



Task Agnostic Continual RL

Massimo Caccia, Jonas Mueller, Taesup Kim, Laurent Charlin, Rasool Fakoor



TL;DR:



Overview

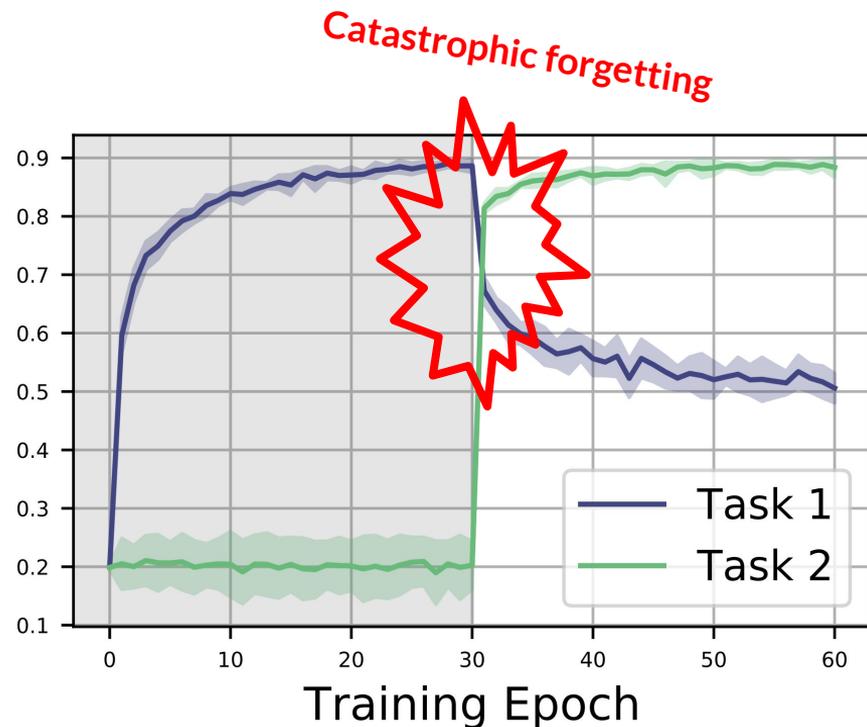


- Background
 - Continual Learning
 - Reinforcement Learning
- Task Agnostic Continual RL
 - Problem Statement
 - Soft upper Bounds: Task-awareness, Multi-task learning
- Methods
 - Backbone RL algo (SAC)
 - Baselines
- Empirical Findings
 - Benchmark
 - Task-agnostic > Task-Aware
 - Continual learning = multi-task learning
 - Hypothesis testing
- Discussion

Background

Continual Learning (CL)

- Accumulating knowledge on **non-stationary** data distributions
- ML/DL can't learn on **changing data distributions** (or tasks)



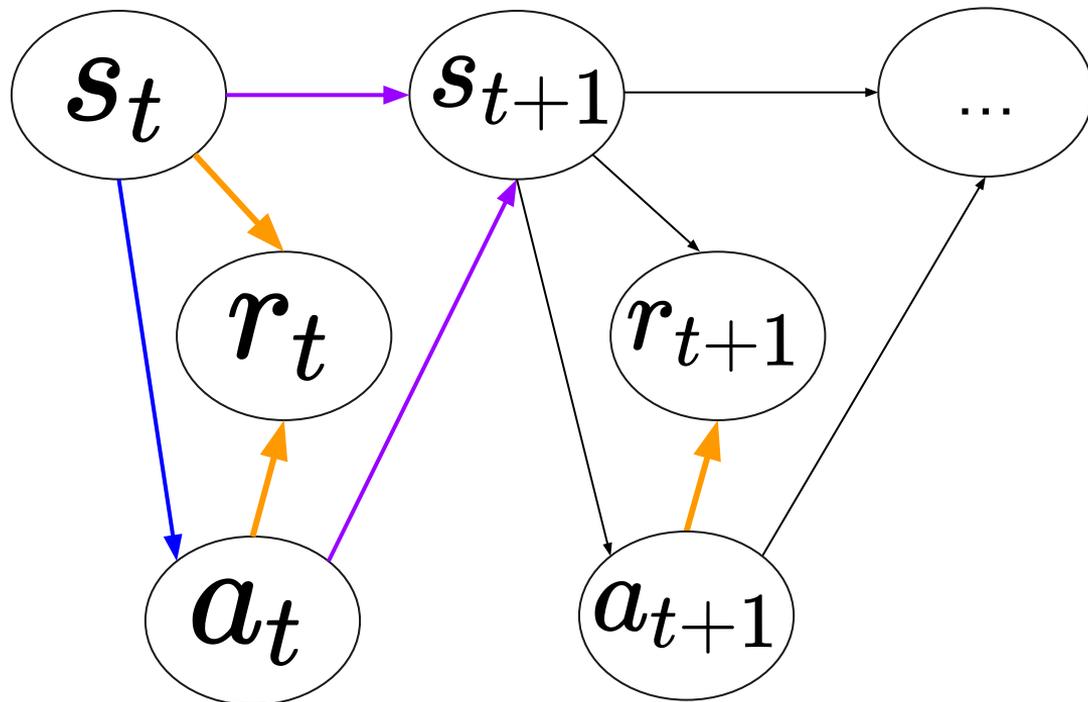
Continual Learning (CL): Why we care



- industry/deployment argument
- Curriculum learning argument
 - Or continually increasing sample/compute efficiency
- Learning autonomously in an open world → AGI

Background

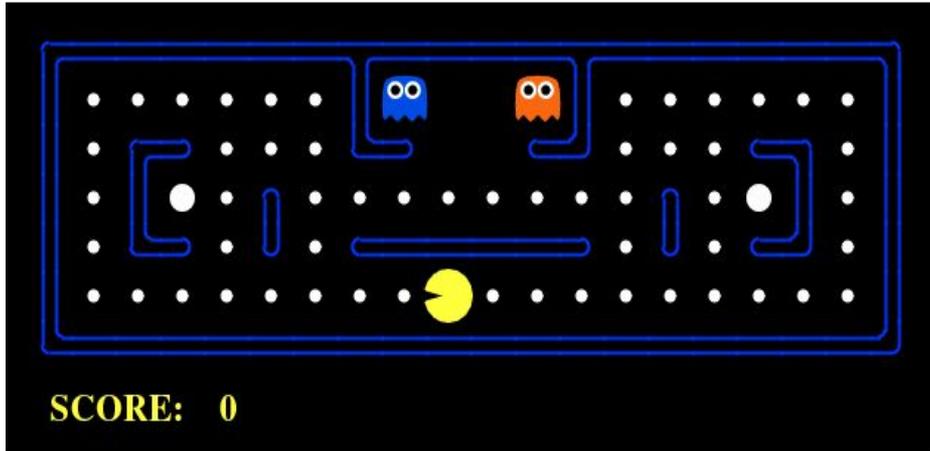
Reinforcement Learning (RL) - MDP



$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E} \sum_{t=0}^{\infty} \gamma^t r_t$$

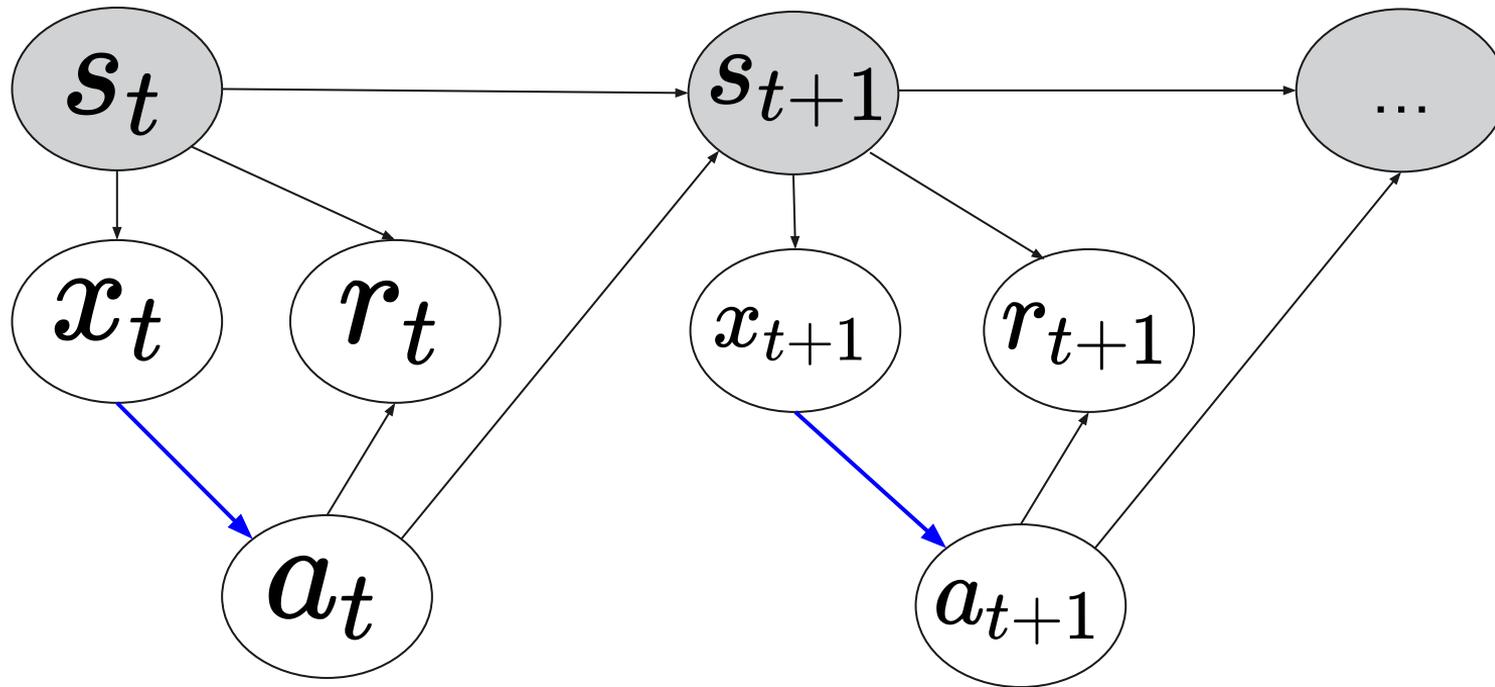
Background

Reinforcement Learning (RL) - MDP



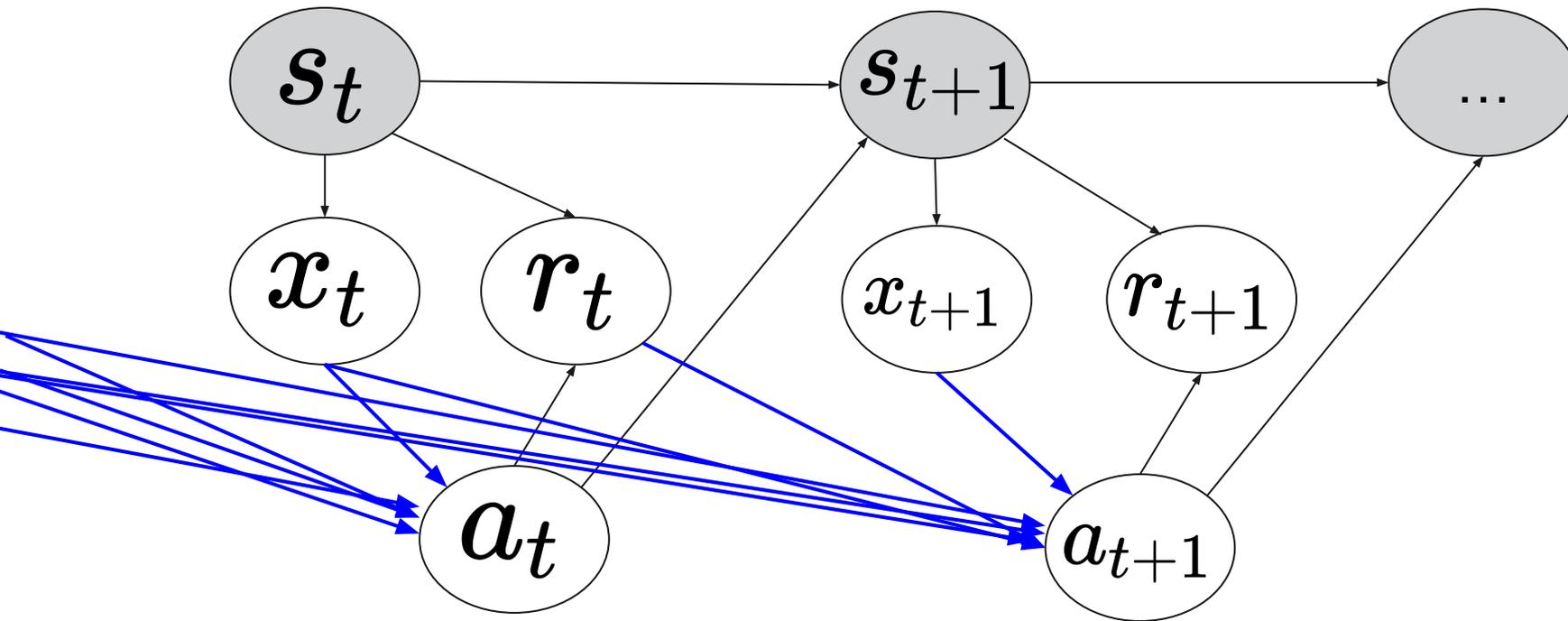
Background

Reinforcement Learning (RL): POMDP



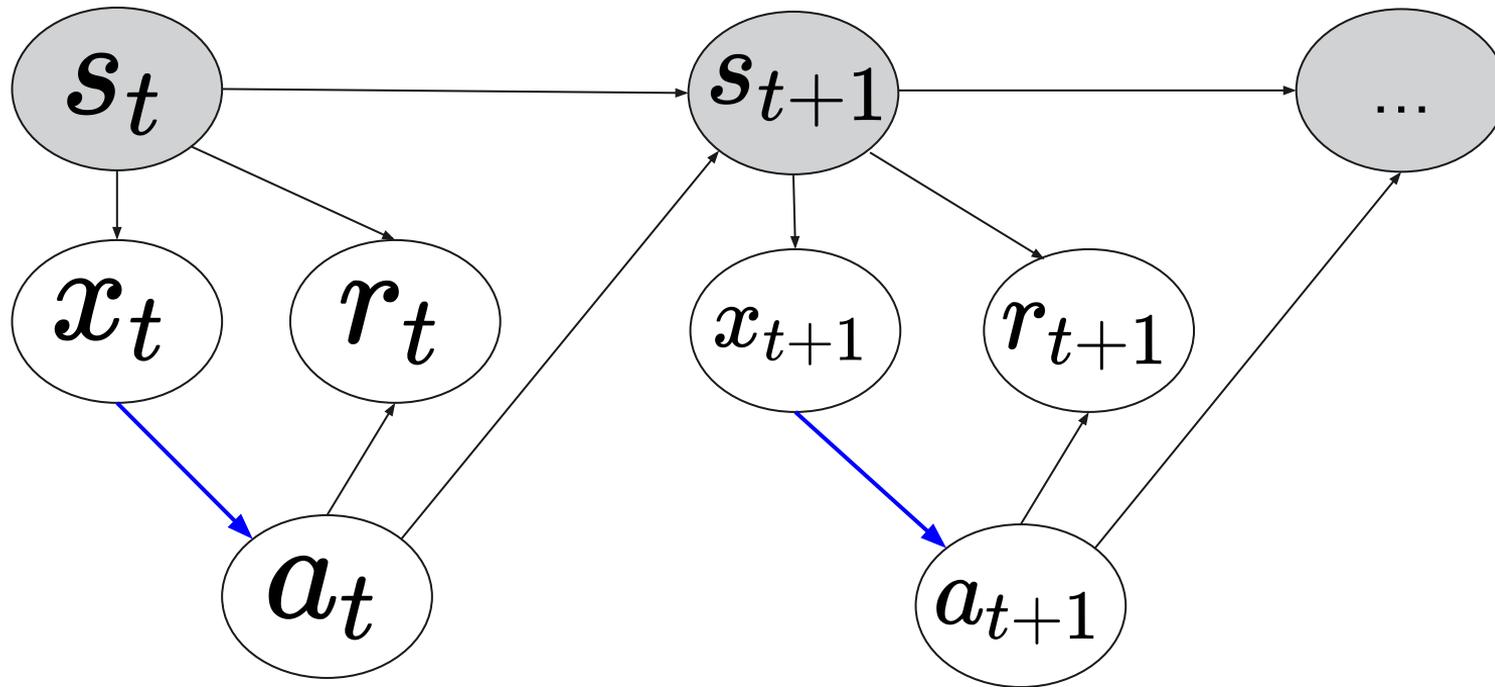
Background

Reinforcement Learning (RL): POMDP



Background

Reinforcement Learning (RL): POMDP



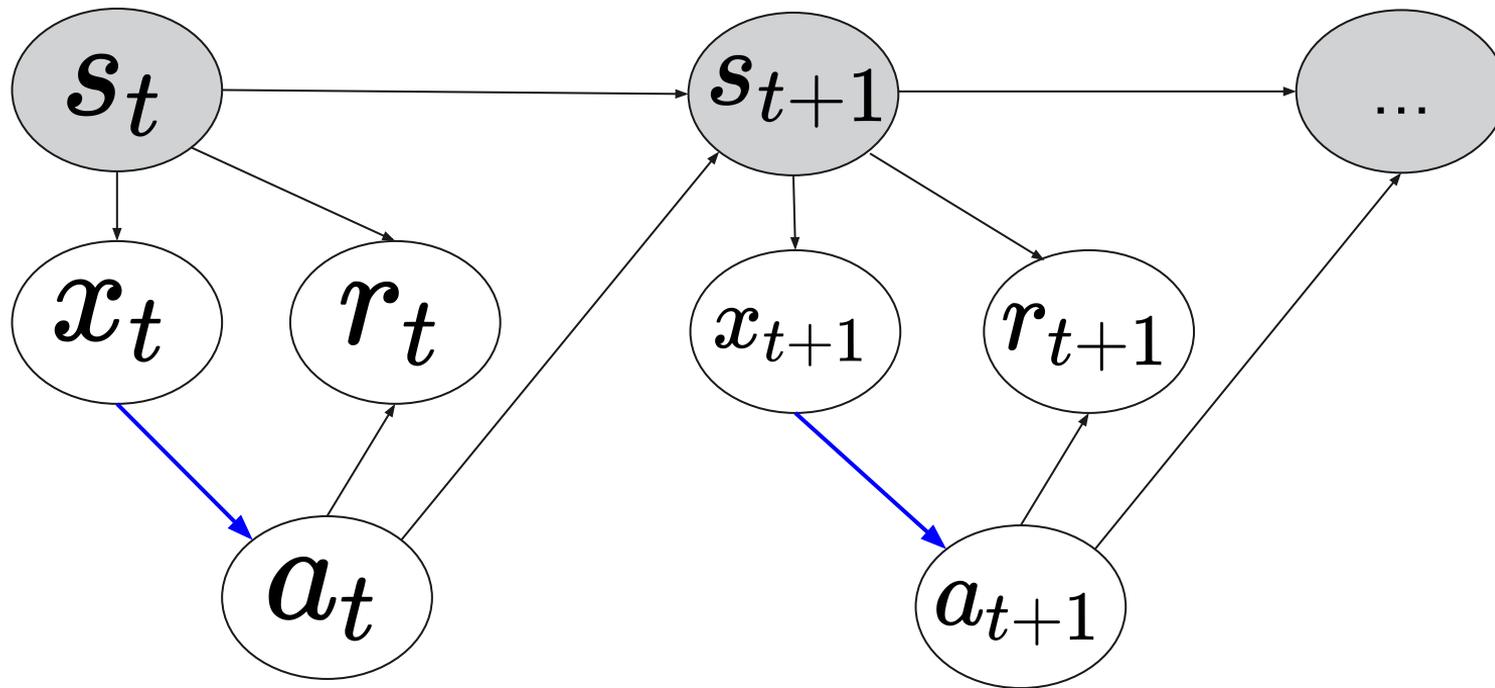
Background

Reinforcement Learning (RL): - POMDP



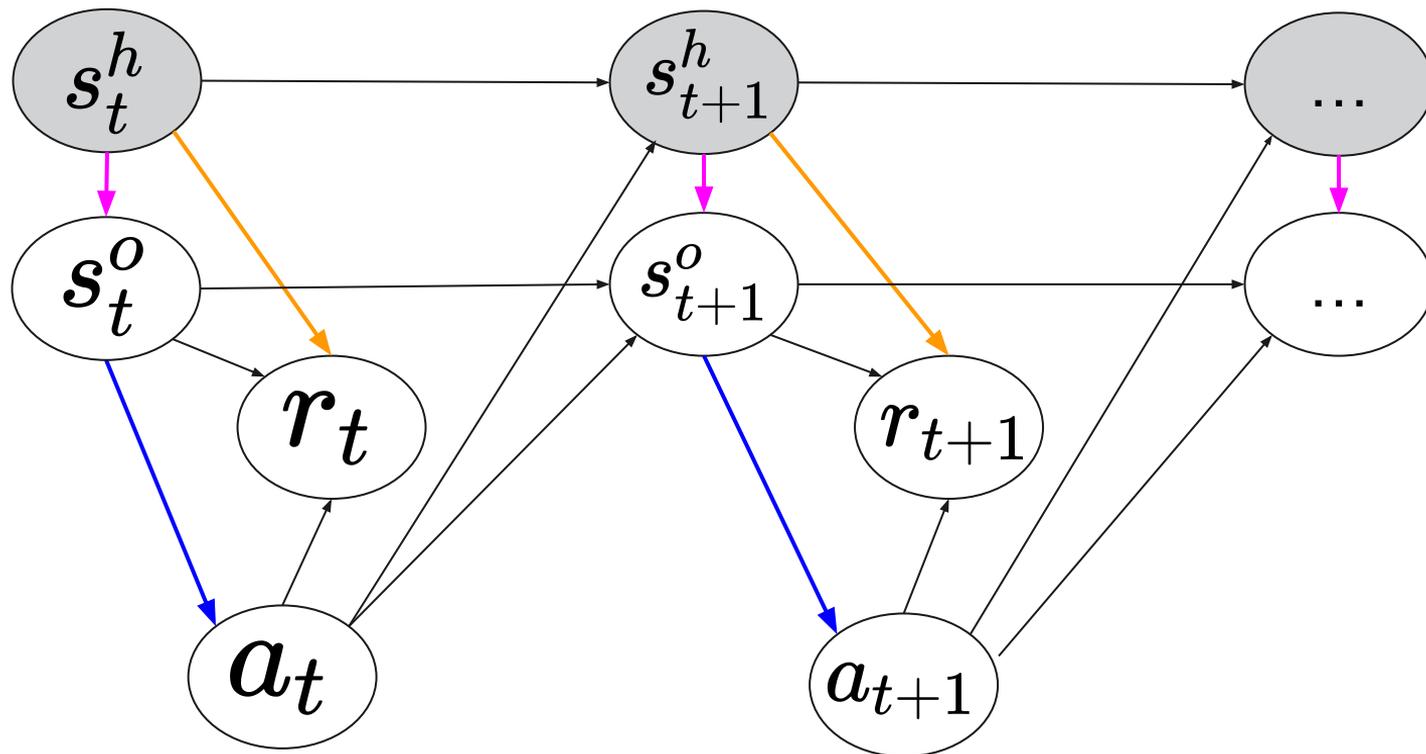
Task-agnostic Continual RL (TACRL)

TACRL as a POMDP special case



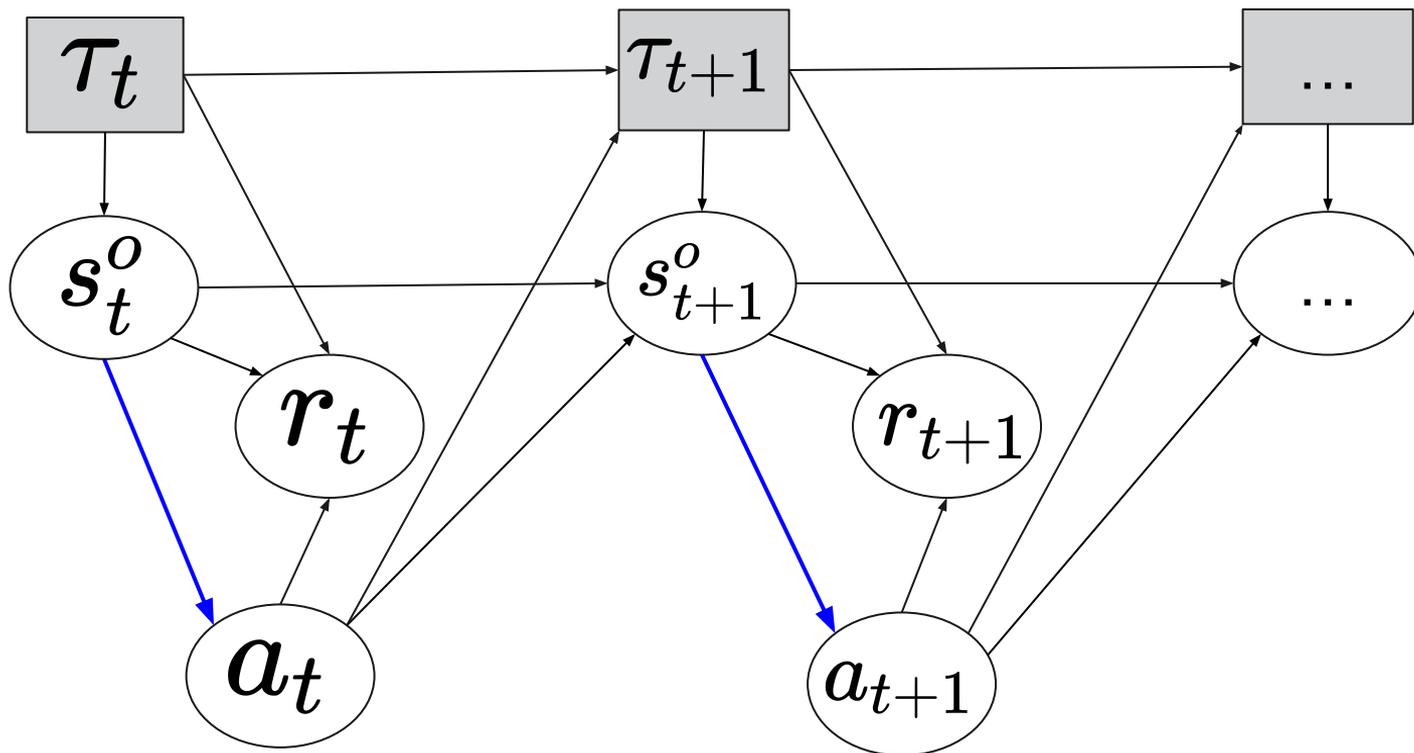
Task-agnostic Continual RL (TACRL)

TACRL as a POMDP special case



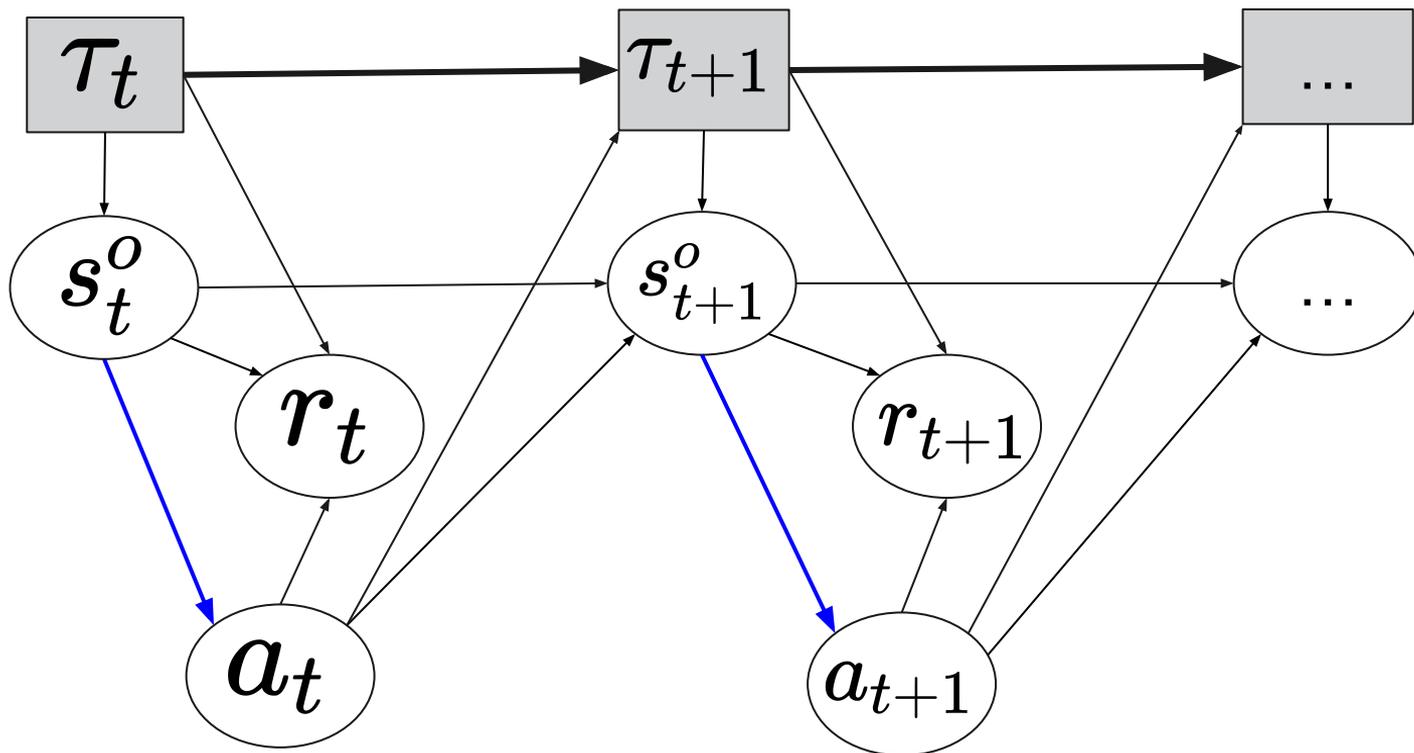
Task-agnostic Continual RL (TACRL)

TACRL as a POMDP special case



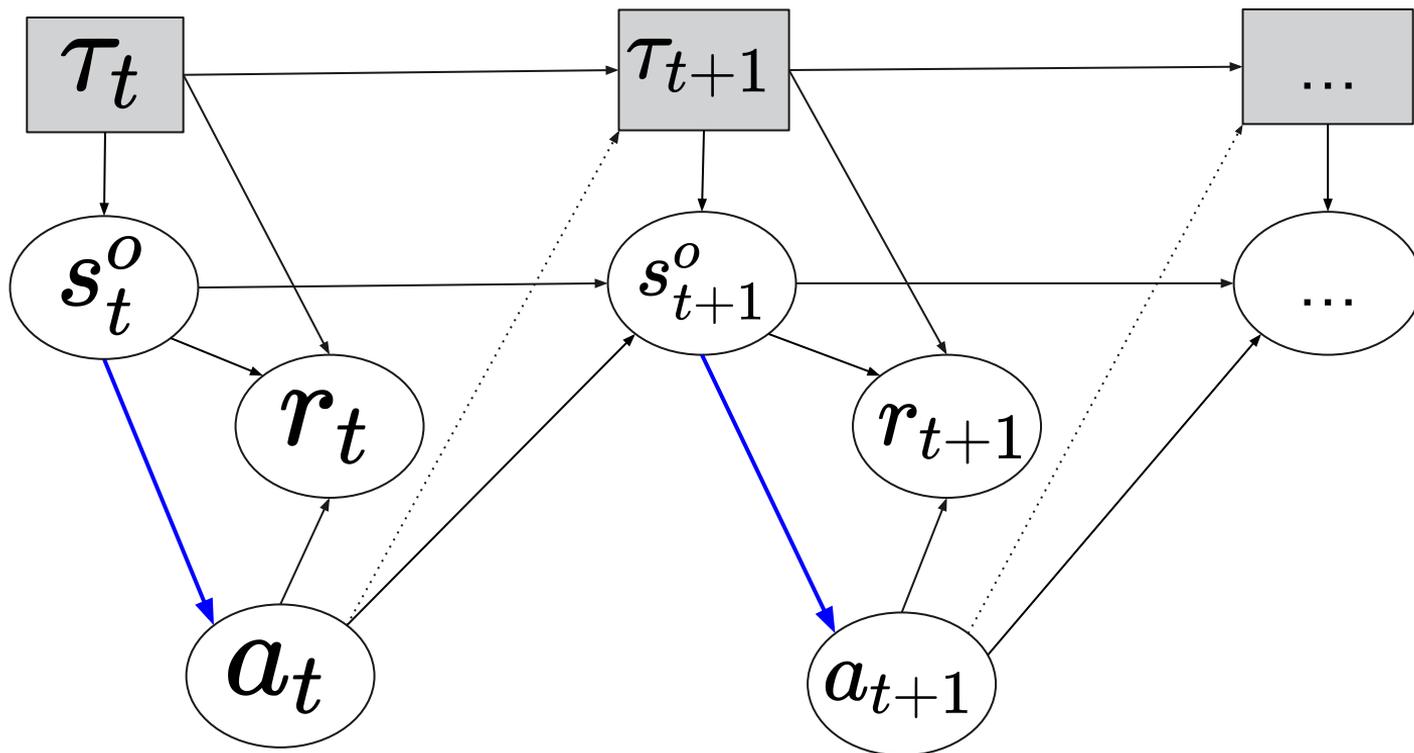
Task-agnostic Continual RL (TACRL)

TACRL as a POMDP special case



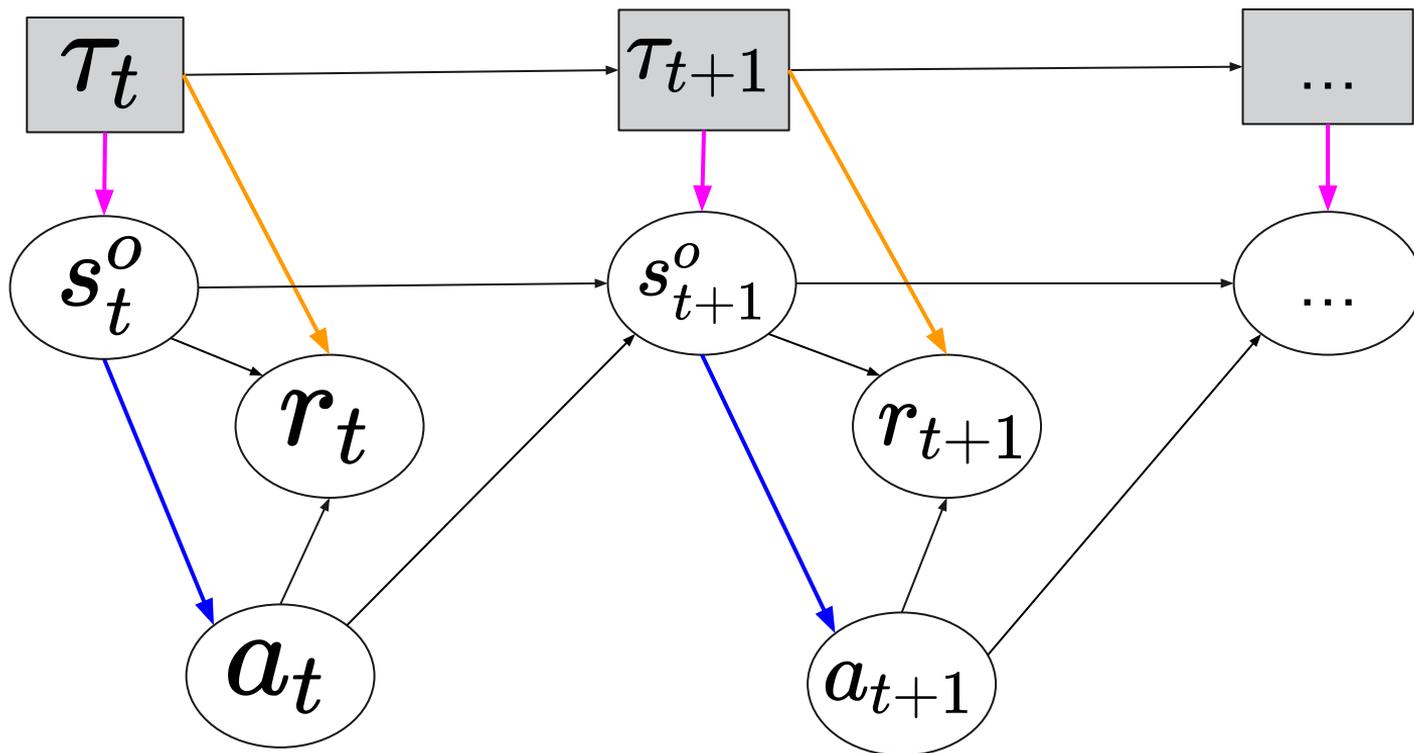
Task-agnostic Continual RL (TACRL)

TACRL as a POMDP special case



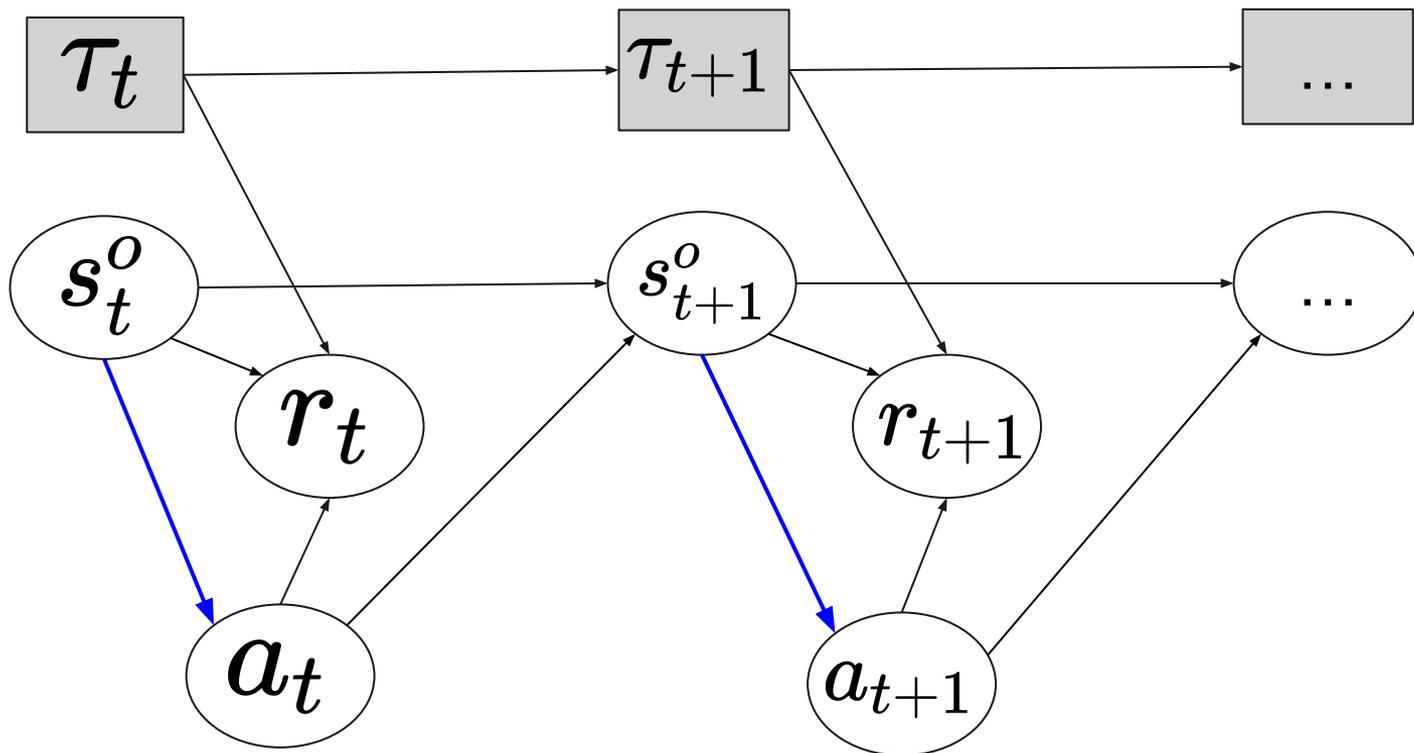
Task-agnostic Continual RL (TACRL)

Task-dependant **reward function** and **dynamics**



Task-agnostic Continual RL (TACRL)

TACRL (in our case)



Problem



POMDP

- + locally-stationary hidden states
- + Passive non-stationarity
- + single hidden state (task)

Task-Agnostic Continual RL

Task-Awareness soft upper Bound



- Provide task label to agent:
 - POMDP \rightarrow MDP

multi-task soft upper Bound



- Train on a stationary data distribution
 - *Catastrophic Forgetting* disappears

Task-Agnostic Continual RL (TACRL)

Related settings



	T	π	Objective	Evaluation
MDP [51]	$p(s_{t+1} s_t, a_t)$	$\pi(a_t s_t)$	$\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t]$	-
POMDP [24]	$p(s_{t+1}^h, s_{t+1}^o s_t^h, s_t^o, a_t)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h} \left[\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] s^h \right]$	-
HM-MDP [9]	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t) p(s_{t+1}^h s_t^h)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h} \left[\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] s^h \right]$	-
Task-agnostic CRL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t) p(s_{t+1}^h s_t^h)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h} \left[\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] s^h \right]$	$\mathbb{E}_{\tilde{s}^h} \left[\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] s^h \right]$
Task-Aware CRL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t) p(s_{t+1}^h s_t^h)$	$\pi(a_t s_t^h, s_t^o)$	$\mathbb{E}_{s^h} \left[\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] s^h \right]$	$\mathbb{E}_{\tilde{s}^h} \left[\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] s^h \right]$
Multi-task RL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t) p(s_{t+1}^h)$	$\pi(a_t s_t^h, s_t^o)$	$\mathbb{E}_{\tilde{s}^h} \left[\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] s^h \right]$	-

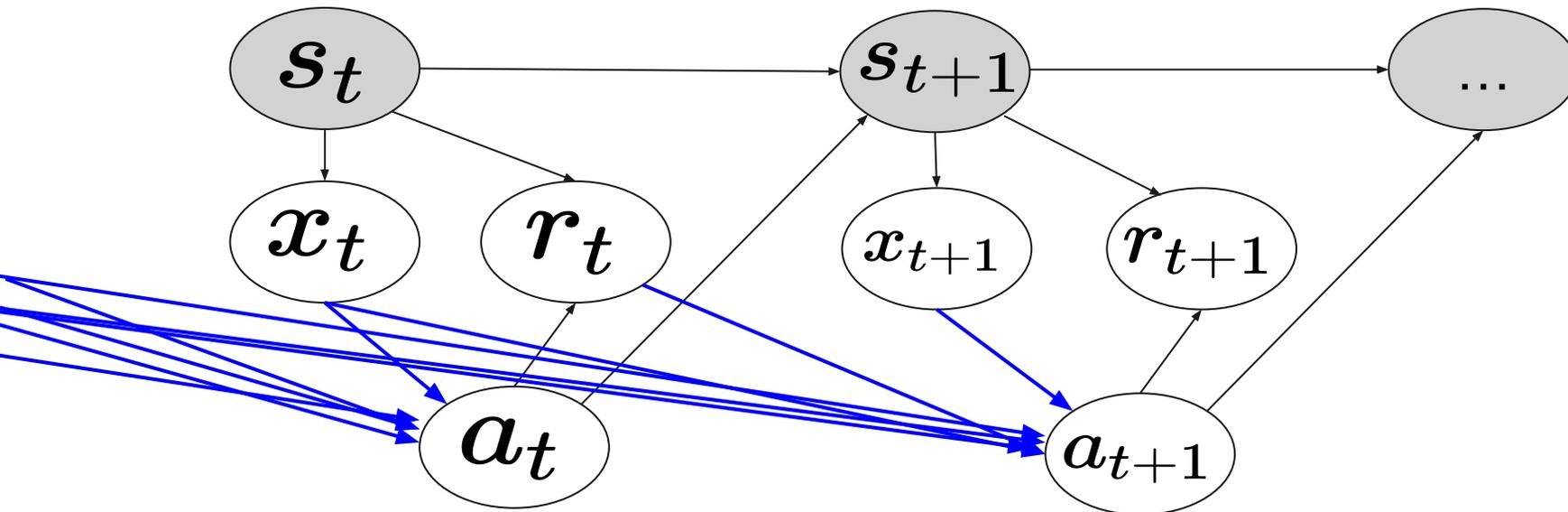
Table 1: Summarizing table of the settings relevant to TACRL. For readability purposes, \tilde{s}^h denotes the stationary distribution of s^h . The Evaluation column is left blank when it is equivalent to the Objective one.

Handling partial observability



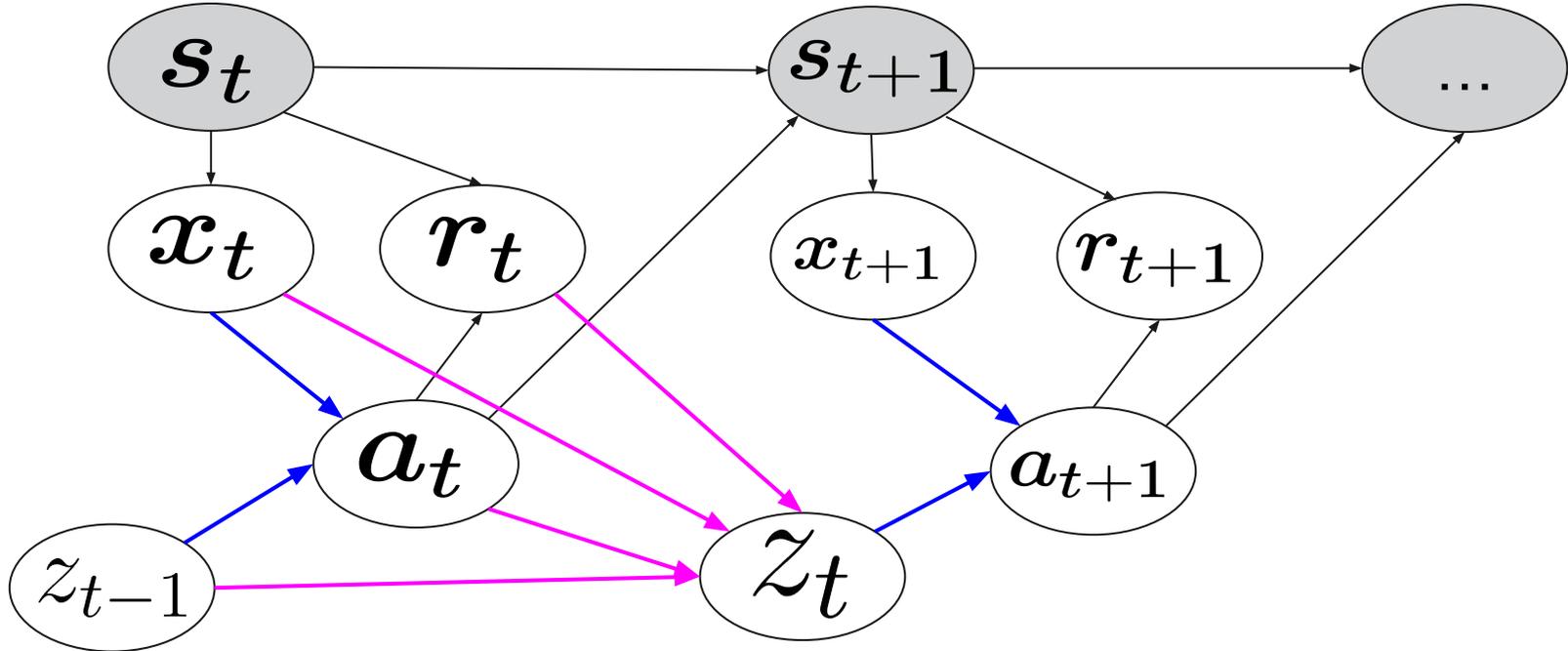
Background

Reinforcement Learning (RL): POMDP



Methods

Handling partial observability w/ a working memory



Replay-based Recurrent RL (3RL)

- Trick to handle partial observability → Working Memory (RNN)
- Trick to handle forgetting → Experience Replay
- RL algo + RNN + Experience Replay =



Other baselines



- Strategies
 - Fine Tuning
 - Experience Replay
 - Multi-task
- Task-Aware Modeling:
 - += Multi-head (MH)
 - += TaskID

SAC



Q-learning: $Q(s, a) \quad \pi(s) = \operatorname{argmax}_a Q(s, a)$

Doesn't scale!

Deep Q-learning (DQN) $Q_\theta(s, a) \quad \pi(s) = \operatorname{argmax}_a Q_\theta(s, a)$

What if my actions are continuous??

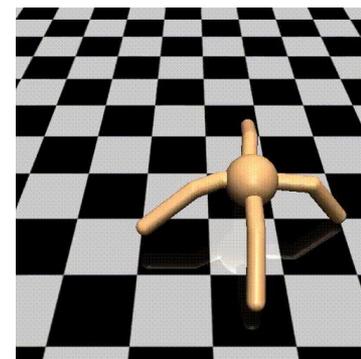
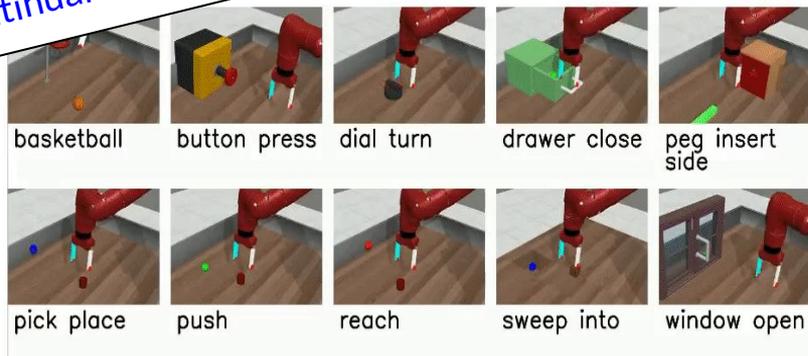
Soft-Actor critic (SAC) $\pi_\phi(s) \approx \operatorname{argmax}_\phi Q_\theta(s, \pi_\phi(s))$

Empirical Findings

Benchmark



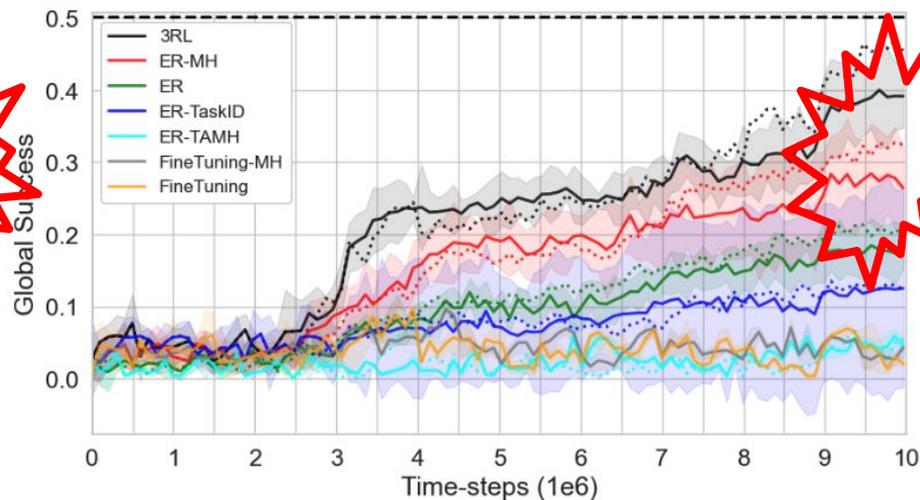
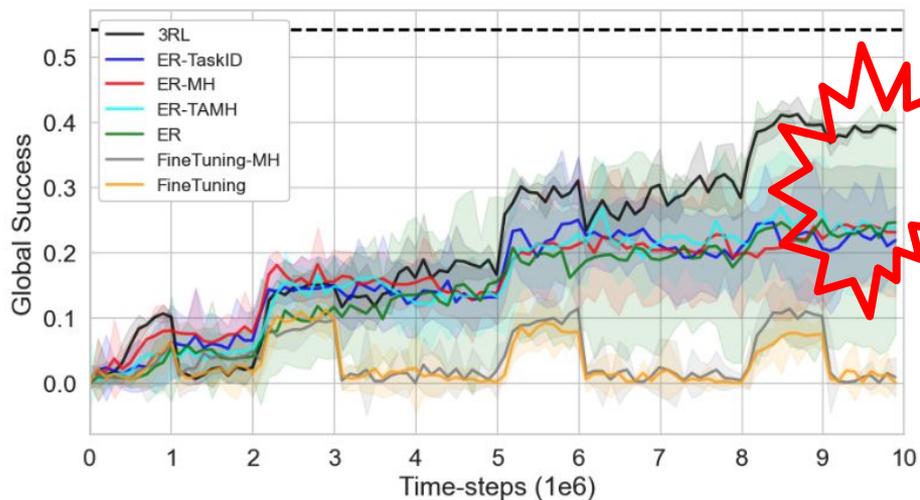
Continual World



Task similarity

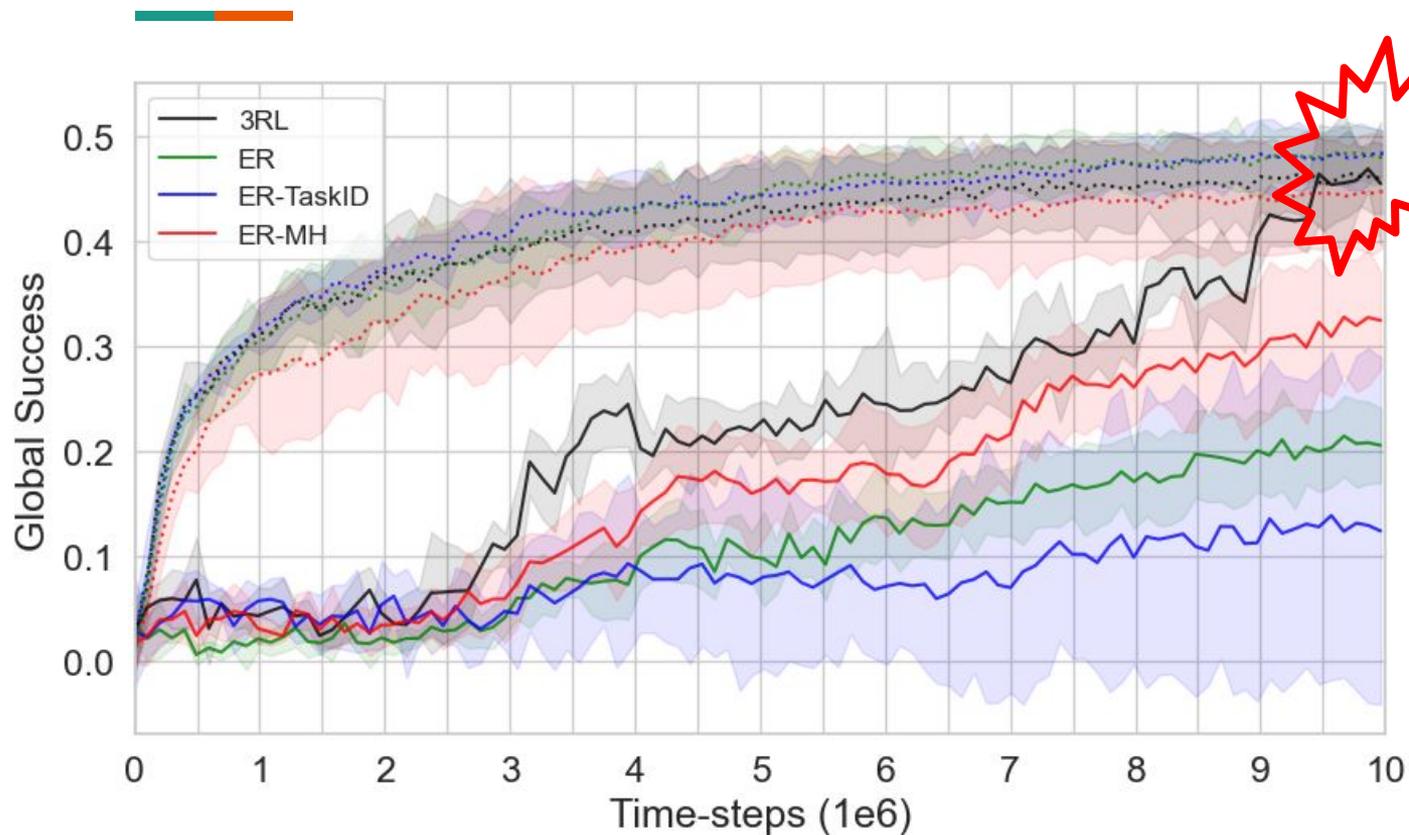
Empirical Findings

Task-agnostic > Task-Aware



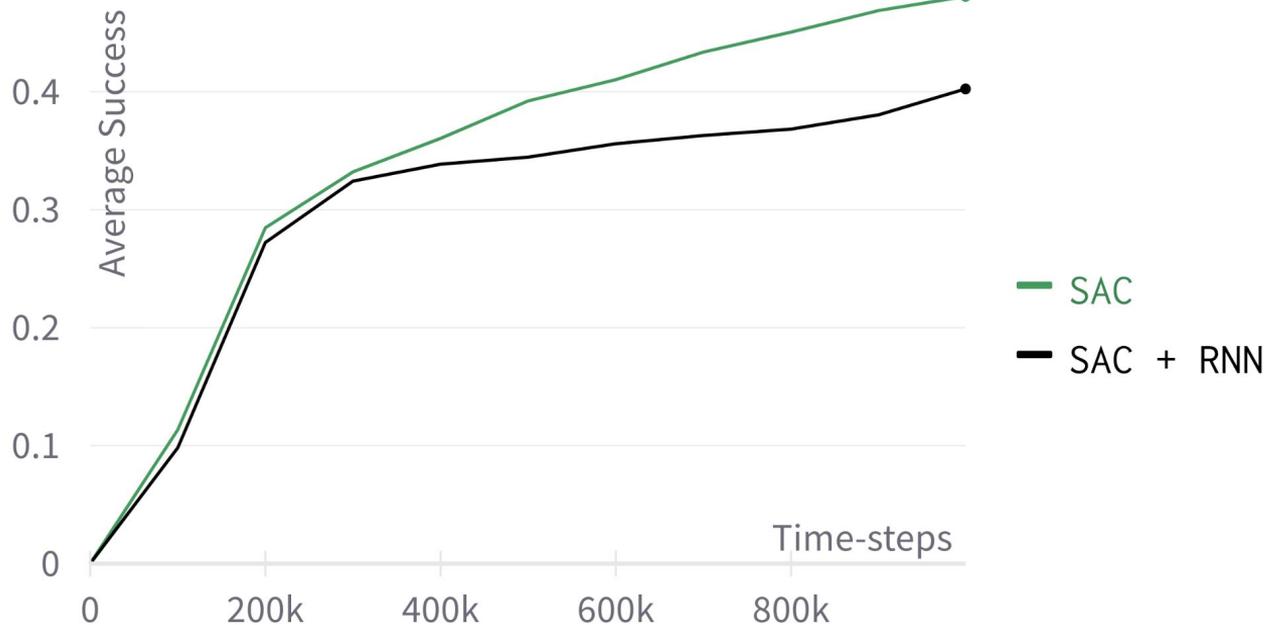
Empirical Findings

CL = MTL



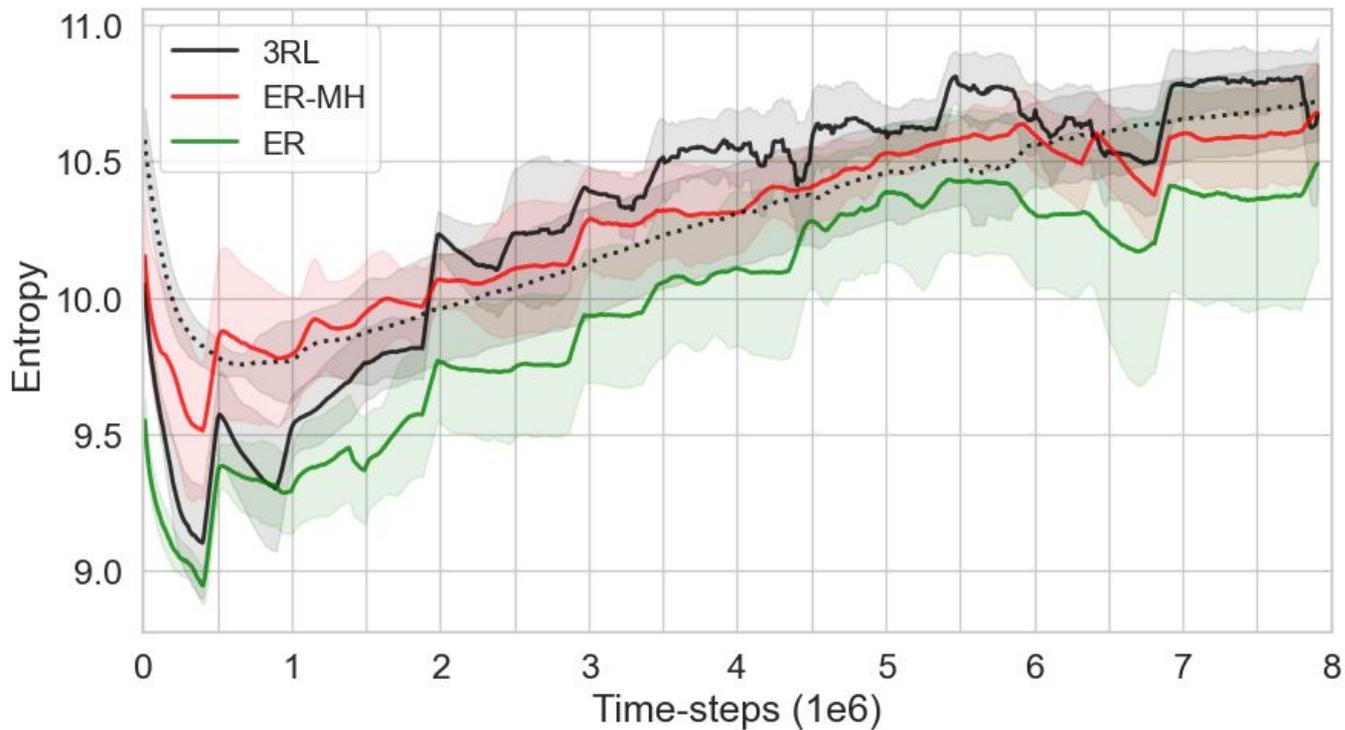
Empirical Findings

Hypothesis #1: rnn individually improves single-task performance



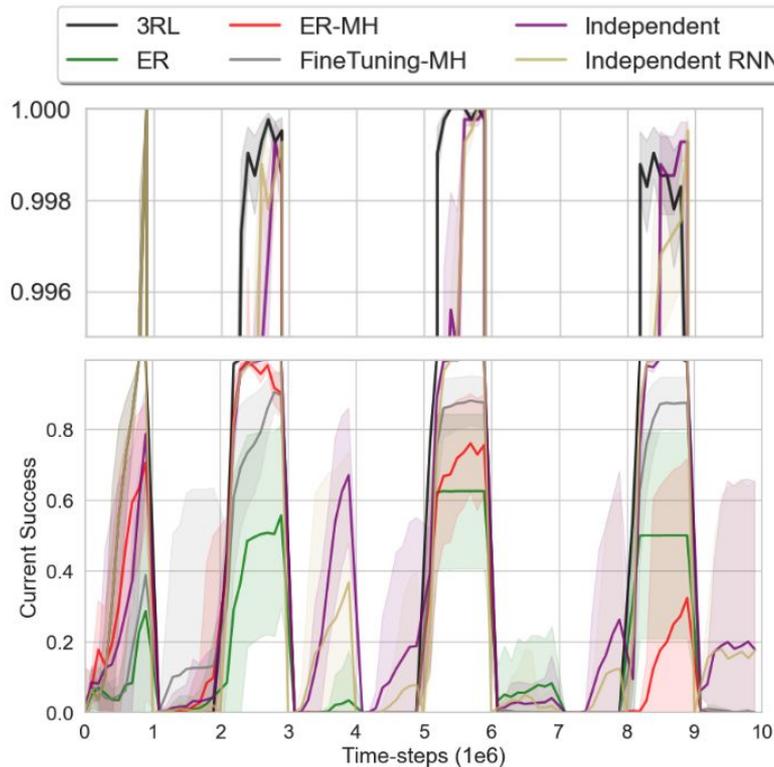
Empirical Findings

Hypothesis #2: increases parameter stability



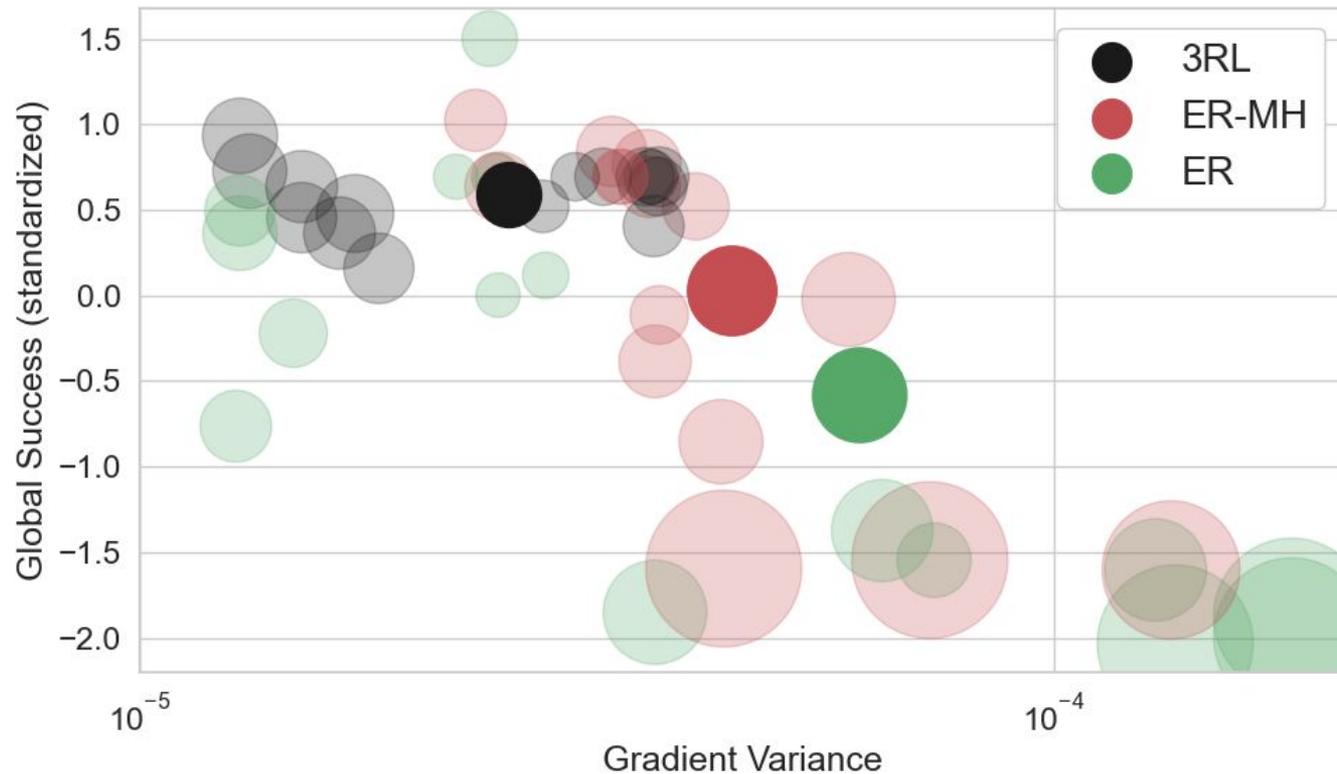
Empirical Findings

Hypothesis #3: rnn correctly places the new tasks in the context for previous ones



Empirical Findings

Hypothesis #3: rnn correctly places the new tasks in the context for previous ones

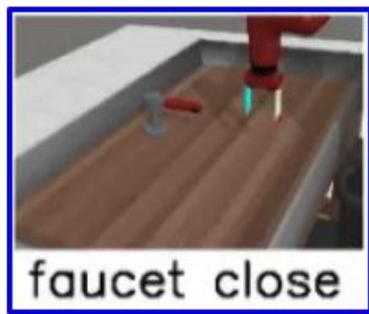


Empirical Findings

Qualitative analysis of the representations



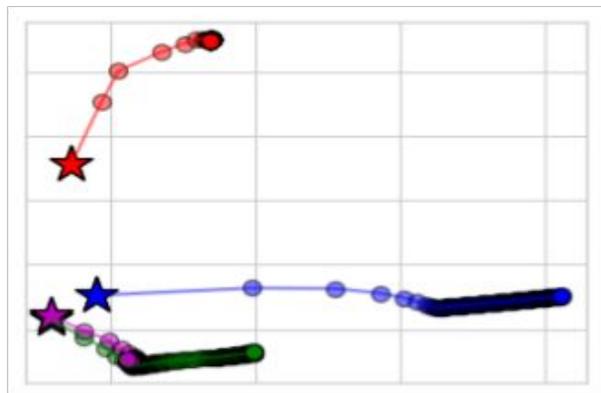
...



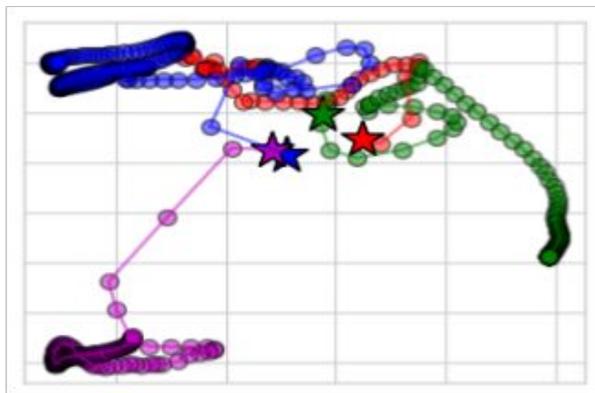
...



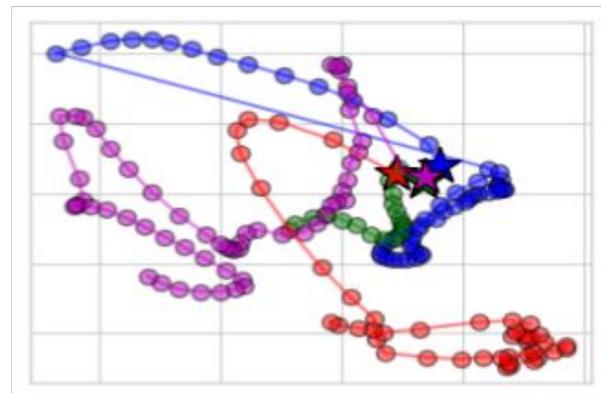
...



Before training



After one task



End of training



- Implications for CL research
- CSL vs CRL