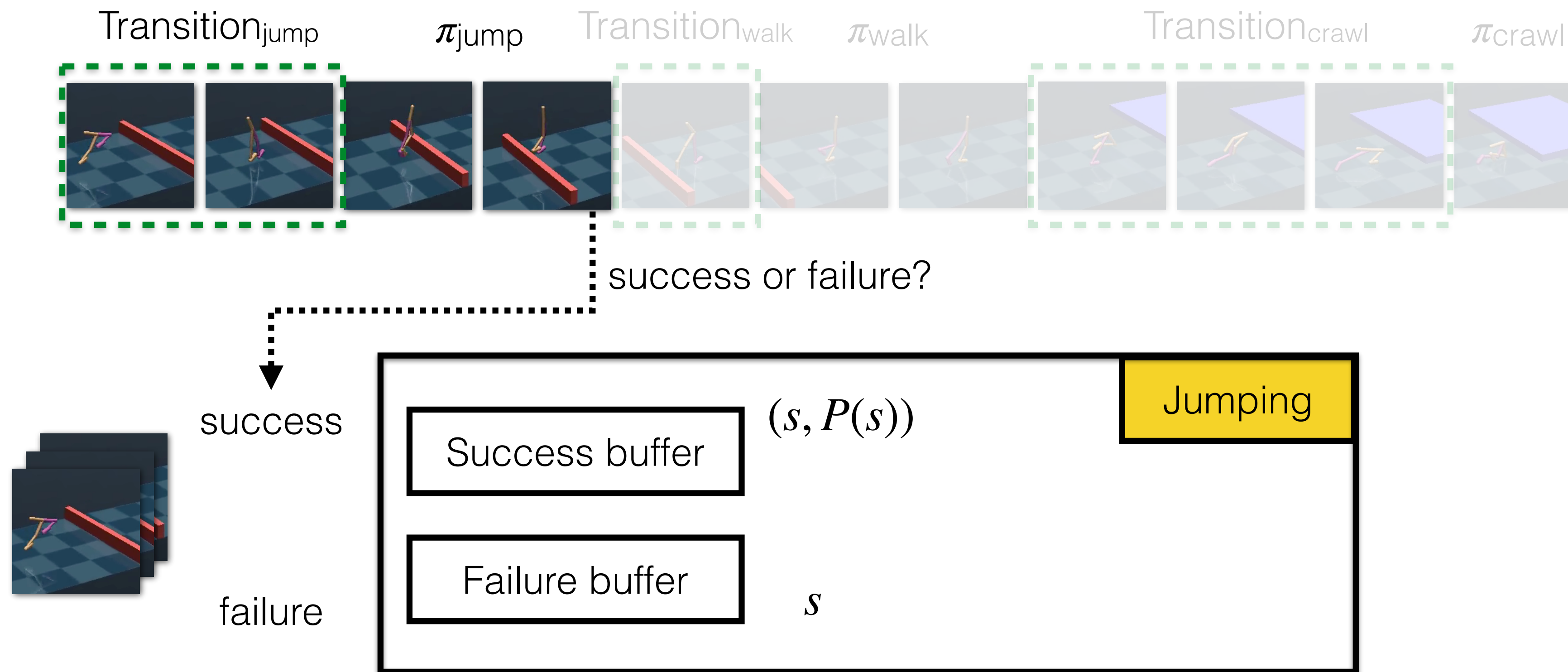
and provide *proximity reward* every step:  $P(s_{t+1}) - P(s_t)$



# Training Proximity Predictor

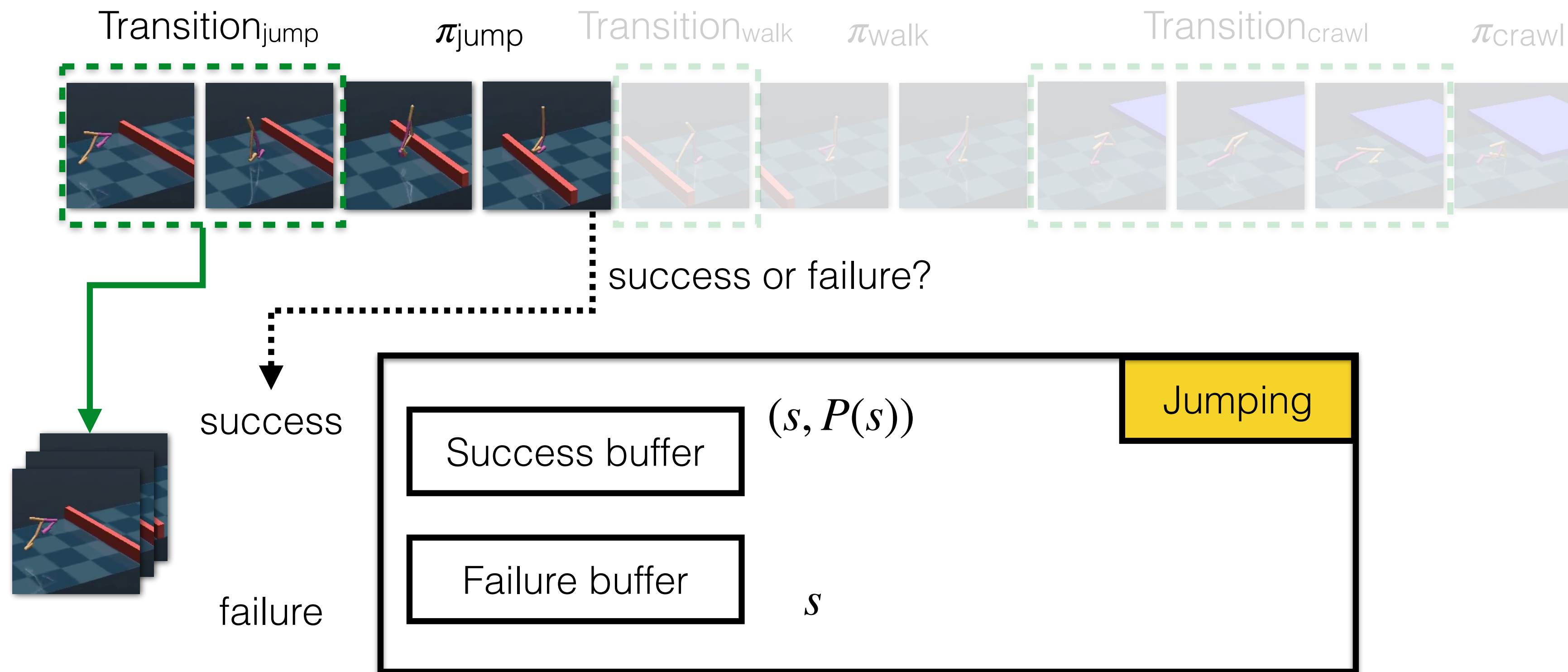


# Training Proximity Predictor



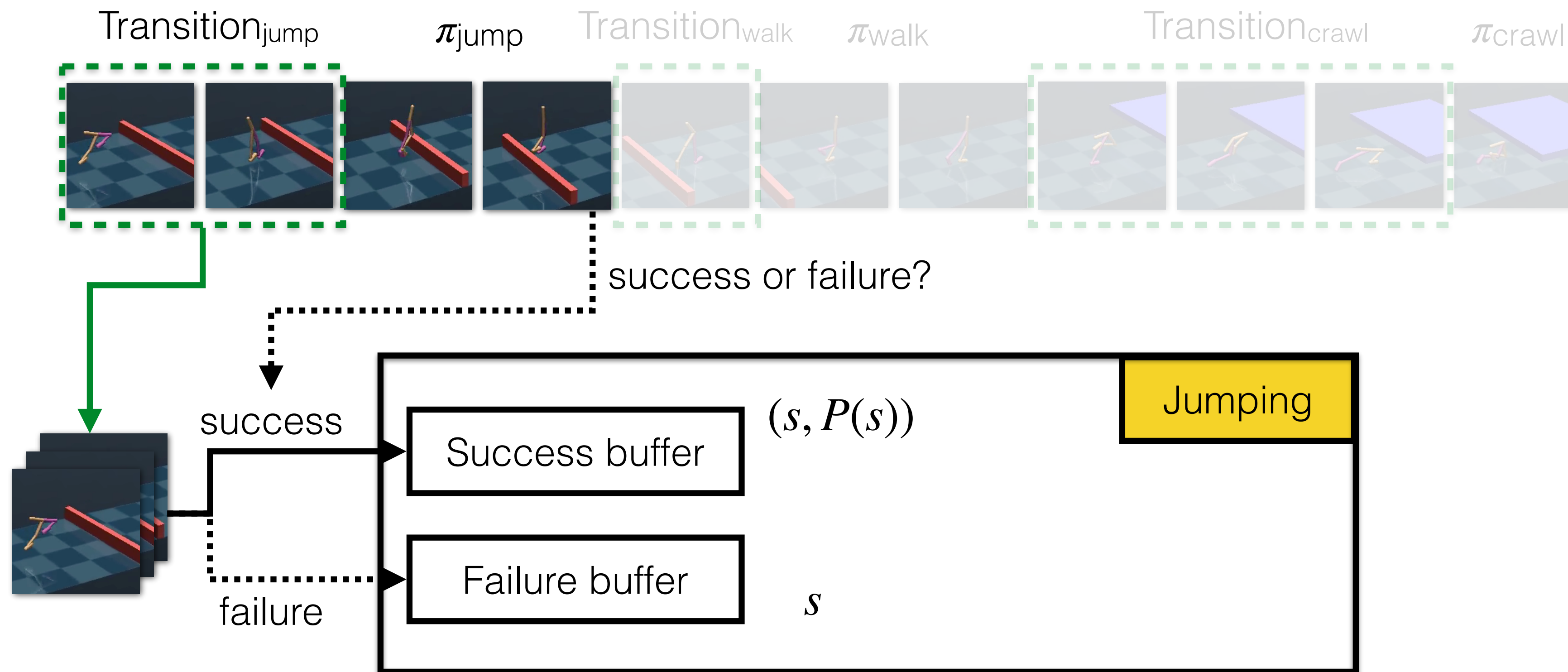
Collect training data for proximity predictors

# Training Proximity Predictor



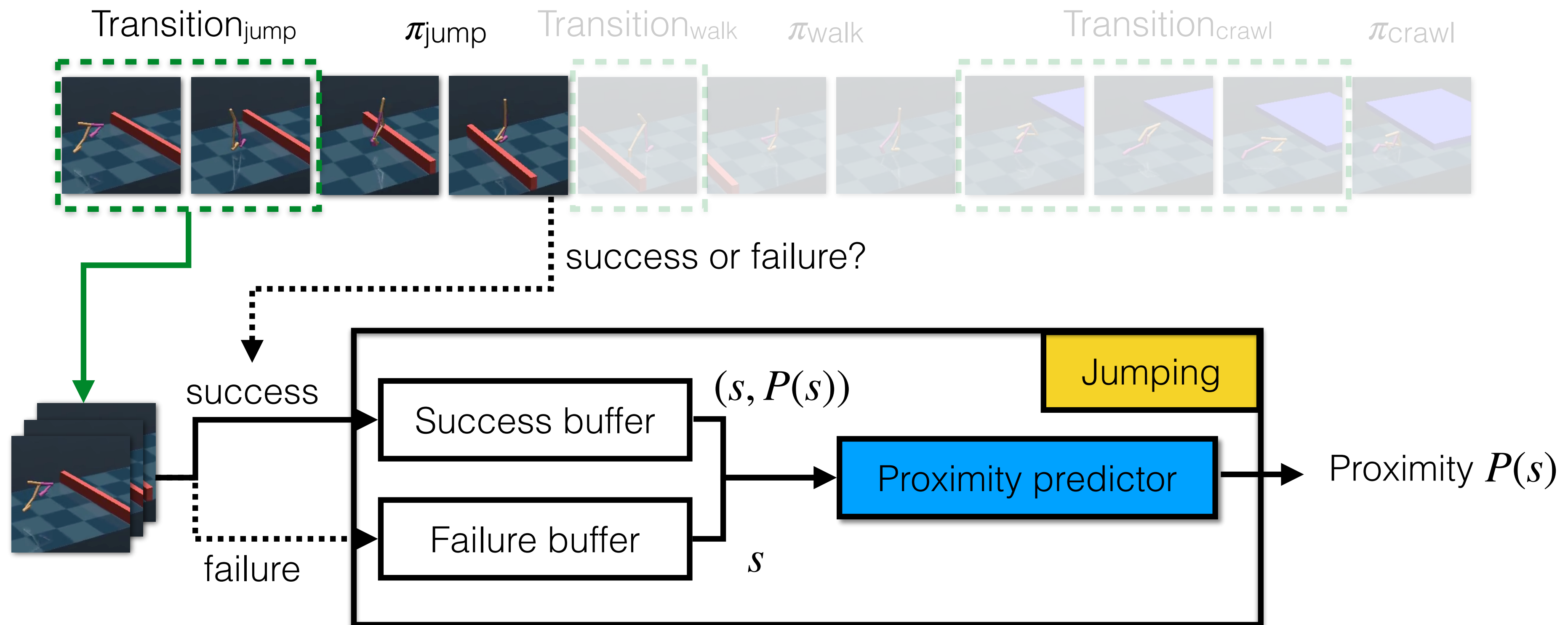
Collect training data for proximity predictors

# Training Proximity Predictor



Collect training data for proximity predictors

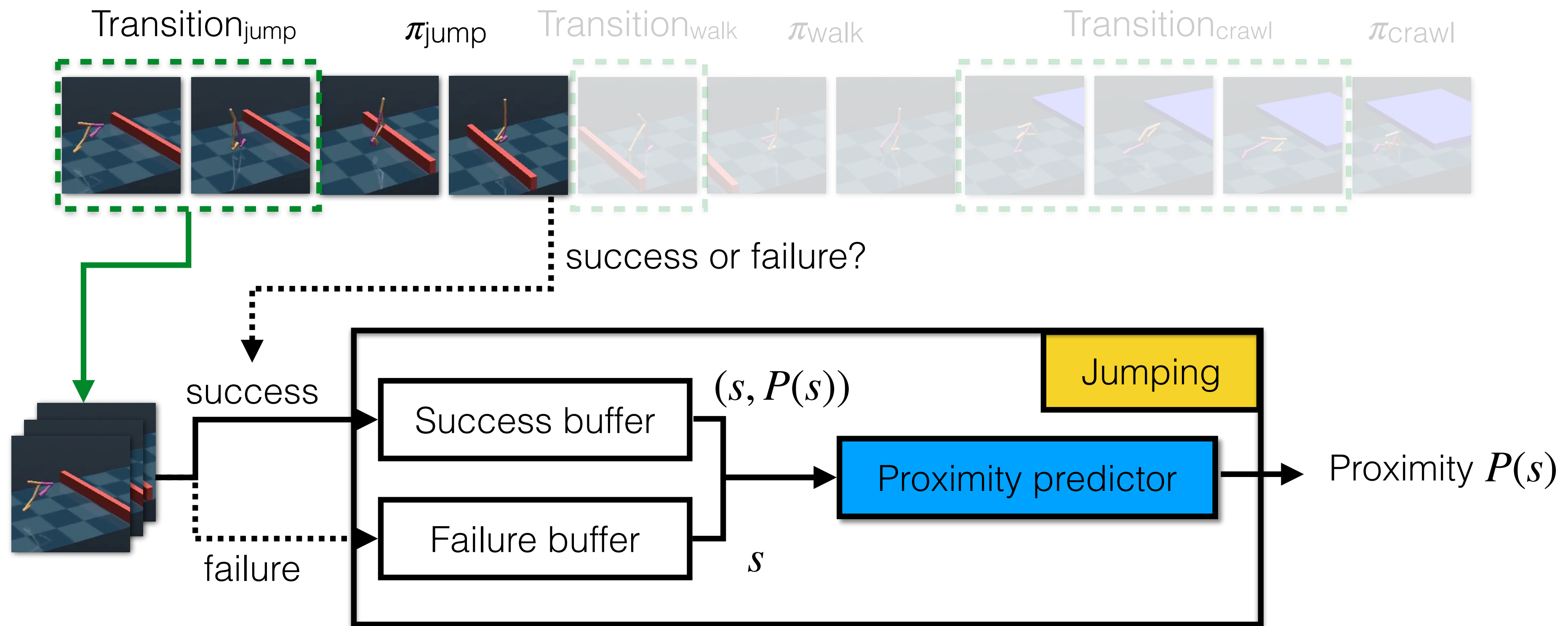
# Training Proximity Predictor



Train proximity predictors

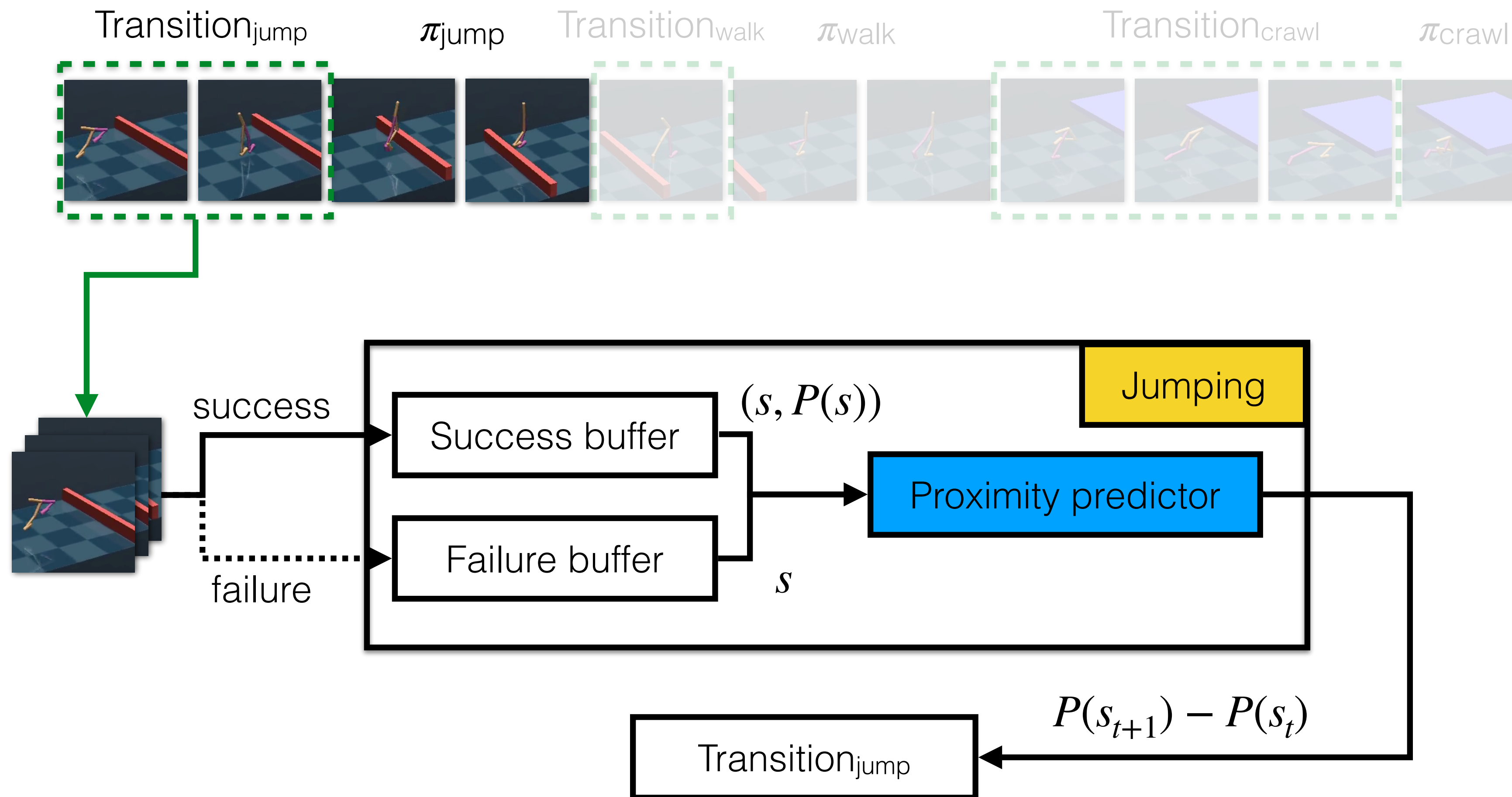


# Training Proximity Predictor



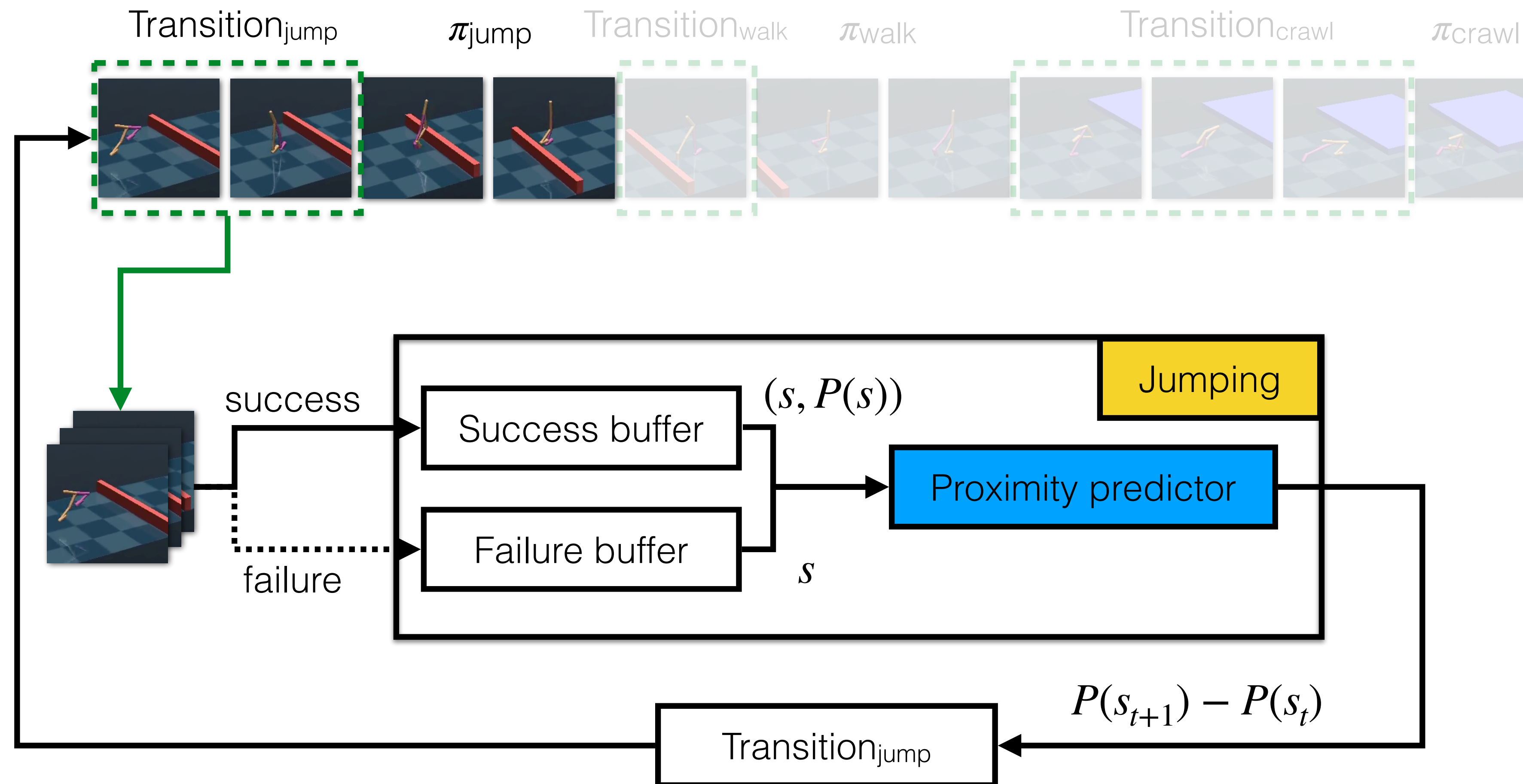
Train proximity predictors

# Training Proximity Predictor



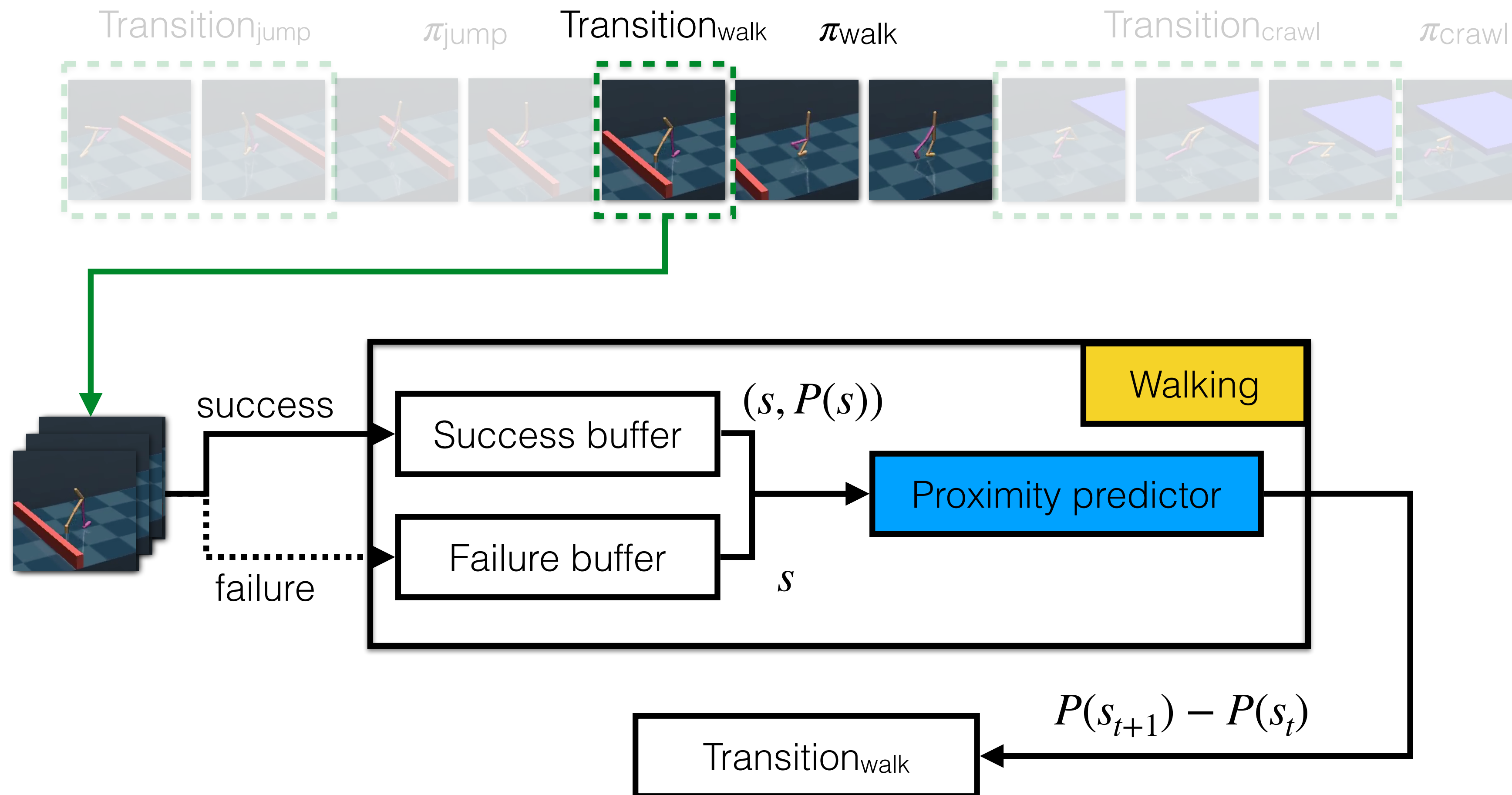
Provide more accurate proximity reward

# Training Proximity Predictor



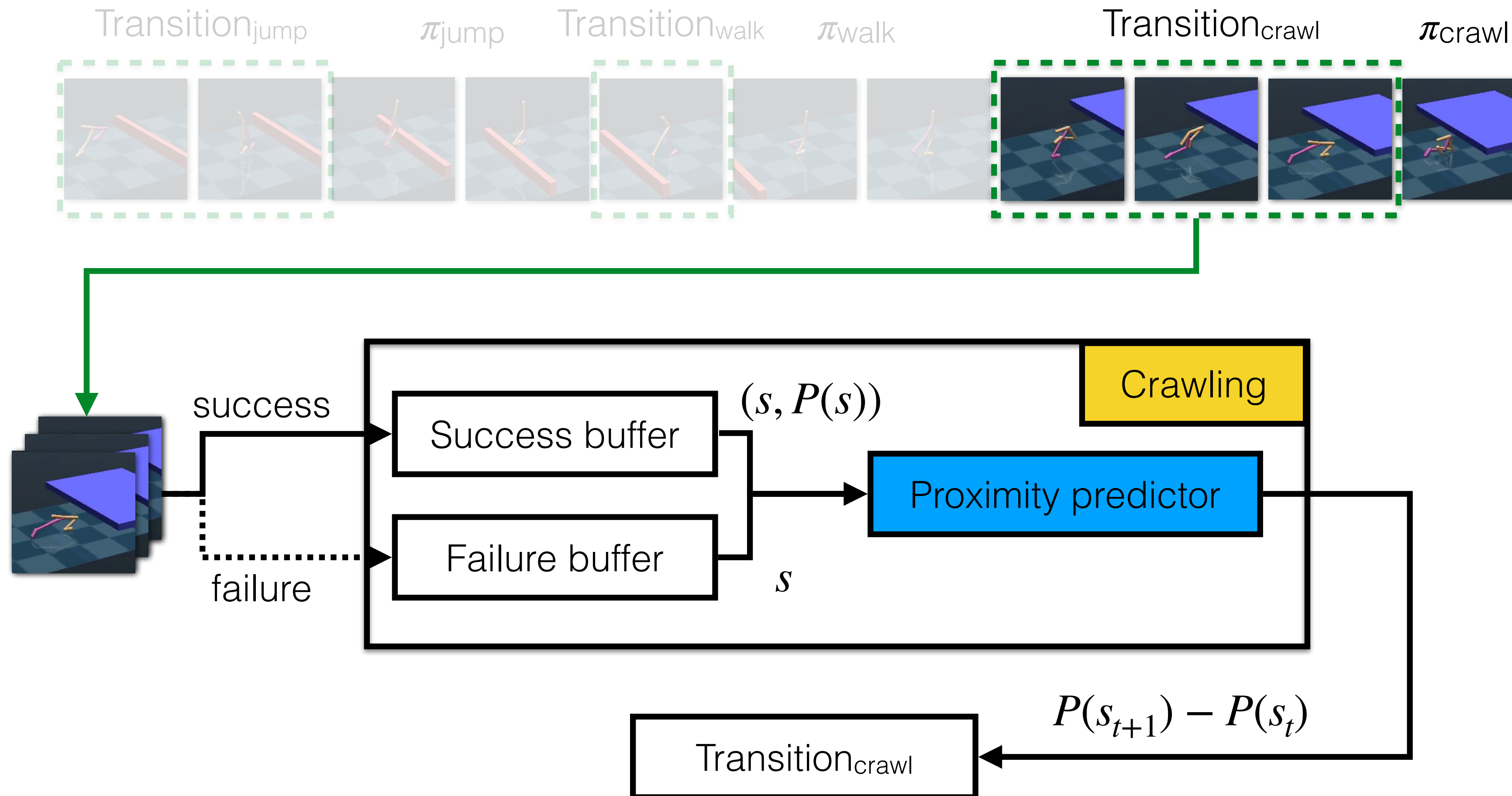
Getter better data with improved policy

# Training Proximity Predictor



Train all transition policies simultaneously

# Training Proximity Predictor



Train all transition policies simultaneously

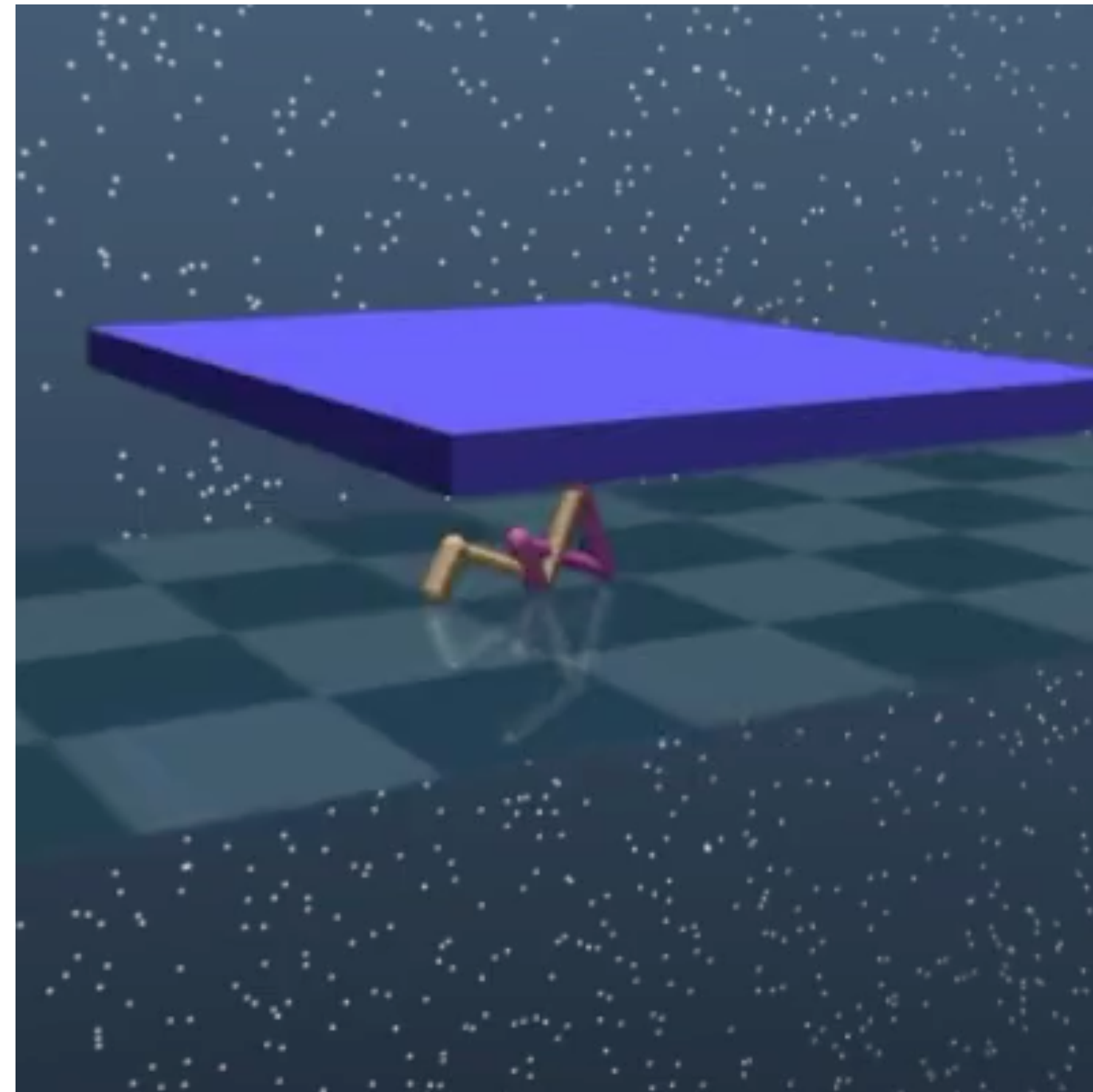


# Obstacle Course

Crawl

***Transition***

Walk

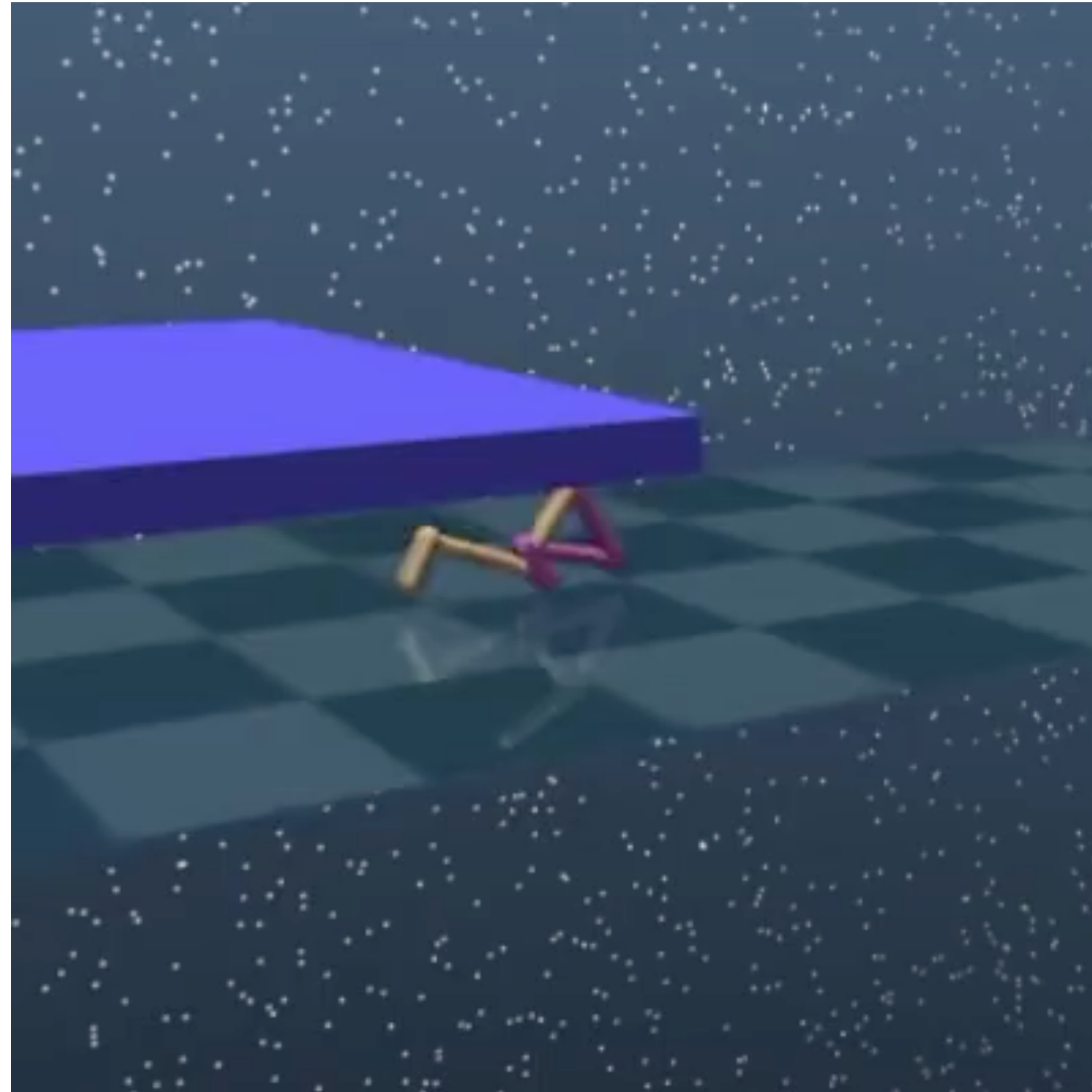


# Obstacle Course

Crawl

***Transition***

Walk

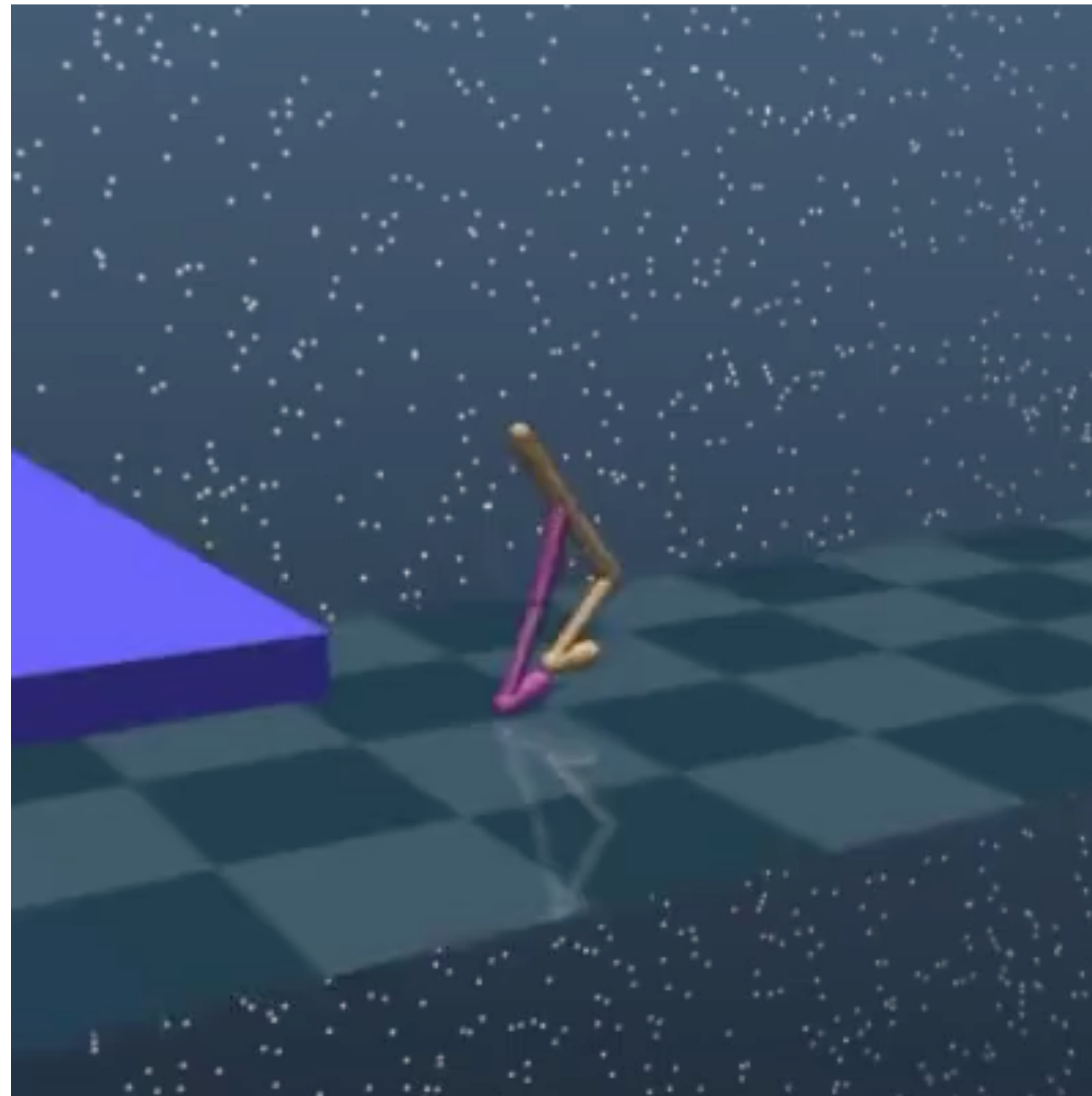


# Obstacle Course

Crawl

***Transition***

Walk

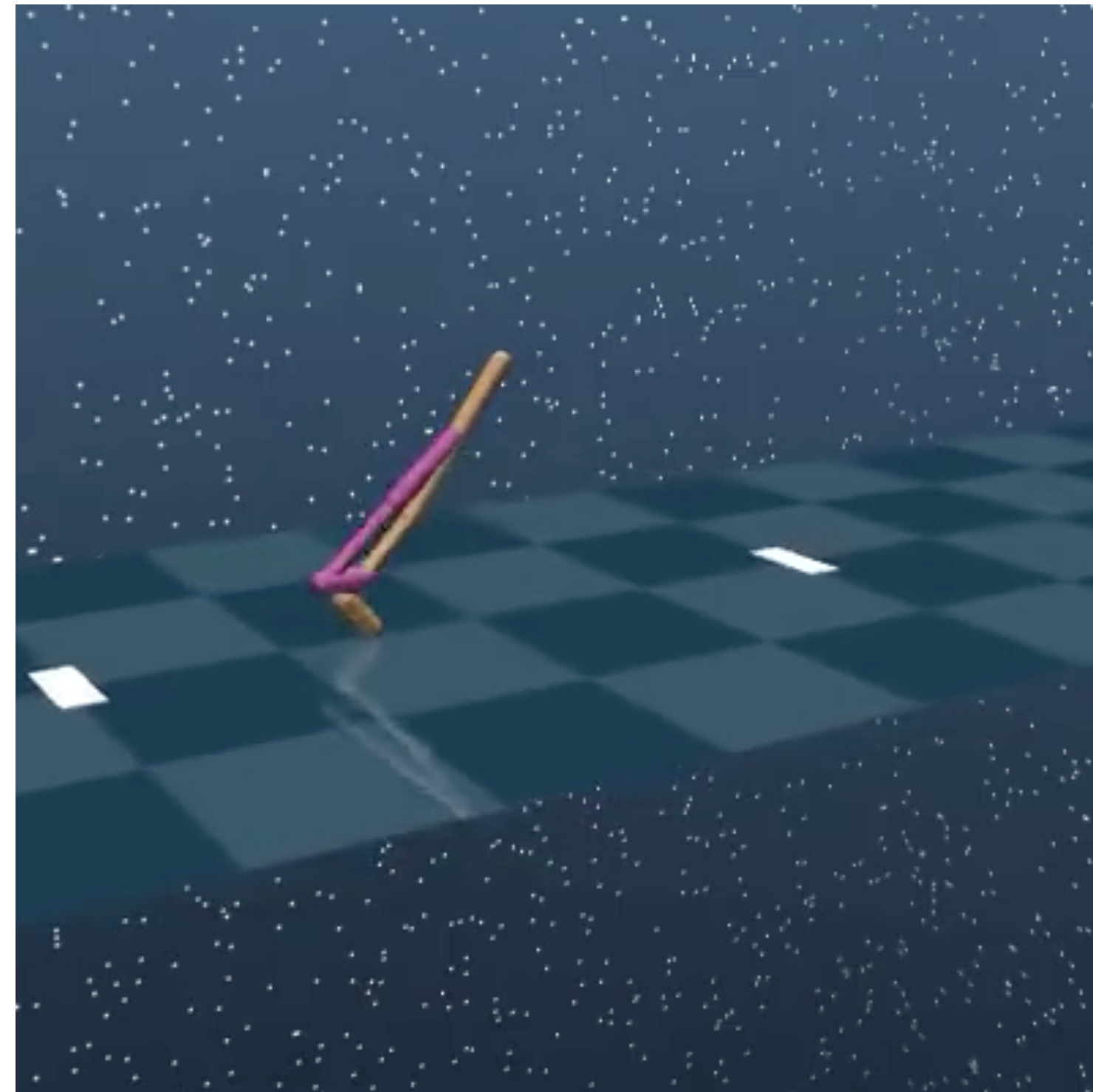


# Walk Forward & Backward

Walk Forward

***Transition***

Walk Backward

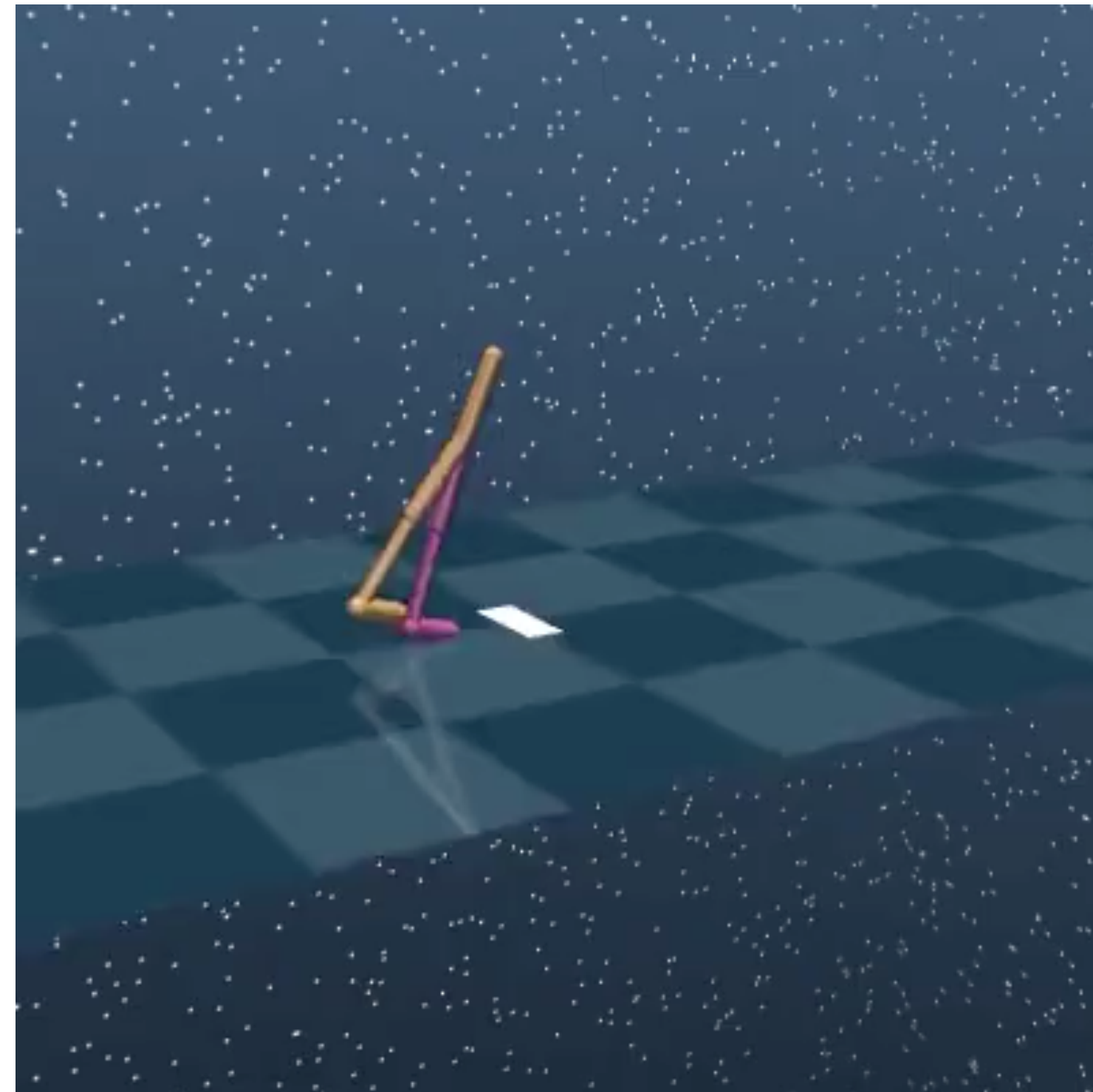


# Walk Forward & Backward

Walk Forward

***Transition***

Walk Backward



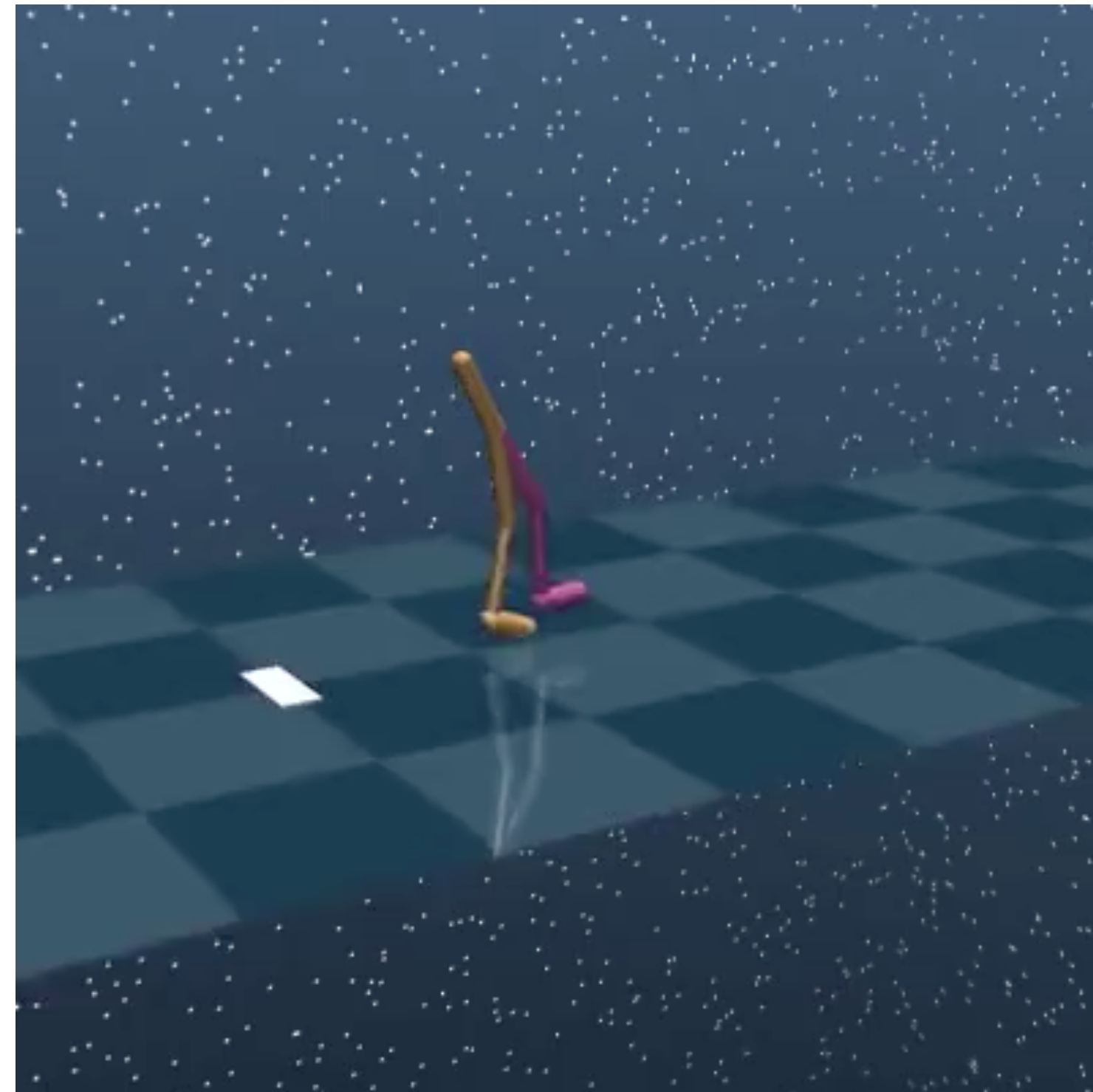


# Walk Forward & Backward

Walk Forward

***Transition***

Walk Backward

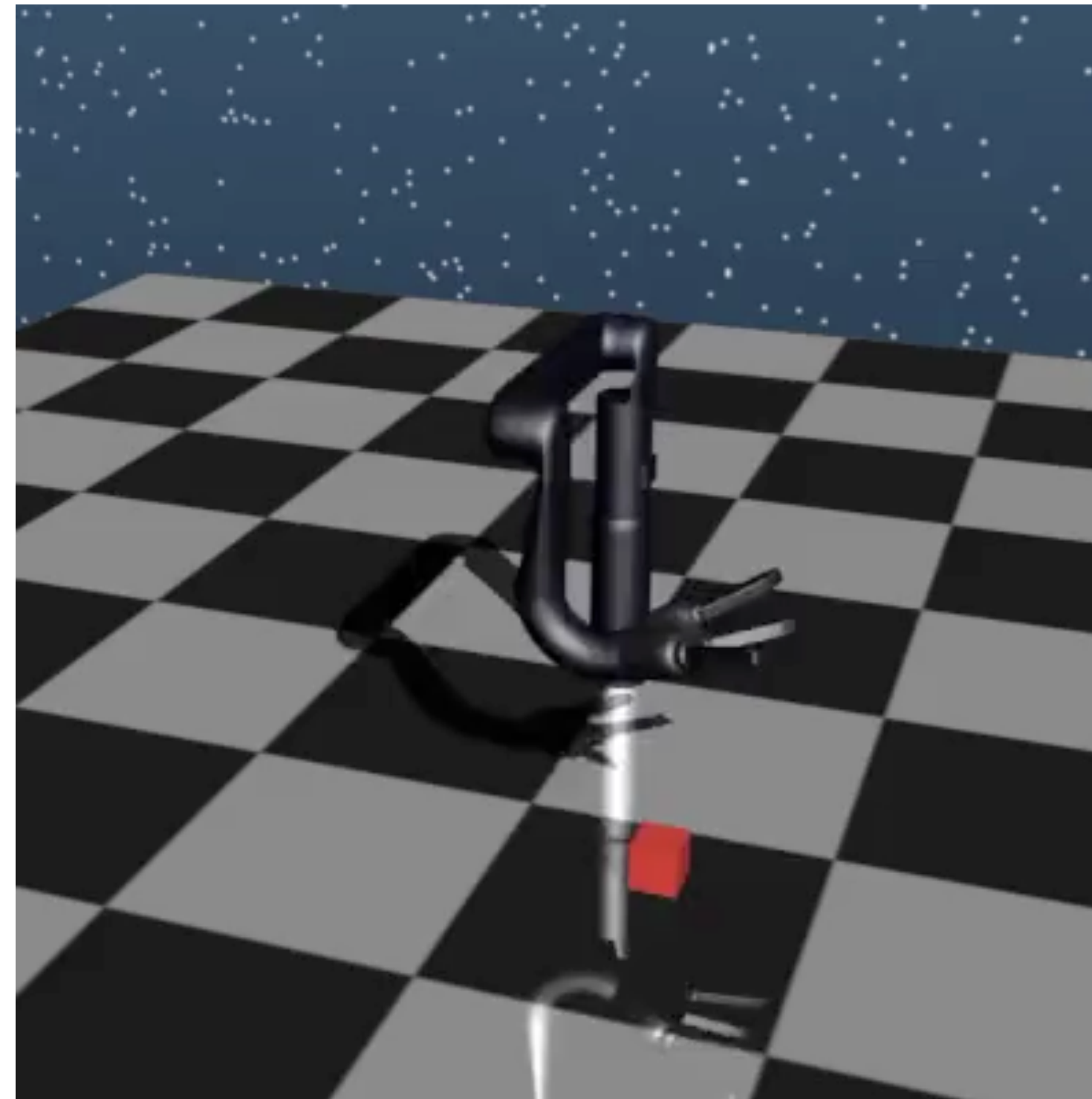


# Repetitive Pick

Pick

***Transition***

Pick

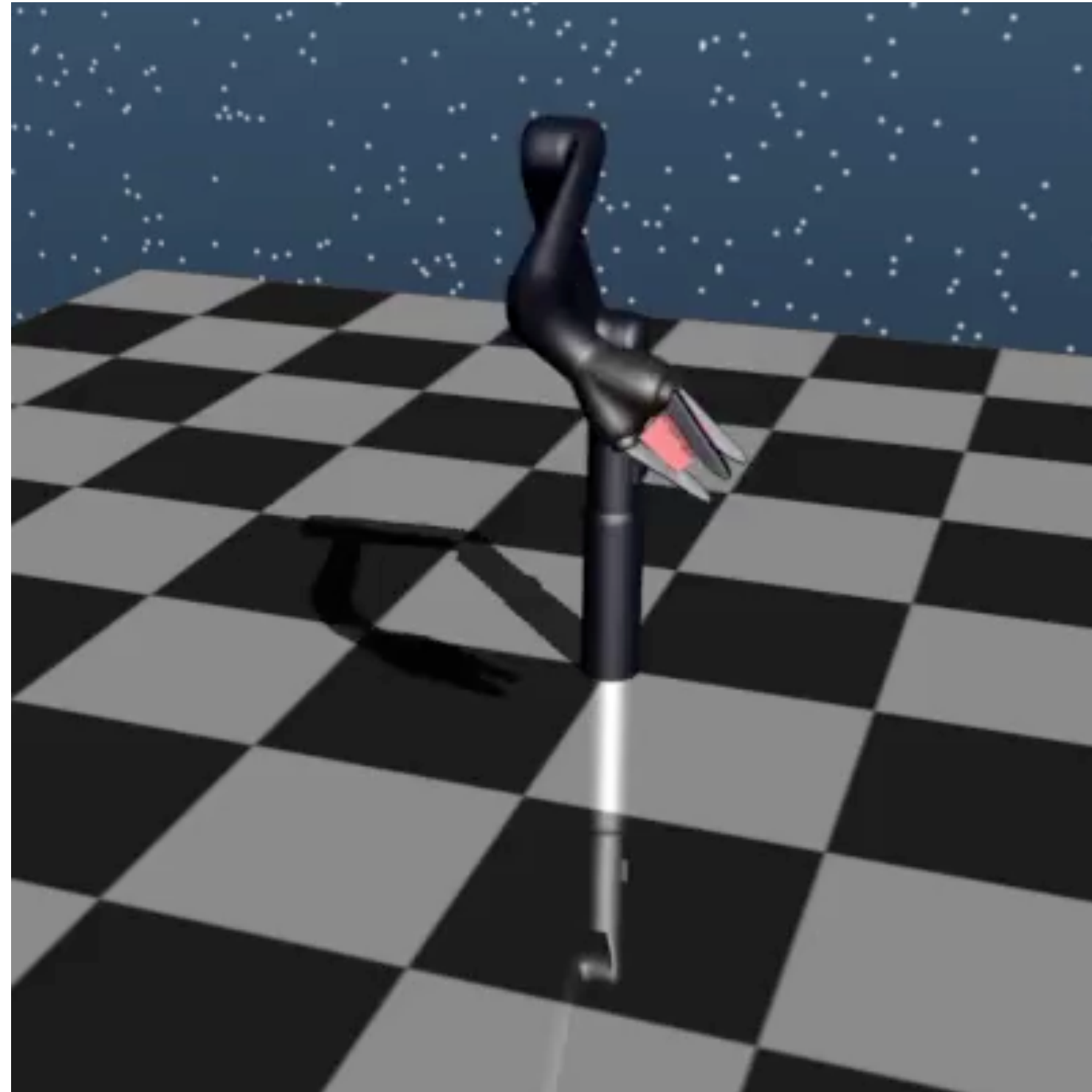


# Repetitive Pick

Pick

***Transition***

Pick

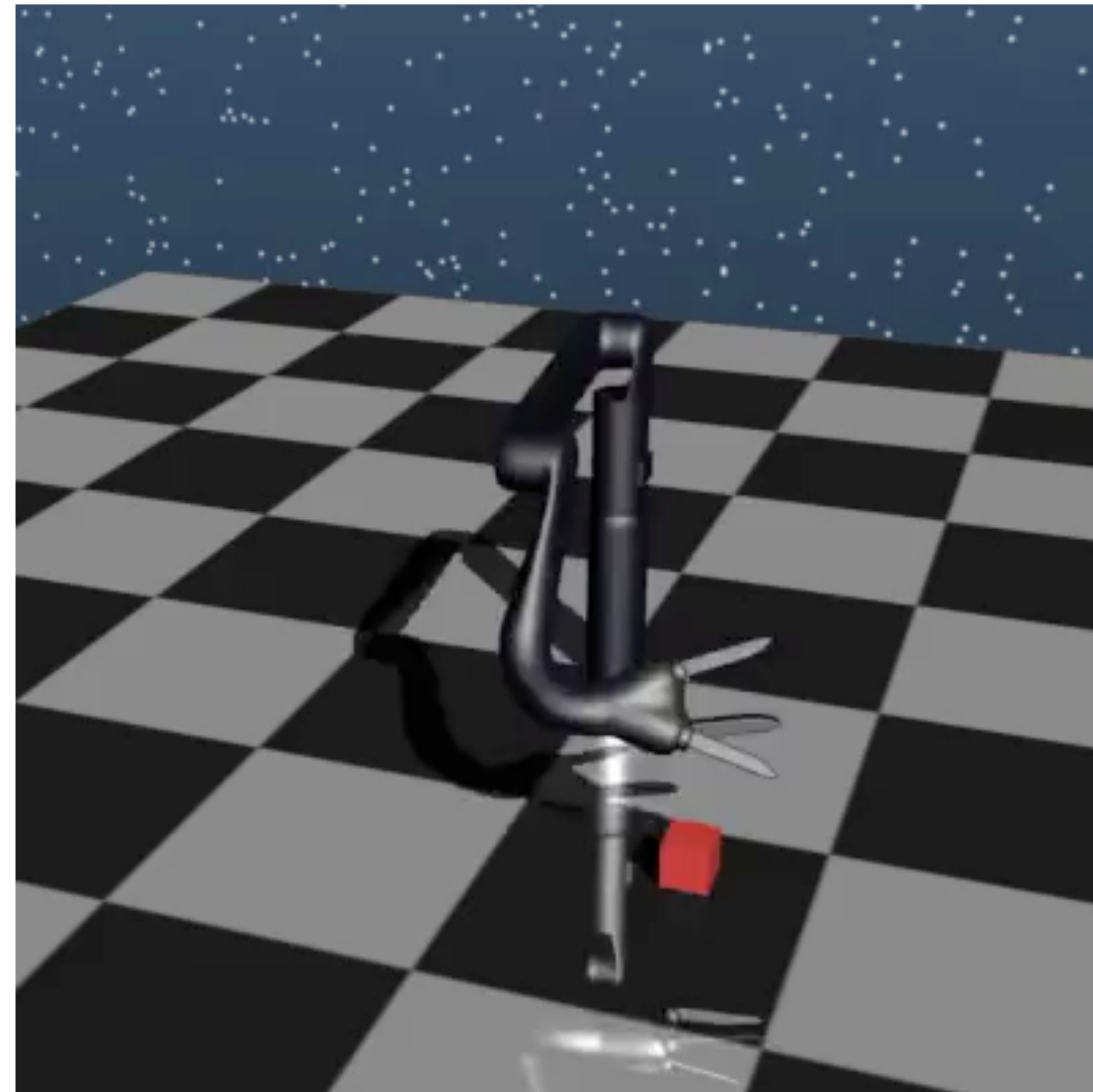


# Repetitive Pick

Pick

***Transition***

Pick

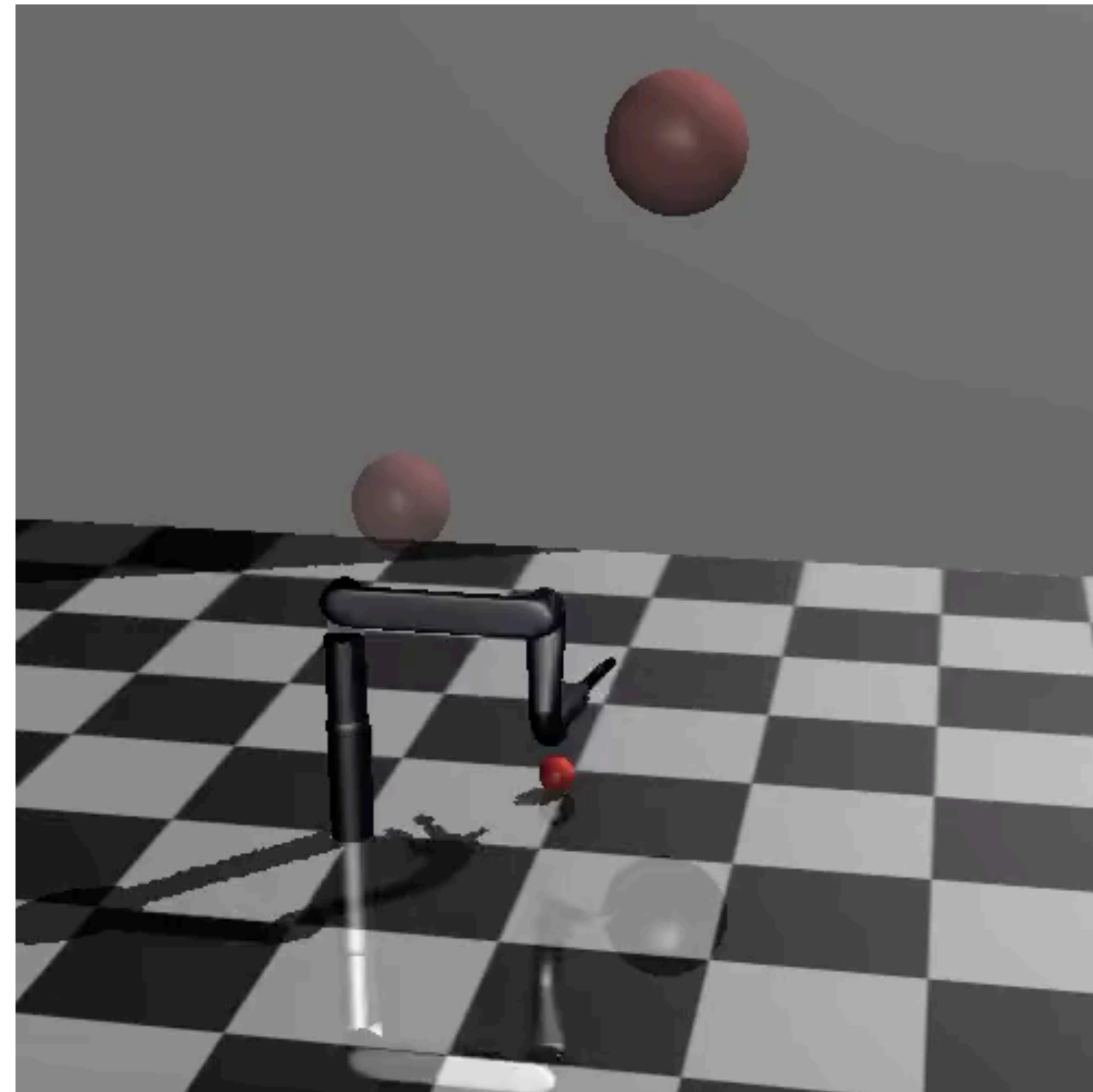


# Toss & Hit

Toss

*Transition*

Hit



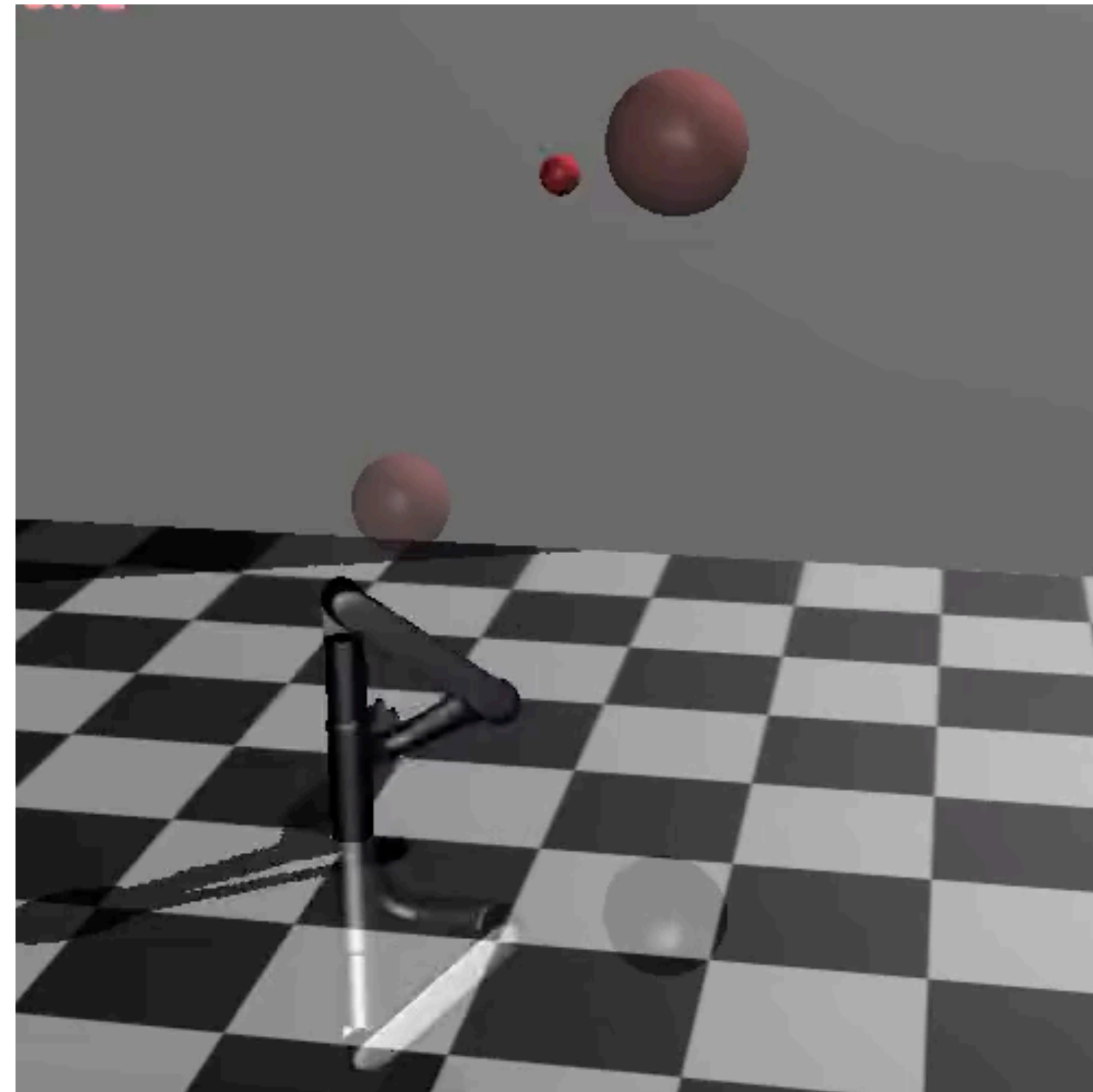


# Toss & Hit

Toss

*Transition*

Hit

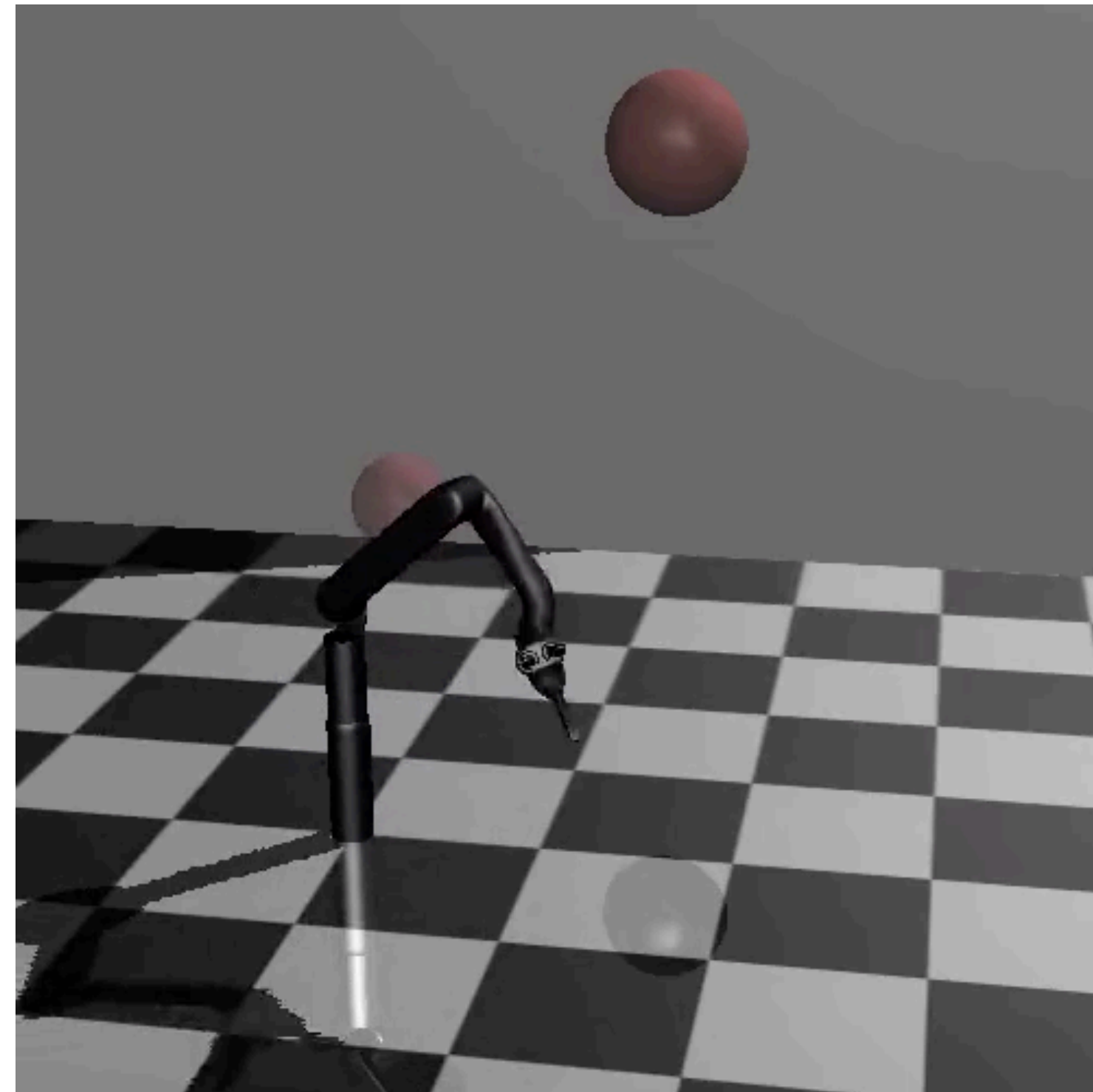


# Toss & Hit

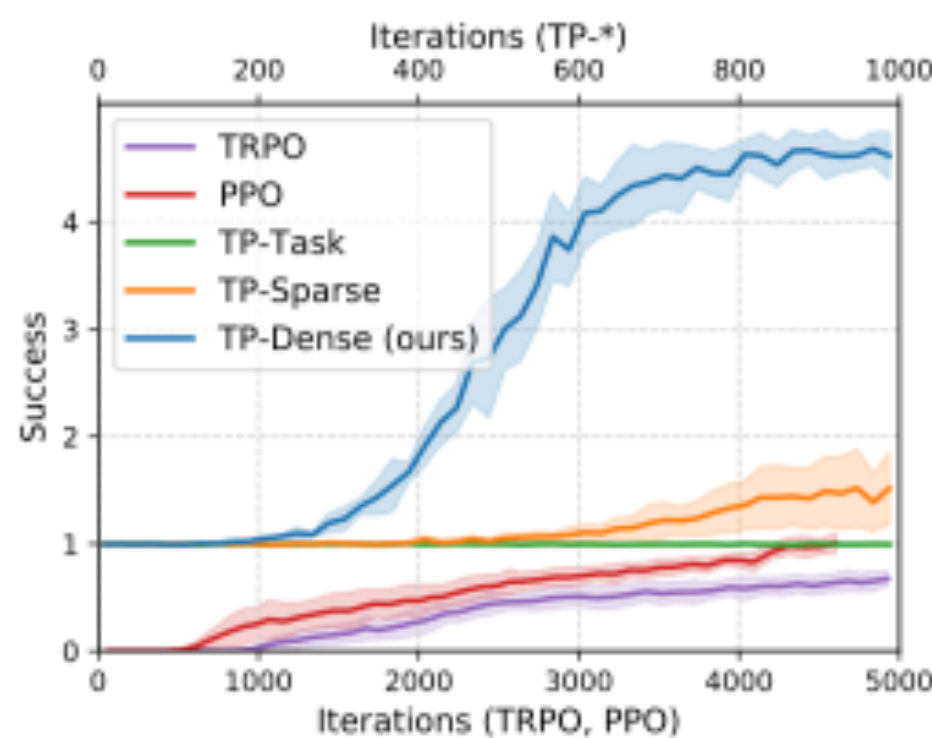
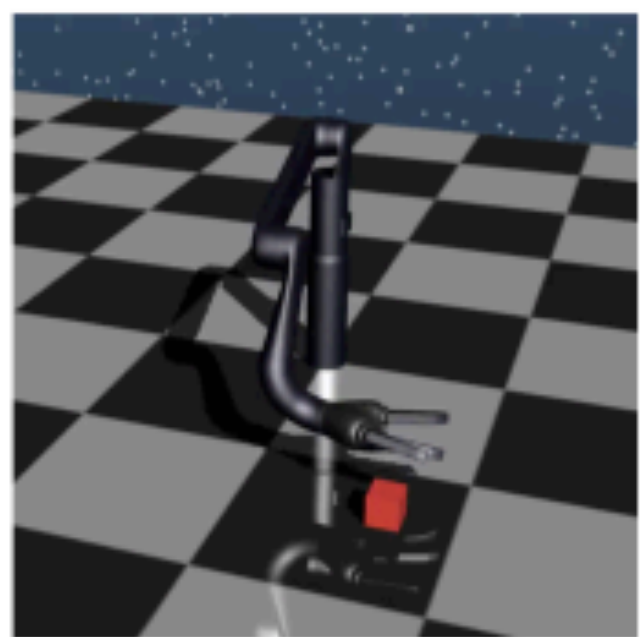
Toss

*Transition*

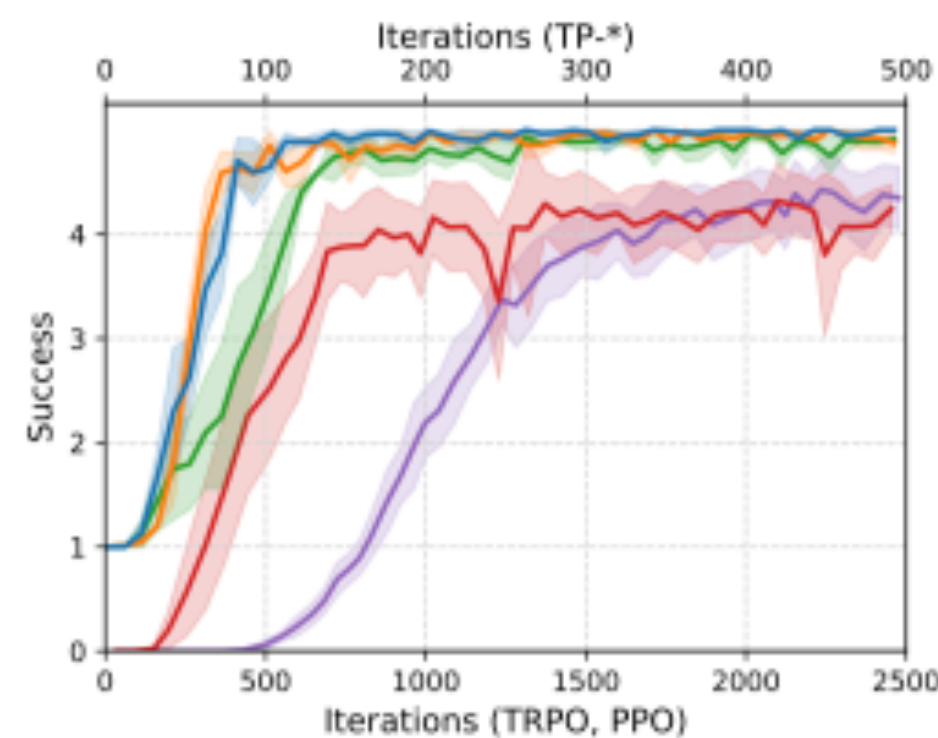
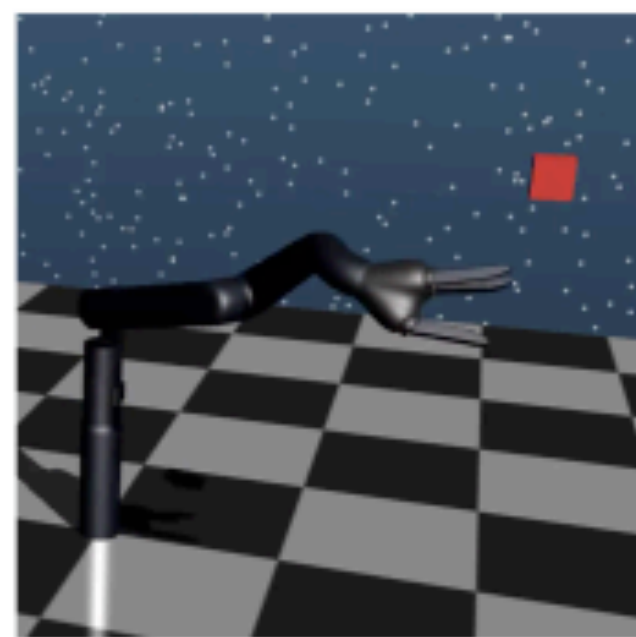
Hit



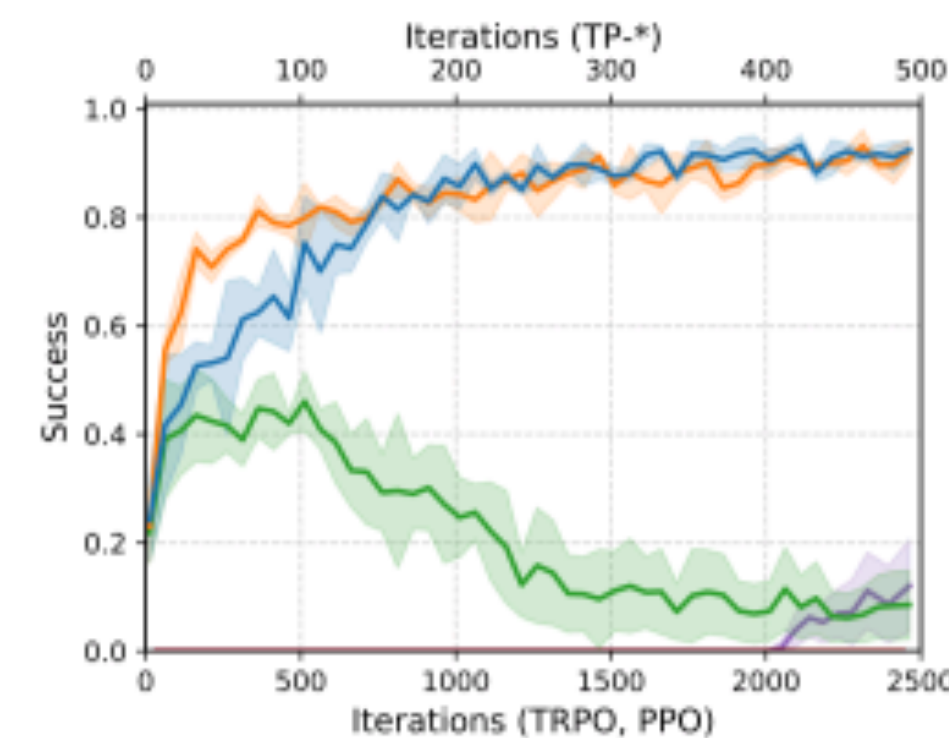
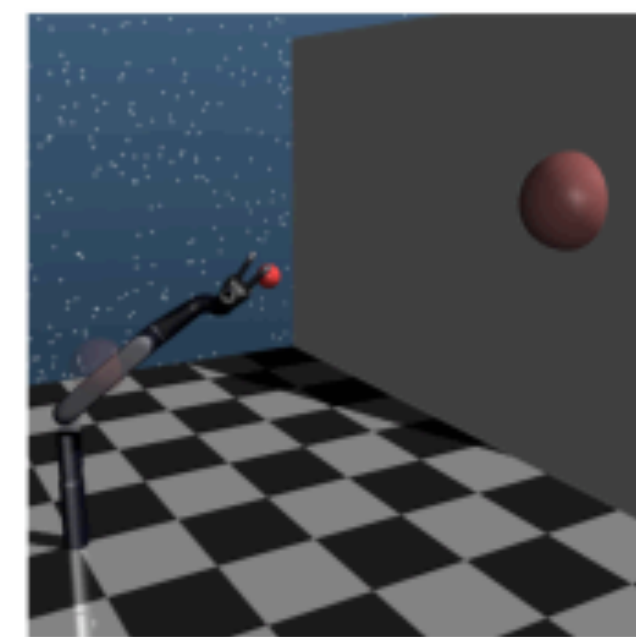
# Quantitative Results



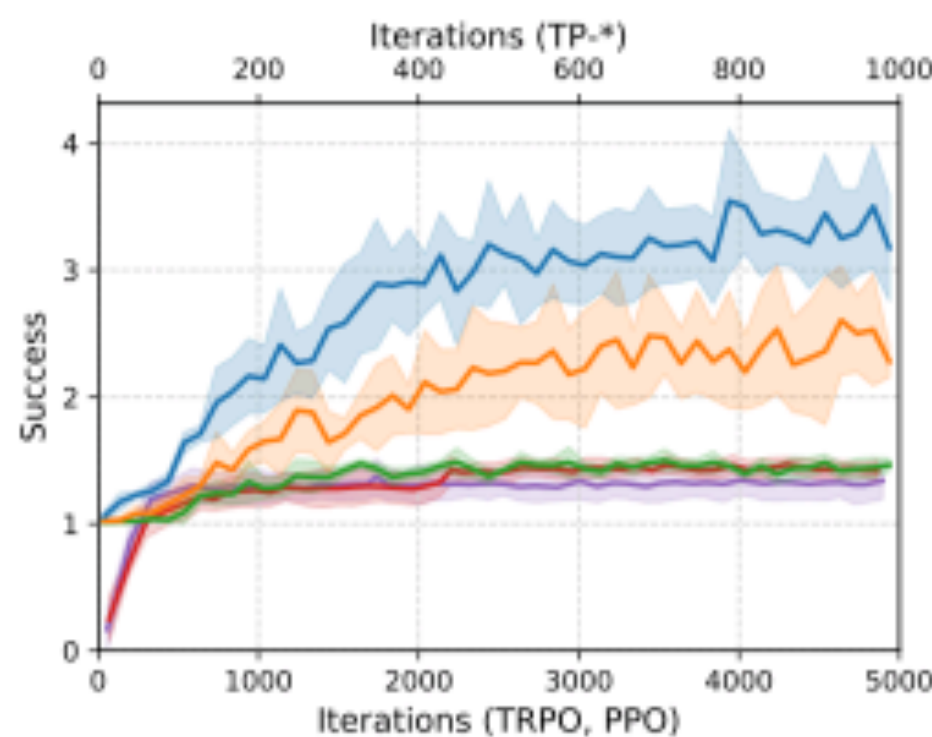
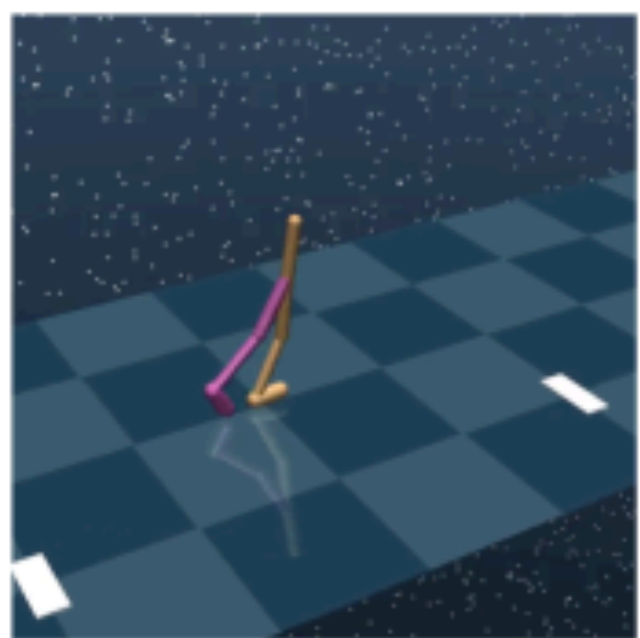
(a) Repetitive picking up



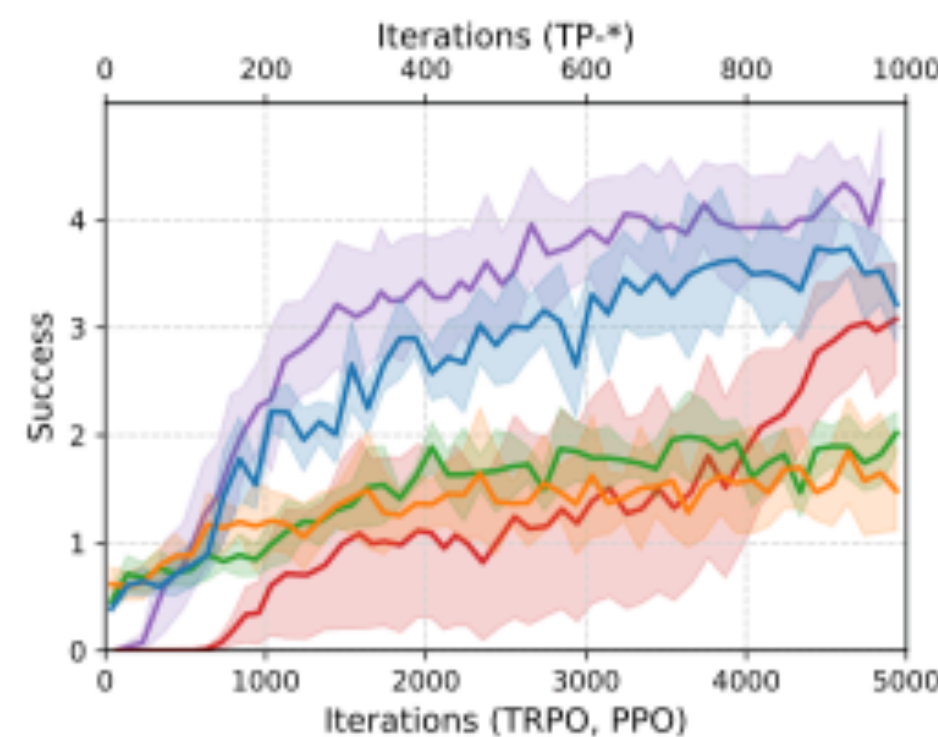
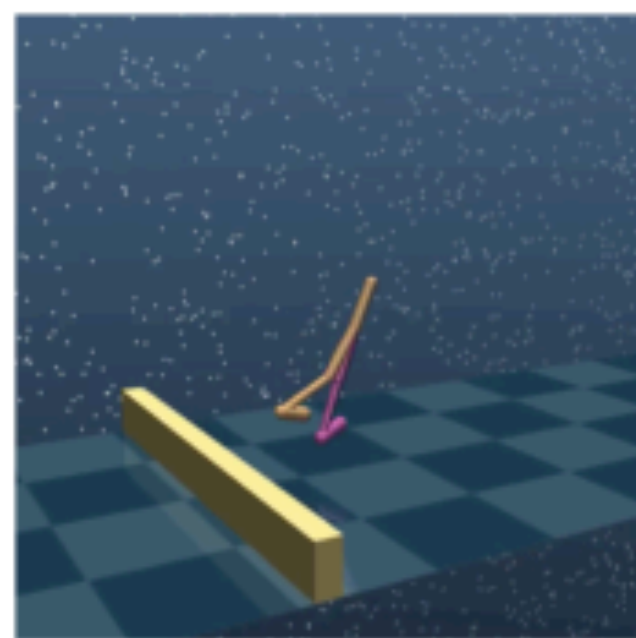
(b) Repetitive catching



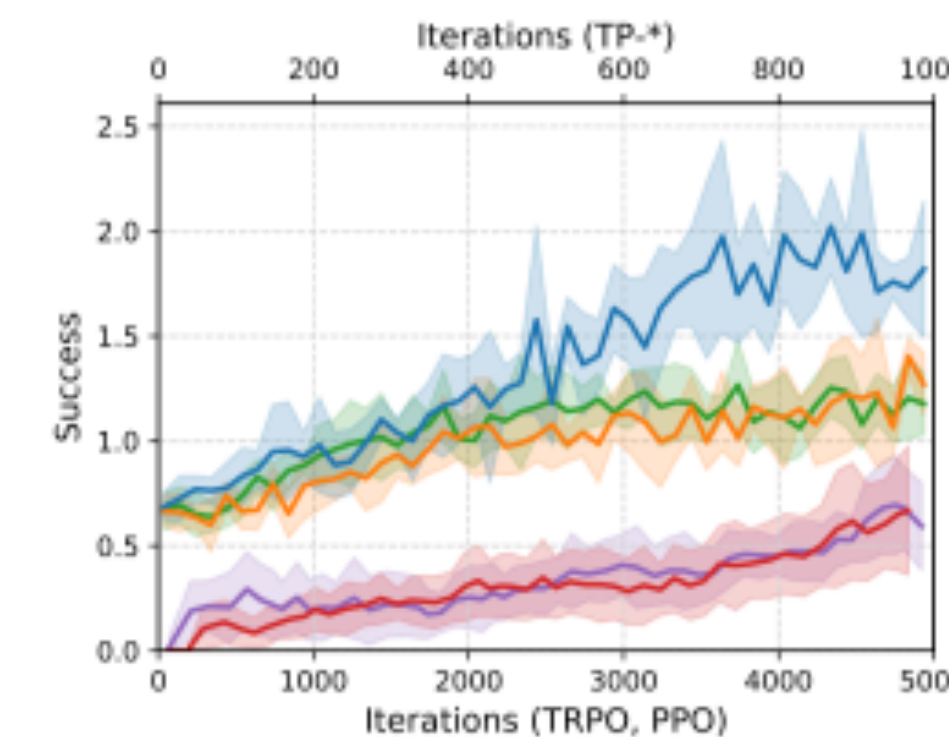
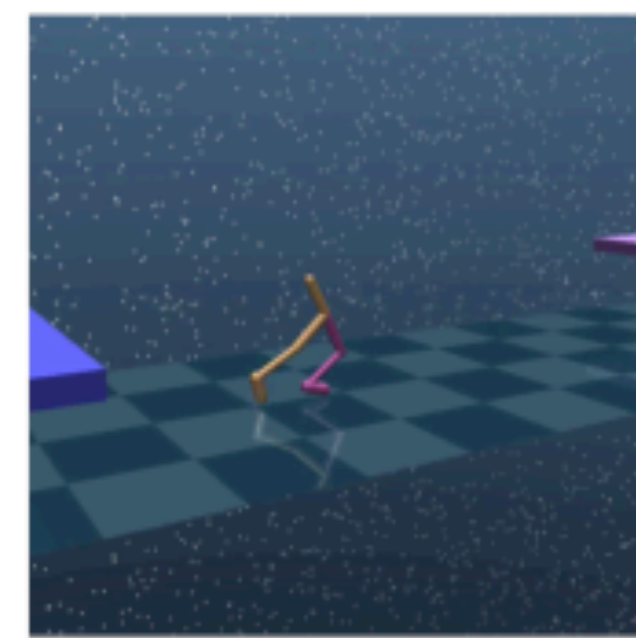
(c) Serve



(d) Patrol



(e) Hurdle



(f) Obstacle course



# Quantitative Results

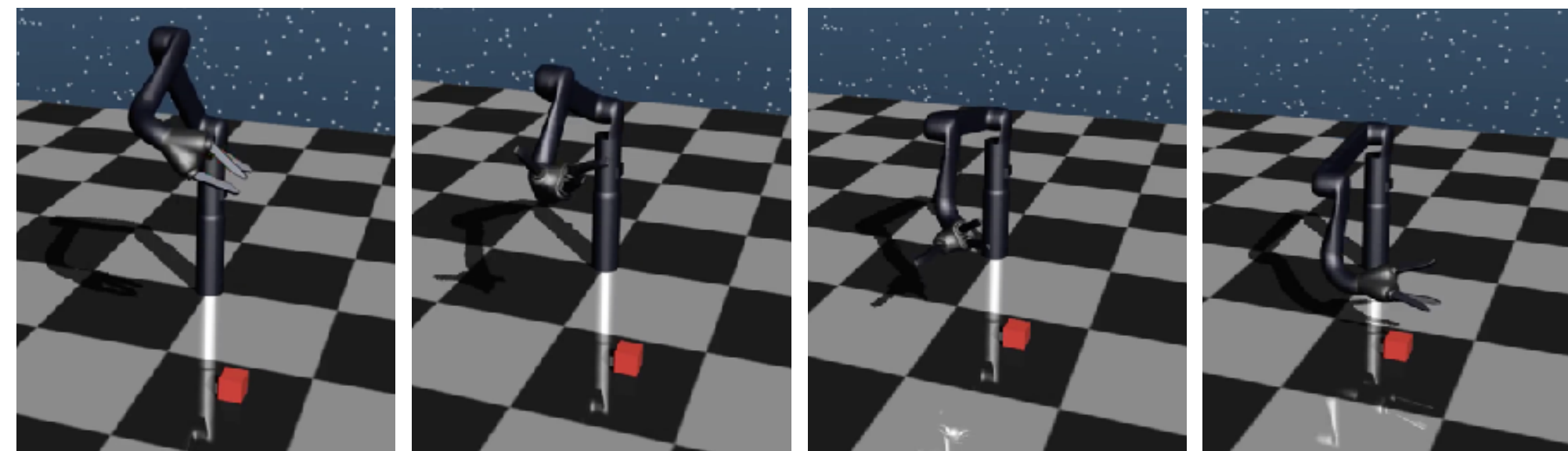
## Manipulation

	Reward	Repetitive picking up	Repetitive catching	Serve
TRPO	dense	$0.69 \pm 0.46$	$4.54 \pm 1.21$	$0.32 \pm 0.47$
PPO	dense	$0.95 \pm 0.53$	$4.26 \pm 1.63$	$0.00 \pm 0.00$
Without TP	sparse	$0.99 \pm 0.08$	$1.00 \pm 0.00$	$0.11 \pm 0.32$
TP-Task	sparse	$0.99 \pm 0.08$	$4.87 \pm 0.58$	$0.05 \pm 0.21$
TP-Sparse	sparse	$1.52 \pm 1.12$	$4.88 \pm 0.59$	<b><math>0.92 \pm 0.27</math></b>
TP-Dense (ours)	sparse	<b><math>4.84 \pm 0.63</math></b>	<b><math>4.97 \pm 0.33</math></b>	<b><math>0.92 \pm 0.27</math></b>

## Locomotion

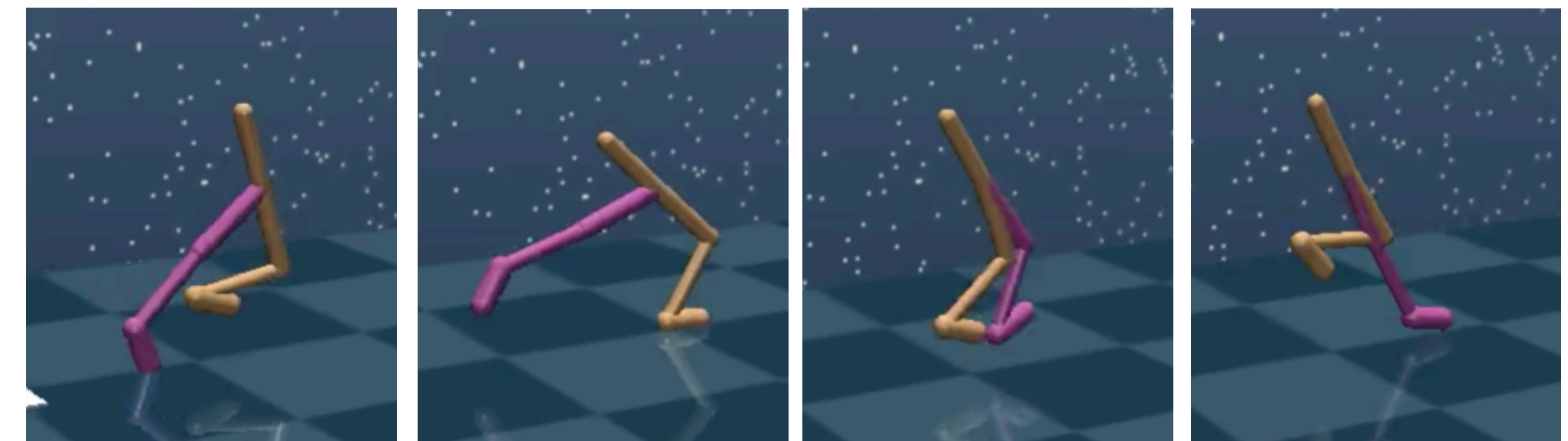
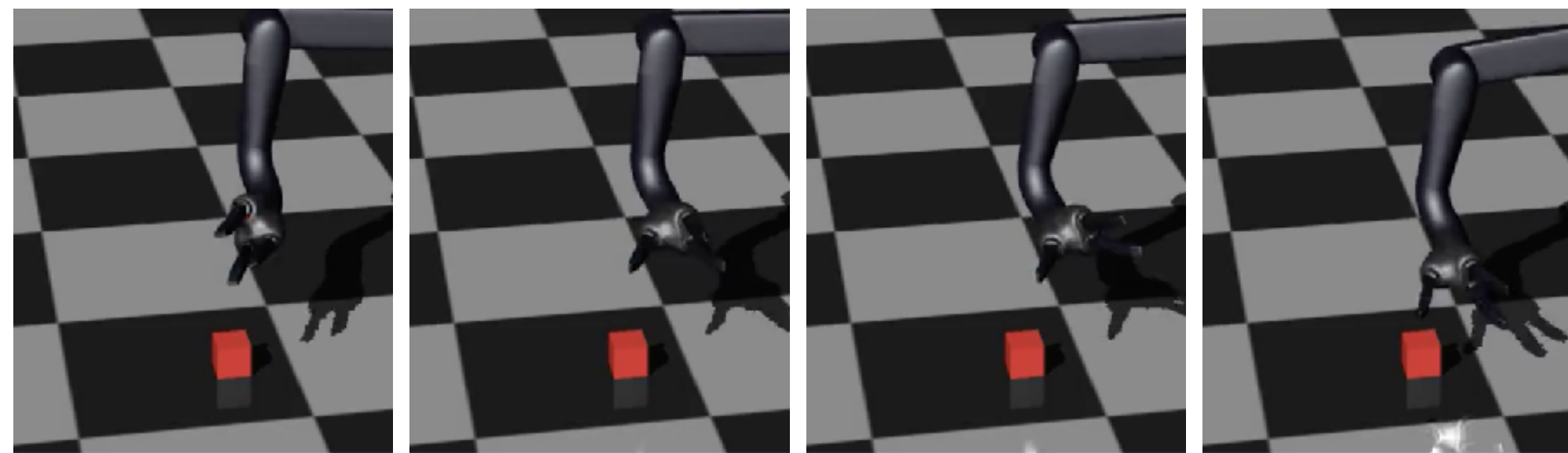
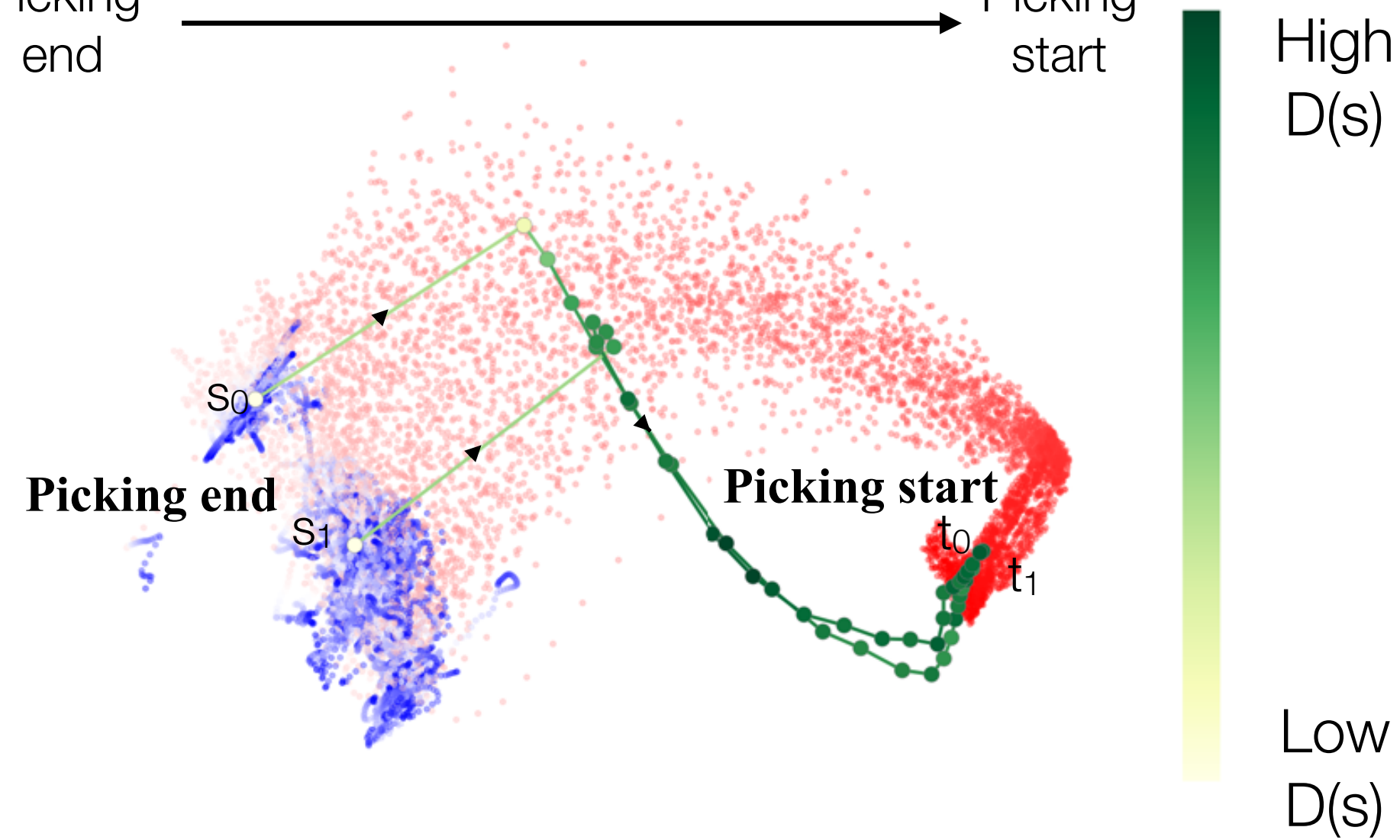
	Reward	Patrol	Hurdle	Obstacle course
TRPO	dense	$1.37 \pm 0.52$	<b><math>4.13 \pm 1.54</math></b>	$0.98 \pm 1.09$
PPO	dense	$1.53 \pm 0.53$	$2.87 \pm 1.92$	$0.85 \pm 1.07$
Without TP	sparse	$1.02 \pm 0.14$	$0.49 \pm 0.75$	$0.72 \pm 0.72$
TP-Task	sparse	$1.69 \pm 0.63$	$1.73 \pm 1.28$	$1.08 \pm 0.78$
TP-Sparse	sparse	$2.51 \pm 1.26$	$1.47 \pm 1.53$	$1.32 \pm 0.99$
TP-Dense (Ours)	sparse	<b><math>3.33 \pm 1.38</math></b>	<b><math>3.14 \pm 1.69^*</math></b>	<b><math>1.90 \pm 1.45</math></b>

# Transition Trajectories



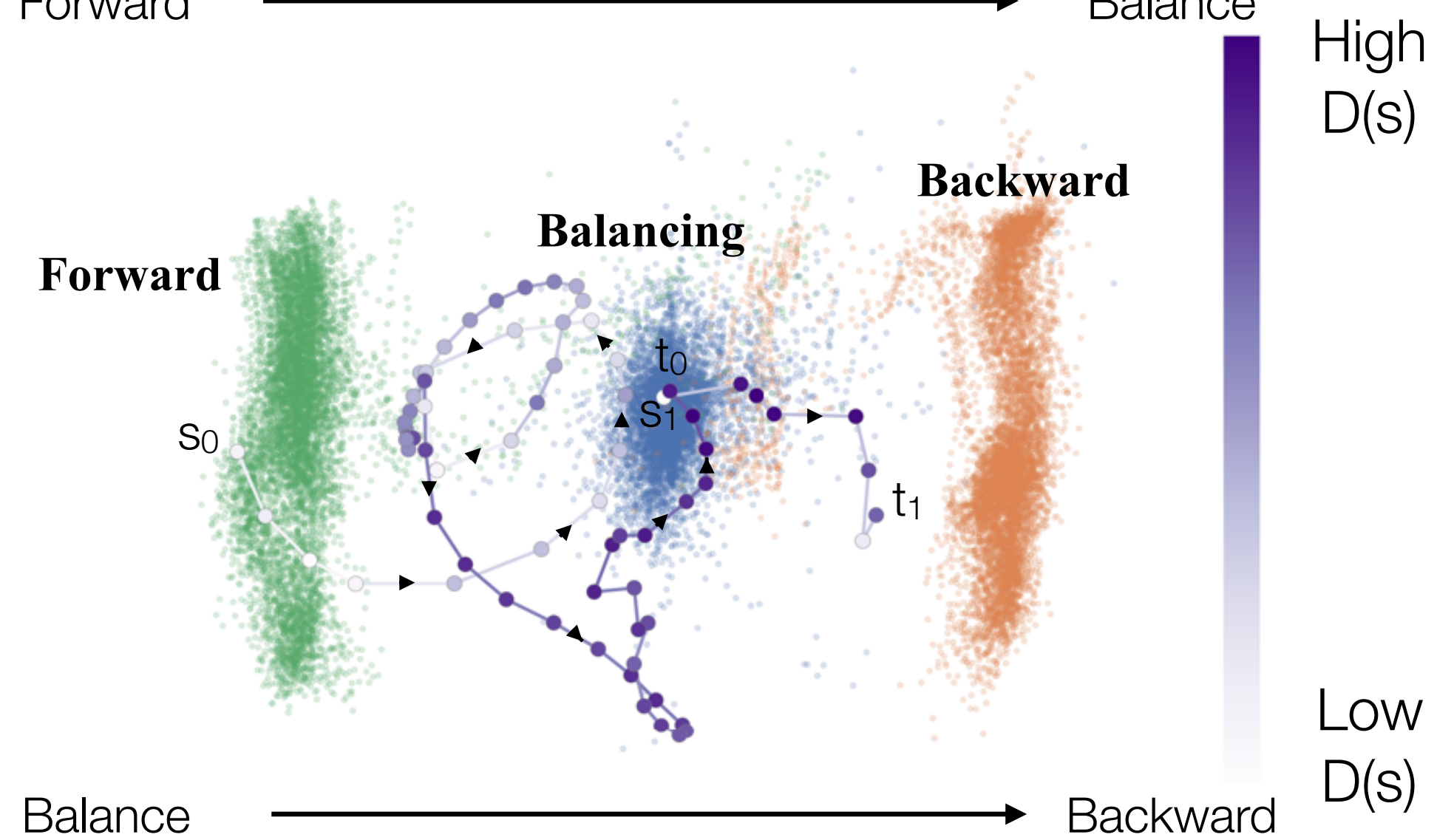
Picking  
end

Picking  
start



Forward

Balance



Balance

Backward



# Summary

We propose to **reuse skills** to compose **complex, long-horizon tasks**.

Naive execution of skills fail since the skills never learned to connect.

**Transition policies** learn to smoothly connect skills.

**Proximity predictors** provide dense reward for efficient training of transition policies.