

Self-Supervised Policy Adaptation during Deployment

Nicklas Hansen, Rishabh Jangir, Yu Sun, Guillem Alenyà,
Pieter Abbeel, Alexei A. Efros, Lerrel Pinto, Xiaolong Wang

ICLR 2021

UC San Diego



Berkeley
UNIVERSITY OF CALIFORNIA

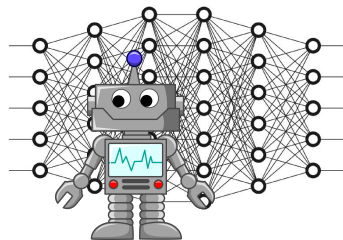
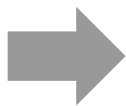


Institut de Robòtica
i Informàtica Industrial

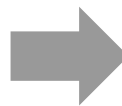




Observation

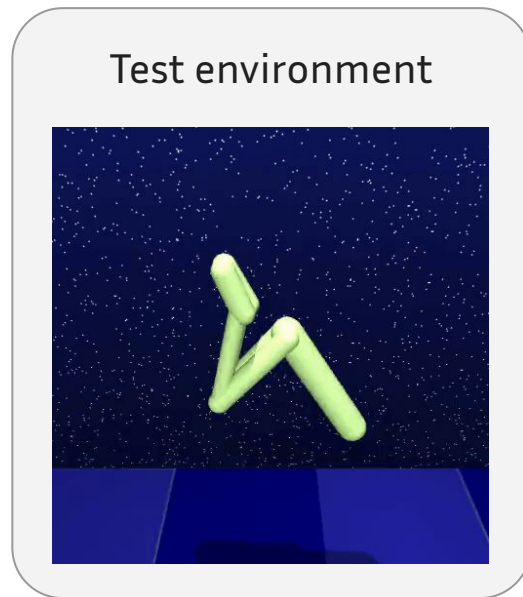
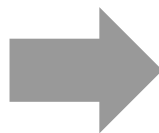
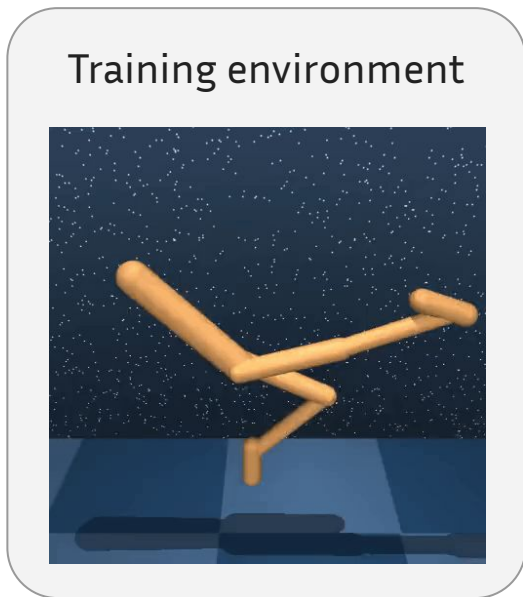


Learned policy

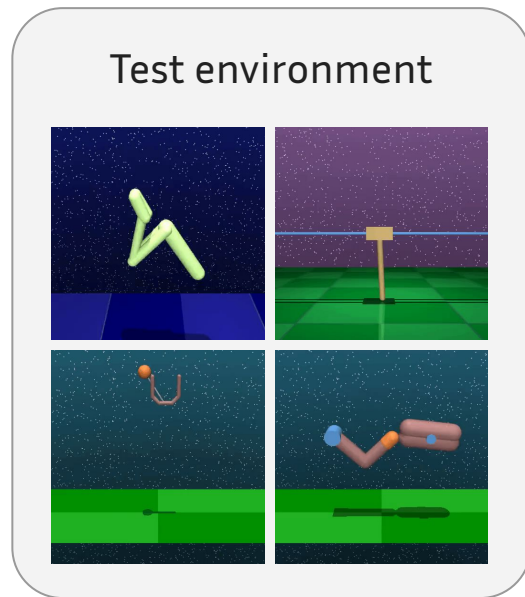
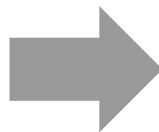
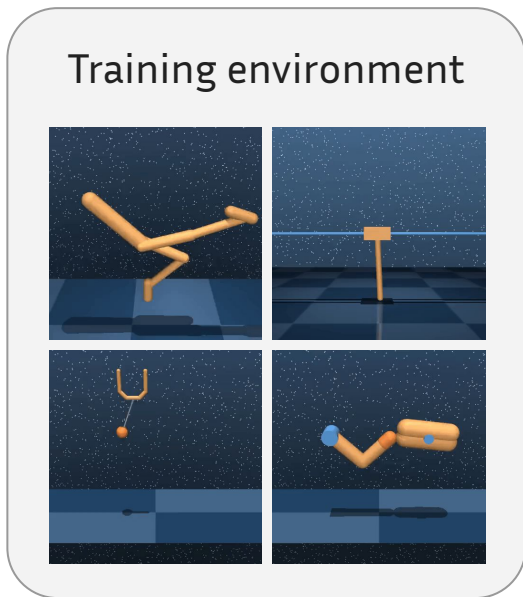


Environment interaction

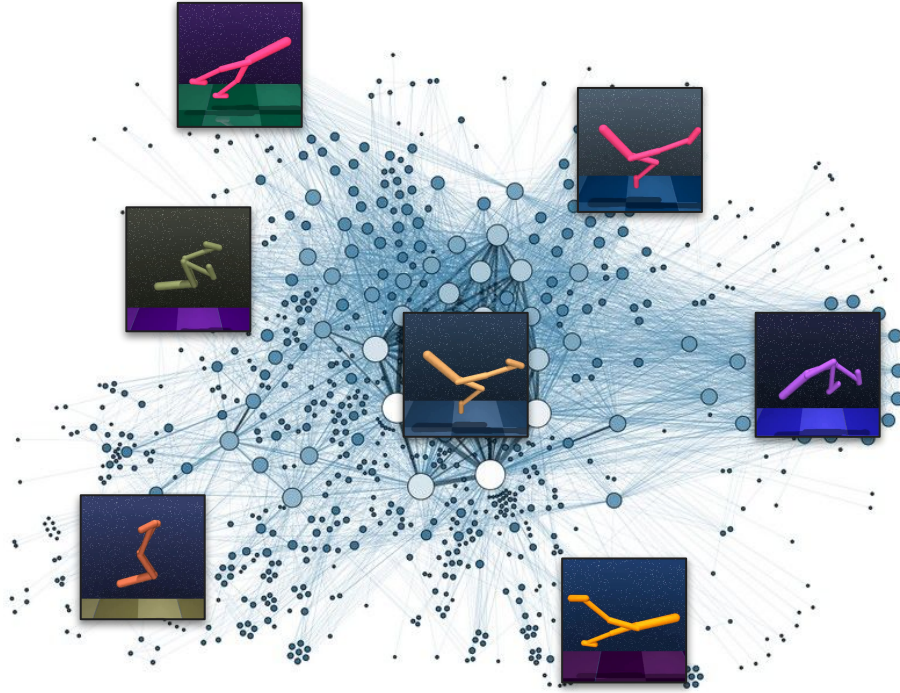
Complex tasks can be solved by **end-to-end** policy learning



Generalization across different environments is **hard**

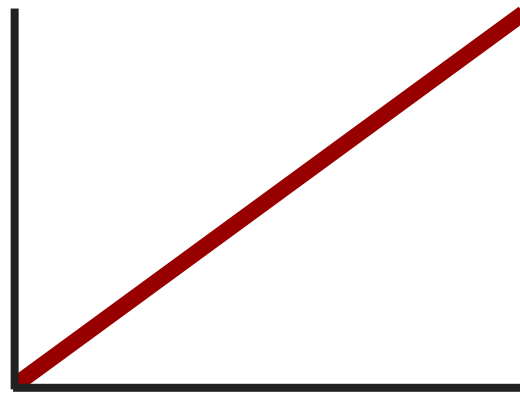


Generalization across different environments is **hard**



We can **randomize** elements that we expect to vary at test-time

Training difficulty



Size of training distribution

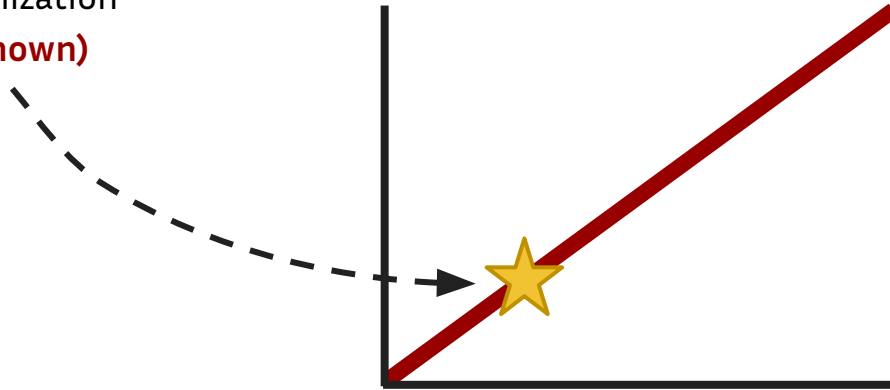
Factors of variation

- Color
- Texture
- Camera pose
- Mass
- Friction
- etc...

Larger training distribution = **harder** problem to solve

optimal level of
randomization
(unknown)

Training difficulty

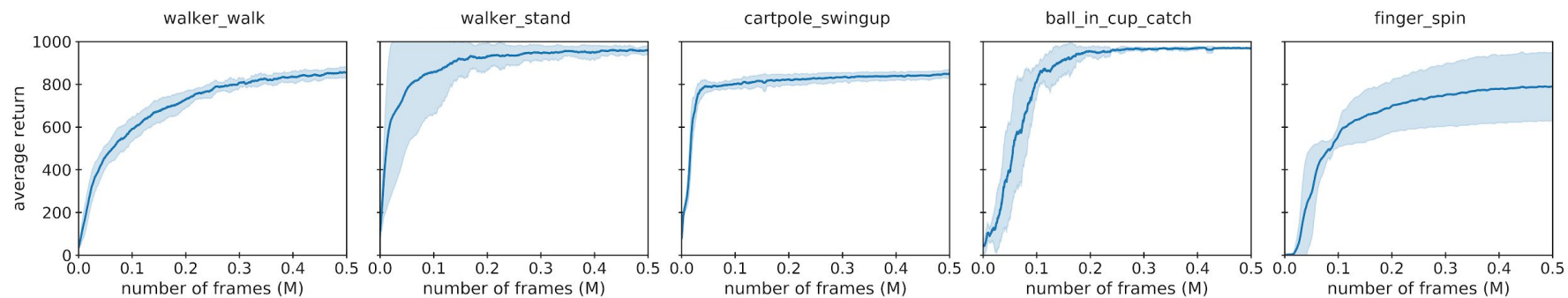


Size of training distribution

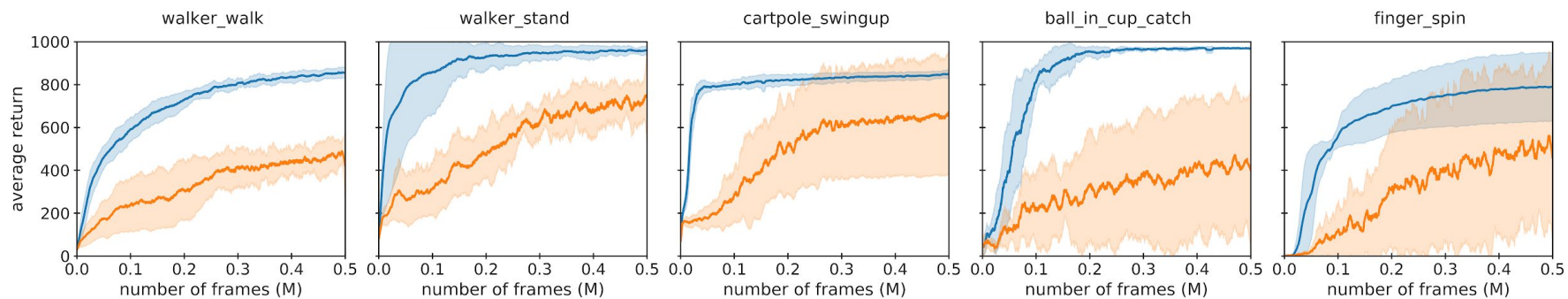
Factors of variation

- Color
- Texture
- Camera pose
- Mass
- Friction
- etc...

Larger training distribution = **harder** problem to solve

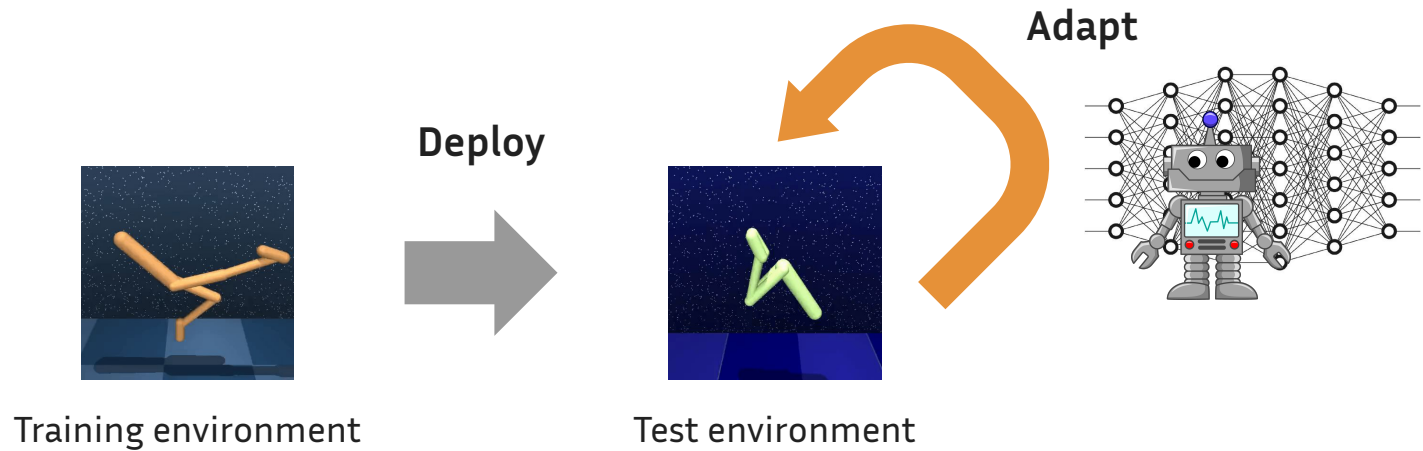


— Single environment



— Single environment

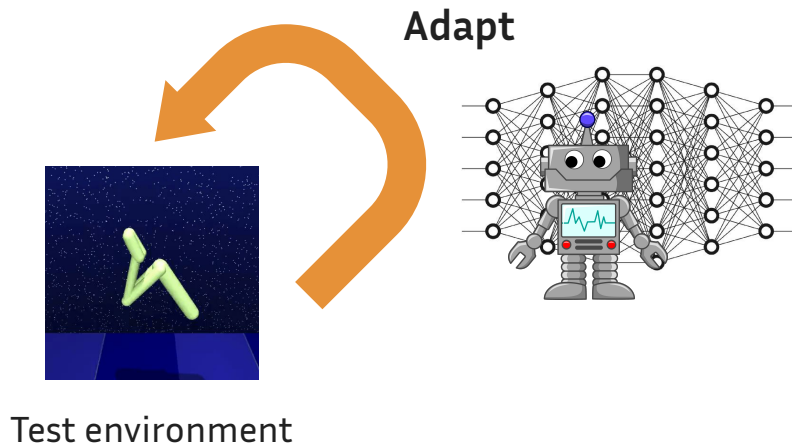
— Domain randomization



Can we instead **adapt** to new environments **during deployment**?

Challenges:

- No data prior to deployment
- Potentially no reward signal



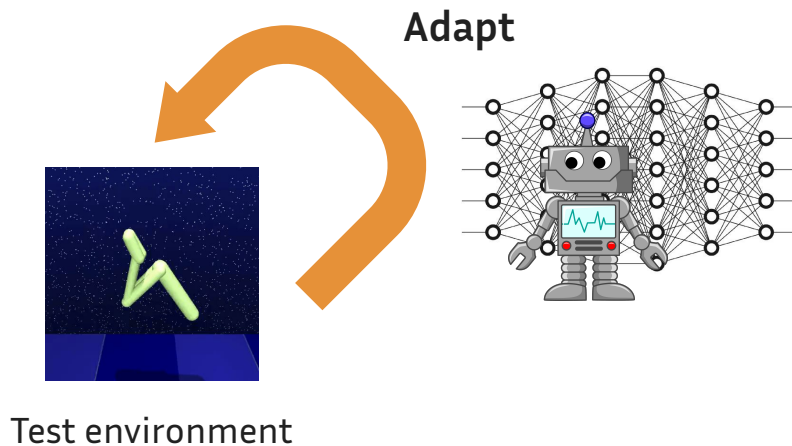
Can we instead **adapt** to new environments **during deployment**?

Challenges:

- No data prior to deployment
- Potentially no reward signal

Solution:

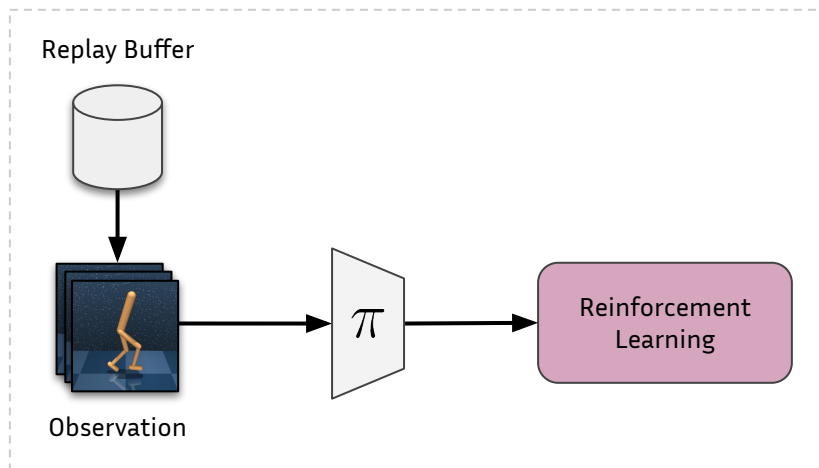
- **Online** adaptation
- **Self-supervised** training signal



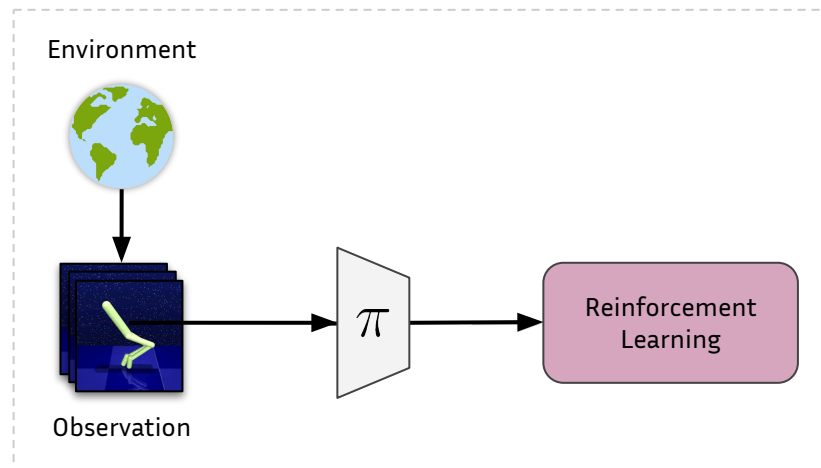
Can we instead **adapt** to new environments **during deployment**?

Algorithmic overview

Training

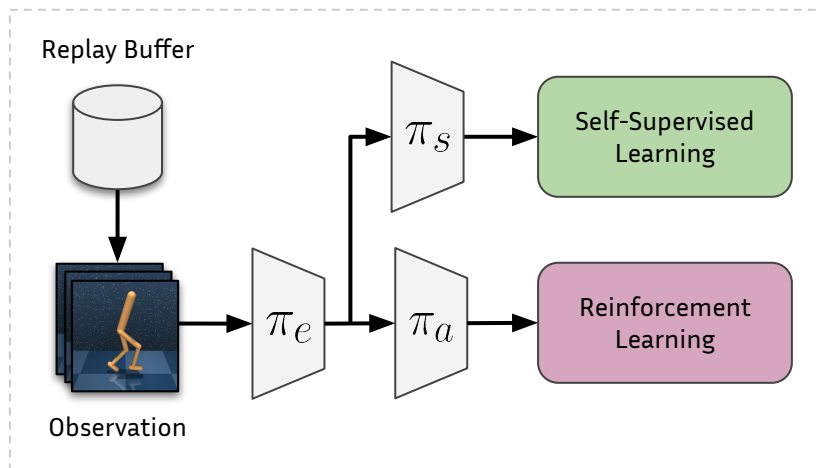


Policy Adaptation during Deployment

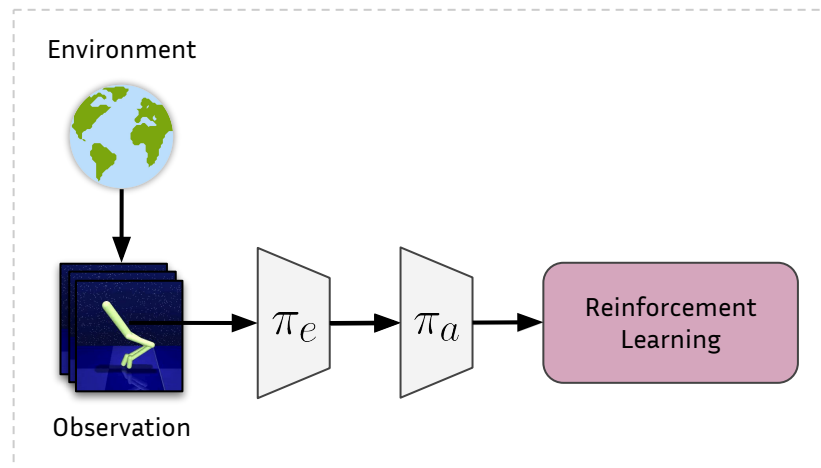


Algorithmic overview

Training

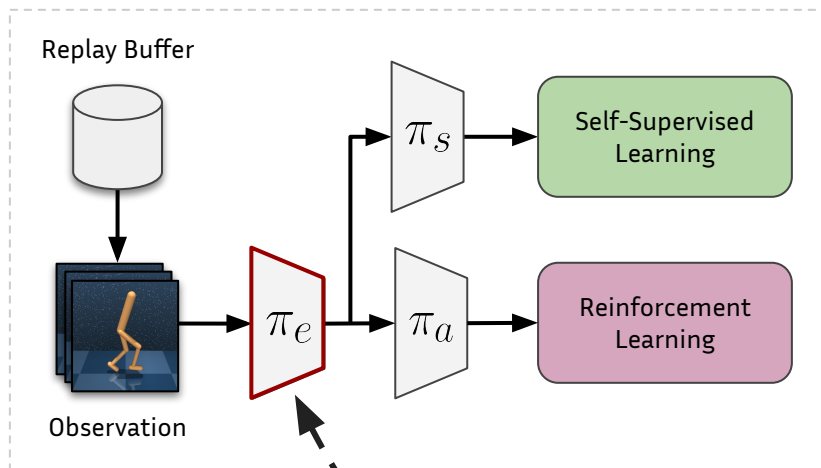


Policy Adaptation during Deployment

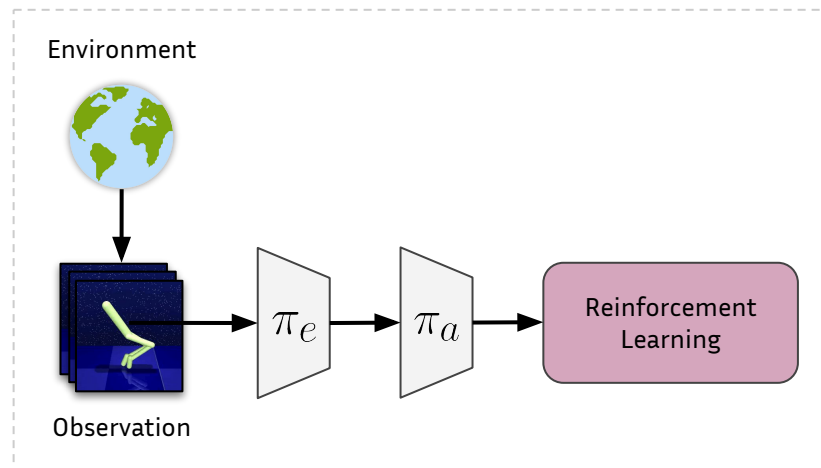


Algorithmic overview

Training



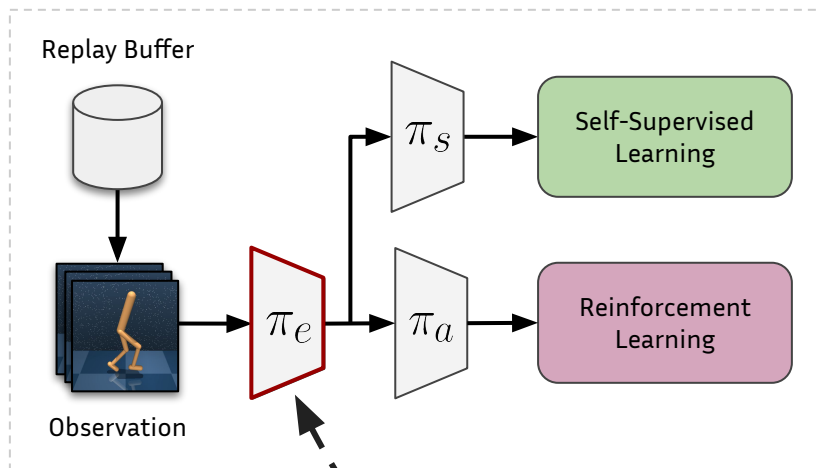
Policy Adaptation during Deployment



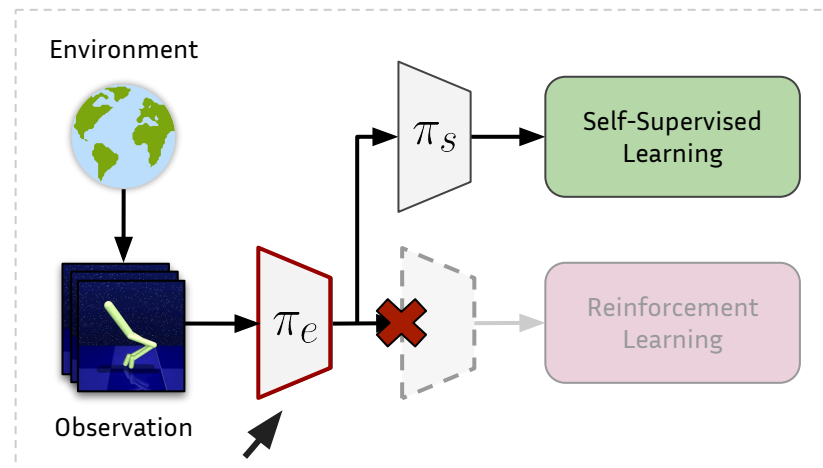
shared encoder

Algorithmic overview

Training



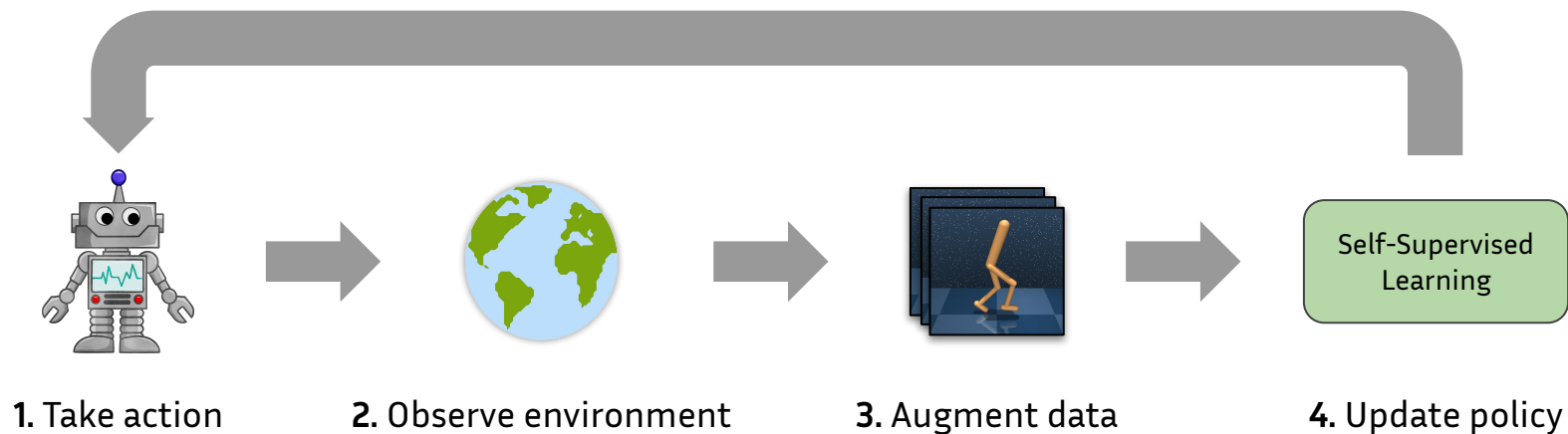
Policy Adaptation during Deployment



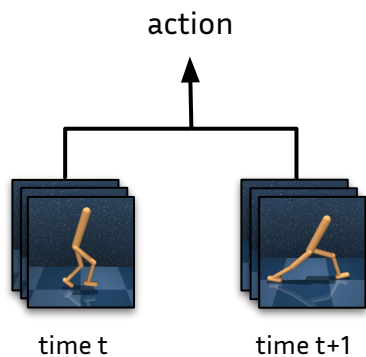
shared encoder

Algorithmic overview

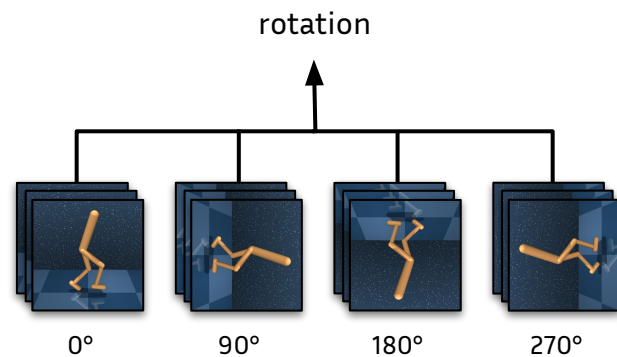
next step



Algorithmic overview

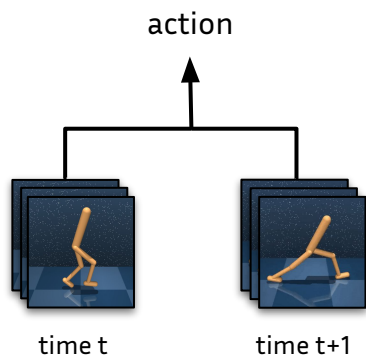


Inverse dynamics model

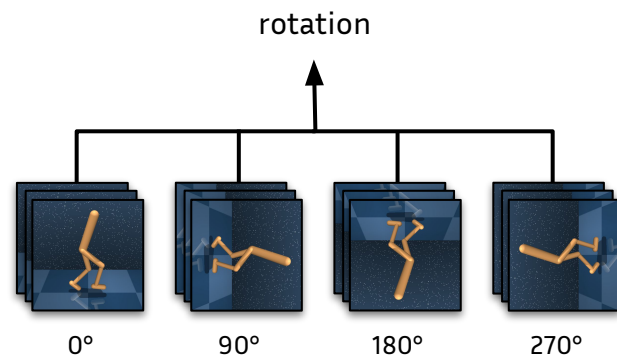


Rotation prediction

Algorithmic overview



Inverse dynamics model

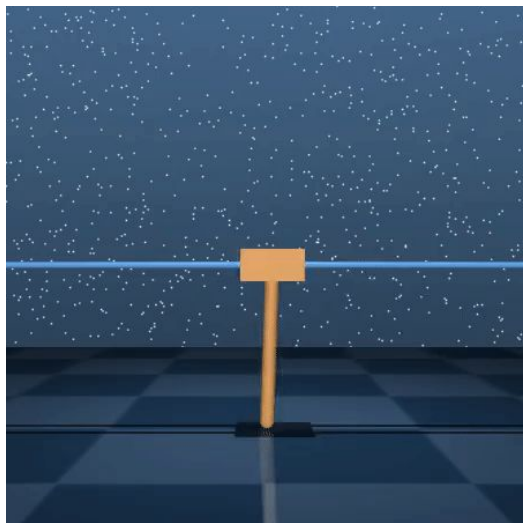


Rotation prediction

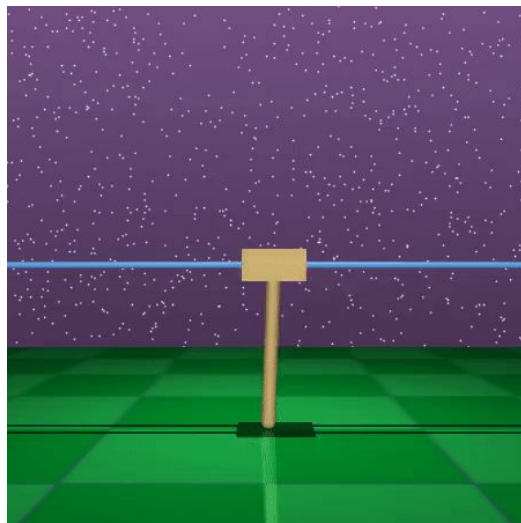
Choice of self-supervised task is **task-dependent**



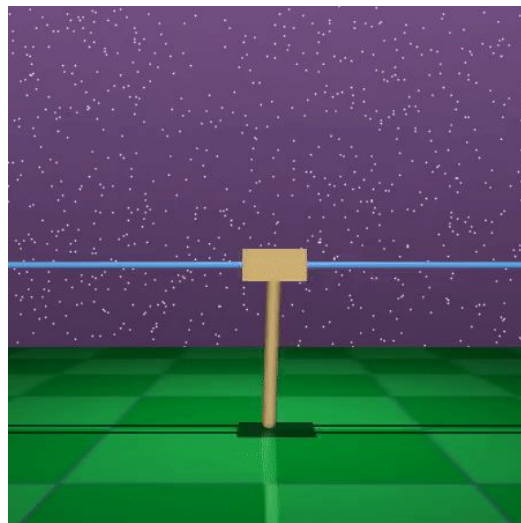
We evaluate generalization to **100+** environments



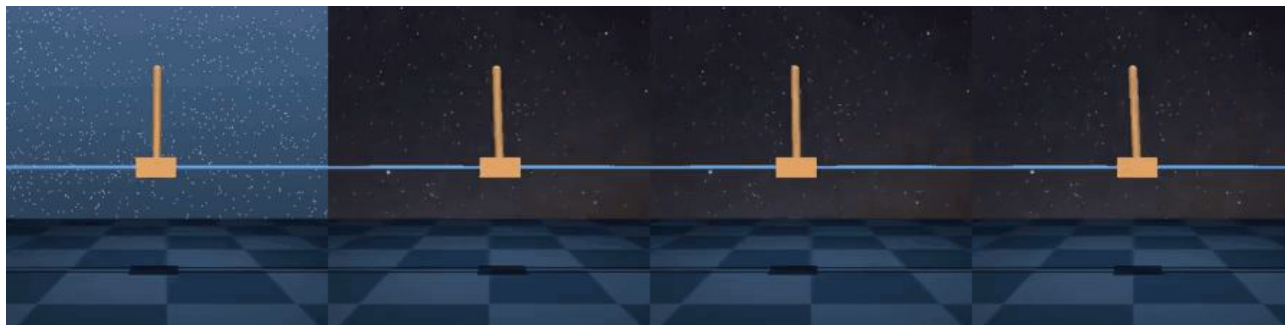
Training environment



Direct transfer



PAD (ours)



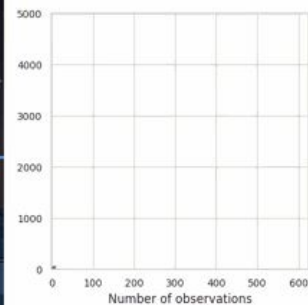
Training environment

Direct transfer

CURL

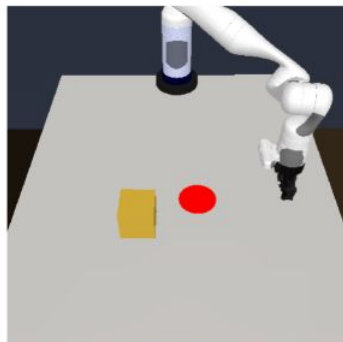
PAD (ours)

Cumulative reward

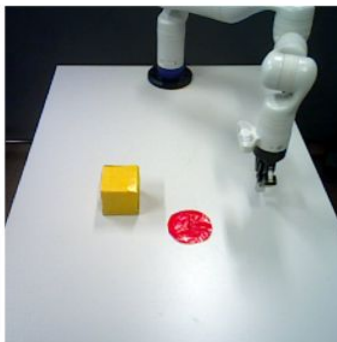




Adapting to the real world

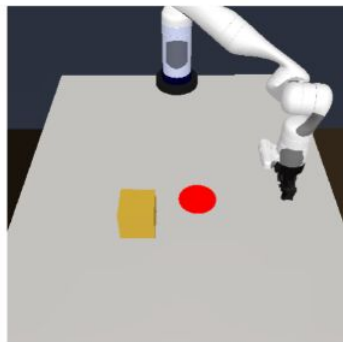


Simulation

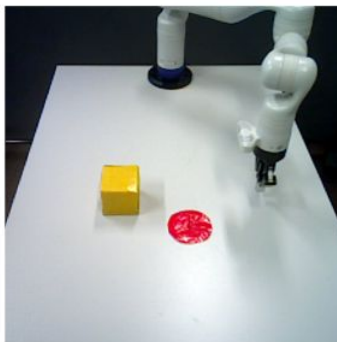


Default transfer

Adapting to the real world



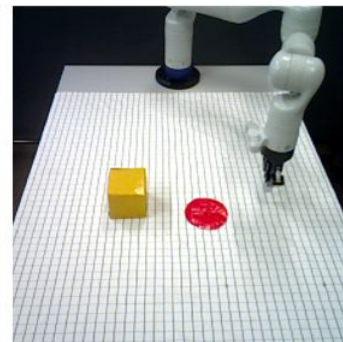
Simulation



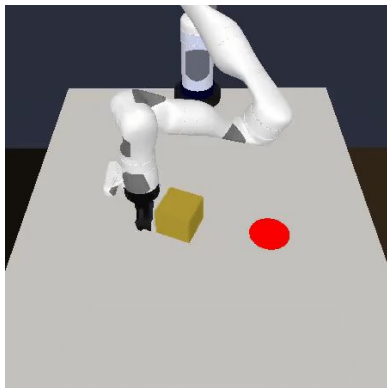
Default transfer



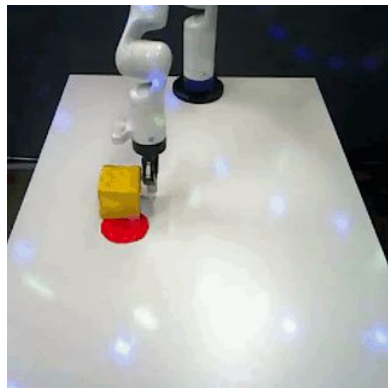
Disco lights



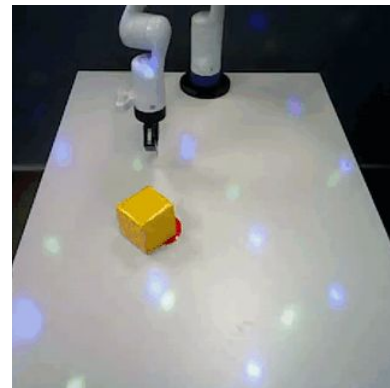
Tablecloth



Training environment

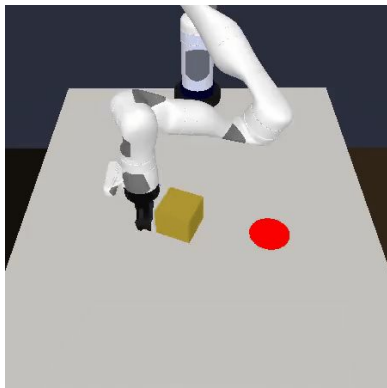


Direct transfer

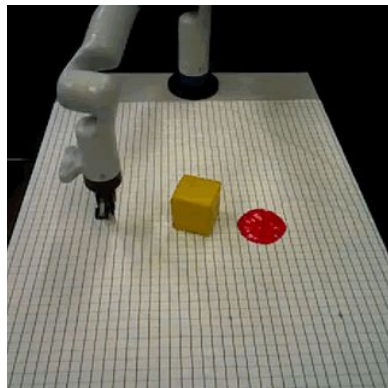


PAD (ours)

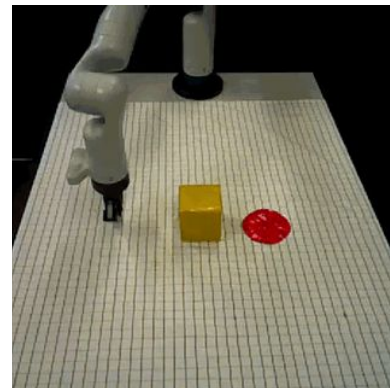
Our method improves generalization in **sim2real**



Training environment

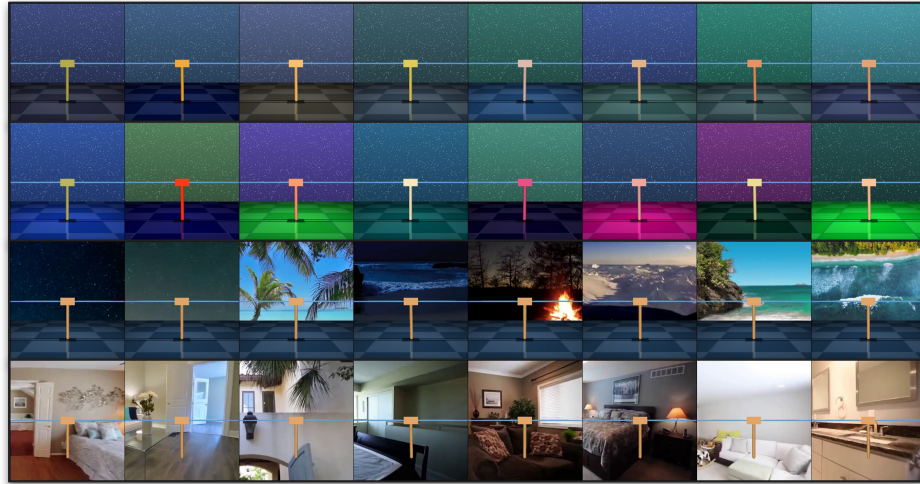


Direct transfer



PAD (ours)

Our method improves generalization in **sim2real**



For more information:

<https://nicklashansen.github.io/PAD>

Self-Supervised Policy Adaptation during Deployment

Nicklas Hansen, Rishabh Jangir, Yu Sun, Guillem Alenyà,
Pieter Abbeel, Alexei A. Efros, Lerrel Pinto, Xiaolong Wang

ICLR 2021

UC San Diego



Berkeley
UNIVERSITY OF CALIFORNIA



Institut de Robòtica
i Informàtica Industrial



