# Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning

Richard S. Sutton, Doina Precup, Satinder Singh

Presenter: Arseny Moskvichev

# Overview

# Motivation for temporal abstraction

# Motivation for temporal abstraction (airport example)

# Motivation for temporal abstraction (airport example)

# Motivation for temporal abstraction (airport example)

# Motivation for temporal abstraction (airport example)



What next?

# Motivation for temporal abstraction (airport example)
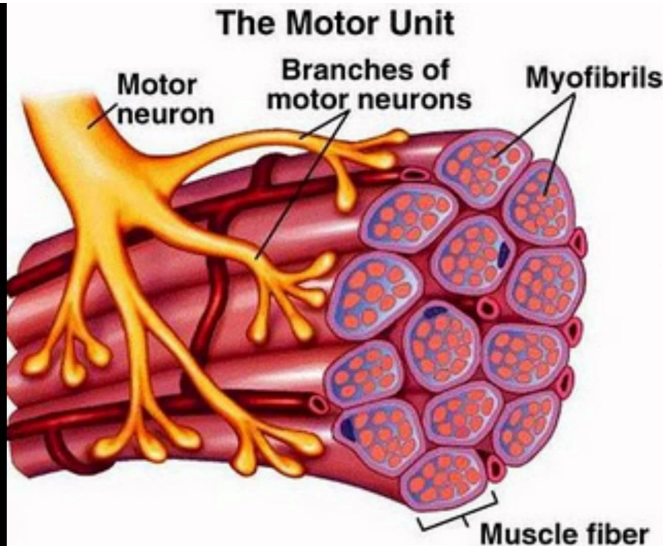
Make a call.

# Motivation for temporal abstraction (airport example)

Which implies getting up and finding your phone

# Motivation for temporal abstraction (airport example)

… Which implies contracting your muscles in
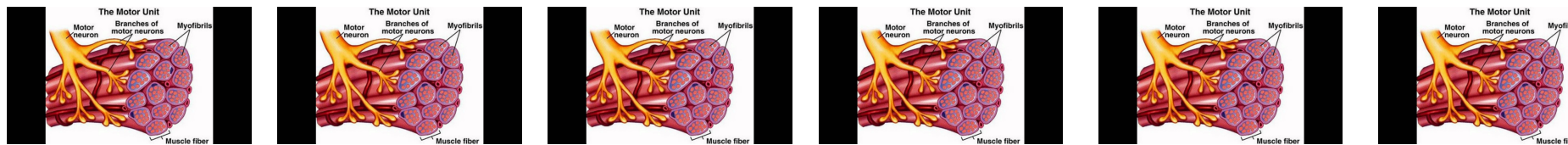a specific sequence



The Motor Unit

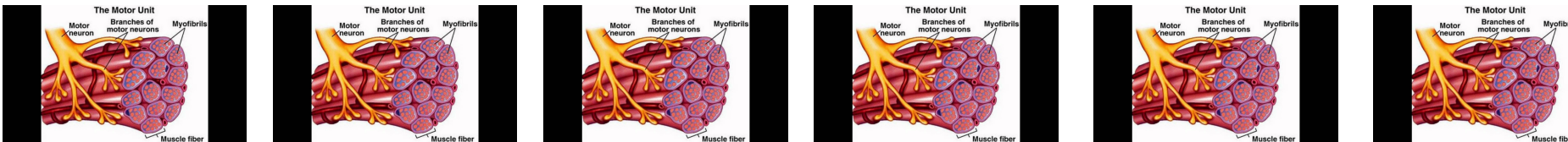# Motivation for temporal abstraction (airport example)

… etc.

Now imagine for a second, how a trip to an airport would look like if we had no temporal abstraction mechanisms.

# Motivation for temporal abstraction (airport example)

Should I walk?



Or take a taxi?

# MDPs, SMDPs, and Options

Theoretical basis for temporal abstraction

# MDPs

An MDP consists of:

A set of states

A set of actions,

Transition dynamics

$$p_{ss'}^a = \Pr\{s_{t+1} = s' \mid s_t = s, \ a_t = a\}$$

One-step expected reward

$$r_s^a = E\{r_{t+1} \mid s_t = s, \ a_t = a\}$$

# MDPs

$$\pi : \mathcal{S} \times \mathcal{A} \to [0, 1]$$

$$
\begin{aligned}
V^\pi(s) &= E\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots \,\big|\, s_t = s, \pi\} \\
&= E\{r_{t+1} + \gamma V^\pi(s_{t+1}) \,\big|\, s_t = s, \pi\} \\
&= \sum_{a \in \mathcal{A}_s} \pi(s, a) \left[ r_s^a + \gamma \sum_{s'} p_{ss'}^a V^\pi(s') \right],
\end{aligned}
$$

# MDPs

$$V^*(s) = \max_\pi V^\pi(s)$$

$$= \max_{a \in \mathcal{A}_s} E\left\{r_{t+1} + \gamma V^*(s_{t+1}) \,\big|\, s_t = s,\ a_t = a\right\}$$

$$= \max_{a \in \mathcal{A}_s}\left[r_s^a + \gamma \sum_{s'} p_{ss'}^a V^*(s')\right].$$

# Semi MDP

# Options

An option $\langle \mathcal{I}, \pi, \beta \rangle$ consists of three components:

A policy:, $\pi : \mathcal{S} \times \mathcal{A} \to [0, 1]$

A termination condition: $\beta : \mathcal{S}^+ \to [0, 1]$

An initiation set: $\mathcal{I} \subseteq \mathcal{S}$

We say that the option is **available** at some state, if this state belongs to the option's initiation set.

# Semi-markov Options

Like options, but decisions could depend on history since initiation,

$s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}, \ldots, r_\tau, s_\tau$, denoted $h_{t\tau}$

Therefore, now we need to re-define policy and termination conditions:

$\pi : \Omega \times \mathcal{A} \to [0, 1]$ and $\beta : \Omega \to [0, 1]$

# Semi-markov Options

Semi markov options naturally arise when:

# Semi-markov Options

Semi markov options naturally arise when:

- An option needs to terminate after being executed for some time, regardless of the resulting state
- Options have access to a more fine-grained representation of the state that the higher-level policy

# Options. More definitions and details

$\mathcal{O}_s$ - like action sets, but for options
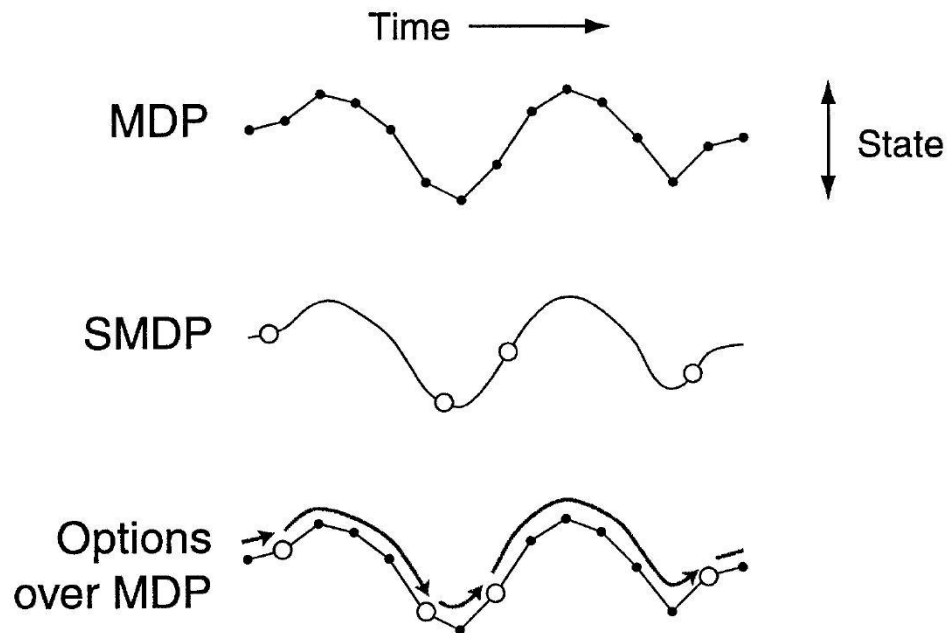
Viewing simple actions as single-step options

Composing options (note: composition of Markov options is semi-Markov)

Policies over options: $\mu : \mathcal{S} \times \mathcal{O} \rightarrow [0, 1]$

Associated flat policy: $\pi = flat(\mu)$

Note: even a policy composed of markov options is likely to be only semi-Markov.

# MDPs, SMDPs, and Options (summary)

# Relationship between Options and SMDPs

**Theorem 1** (MDP + Options = SMDP). *For any MDP, and any set of options defined on that MDP, the decision process that selects only among those options, executing each to termination, is an SMDP.*

# Rewriting Bellman equations with options

A few more definitions:

$$r_s^o = E\{r_{t+1} + \gamma r_{t+2} + \cdots + \gamma^{k-1} r_{t+k} \,|\, \mathcal{E}(o, s, t)\}$$

$$p_{ss'}^o = \sum_{k=1}^{\infty} p(s', k)\gamma^k,$$

Where $p(s', k)$ is the probability that option terminates in a state s' after k steps.

# Rewriting Bellman equations with options

Using these definitions, we can finally write

$$V^{\mu}(s) = E\left\{ r_{t+1} + \cdots + \gamma^{k-1} r_{t+k} + \gamma^{k} V^{\mu}(s_{t+k}) \,\middle|\, \mathcal{E}(\mu, s, t) \right\}$$

(where $k$ is the duration of the first option selected by $\mu$)

$$= \sum_{o \in \mathcal{O}_s} \mu(s, o) \left[ r_s^o + \sum_{s'} p_{ss'}^o V^{\mu}(s') \right],$$

# Rewriting Bellman equations with options

$$Q^{\mu}(s, o) = E\{r_{t+1} + \cdots + \gamma^k \ {}^1 r_{t+k} + \gamma^k V^{\mu}(s_{t+k}) \,\big|\, \mathcal{E}(o, s, t)\},$$

$$= E\Big\{r_{t+1} + \cdots + \gamma^k \ {}^1 r_{t+k}$$

$$+ \gamma^k \sum_{o' \in \mathcal{O}_s} \mu(s_{t+k}, o') Q^{\mu}(s_{t+k}, o') \,\Big|\, \mathcal{E}(o, s, t)\Big\}$$

$$= r_s^o + \sum_{s'} p_{ss'}^o \sum_{o' \in \mathcal{O}_{s'}} \mu(s', o') Q^{\mu}(s', o').$$

# Rewriting Bellman equations with options

Similarly we can rewrite the bellman optimality equations.

$$V_{\mathcal{O}}^*(s) \overset{\text{def}}{=} \max_{\mu \in \Pi(\mathcal{O})} V^\mu(s)$$

$$= \max_{o \in \mathcal{O}_s} E\{r_{t+1} + \cdots + \gamma^{k\ -1} r_{t+k} + \gamma^k V_{\mathcal{O}}^*(s_{t+k}) \,|\, \mathcal{E}(o, s, t)\}$$

$$\text{(where } k \text{ is the duration of } o \text{ when taken in } s)$$

$$= \max_{o \in \mathcal{O}_s} \left[ r_s^o + \sum_{s'} p_{ss'}^o V_{\mathcal{O}}^*(s') \right]$$

$$= \max_{o \in \mathcal{O}_s} E\{r + \gamma^k V_{\mathcal{O}}^*(s') \,|\, \mathcal{E}(o, s)\},$$

# Rewriting Bellman equations with options

Similarly we can rewrite the bellman optimality equations.

$$Q_{\mathcal{O}}^*(s, o) \stackrel{\text{def}}{=} \max_{\mu \in \Pi(\mathcal{O})} Q^\mu(s, o)$$

$$= E\left\{r_{t+1} + \cdots + \gamma^{k-1} r_{t+k} + \gamma^k V_{\mathcal{O}}^*(s_{t+k}) \,\middle|\, \mathcal{E}(o, s, t)\right\}$$

(where $k$ is the duration of $o$ from $s$)

$$= E\left\{r_{t+1} + \cdots + \gamma^{k-1} r_{t+k} + \gamma^k \max_{o' \in \mathcal{O}_{s_{t+k}}} Q_{\mathcal{O}}^*(s_{t+k}, o') \,\middle|\, \mathcal{E}(o, s, t)\right\},$$

$$= r_s^o + \sum_{s'} p_{ss'}^o \max_{o' \in \mathcal{O}_{s'}} Q_{\mathcal{O}}^*(s', o')$$

$$= E\left\{r + \gamma^k \max_{o' \in \mathcal{O}_{s'}} Q_{\mathcal{O}}^*(s', o') \,\middle|\, \mathcal{E}(o, s)\right\},$$

# Why is it all useful?

Because we can now apply conventional DP algorithms to find policies over options.
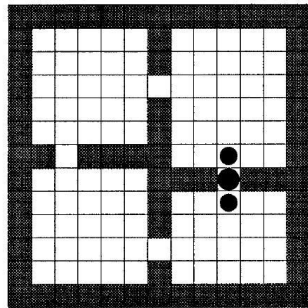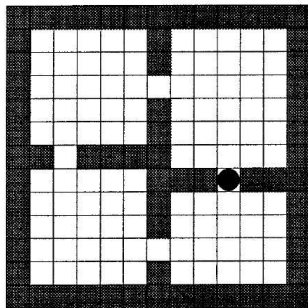
# Example



4 stochastic
primitive actions

up
left ← → right
down

Fail 33%
of the time

8 multi-step options
(to each room's 2 hallways)
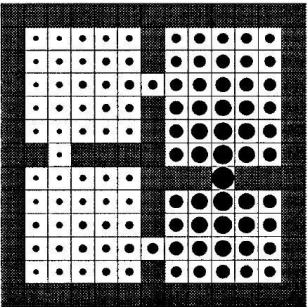
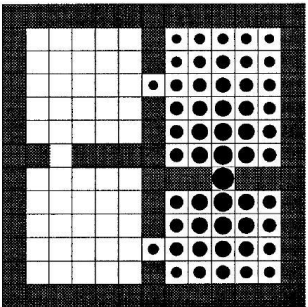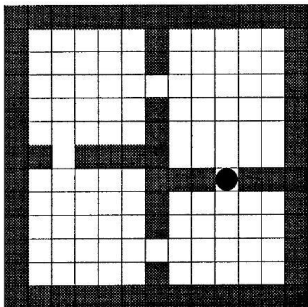# Example (option)



Target Hallway

# Example (hallway goal)

Primitive
options
$\mathcal{O}=\mathcal{A}$

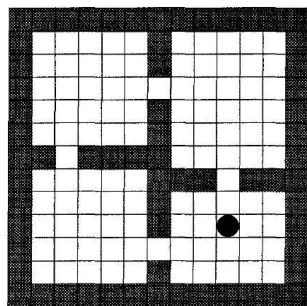Hallway
options
$\mathcal{O}=\mathcal{H}$



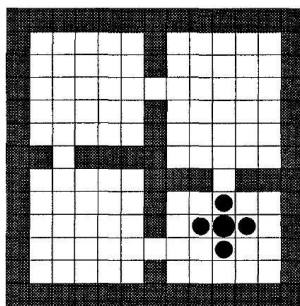Initial Values          Iteration #1          Iteration #2

# Example (non-hallway goal)
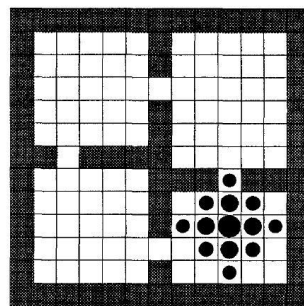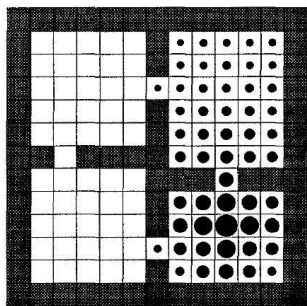
Primitive
and
hallway
options
$\mathcal{O}=\mathcal{A}\cup\mathcal{H}$
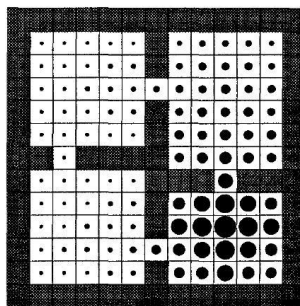


Initial values

Iteration #1

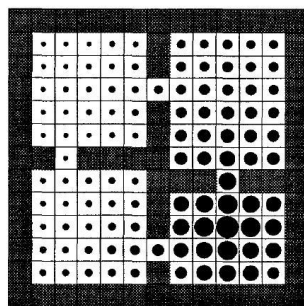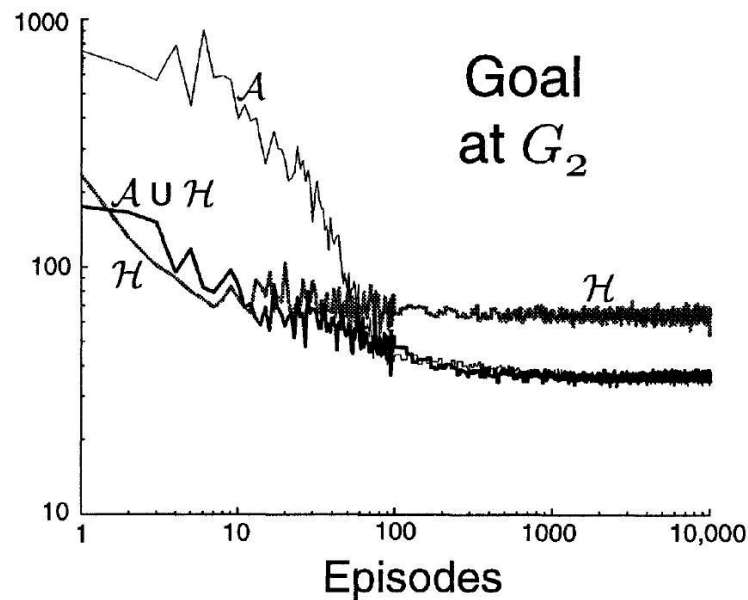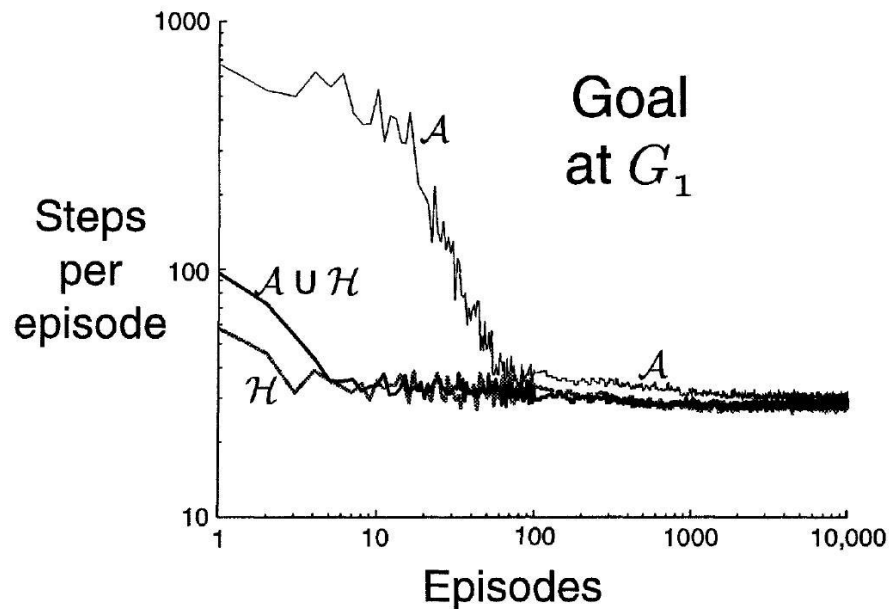Iteration #2

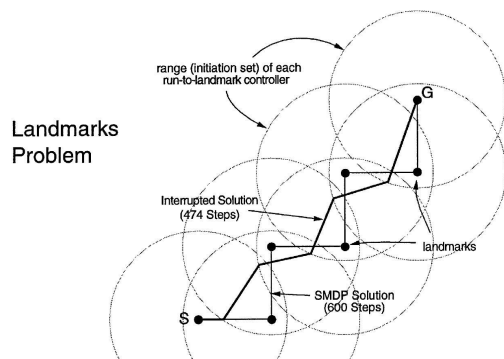Iteration #3

Iteration #4

Iteration #5

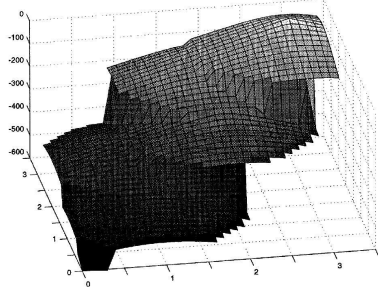# Learning

# Option-specific methods: Interruption

**Theorem 2** (Interruption). *For any MDP, any set of options $\mathcal{O}$, and any Markov policy $\mu : \mathcal{S} \times \mathcal{O} \to [0, 1]$, define a new set of options, $\mathcal{O}'$, with a one-to-one mapping between the two option sets as follows: for every $o = \langle \mathcal{I}, \pi, \beta \rangle \in \mathcal{O}$ we define a corresponding $o' = \langle \mathcal{I}, \pi, \beta' \rangle \in \mathcal{O}'$, where $\beta' = \beta$ except that for any history $h$ that ends in state $s$ and in which $Q^\mu(h, o) < V^\mu(s)$, we may choose to set $\beta'(h) = 1$. Any histories whose termination conditions are changed in this way are called* interrupted histories. *Let the interrupted policy $\mu'$ be such that for all $s \in \mathcal{S}$, and for all $o' \in \mathcal{O}'$, $\mu'(s, o') = \mu(s, o)$, where $o$ is the option in $\mathcal{O}$ corresponding to $o'$. Then*

(i) $V^{\mu'}(s) \geqslant V^\mu(s)$ *for all $s \in \mathcal{S}$.*

(ii) *If from state $s \in \mathcal{S}$ there is a non-zero probability of encountering an interrupted history upon initiating $\mu'$ in $s$, then $V^{\mu'}(s) > V^\mu(s)$.*
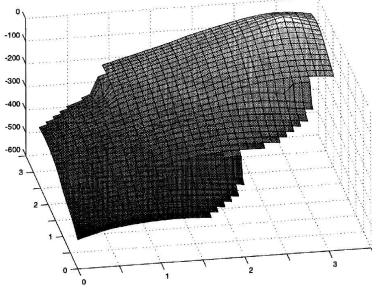
# Option-specific methods: Interruption



Landmarks Problem

range (initiation set) of each run-to-landmark controller

G

Interrupted Solution (474 Steps)

landmarks

SMDP Solution (600 Steps)

S

$V_{\mathcal{O}}^*$- SMDP Value Function

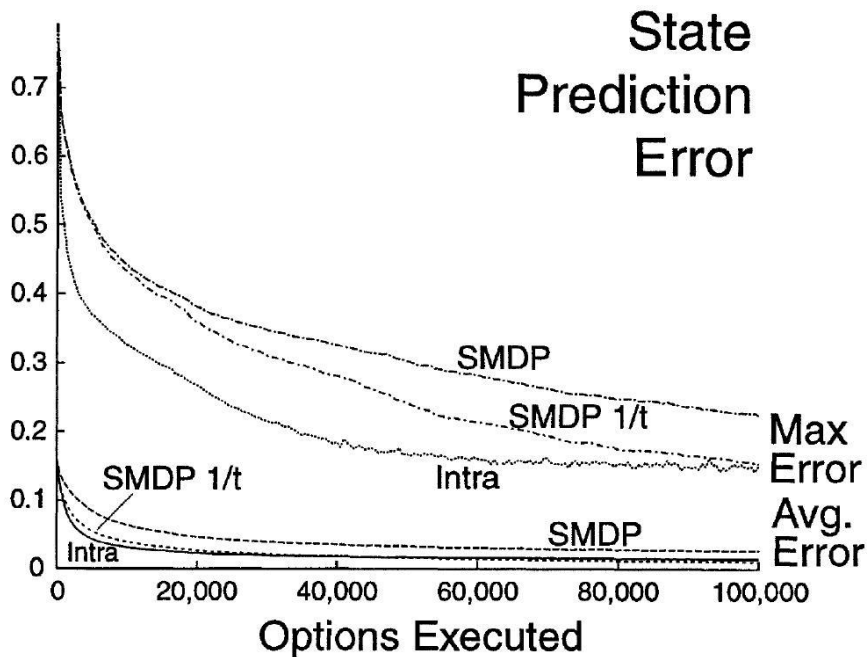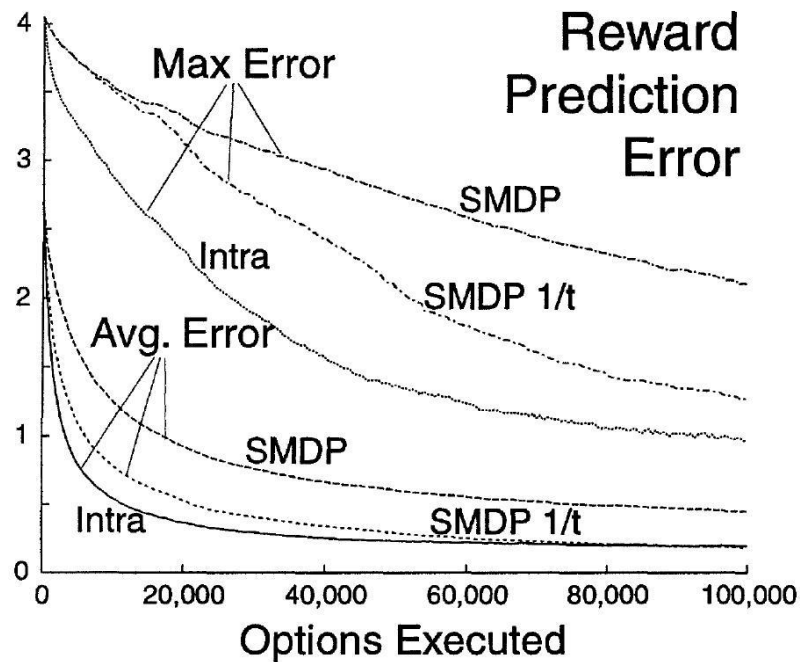$V^{\mu'}$- Values with Interruption

Option-specific methods:

Intra-option **model** learning and
Intra-option **value** learning
Subgoals for learning options

# Intra-option model learning:

# Intra-option value learning:

**Theorem 3** (Convergence of intra-option Q-learning). *For any set of Markov options, $\mathcal{O}$, with deterministic policies, one-step intra-option Q-learning converges with probability 1 to the optimal Q-values, $Q_{\mathcal{O}}^*$, for every option regardless of what options are executed during learning, provided that every action gets executed in every state infinitely often.*

# Subgoals for learning options:



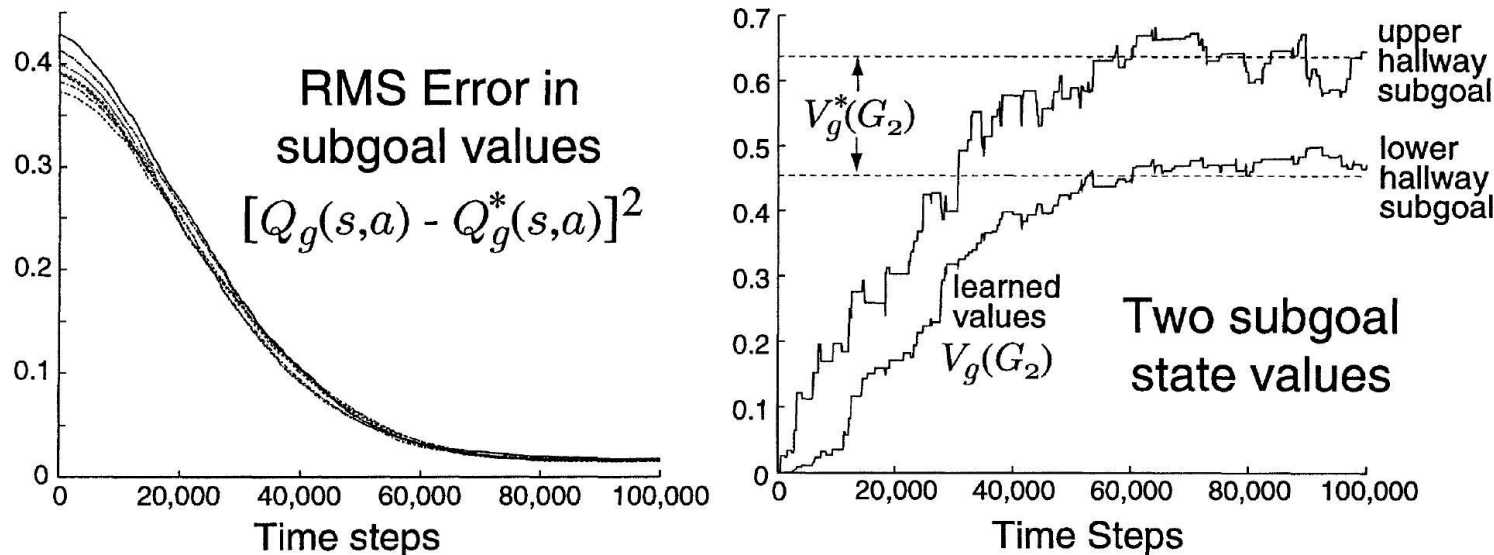RMS Error in subgoal values $[Q_g(s,a) - Q_g^*(s,a)]^2$

Two subgoal state values

Fig. 11. Learning subgoal-achieving hallway options under random behavior. Shown on the left is the error between $Q_g(s,a)$ and $Q_g^*(s,a)$ averaged over $s \in \mathcal{I}$, $a \in \mathcal{A}$, and 30 repetitions. The right panel shows the learned state values (maximum over action values) for two options at state $G_2$ approaching their correct values.

# Summary

Options are temporally extended actions with initiation sets and termination conditions. Options defined over an MDP yield SMDPs and all SMDP methods and results apply.

Options compare favorably to SMDPs, since they also allow to leverage our knowledge about the structure of a temporally-extended actions. In particular, option-specific methods include:

- Interrupting options

- Intra-option model/value learning

- Learning with subgoals