

# The Online Discovery Problem and Its Application to Lifelong Reinforcement Learning

**Emma Brunskill**

Carnegie Mellon University

**Lihong Li**

Microsoft Research

**Multi-disciplinary Conference on Reinforcement Learning and Decision Making 2015**

Edmonton, AB, Canada

**Full version to be available on arXiv**

# Lifelong Learning Example: Intelligent Tutoring Systems



Alice



# Lifelong Learning Example: Intelligent Tutoring Systems



Alice



State = (courses taken, skills mastered, grades, ...)

# Lifelong Learning Example: Intelligent Tutoring Systems

Action  $\in$  { test skill,  
teach new concept,  
review old lectures,  
... }



Alice



State = (courses taken, skills mastered, grades, ...)

# Lifelong Learning Example: Intelligent Tutoring Systems

Action  $\in$  { test skill,  
teach new concept,  
review old lectures,  
... }



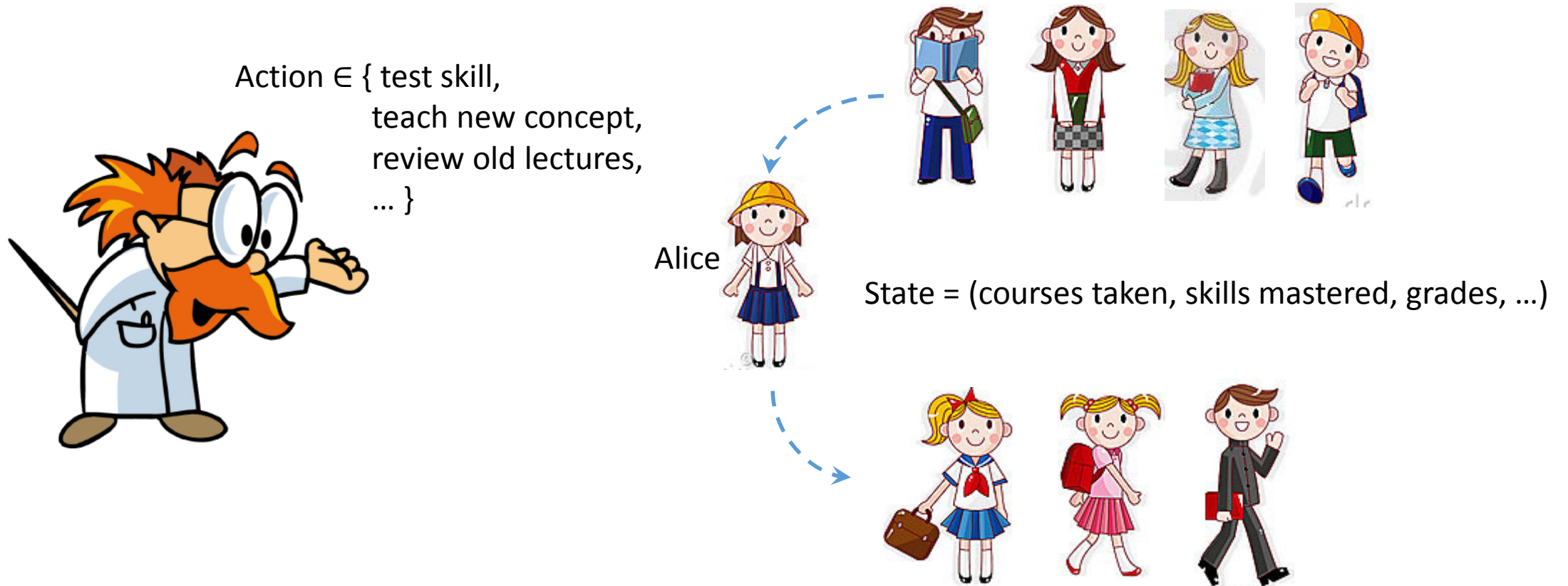
Alice



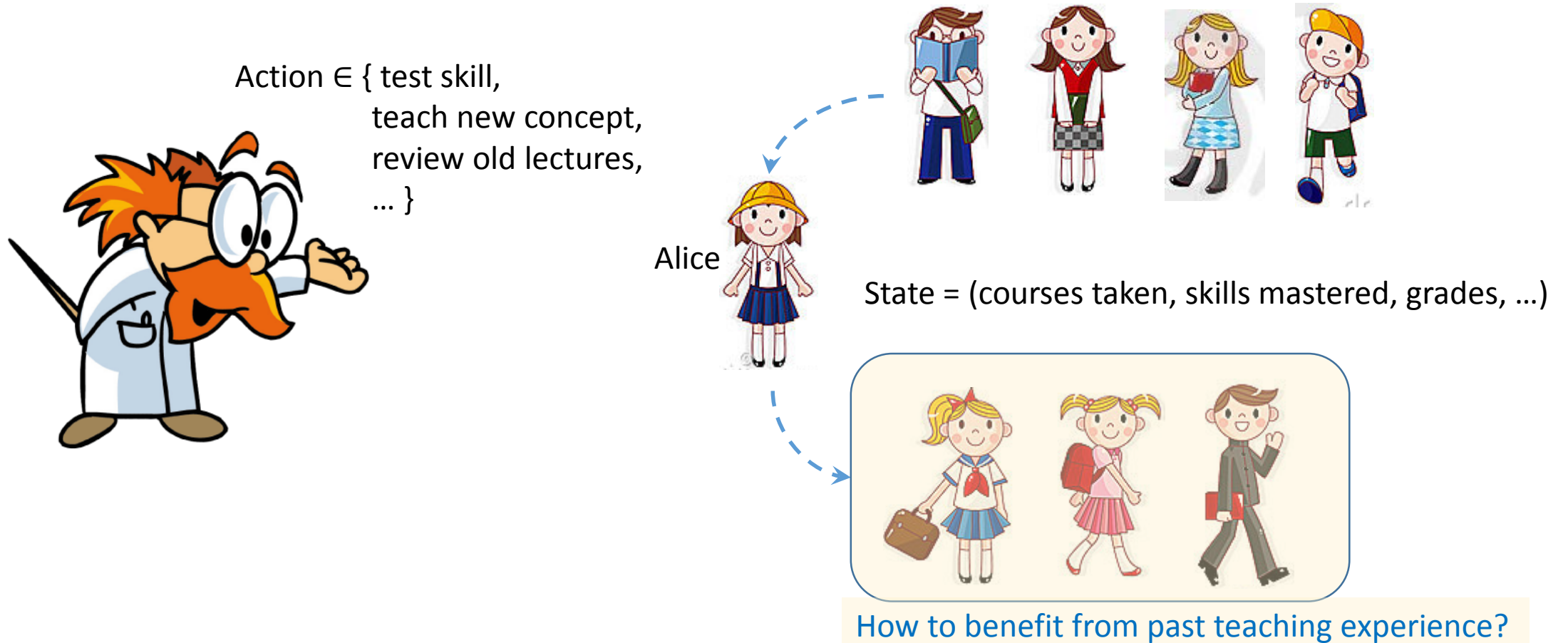
State = (courses taken, skills mastered, grades, ...)



# Lifelong Learning Example: Intelligent Tutoring Systems



# Lifelong Learning Example: Intelligent Tutoring Systems



# Lifelong Learning Example: Intelligent Tutoring Systems

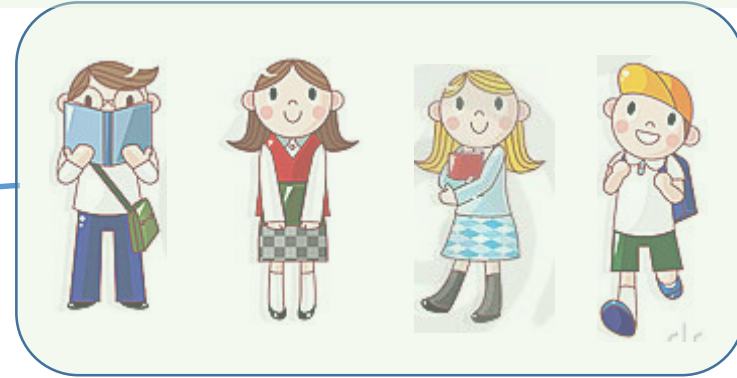


Action  $\in$  { test skill,  
teach new concept,  
review old lectures,  
... }

Alice



How to teach Alice to benefit future students?



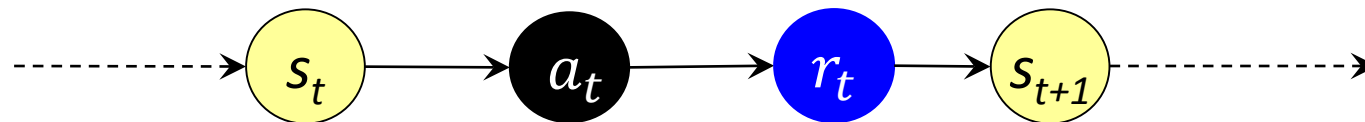
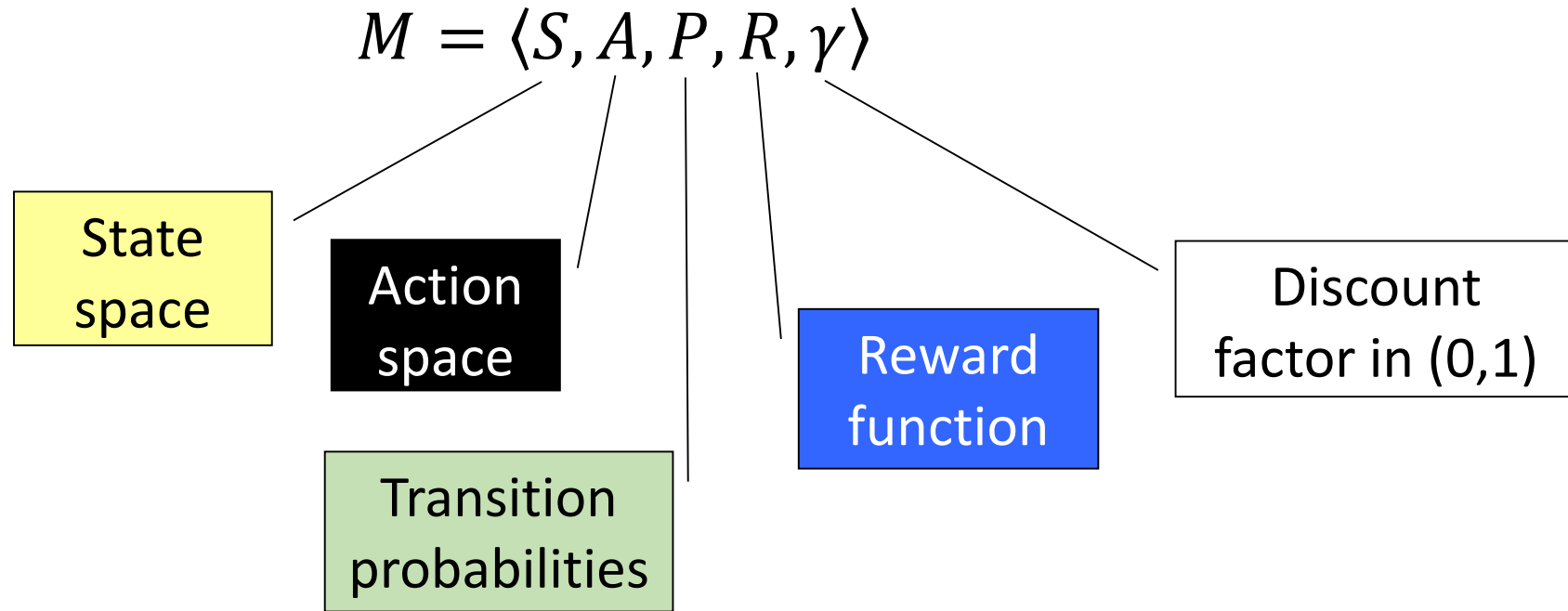
State = (courses taken, skills mastered, grades, ...)



How to benefit from past teaching experience?



# Task as Finite Markov Decision Process (MDP)



$$E[r_t] = R(s_t, a_t)$$

$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

# A Class of Lifelong RL Problems

- Given (known):  $S$  (finite),  $A$  (finite),  $\gamma \in (0,1)$
- Unknown:  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$   
 $\forall M \in \mathbf{M}, M = \langle S, A, P_M, R_M, \gamma \rangle$

For  $t = 1, 2, \dots, T$

- Environment chooses an unknown  $M_t \in \mathbf{M}$
- Agent acts in  $M_t$  for  $H$  steps

Note: Many previous works on LLRL  
with different setups

# A Class of Lifelong RL Problems

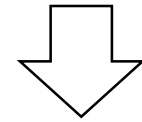
- Given (known):  $S$  (finite),  $A$  (finite),  $\gamma \in (0,1)$
- Unknown:  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$   
 $\forall M \in \mathbf{M}, M = \langle S, A, P_M, R_M, \gamma \rangle$

For  $t = 1, 2, \dots, T$

- Environment chooses an unknown  $M_t \in \mathbf{M}$
- Agent acts in  $M_t$  for  $H$  steps

Note: Many previous works on LLRL  
with different setups

Finite  $S$  and  $A$



Finitely many MDPs with  
"large" model differences

# A Class of Lifelong RL Problems

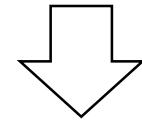
- Given (known):  $S$  (finite),  $A$  (finite),  $\gamma \in (0,1)$
- Unknown:  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$   
 $\forall M \in \mathbf{M}, M = \langle S, A, P_M, R_M, \gamma \rangle$

For  $t = 1, 2, \dots, T$

- Environment chooses an unknown  $M_t \in \mathbf{M}$
- Agent acts in  $M_t$  for  $H$  steps

Note: Many previous works on LLRL  
with different setups

Finite  $S$  and  $A$



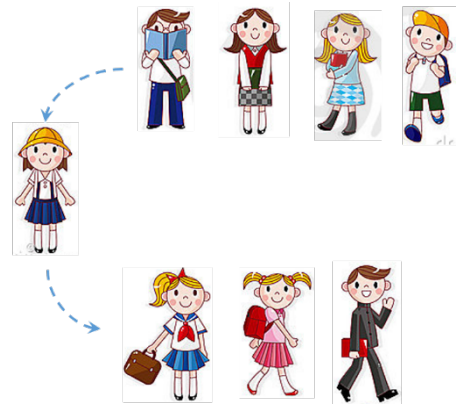
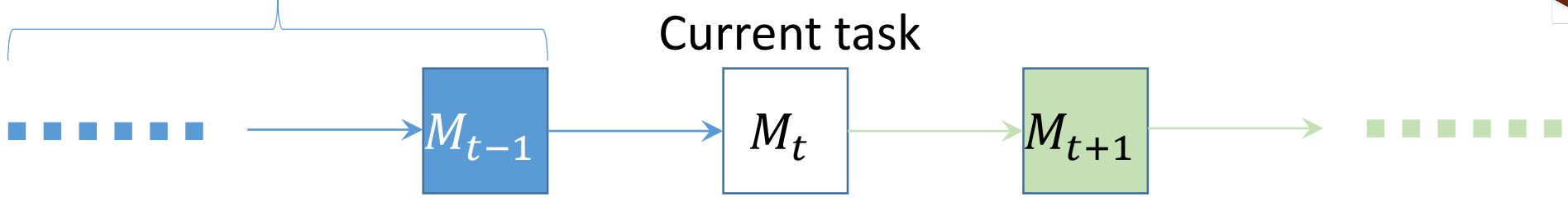
Finitely many MDPs with  
“large” model differences

## Examples

- Student types w/ varying learning rates [Liu&Koedinger]
- User types in human robot interaction [Nikolaidis et al.]
- User goal recognition for task assistance [Fern et al.]

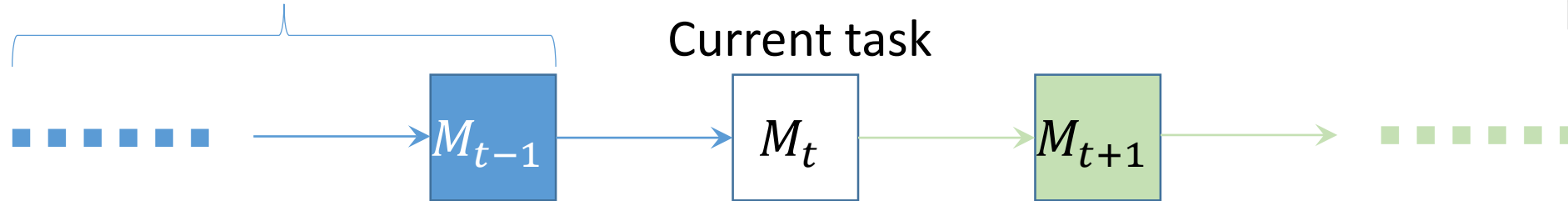
# Two Kinds of Exploration

$\hat{M}$ : set of discovered types before  $t$



# Two Kinds of Exploration

$\hat{\mathbf{M}}$ : set of discovered types before  $t$



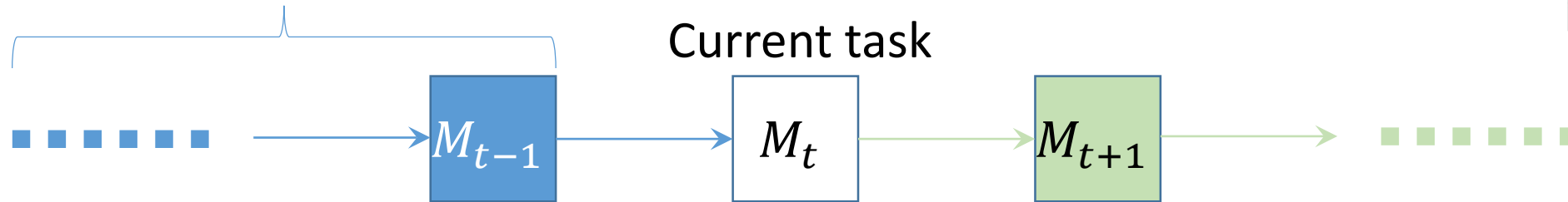
When in  $M_t$ ...

- Within-task learning:

- Cross-task knowledge transfer:

# Two Kinds of Exploration

$\hat{\mathbf{M}}$ : set of discovered types before  $t$



When in  $M_t$ ...

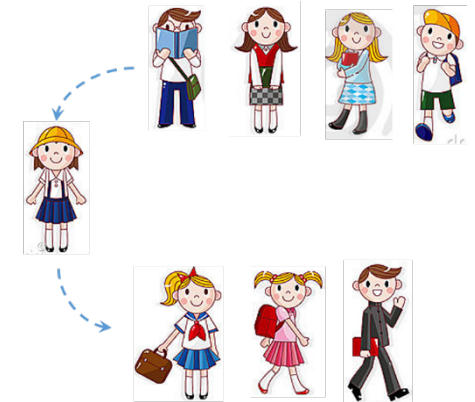
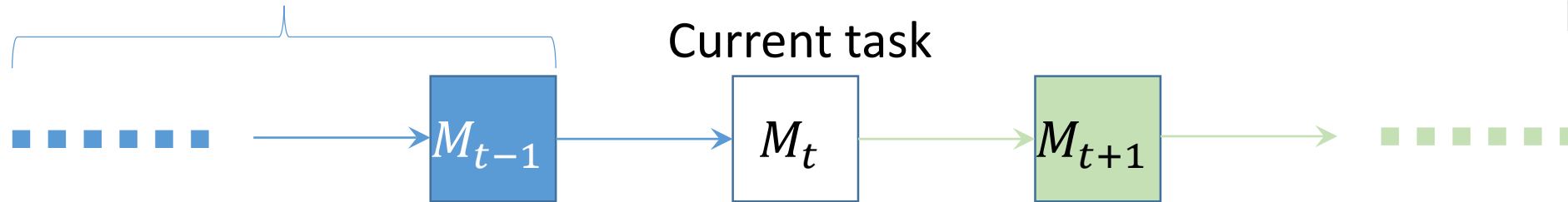
○ Within-task learning:

- Goal: maximize reward in  $M_t$
- Explore **promising** states in  $M_t$  until policy is  $\epsilon$ -optimal

○ Cross-task knowledge transfer:

# Two Kinds of Exploration

$\hat{\mathbf{M}}$ : set of discovered types before  $t$



When in  $M_t$ ...

○ Within-task learning:

- Goal: maximize reward in  $M_t$
- Explore **promising** states in  $M_t$  until policy is  $\epsilon$ -optimal

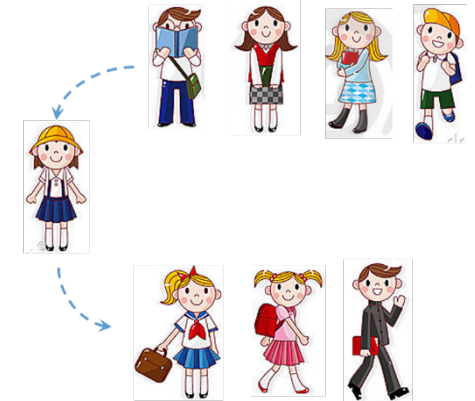
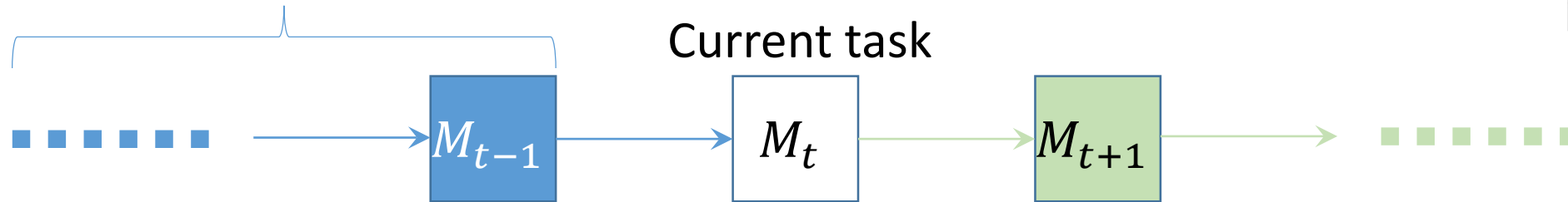
○ Cross-task knowledge transfer:

- Goal: maximize reward in  $M_{t+1}, \dots$  w/ transferable info.
- Explore **possibly all** states in  $M_t$  to discover novel types



# Two Kinds of Exploration

$\hat{\mathbf{M}}$ : set of discovered types before  $t$



When in  $M_t$ ...

○ Within-task learning:

- Goal: maximize reward in  $M_t$
- Explore **promising** states in  $M_t$  until policy is  $\epsilon$ -optimal

○ Cross-task knowledge transfer:

- Goal: maximize reward in  $M_{t+1}, \dots$  w/ transferable info.
- Explore **possibly all** states in  $M_t$  to discover novel types

Cross-task E/E tradeoff  
over  
within-task E/E tradeoff

# The *Online Discovery* Problem: Abstraction of Cross-task Exploration

**Environment** has an unknown set  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$

**Agent** starts with  $\widehat{\mathbf{M}} = \emptyset$

For  $t = 1, 2, \dots, T$

- $M_t \in \mathbf{M}$
- **Agent** chooses to explore ( $A_t = 1$ ) or exploit ( $A_t = 0$ )
  - If  $A_t = 1$ ,  $\widehat{\mathbf{M}} \leftarrow \widehat{\mathbf{M}} \cup \{M_t\}$
- Loss to agent

**Agent** aims to minimize total loss

# The *Online Discovery* Problem: Abstraction of Cross-task Exploration

## **MM**

**Environment** has an unknown set  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$

**Agent** starts with  $\widehat{\mathbf{M}} = \emptyset$

For  $t = 1, 2, \dots, T$

- **Environment** selects  $M_t \in \mathbf{M}$
- **Agent** chooses to explore ( $A_t = 1$ ) or exploit ( $A_t = 0$ )
  - If  $A_t = 1$ ,  $\widehat{\mathbf{M}} \leftarrow \widehat{\mathbf{M}} \cup \{M_t\}$
- Loss to agent

**Agent** aims to minimize total loss

# The *Online Discovery* Problem: Abstraction of Cross-task Exploration

1,  $\mathbf{M} \leftarrow \mathbf{M} \cup \{M_t\}$

1) or exploit ( $A_t = 0$ )

$\hat{\mathbf{M}}$

**Environment** has an unknown set  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$

**Agent** starts with  $\hat{\mathbf{M}} = \emptyset$

For  $t = 1, 2, \dots, T$

- If  $A_t = 1$ ,  $\hat{\mathbf{M}} \leftarrow \hat{\mathbf{M}} \cup \{M_t\}$
- **Agent** chooses to explore ( $A_t = 1$ ) or exploit ( $A_t = 0$ )
  - If  $A_t = 1$ ,  $\hat{\mathbf{M}} \leftarrow \hat{\mathbf{M}} \cup \{M_t\}$
- Loss to agent

# The *Online Discovery* Problem: Abstraction of Cross-task Exploration

$$1, \mathbf{M} \leftarrow \mathbf{M} \cup \{M_t\}$$

1) or exploit ( $A_t = 0$ )

$\mathbf{M}$

**Environment** has an unknown set  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$

**Agent** starts with  $\hat{\mathbf{M}} = \emptyset$

For  $t = 1, 2, \dots, T$

- Loss to agent

- Agent** chooses

- If  $A_t = 1$ ,  $\mathbf{M} \leftarrow \mathbf{M} \cup \{M_t\}$

- Loss to agent

	$M_t \in \hat{\mathbf{M}}$	$M_t \notin \hat{\mathbf{M}}$
$A_t = 0$	$\rho_0$	$\rho_3$
$A_t = 1$	$\rho_1$	$\rho_2$

$$(\rho_0 \ll \rho_1 \leq \rho_2 \ll \rho_3)$$

1) or exploit ( $A_t = 0$ )

# The *Online Discovery* Problem: Abstraction of Cross-task Exploration

$$1, \mathbf{M} \leftarrow \mathbf{M} \cup \{M_t\}$$

1) or exploit ( $A_t = 0$ )

$\mathbf{M}$

**Environment** has an unknown set  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$

**Agent** starts with  $\hat{\mathbf{M}} = \emptyset$

For  $t = 1, 2, \dots, T$

- Loss to agent

- Agent** chooses

- If  $A_t = 1$ ,  $\mathbf{M} \leftarrow \mathbf{M} \cup \{M_t\}$

	$M_t \in \hat{\mathbf{M}}$	$M_t \notin \hat{\mathbf{M}}$
$A_t = 0$	$\rho_0$	$\rho_3$
$A_t = 1$	$\rho_1$	$\rho_2$

$$(\rho_0 \ll \rho_1 \leq \rho_2 \ll \rho_3)$$

1) or exploit ( $A_t = 0$ )

**Agent** aims to minimize total loss

# The *Online Discovery* Problem: Abstraction of Cross-task Exploration

$$1, \mathbf{M} \leftarrow \mathbf{M} \cup \{M_t\}$$

1) or exploit ( $A_t = 0$ )

$\mathbf{M}$

**Environment** has an unknown set  $\mathbf{M} = \{M^1, M^2, \dots, M^C\}$

**Agent** starts with  $\hat{\mathbf{M}} = \emptyset$

For  $t = 1, 2, \dots, T$

- Loss to agent

- Agent** chooses

- If  $A_t = 1$ ,  $\mathbf{M} \leftarrow \mathbf{M} \cup \{M_t\}$

	$M_t \in \hat{\mathbf{M}}$	$M_t \notin \hat{\mathbf{M}}$
$A_t = 0$	$\rho_0$	$\rho_3$
$A_t = 1$	$\rho_1$	$\rho_2$

1) or exploit ( $A_t = 0$ )  
 $\rho_0$ : successful transfer  
 $\rho_3$ : negative transfer

**Agent** aims to minimize total loss

# Explore-First Algorithm

Stochastic assumptions:

$$M_t \sim \mu \text{ i.i.d. over } \mathbf{M}, \text{ and } \mu_m := \min_{M \in \mathbf{M}} \mu(M)$$

- Action selection

$$A_t = \begin{cases} 1 & \text{if } t \leq E & \text{(Exploration phase)} \\ 0 & \text{otherwise} & \text{(Exploitation phase)} \end{cases}$$

- $\text{AverageLoss} = O(\mu_m^{-1} \log(C \mu_m T))$ , then

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{T \cdot \mu_m} \log \left( \frac{TC \mu_m \rho_3}{\rho_1} \right)$$



# Explore-First Algorithm

$$\leq \text{OptLoss} + \frac{1}{T \cdot \mu_m} \log \left( \frac{TC \mu_m \rho_3}{\rho_1} \right)$$

Stochastic assumptions:  $M_t \sim \mu$  i.i.d. over  $\mathbf{M}$ , and  $\mu_m := \min_{M \in \mathbf{M}} \mu(M)$

- Action selection

$$A_t = \begin{cases} 1 & \text{if } t \leq E \quad (\text{Exploration phase}) \\ 0 & \text{otherwise} \quad (\text{Exploitation phase}) \end{cases}$$

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{T \cdot \mu_m} \log \left( \frac{TC \mu_m \rho_3}{\rho_1} \right)$$

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{T \cdot \mu_m} \log \left( \frac{TC \mu_m \rho_3}{\rho_1} \right)$$

# Forced-Exploration Algorithm

No stochastic assumption ( $M_t$  can even be generated **adversarially!**)

- $\eta_1 \geq \eta_2 \geq \dots \geq \eta_T > 0$

- Algorithm chooses action

$$A_t \sim \text{Bernoulli}(\eta_t)$$

- **Theorem**: If choose  $\eta_t = 1/\sqrt{t}$ , then

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{\sqrt{T}} (2\rho_1 + C\rho_3)$$

# Forced-Exploration Algorithm

Bernoulli  $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$

$$\eta_1 \geq \eta_2 \geq \dots \geq \eta_T > 0$$

No stochastic assumption ( $M_t$  can even be generated **adversarially!**)

- Algorithm chooses action

$$A_t \sim \text{Bernoulli}(\eta_t)$$

$$A_t \sim \text{Bernoulli}(\eta_t)$$

- Theorem:** If choose  $\eta_t = 1/\sqrt{t}$ , then

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{\sqrt{T}} (2\rho_1 + C\rho_3)$$

# Forced-Exploration Algorithm

$$\text{OptLoss} + \frac{1}{\sqrt{T}} (2\rho_1 + C\rho_3)$$

1/ t t t t, then

Bernoulli  $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$   $\eta_t$

$$\eta_1 \geq \eta_2 \geq \dots \geq \eta_T > 0$$

No stochastic assumption ( $M_t$  can even be generated **adversarially!**)

- Algorithm chooses action

$$A_t \sim \text{Bernoulli}(\eta_t)$$

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{\sqrt{T}} (2\rho_1 + C\rho_3)$$

- Theorem:** If choose  $\eta_t = 1/\sqrt{t}$ , then

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{\sqrt{T}} (2\rho_1 + C\rho_3)$$

# Forced-Exploration Algorithm

$$\text{OptLoss} + \frac{1}{\sqrt{T}} (2\rho_1 + C\rho_3)$$

1/ t t t t, then

Bernoulli  $\eta_t$

$$\eta_1 \geq \eta_2 \geq \dots \geq \eta_T > 0$$

No stochastic assumption ( $M_t$  can even be generated **adversarially!**)

- Algorithm chooses action

$$A_t \sim \text{Bernoulli}(\eta_t)$$

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{\sqrt{T}} (2\rho_1 + C\rho_3)$$

- Theorem:** If choose  $\eta_t = 1/\sqrt{t}$ , then

$$\text{AverageLoss} \leq \text{OptLoss} + \frac{1}{\sqrt{T}} (2\rho_1 + C\rho_3)$$

We have an  $\Omega\left(\frac{1}{\sqrt{T}}\right)$  lower bound

**→ Forced-Exploration is essentially optimal**

# A Lifelong RL Algorithm based on FE

**Input:**  $S, A, \gamma$

Initia

# A Lifelong RL Algorithm based on FE

**Input:**  $S, A, \gamma$

Initia

# A Lifelong RL Algorithm based on FE

**Input:**  $S, A, \gamma$

Initia



# A Lifelong RL Algorithm based on FE

**Input:**  $S, A, \gamma$

Initia

$= t t t = 1$  if data shows  $M_t$  is novel

# Sample Complexity of Exploration

Sample complexity of algorithm **A** (given  $\epsilon$ ) [Kakade]

Number of steps where  $Q^{\mathbf{A}}_t(s_t, a_t) \leq Q^*(s_t, a_t) - \epsilon$

Measures number of  $\epsilon$ -mistakes made by the algorithm

- \_\_\_\_\_ long enough, with high prob.

$$\text{SampleComplexity(Our Algorithm)} = \tilde{O}\left(\frac{CD}{\Gamma^2} T + SAN\sqrt{T}\right)$$

In contrast, single-task RL's Sample Complexity is  $\Omega(SAT)$

# Sample Complexity of Exploration

Our Algorithm  $O\left(\frac{CD}{\Gamma^2} T + SAN\sqrt{T}\right)$

Sample complexity of algorithm **A** (given  $\epsilon$ ) [Kakade]

Number of steps where  $Q^{\mathbf{A}_t}(s_t, a_t) \leq Q^*(s_t, a_t) - \epsilon$

Measures number of  $\epsilon$ -mistakes made by the algorithm

- **Theorem:** For  $H$  long enough, with high prob.

$$\text{SampleComplexity Our Algorithm} = \tilde{O}\left(\frac{CD}{\Gamma^2} T + SAN\sqrt{T}\right)$$

In contrast, single-task RL's Sample Complexity is  $\Omega(SAT)$

# Sample Complexity of Exploration

Our Algorithm =  $O\left(\frac{CD}{\Gamma^2} T + SAN T\right)$

Sample complexity of algorithm **A** (given  $\epsilon$ ) [Kakade]

Number of steps where  $Q^{\mathbf{A}_t}(s_t, a_t) \leq Q^*(s_t, a_t) - \epsilon$

Measures number of  $\epsilon$ -mistakes made by the algorithm

- **Theorem:** For  $H$  long enough, with high prob.

Sample Complexity Our Algorithm =  $\tilde{O}\left(\frac{CD}{\Gamma^2} T + SAN\sqrt{T}\right)$   
 Asymptotic performance

In contrast, single-task RL's Sample Complexity is  $\Omega(SAT)$

# Sample Complexity of Exploration

*SAT SSAATT SAT*

Our Algorithm =  $O\left(\frac{CD}{\Gamma^2} T + SAN\sqrt{T}\right)$   
*CCDD CD  $\Gamma^2$   $\Gamma^2$   $\Gamma^2$   $\Gamma^2$   $\Gamma^2$   $\Gamma^2$   $CD$   $\Gamma^2$   $TT+SSAANN$   $T$   $TTT$   $T$   $C$   
*D  $\Gamma^2$   $T+SAN$   $T$**

Sample complexity of algorithm **A** (given  $\epsilon$ ) [Kakade]

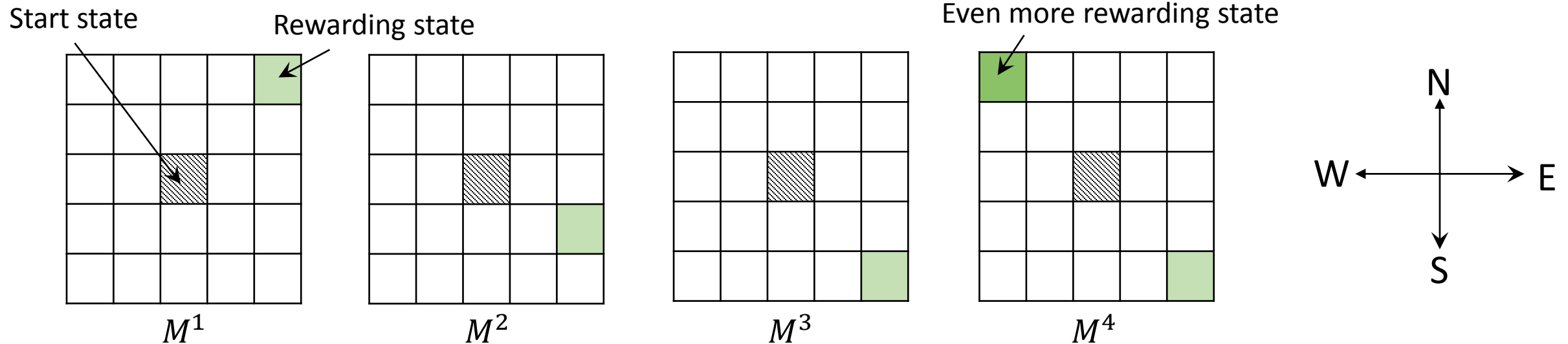
Number of steps where  $Q^{\mathbf{A}_t}(s_t, a_t) \leq Q^*(s_t, a_t) - \epsilon$

Measures number of  $\epsilon$ -mistakes made by algorithm

- **Theorem:** For  $H$  long enough, with high prob.

$$\text{SampleComplexity Our Algorithm} = \tilde{O}\left(\frac{CD}{\Gamma^2} T + SAN\sqrt{T}\right)$$

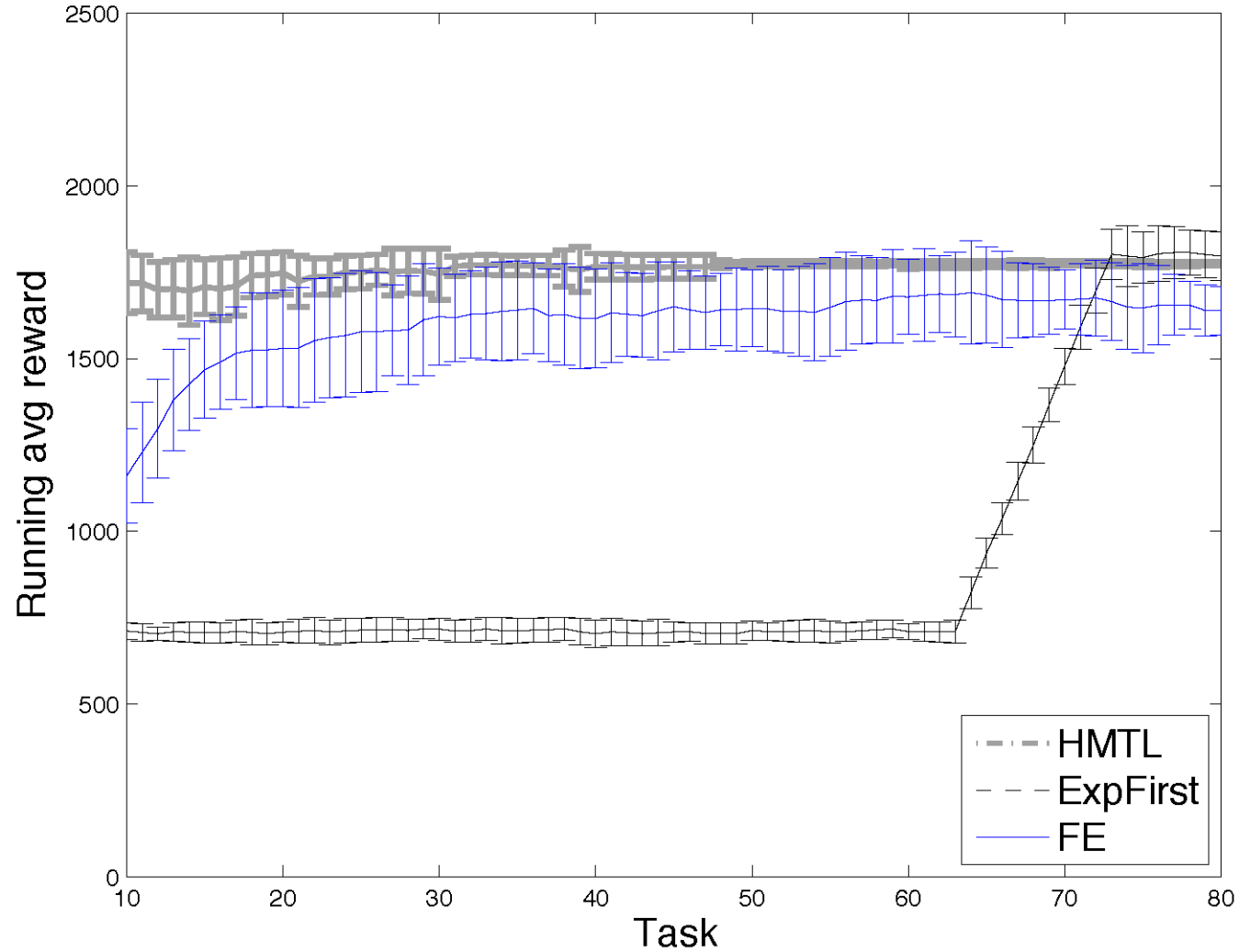
# Experiment



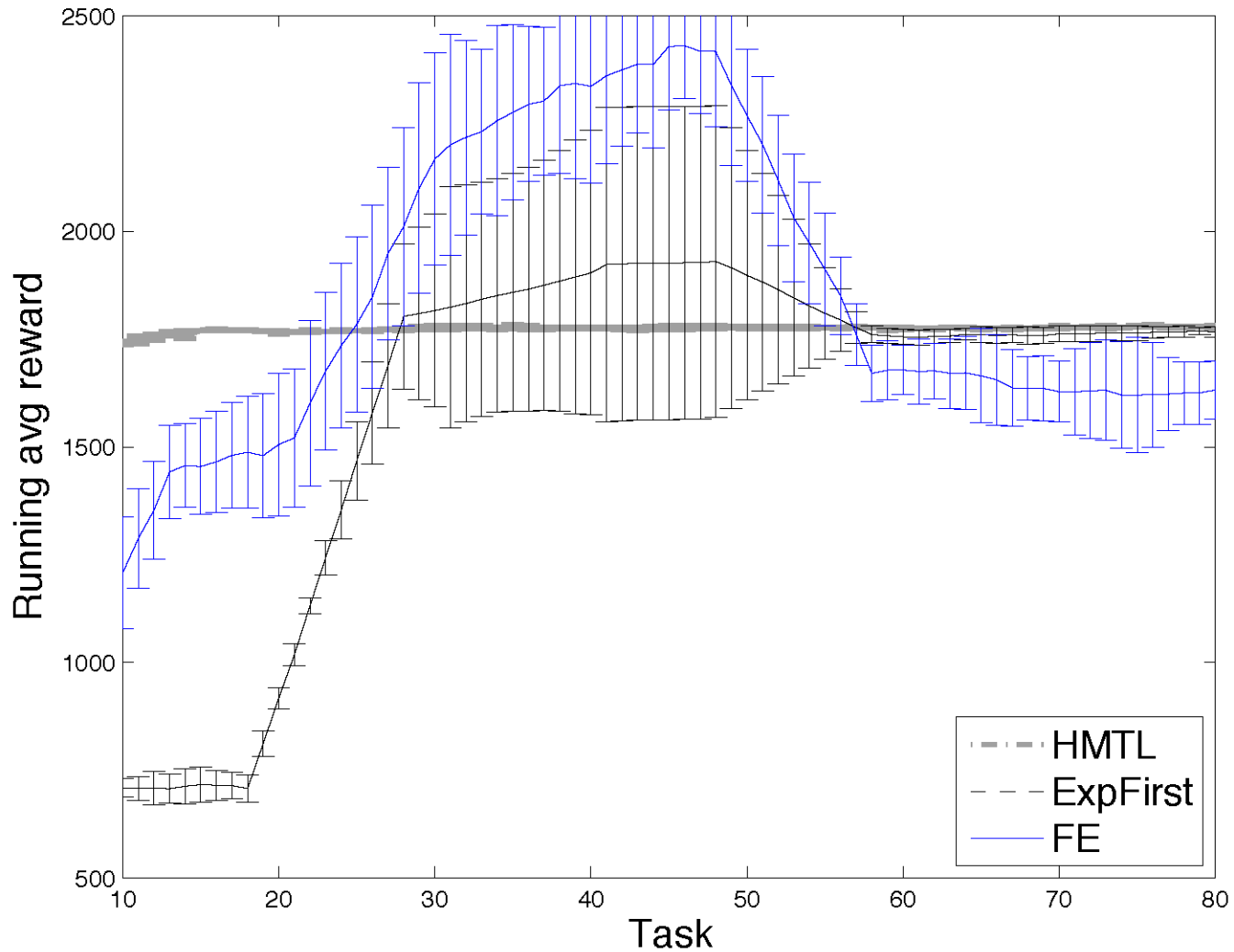
- 4 possible MDPs with
- noisy state transitions
  - different rewarding states

- Algorithms for comparison
- Forced-exploration [this work]
  - Explore-first [Brunskill-Li]
  - Hierarchical Multi-task Learning [Wilson et al.]

# Stochastic Setting with small $\mu_m$



# Adversarial Setting with Changing Distribution





# Conclusions

- Two kinds of exploration needed in LLRL
- Online discovery problem as abstraction for cross-task exploration
- A new lifelong RL algorithm based on optimal ODP algorithm
  - Provably sample complexity better than single-task RL
  - Proof-of-concept experiments demonstrating desired behavior

## Future work

- Function approximation
- Use of prior information