Seminar in Deep Reinforcement Learning
# The Path to Continual Learning Curriculum Learning
Ramon Witschi, ETH Computer Science MSc, 19.05.2020

# What is a Curriculum?

O -1 +

− × ÷

∞ ⟺ □

$\partial$ x x² x³

√ ∫ ∇

< ∑

0 -1 +

− × ÷

x x² x³

√ < ∑

∞ ⇔ □

∂ ∫ ∇

$$0 \quad -1 \quad +$$

$$- \quad \times \quad \div$$

$$\Rightarrow$$

$$x \quad x^2 \quad x^3$$

$$\sqrt{} \quad < \quad \Sigma$$

$$\infty \quad \Leftrightarrow \quad \Box$$

$$\partial \quad \int \quad \nabla$$

# Curriculum over Training Data!

Curriculum Learning (ICML, 2009) – Bengio et al.

1   Start with simple examples

2   Gradually add more difficult ones

3   Arrive at target training distribution

Curriculum Learning (ICML, 2009) – Bengio et al.

1    Start with simple examples

2    Gradually add more difficult ones

3    Arrive at target training distribution

Curriculum Learning (ICML, 2009) – Bengio et al.

1     Start with simple examples

2     Gradually add more difficult ones

3     Arrive at target training distribution

Curriculum Learning (ICML, 2009) – Bengio et al.

# Empirical Results
# Faster Training & sometimes higher Test Scores

Curriculum Learning (ICML, 2009) – Bengio et al.

Faster Training proven on Linear Regression (Convex Optimization) ☺

Curriculum Learning by Transfer Learning: Theory and Experiments with Deep Networks – Weinshall, Cohen & Amir

# Curriculum Learning meets Reinforcement Learning

Intrinsic Motivation and Automatic Curricula via Asymmetric Self-Play

Sukhbaatar et al.

Reverse Curriculum Generation for Reinforcement Learning

Florensa et al.

Mix & Match – Agent Curricula for Reinforcement Learning

Czarnecki et al.

# Curriculum Learning meets Reinforcement Learning

Intrinsic Motivation and Automatic Curricula via Asymmetric Self-Play

Sukhbaatar et al.

Reverse Curriculum Generation for Reinforcement Learning

Florensa et al.

Mix & Match – Agent Curricula for Reinforcement Learning

Czarnecki et al.

# Model-Free Reinforcement Learning

## Sample Inefficient ☹

Jointly learn Environment and optimize for Reward

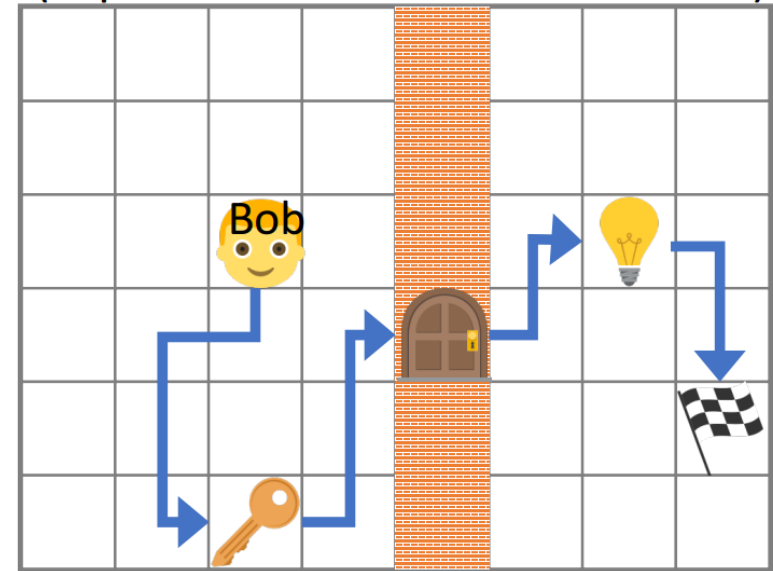"Unsupervised" Exploration!

# Framework



Self Play Episode (no supervision -- internal reward only)

Target Task Episode (supervision from external reward)
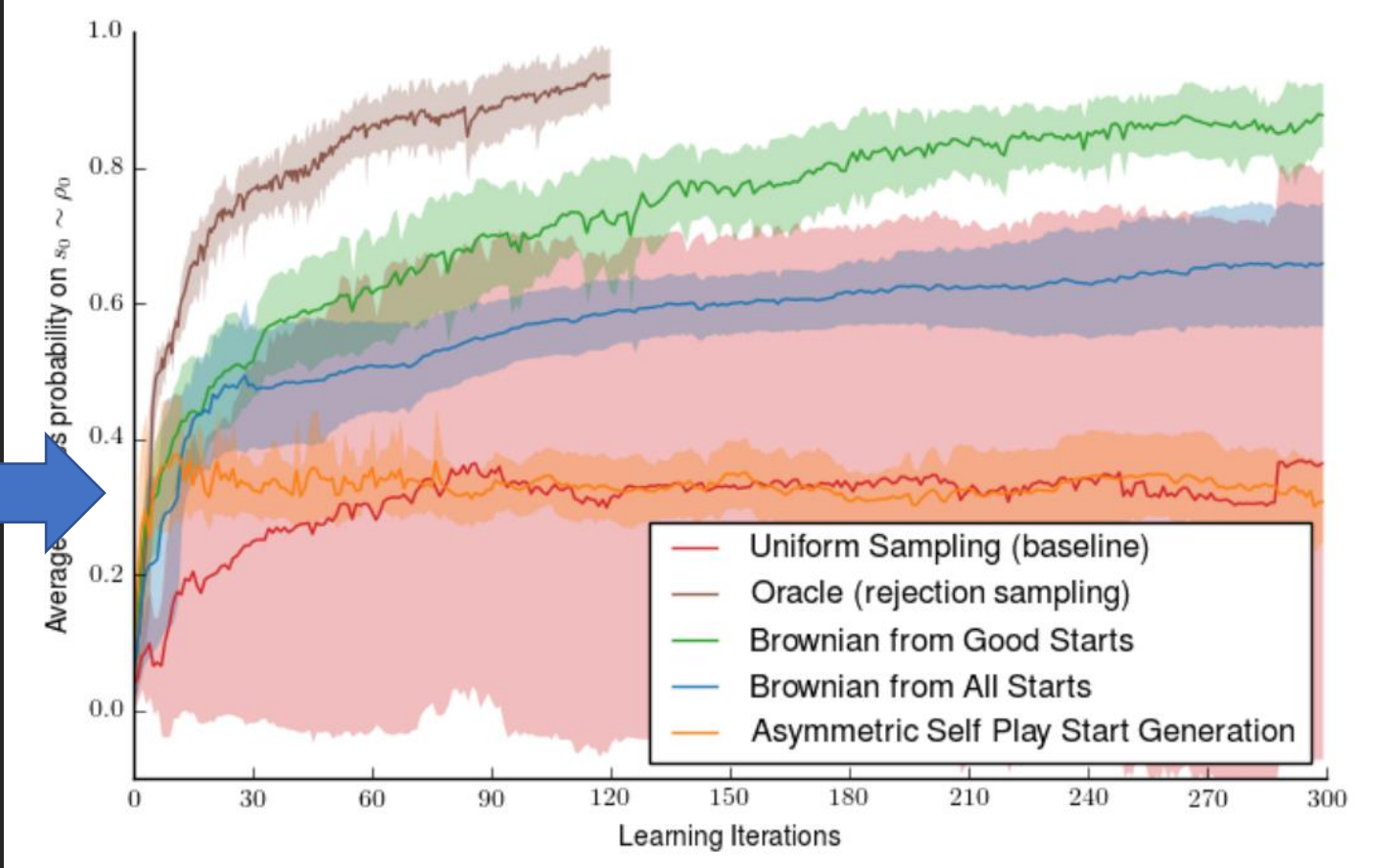
Alice's turn

Bob's turn

Bob applied to target task

# Internal Reward Structure

$$R_A = \max(0, t_B - t_A)$$

$$\xrightarrow{\text{\textit{Bob fast}} \atop \text{\textit{or}} \atop \textcolor{red}{\textbf{\textit{too}}} \textit{ slow}} 0 \quad ☹$$

# Internal Reward Structure
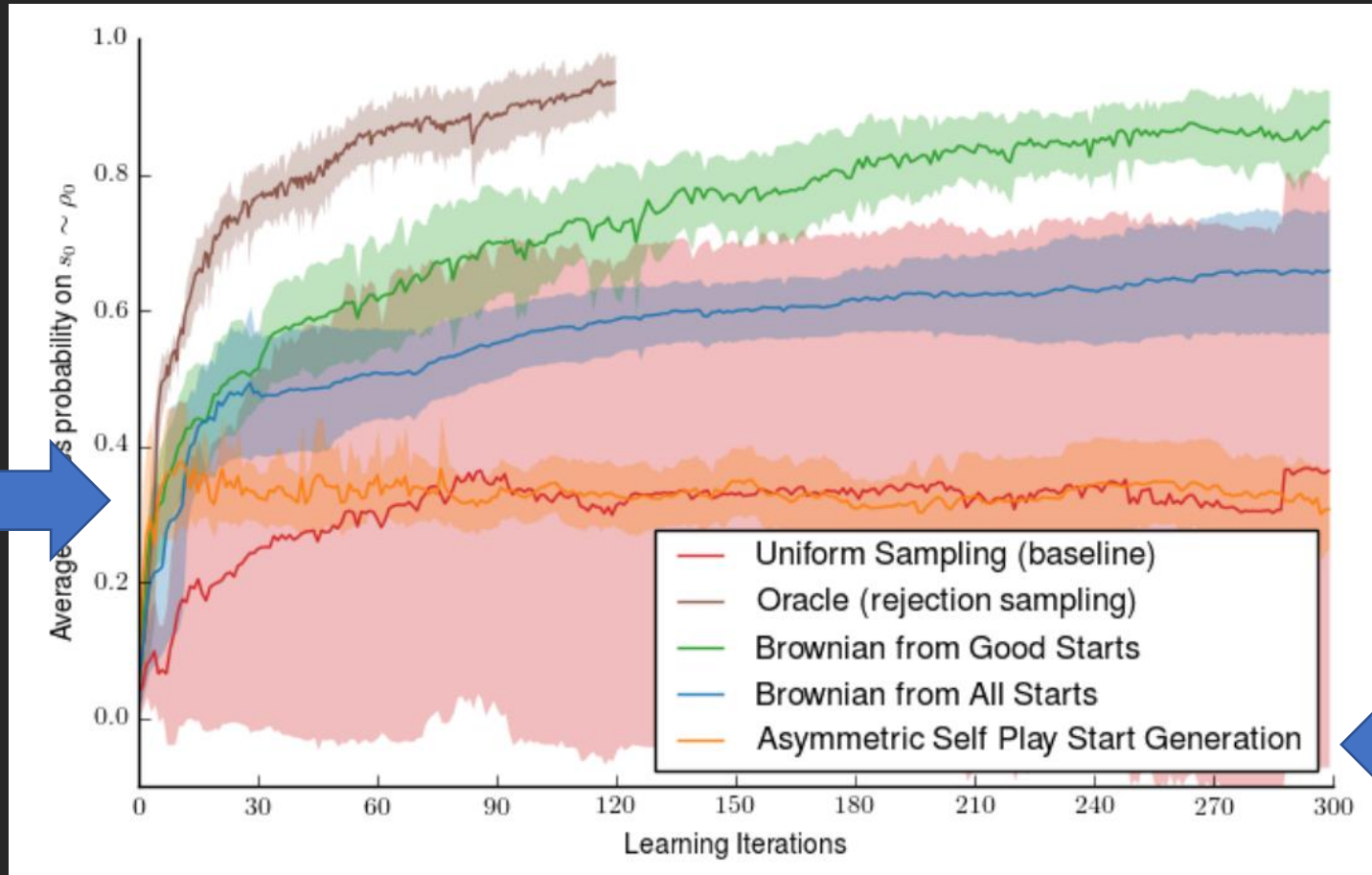


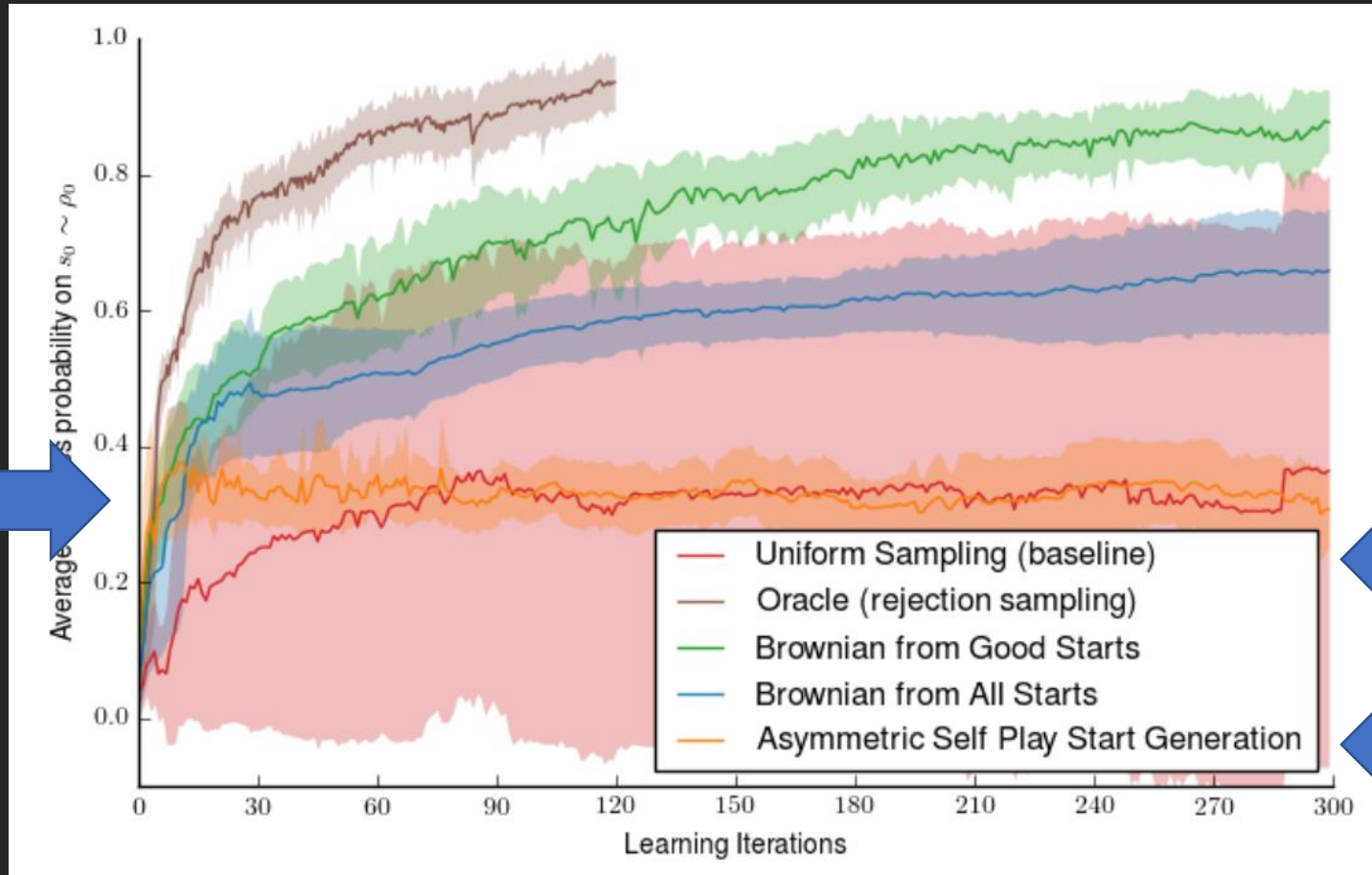$$R_A = \max(0, t_B - t_A) \implies 0$$

*ob just or too slow*

# Curriculum Learning meets Reinforcement Learning

Intrinsic Motivation and Automatic Curricula via Asymmetric Self-Play

Sukhbaatar et al.
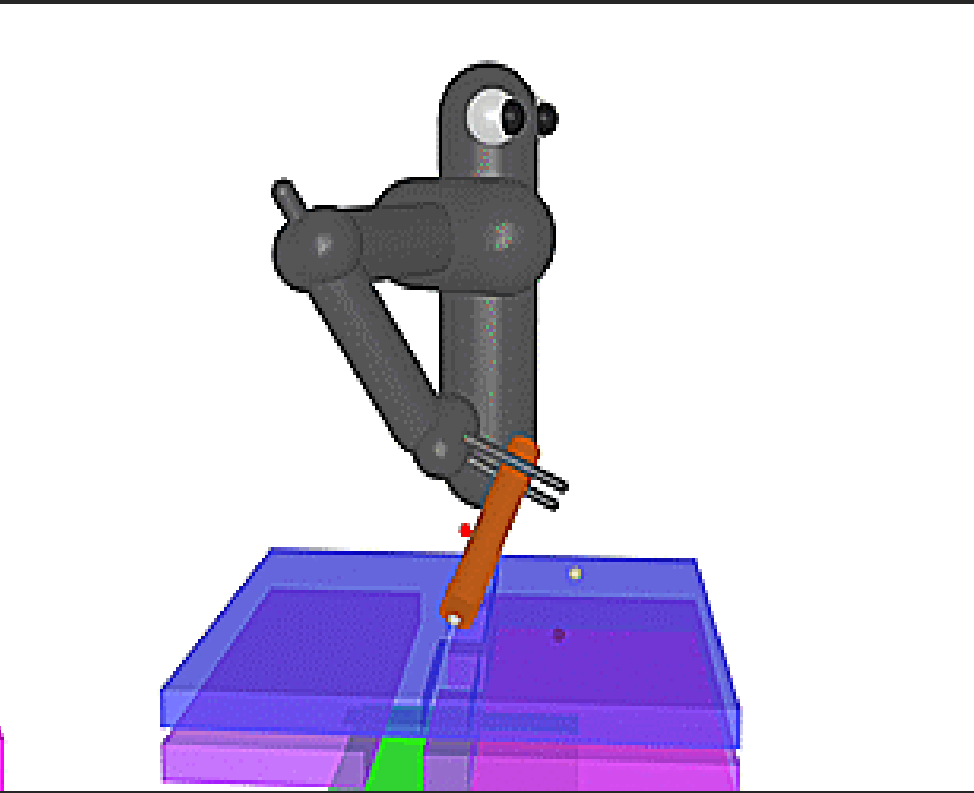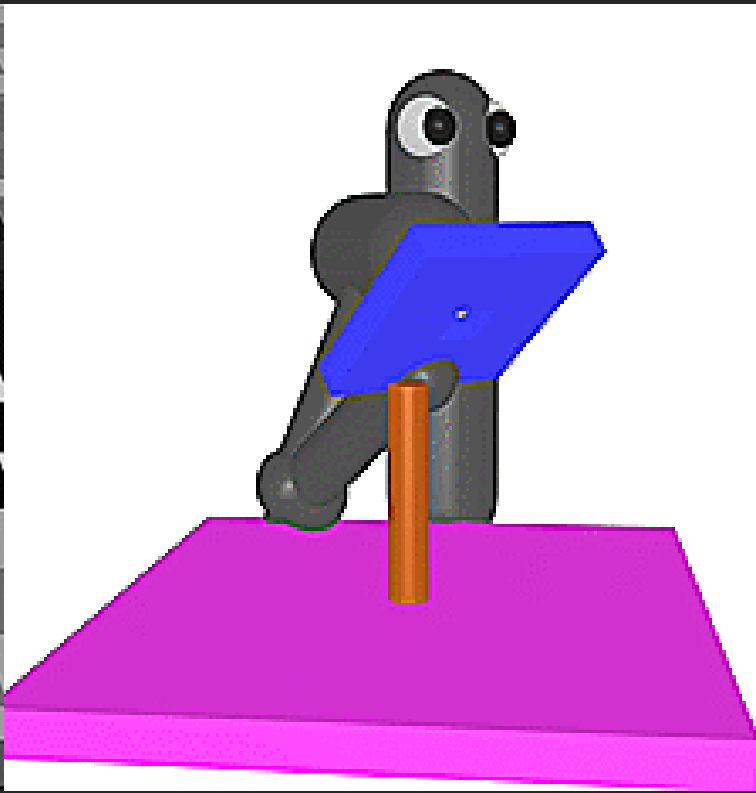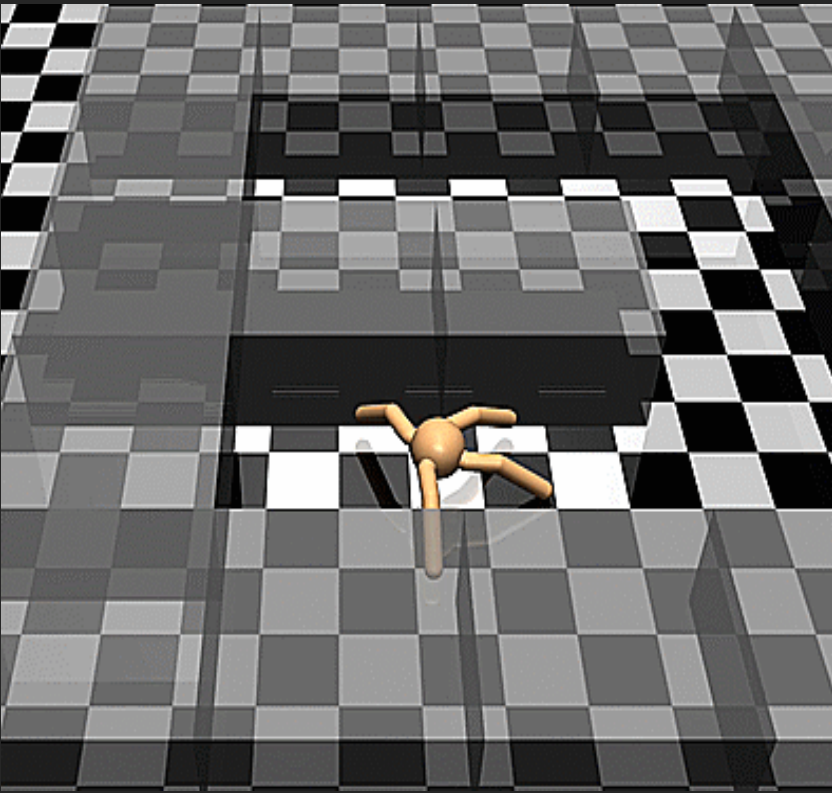
Reverse Curriculum Generation for Reinforcement Learning

Florensa et al.

Mix & Match – Agent Curricula for Reinforcement Learning

Czarnecki et al.

# Goal-Oriented Target Tasks

# Goal-Oriented Target Tasks

Binary Reward Signal  ☹

# Goal-Oriented Target Tasks

Binary Reward Signal ☹

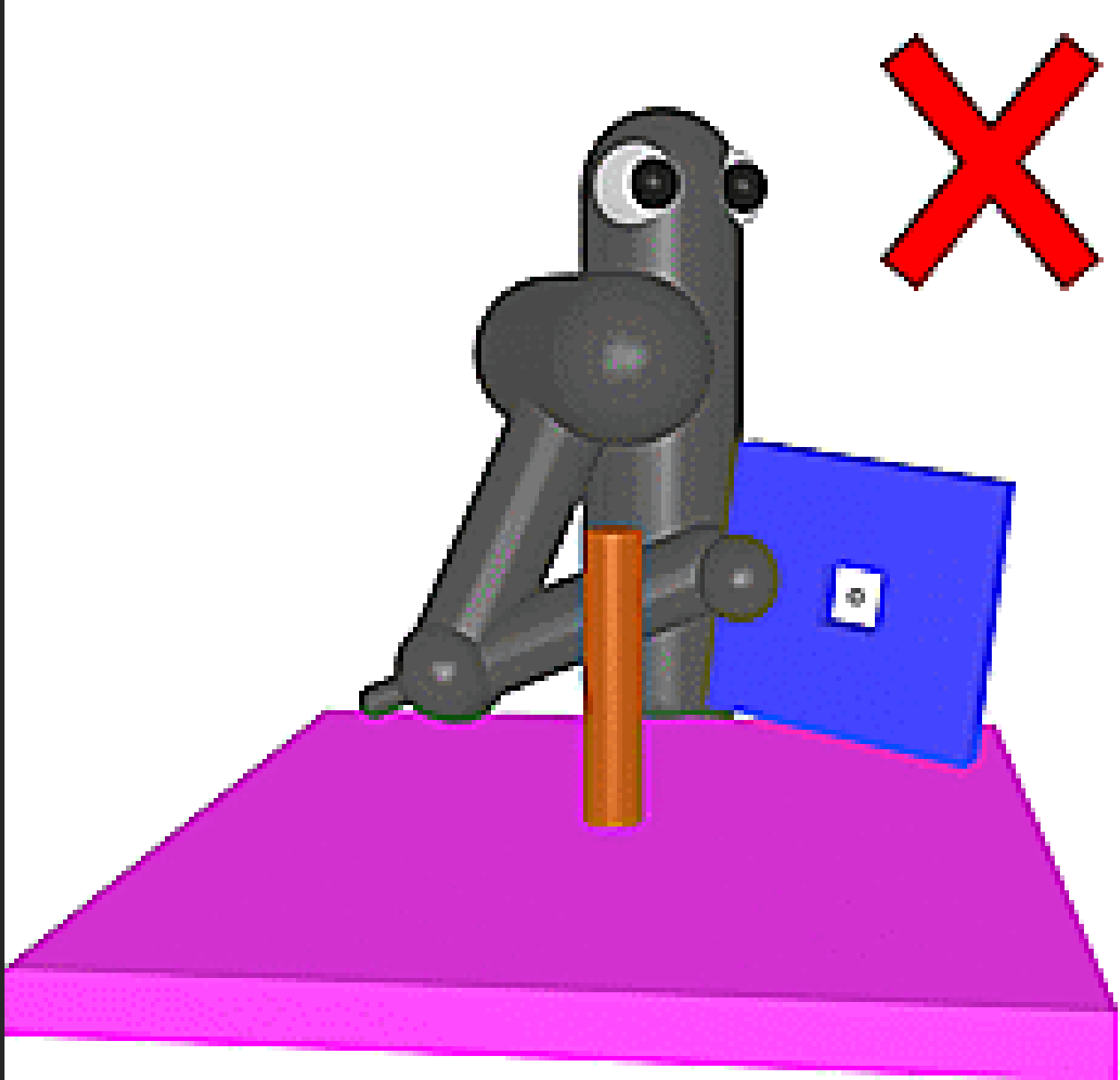+ Model-Free Reinforcement Learning ☹

# Goal-Oriented Target Tasks

Binary Reward Signal ☹
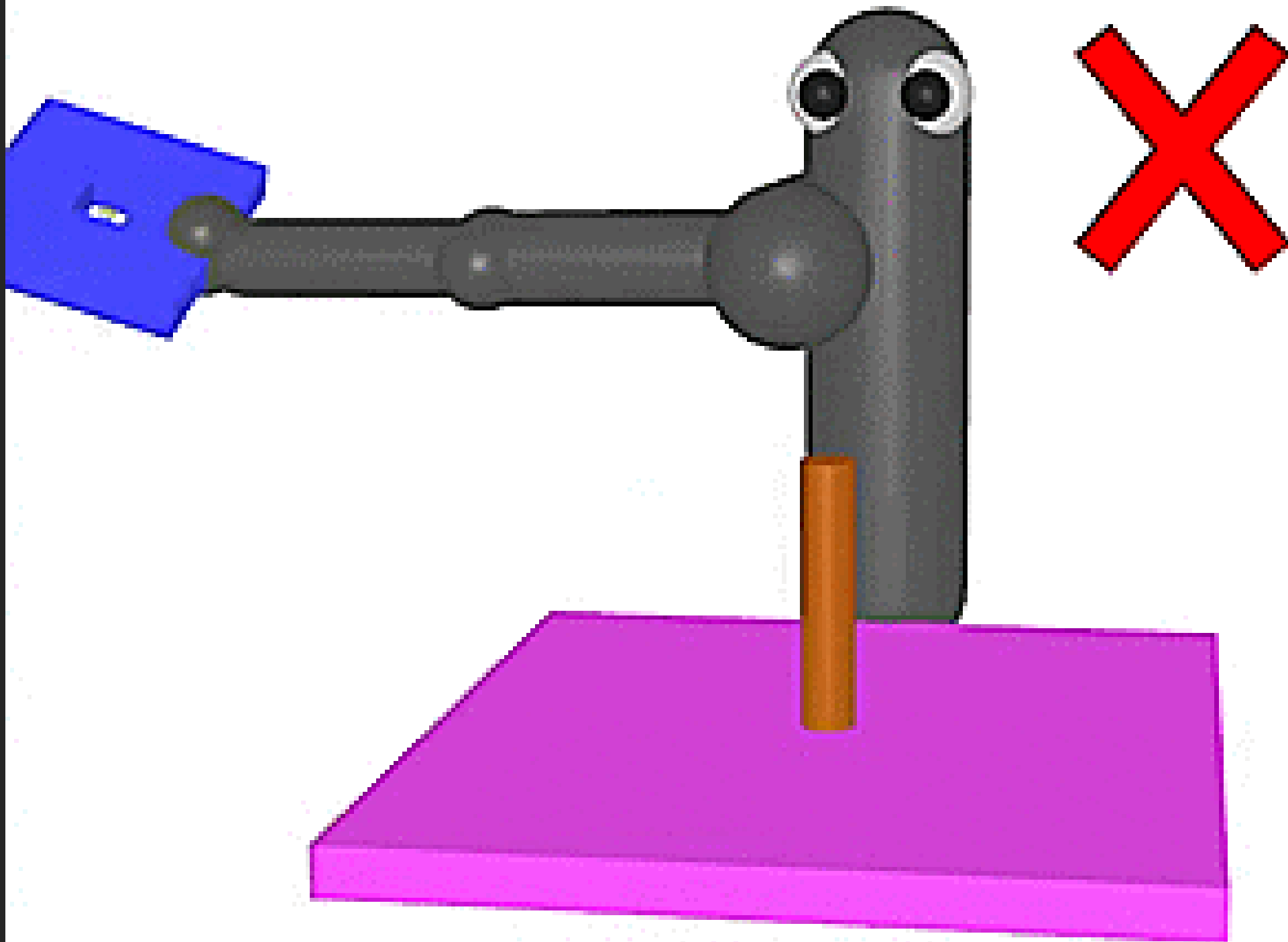
\+ Model-Free Reinforcement Learning ☹

= (╯°□°）╯︵ ┻━┻

# How do We Train the Agent?

# Random Sampling of Starting States?
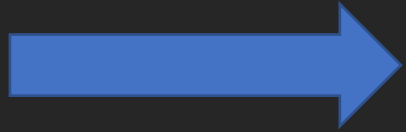
# Add Regularization Term?

# What's the Trick?

Easy to Win, if you Start at the Goal!

# Reverse Curriculum

1    Start almost there

2    Start increasingly further away
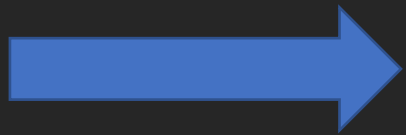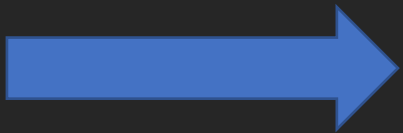
3    Profit from work already done

# Reverse Curriculum

1    Start almost there

2    Start increasingly further away

3    Profit from work already done

# Reverse Curriculum

1    Start almost there

2    Start increasingly further away

3    Profit from work already done

**Automatically** creates a Curriculum over Start States! 🎉

# States of Intermediate Difficulty (SoIDs)

1. States Close to $s^g$ may be good Start States 💡

2. Random Walk in State-Space ☹

3. Brownian Motion in Action-Space ☺

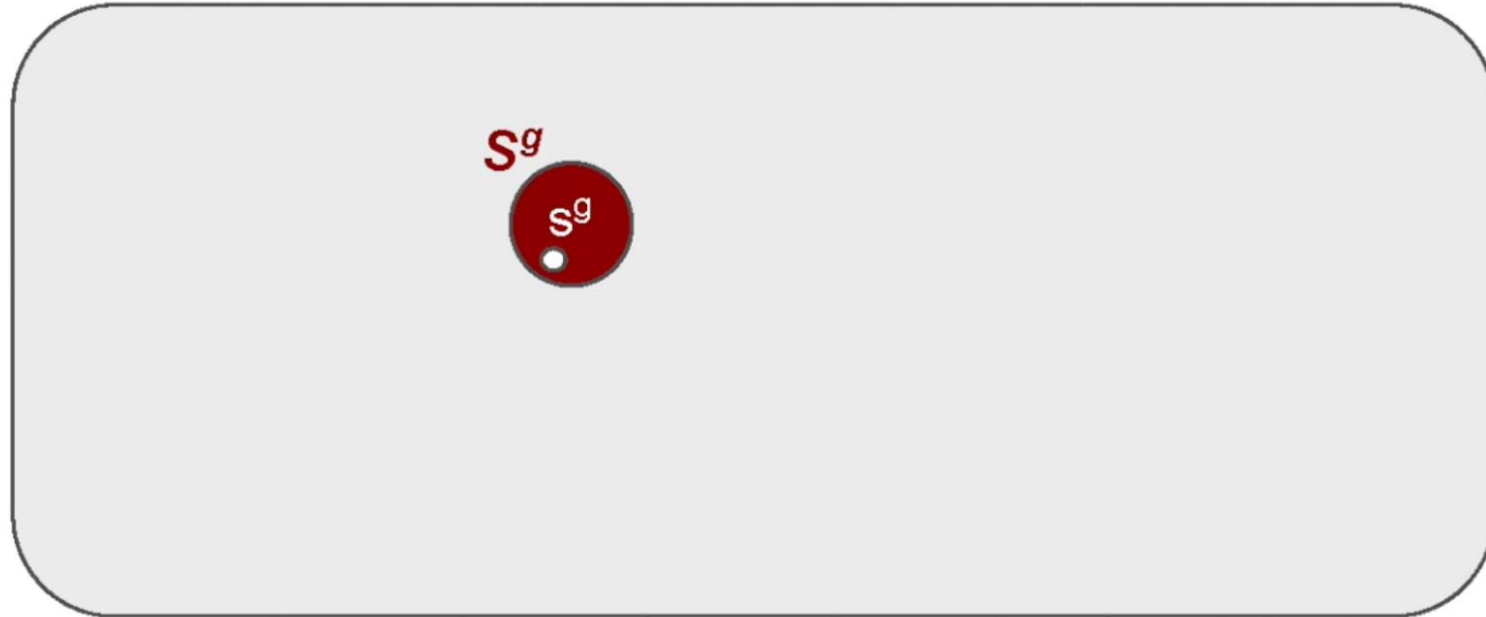# States of Intermediate Difficulty (SoIDs)

1. States Close to $s^g$ may be good Start States 💡

⮕ 2. Random Walk in State-Space ☹

3. Brownian Motion in Action-Space ☺

# States of Intermediate Difficulty (SoIDs)

1. States Close to $s^g$ may be good Start States 💡

2. Random Walk in State-Space ☹
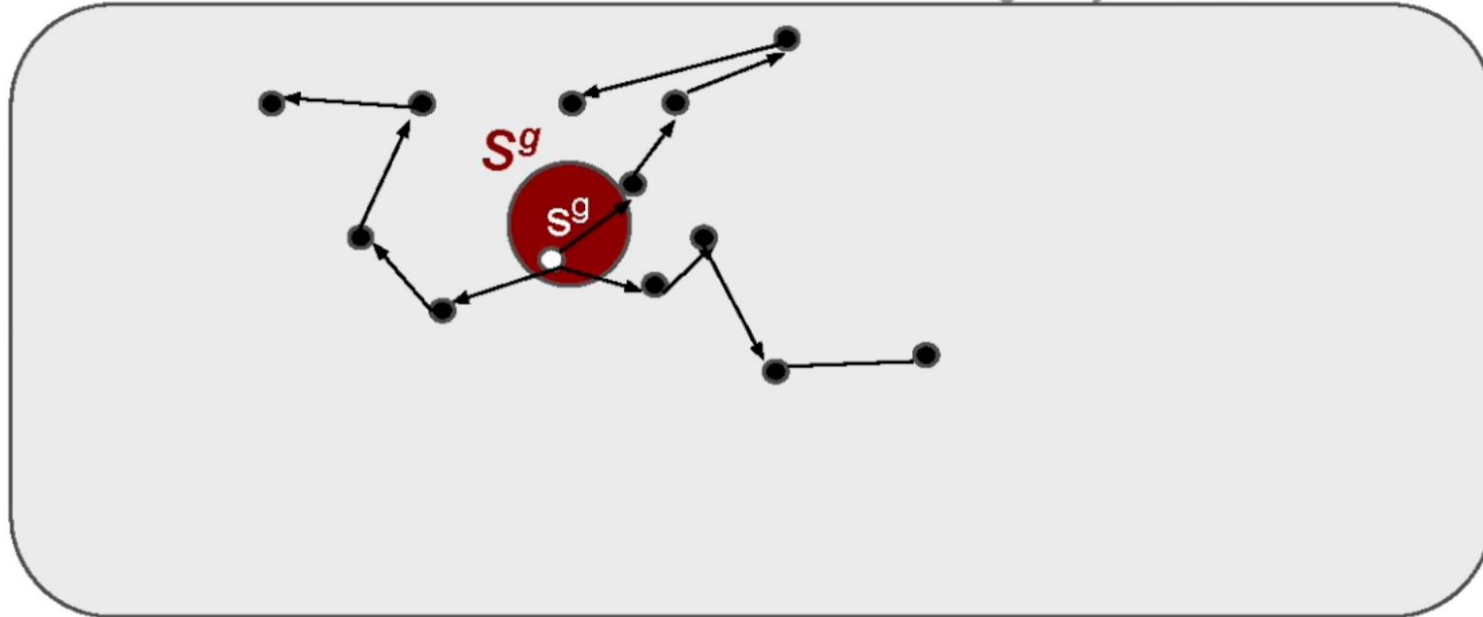
➡ 3. Brownian Motion in Action-Space ☺

**$S^g$**: goal states we want to reach from everywhere.

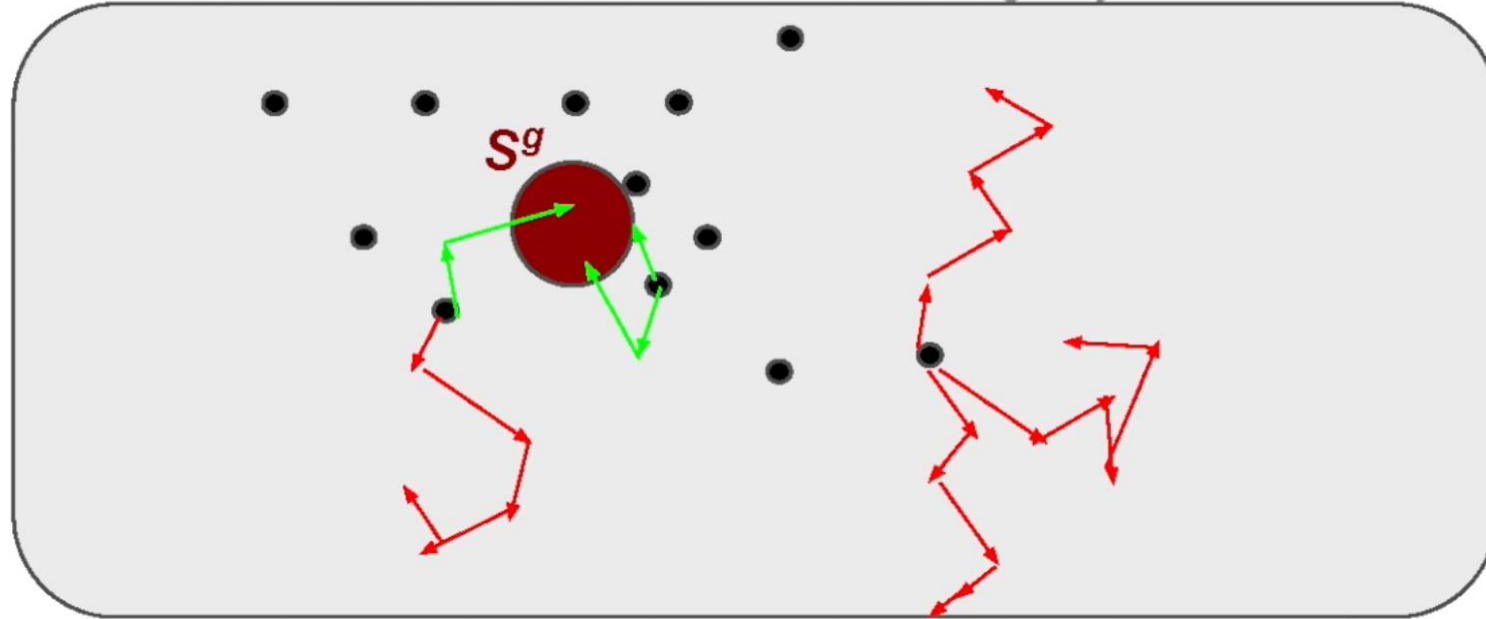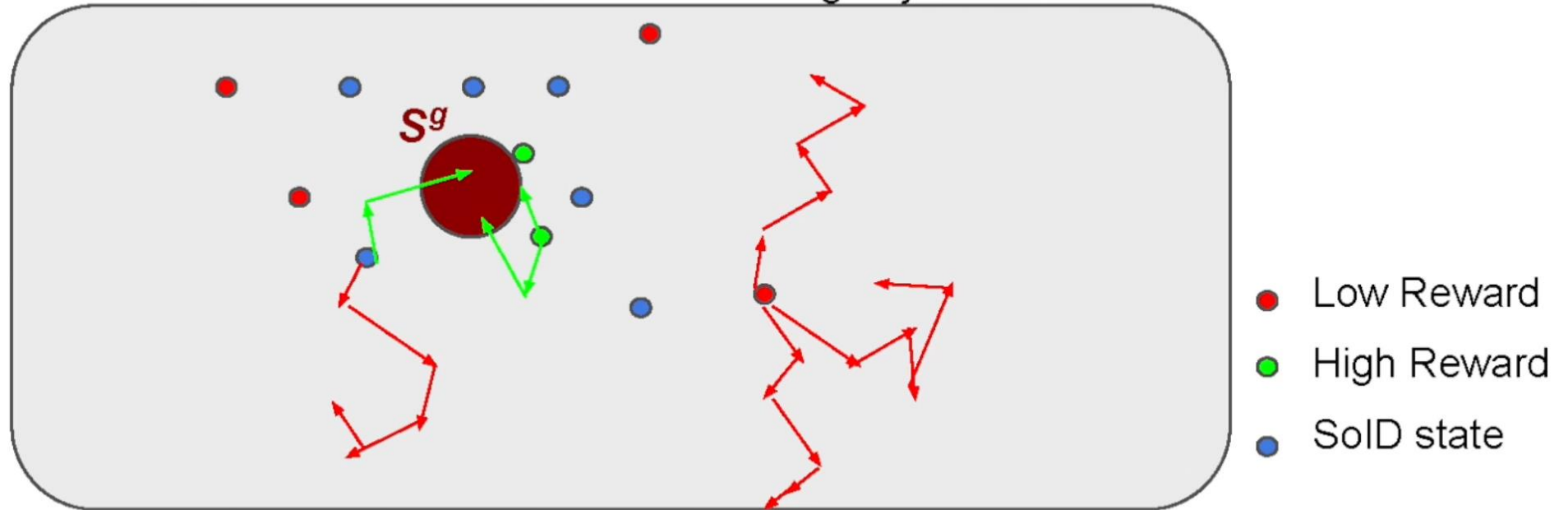**$s^g$**: one goal state is provided

**Iteration 1:**
- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- Label and filter starts based on training trajectories

**Iteration 1:**
- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- Label and filter starts based on training trajectories
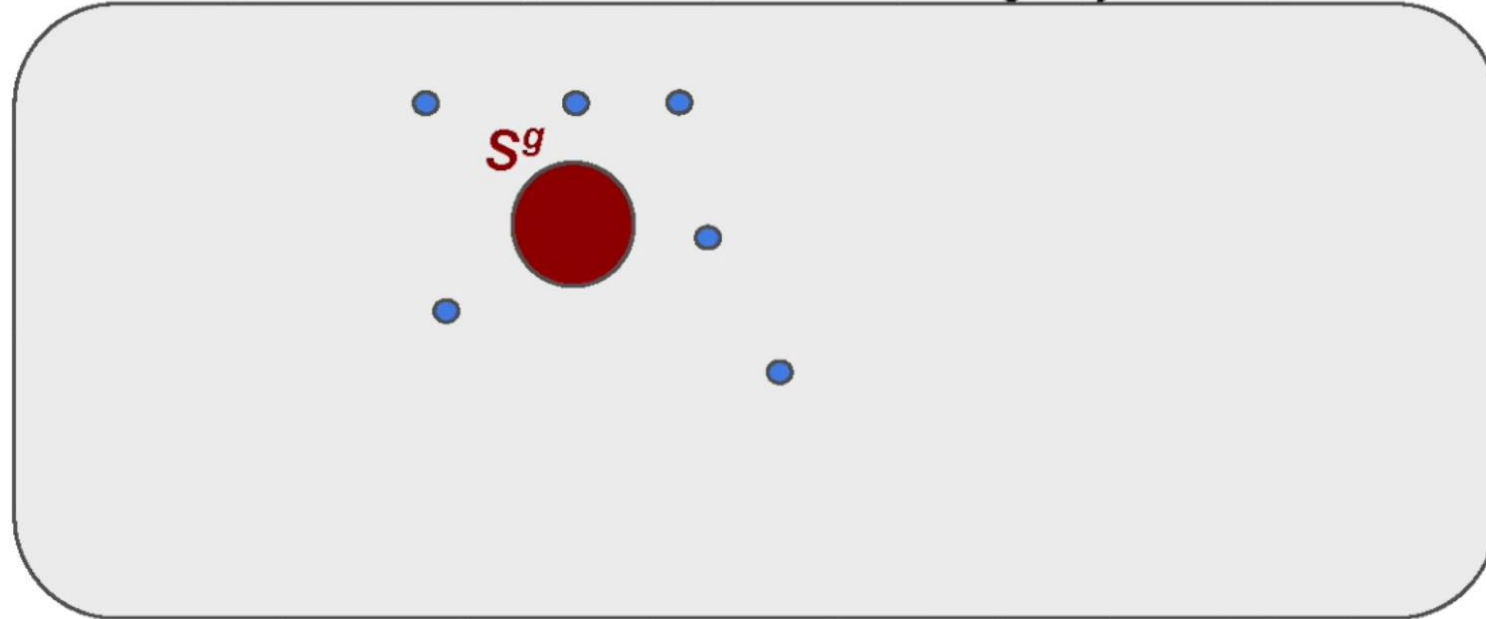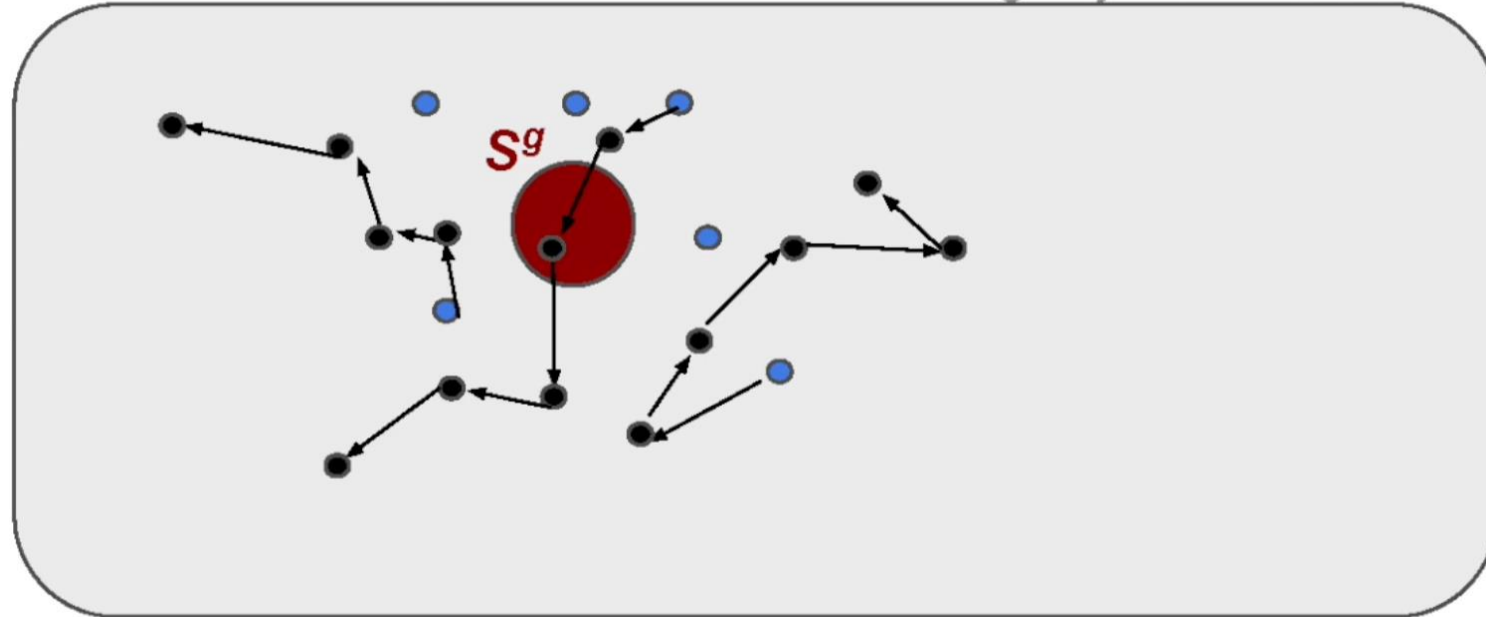
$S^g$

# Iteration 1:
- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- **Label and filter starts based on training trajectories**
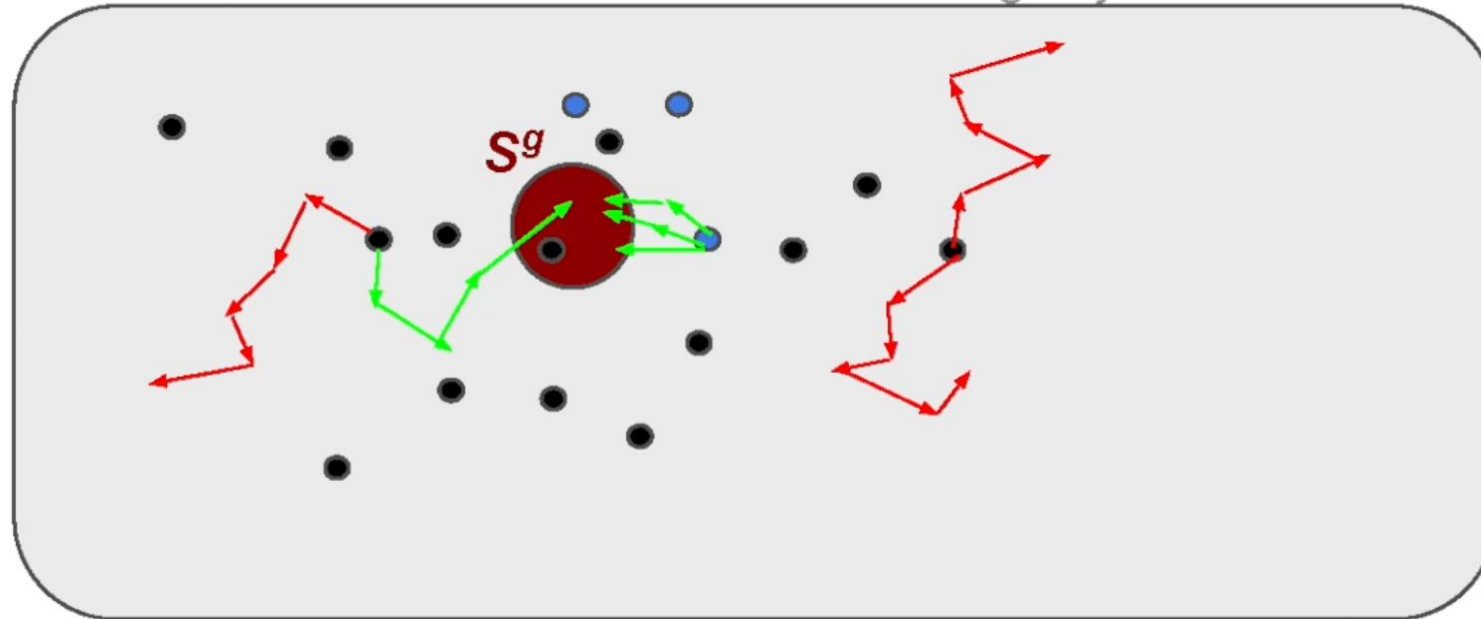


Low Reward

High Reward

SoID state

**Iteration 2:**
- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- Label and filter starts based on training trajectories

$S^g$

Low Reward

High Reward

SolD state

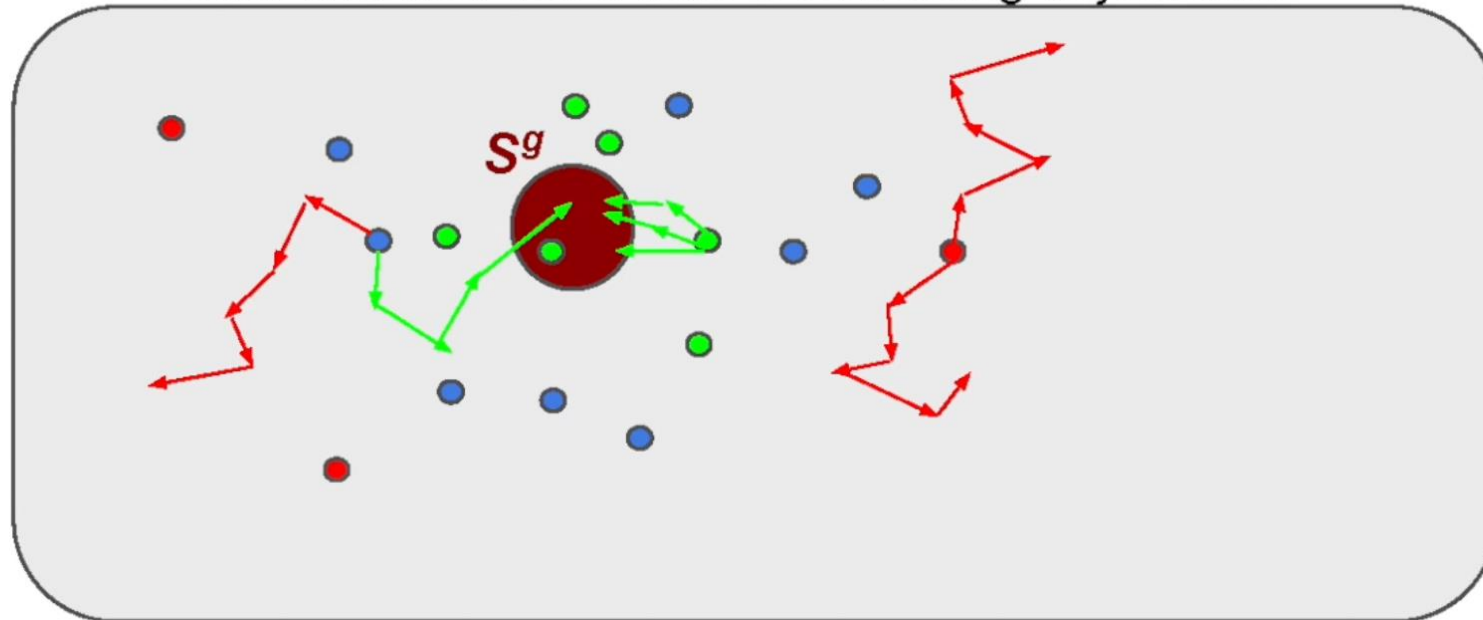**Iteration 2:**
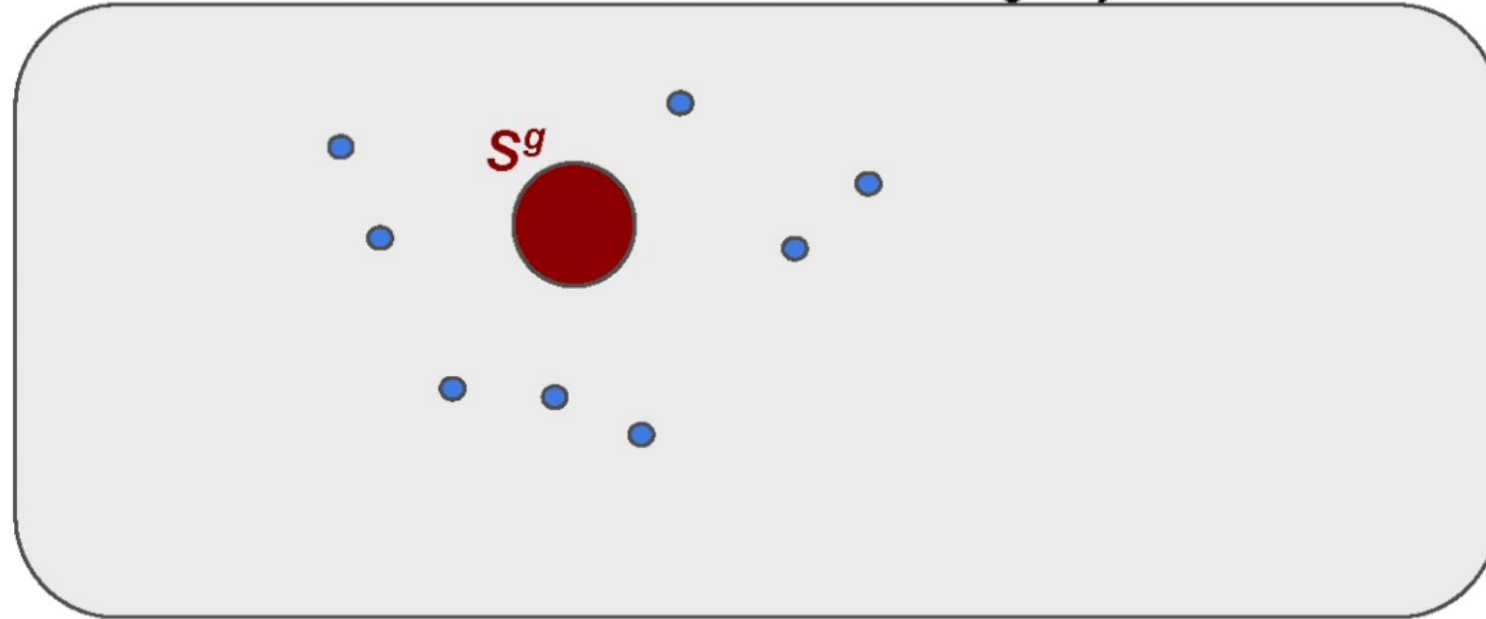- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- Label and filter starts based on training trajectories

$S^g$

Low Reward

High Reward

SoID state

# Iteration 2:
- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- **Label and filter starts based on training trajectories**
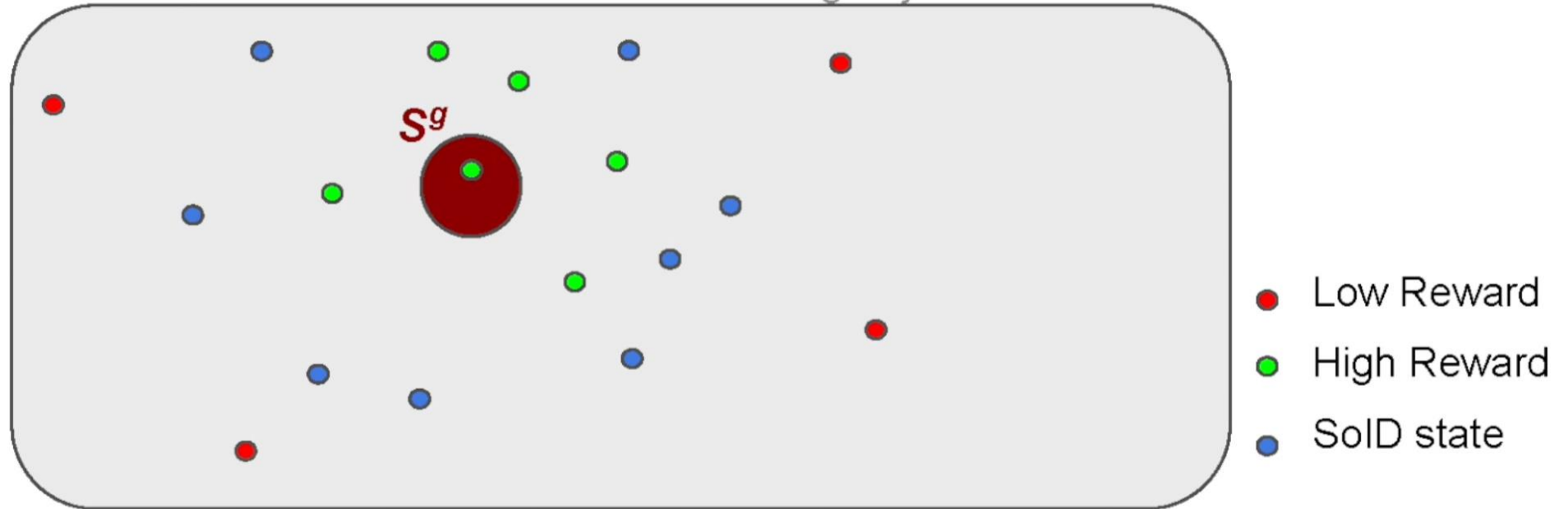
**Iteration 2:**
- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- Label and filter starts based on training trajectories

$S^g$

Low Reward
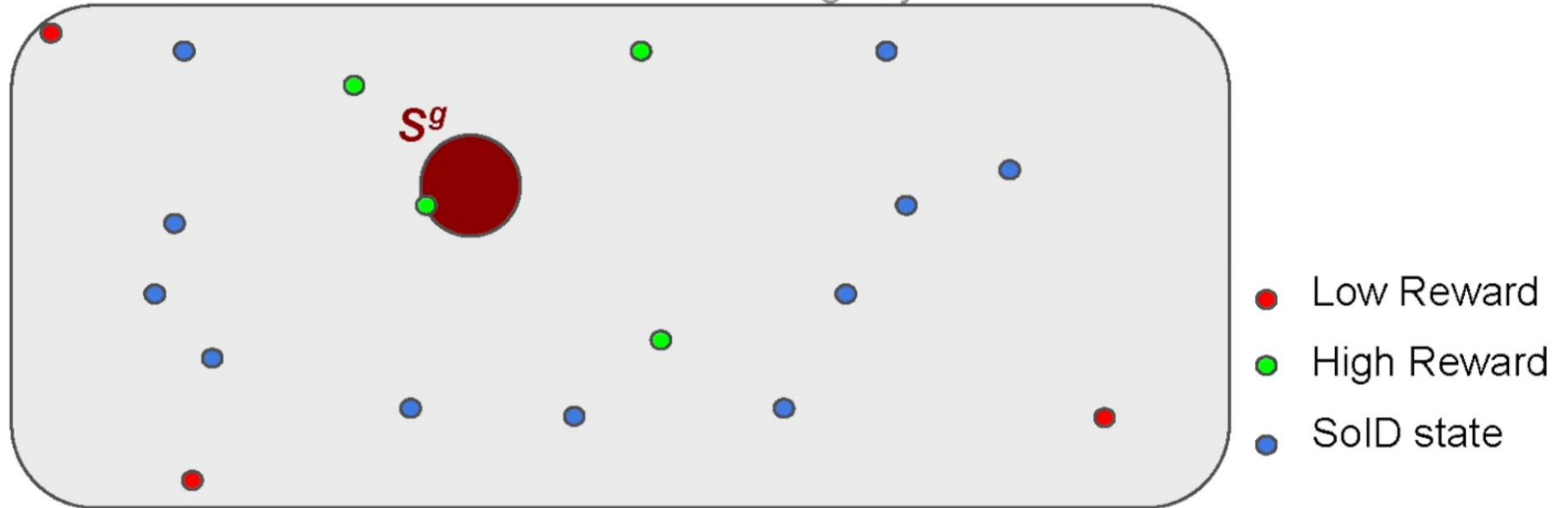
High Reward

SolD state

**Iteration 3:**
- Run Brownian motion
- Obtain trajectories from collected starts to train policy
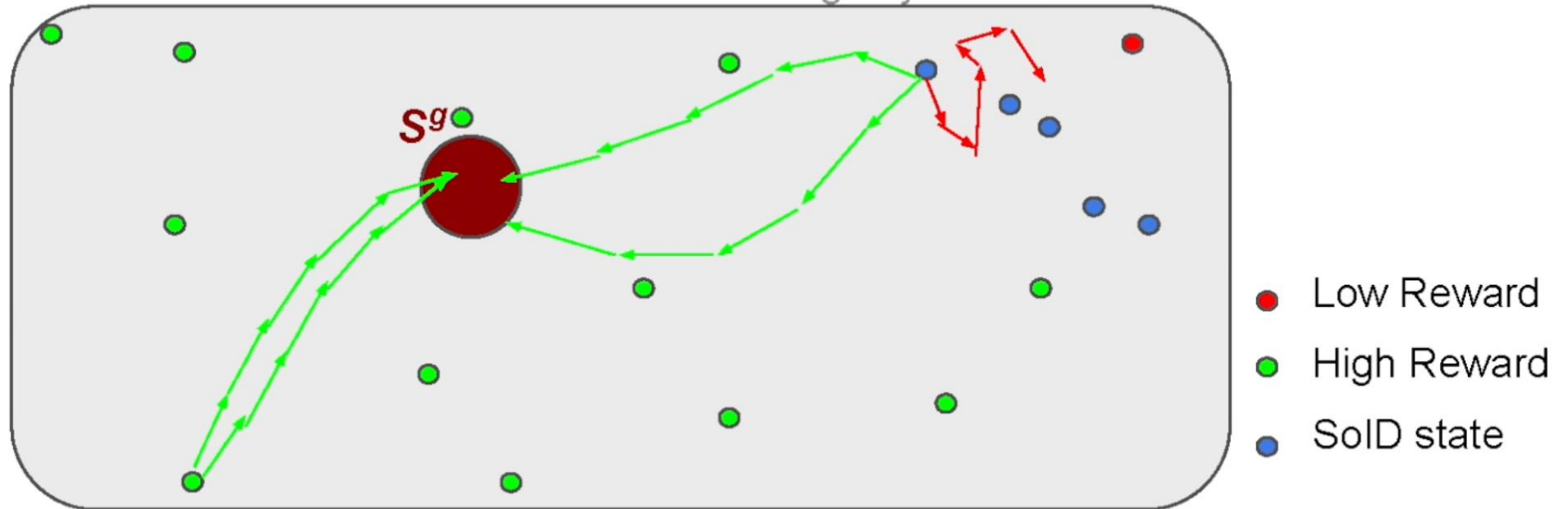- Label and filter starts based on training trajectories

$S^g$

Low Reward

High Reward

SoID state

**Iteration 4:**
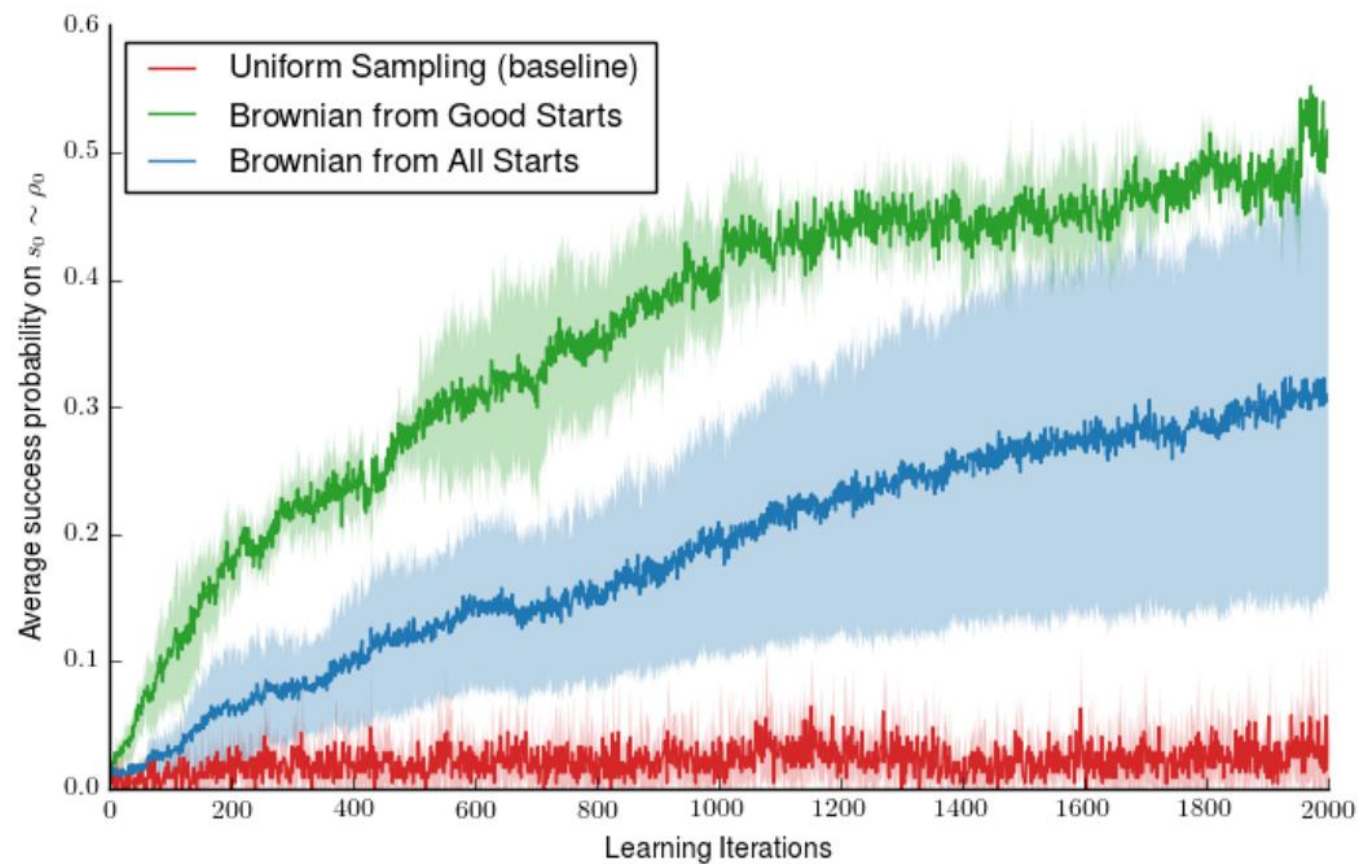- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- Label and filter starts based on training trajectories

$S^g$

Low Reward

High Reward

SoID state

# Iteration 5:
- Run Brownian motion
- Obtain trajectories from collected starts to train policy
- Label and filter starts based on training trajectories



Low Reward
High Reward
SoID state

(d) Key insertion task

# Curriculum Learning meets Reinforcement Learning

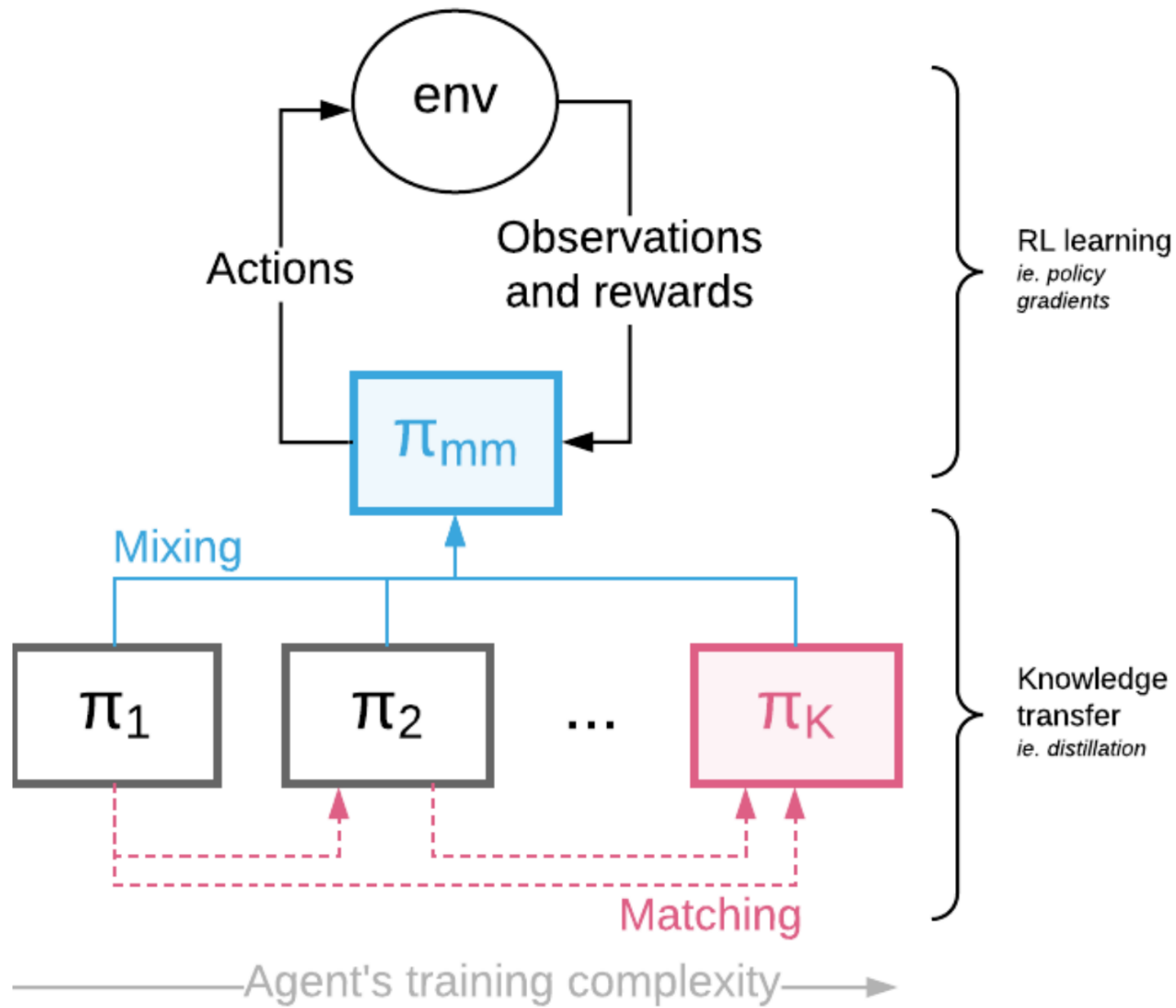Intrinsic Motivation and Automatic Curricula via Asymmetric Self-Play

Sukhbaatar et al.
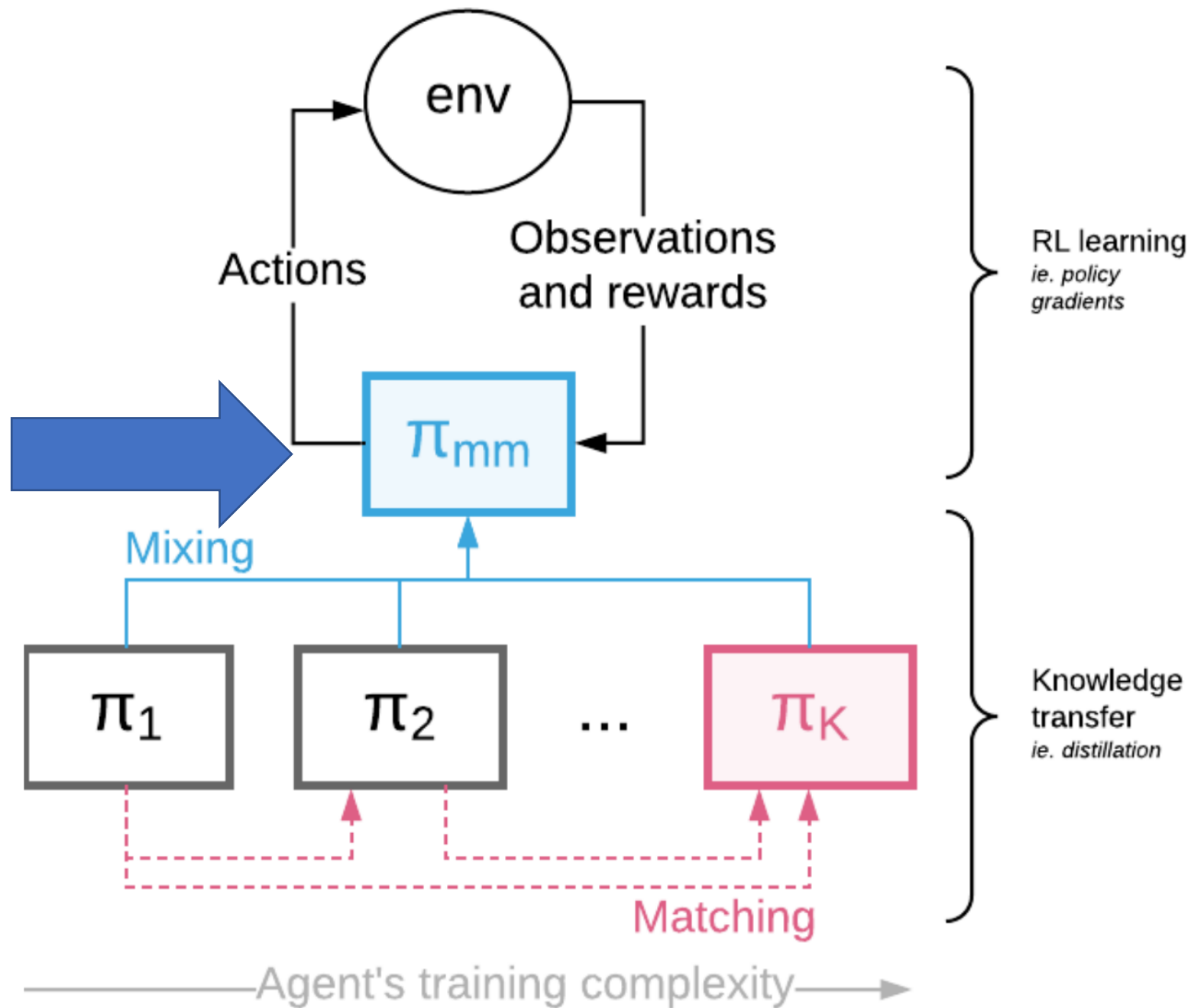
Reverse Curriculum Generation for Reinforcement Learning

Florensa et al.

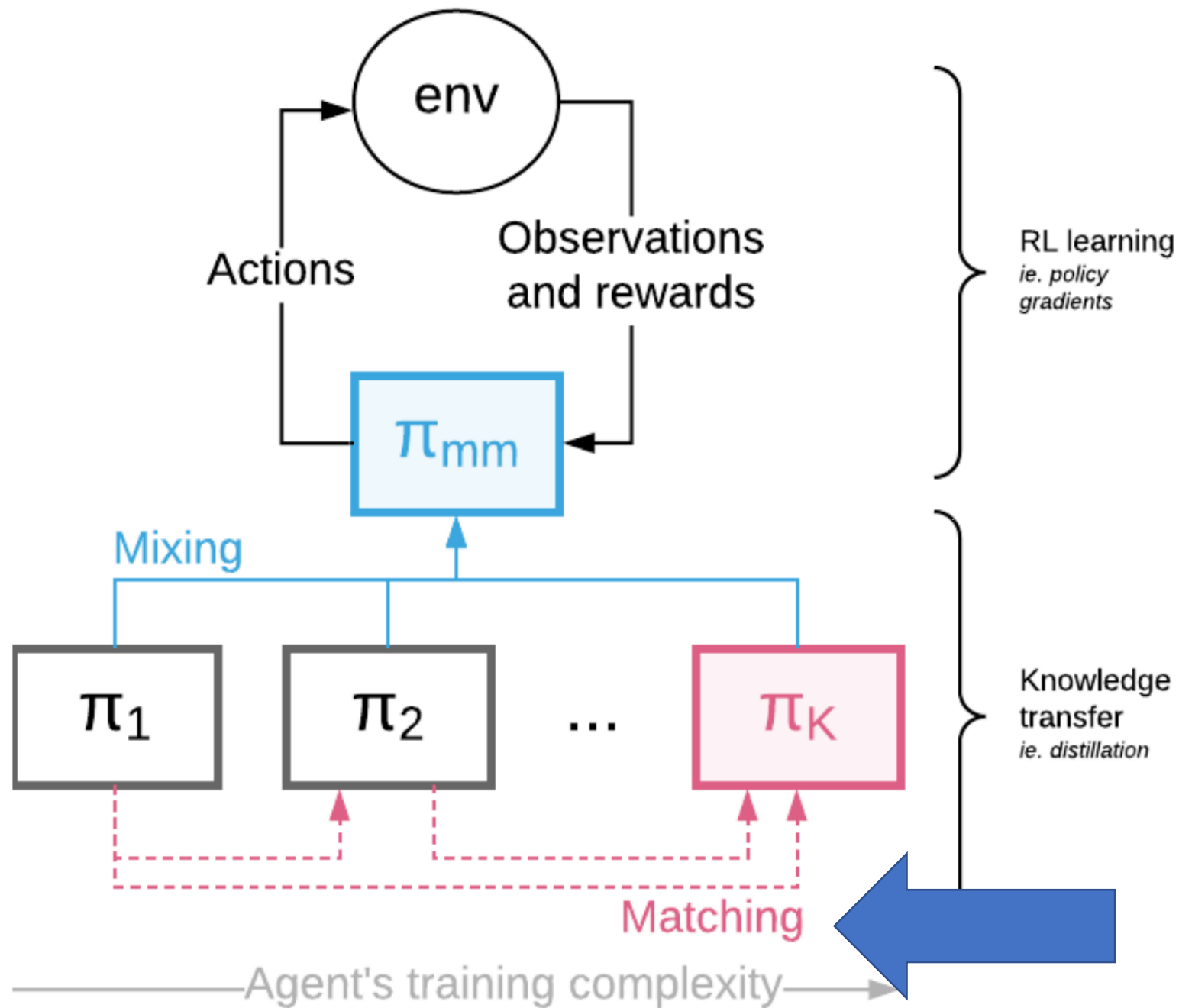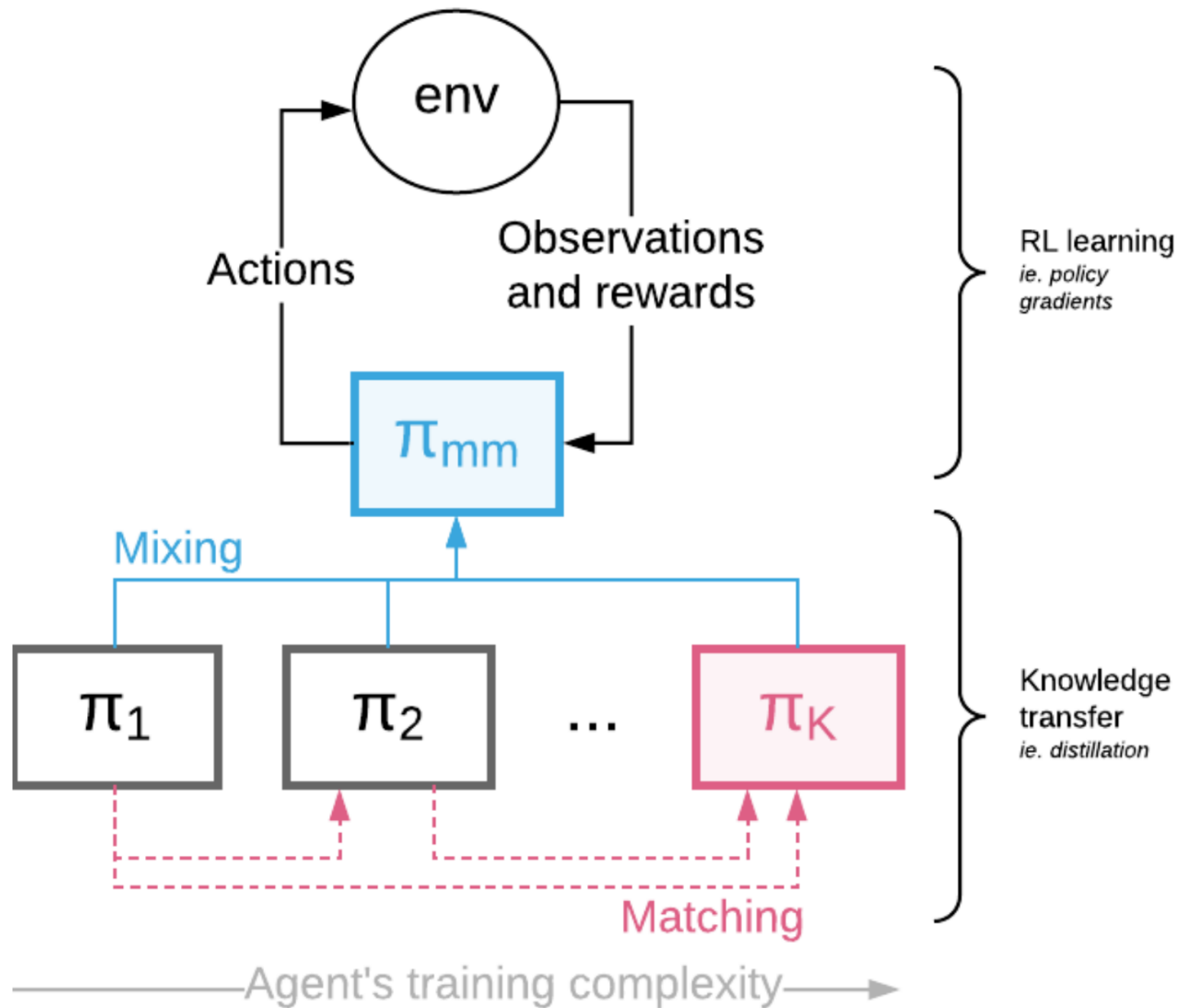Mix & Match – Agent Curricula for Reinforcement Learning

Czarnecki et al.

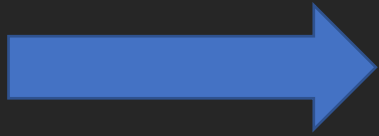# Rethinking the Notion of Curriculum

Curriculum not Automatic!

# What's the Difficulty of an Agent?

Agents **are** Neural Networks!
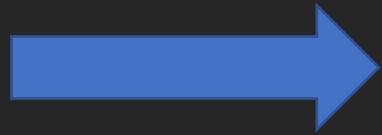
*for all practical purposes

# Architectural Components
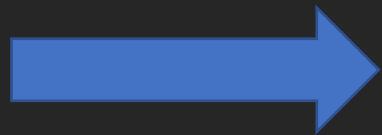**or**
# Performable Actions
**or**
# Jointly-Learnable Tasks
**and**
# Training Iterations

# Architectural Components

**or**

# Performable Actions

**or**

# Jointly-Learnable Tasks

**and**

# Training Iterations

# Architectural Components
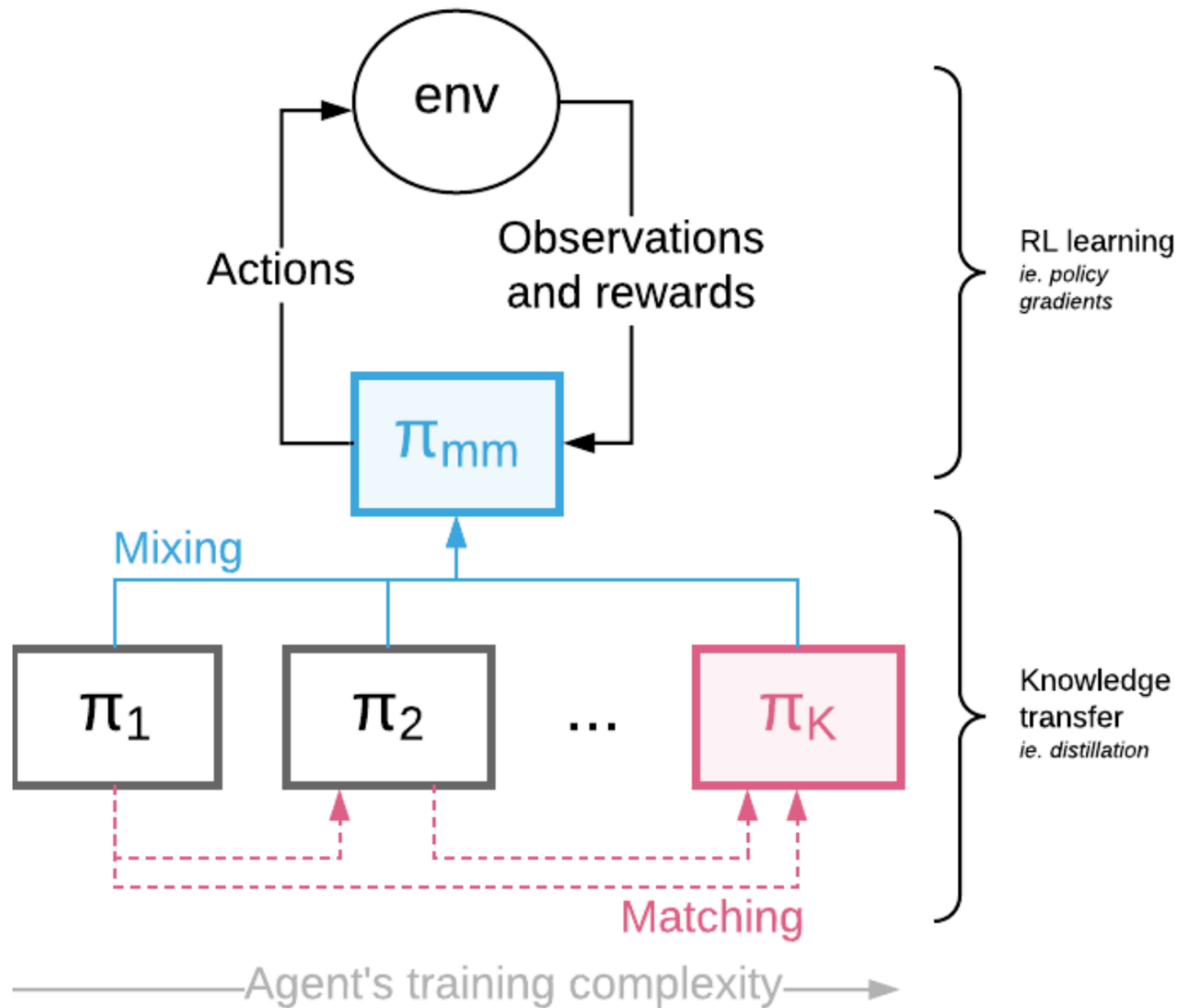
**or**

# Performable Actions
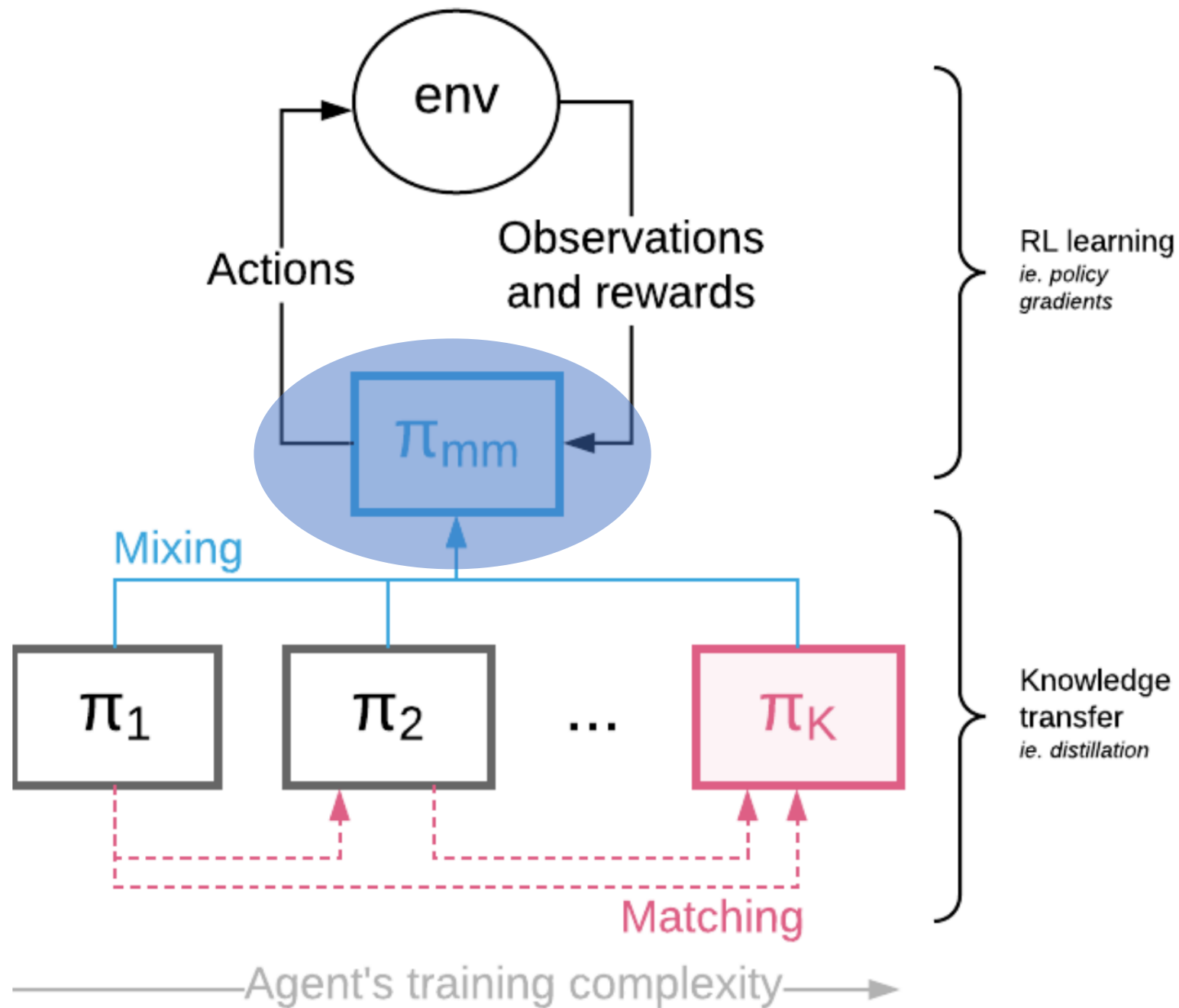
**or**

# Jointly-Learnable Tasks

**and**

# Training Iterations

Difficulty ✓

# Scheduler: Tune Mixture Parameter $\alpha$

# Scheduler: Tune Mixture Parameter $\alpha$

Could use hand crafted scheduler ☹

# Scheduler: Tune Mixture Parameter $\alpha$

Could use hand crafted scheduler ☹

Could use naive hyperparameter tuning ☹

# Scheduler: Tune Mixture Parameter $\alpha$

Could use hand crafted scheduler ☹

Could use naive hyperparameter tuning ☹

Population Based Training ☺

# Population Based Training

# Population Based Training
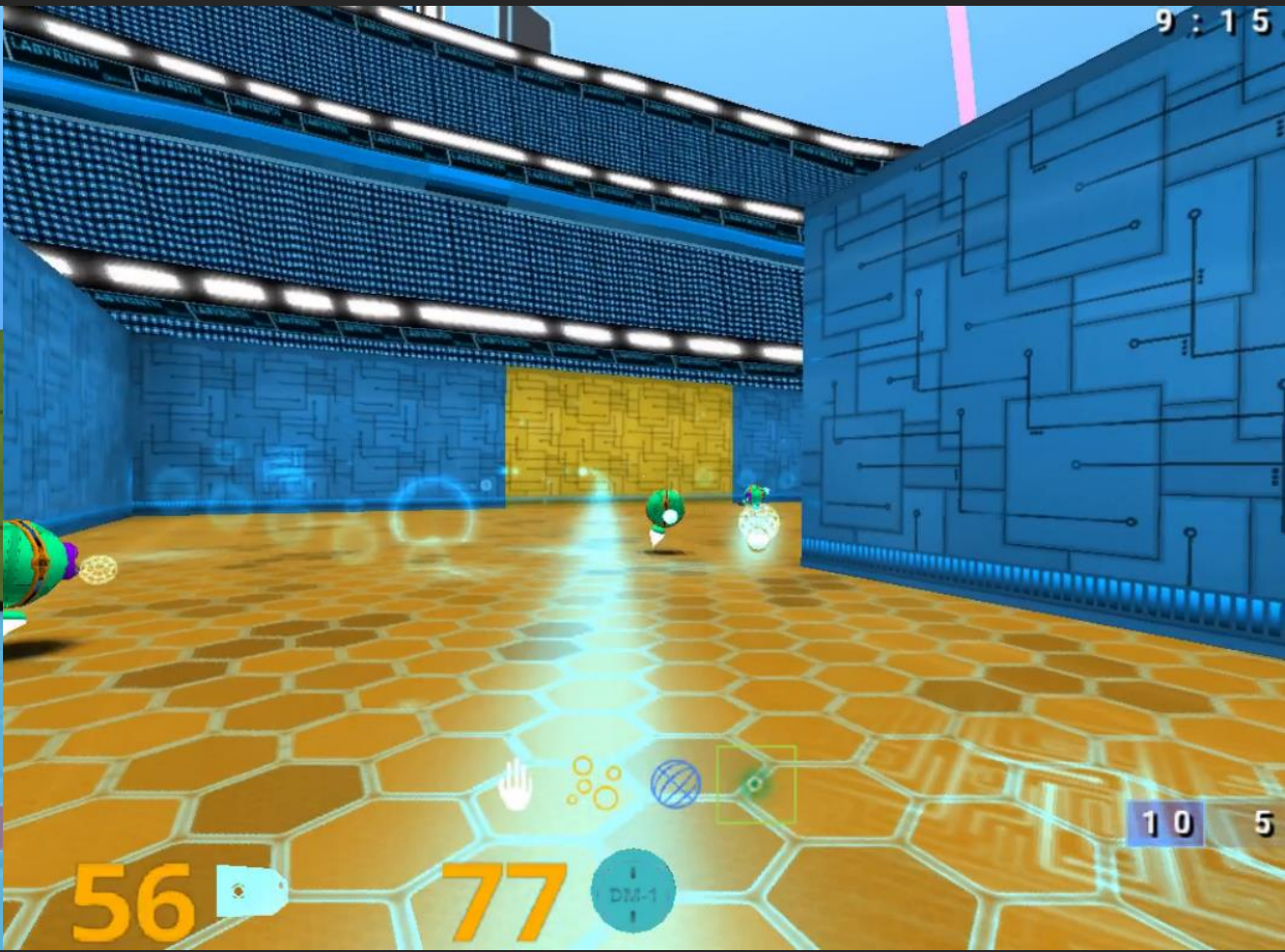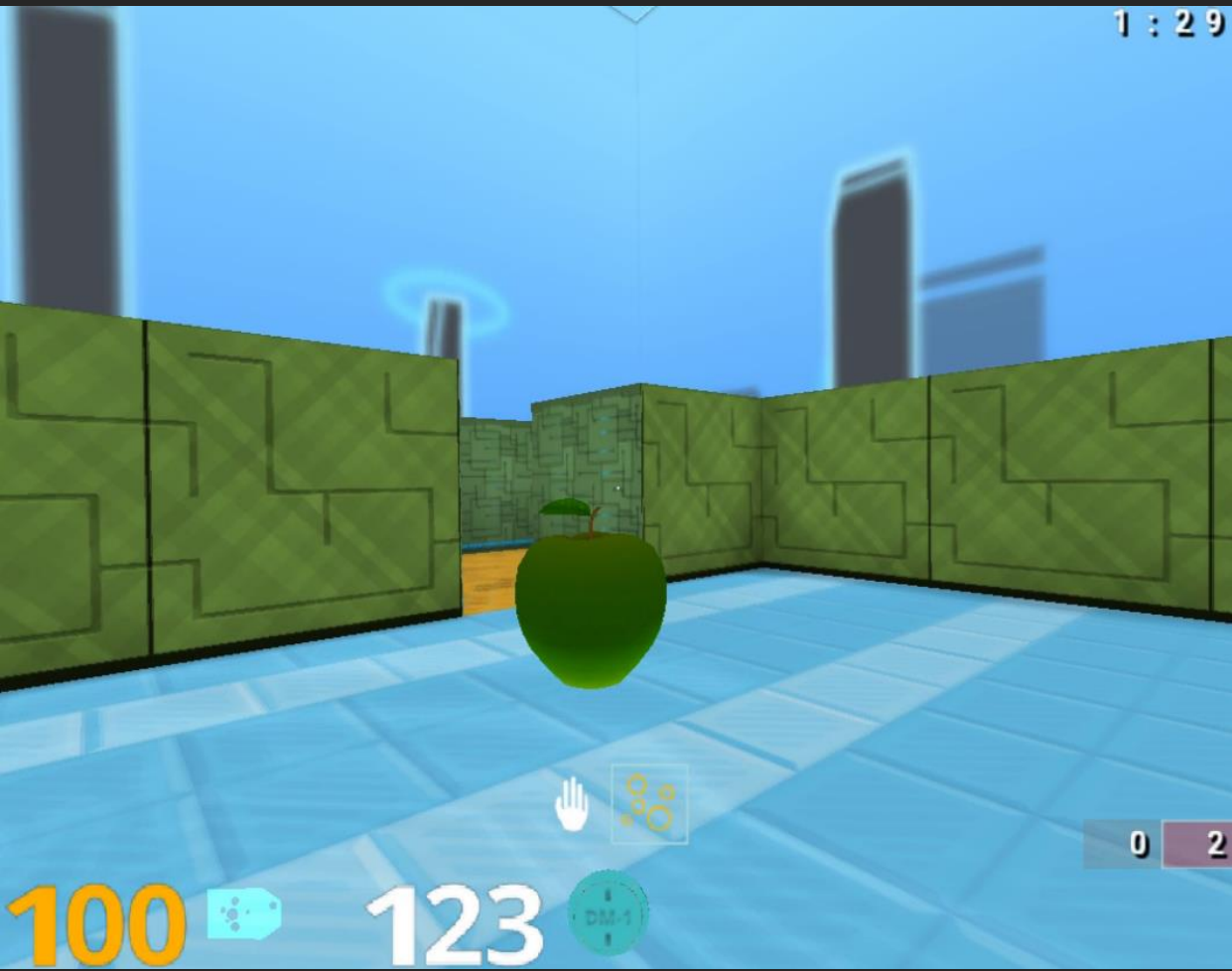
1       Tuning several mixture agents in parallel

# Population Based Training

1  Tuning several mixture agents in parallel
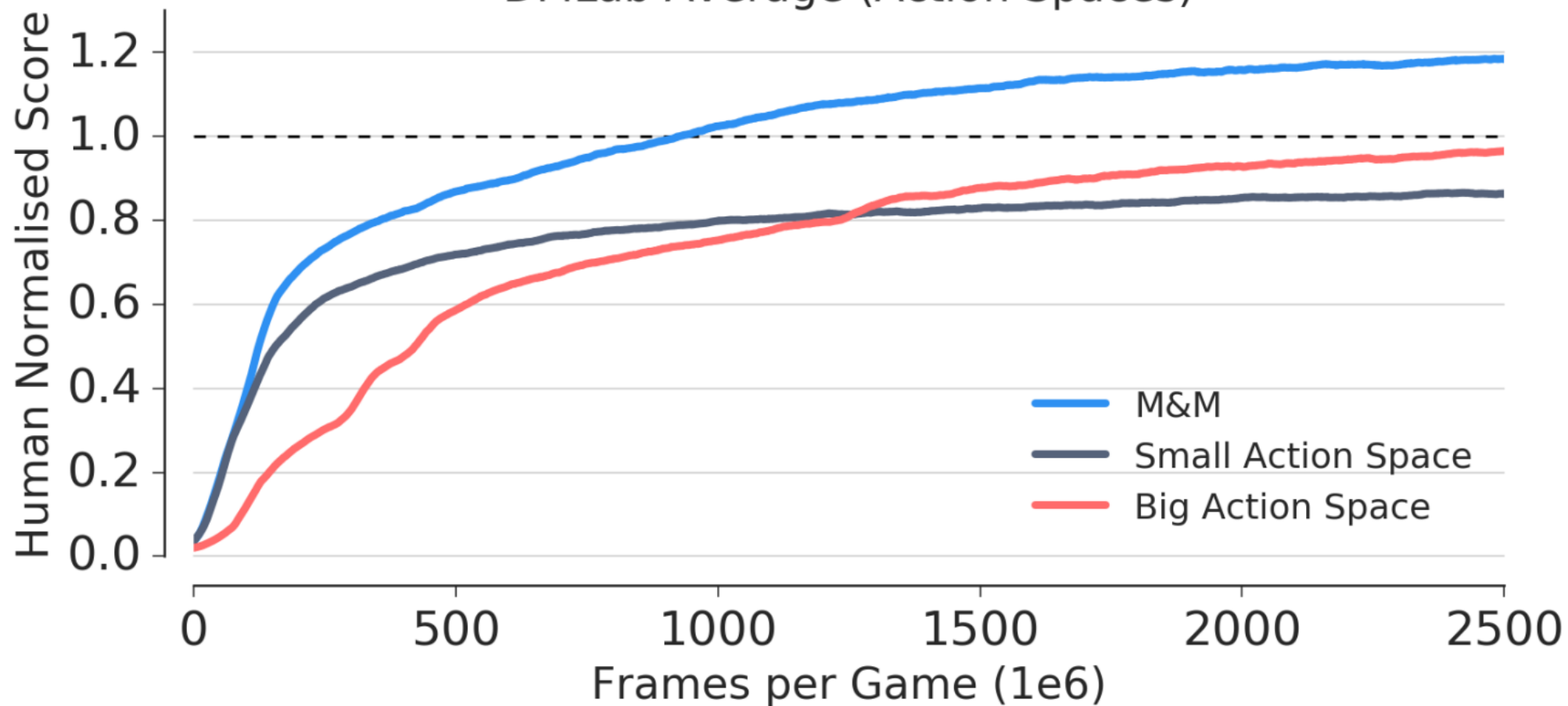
2  Agent A periodically communicates with some B

# Population Based Training

1      Tuning several mixture agents in parallel

2      Agent A periodically communicates with some B

3      Badly performing: Copy weights and
                                hyperparameters ($\alpha$)

# Explore Search Space
## with badly performing Agents

DMLab Average (Action Spaces)

Curriculum Learning
Is Here to Stay! ☺

Yes, Mr. Frodo. It's over now.

Yes, Mr. Frodo. It's over now.

# Take Care!