

Reinforcement Learning

Abstraction and Hierarchy in RL

Subramanian Ramamoorthy
School of Informatics

3 March, 2017

On the Organization of Behaviour

Consider day to day tasks such as opening a bottle or picking up your coat from the stand

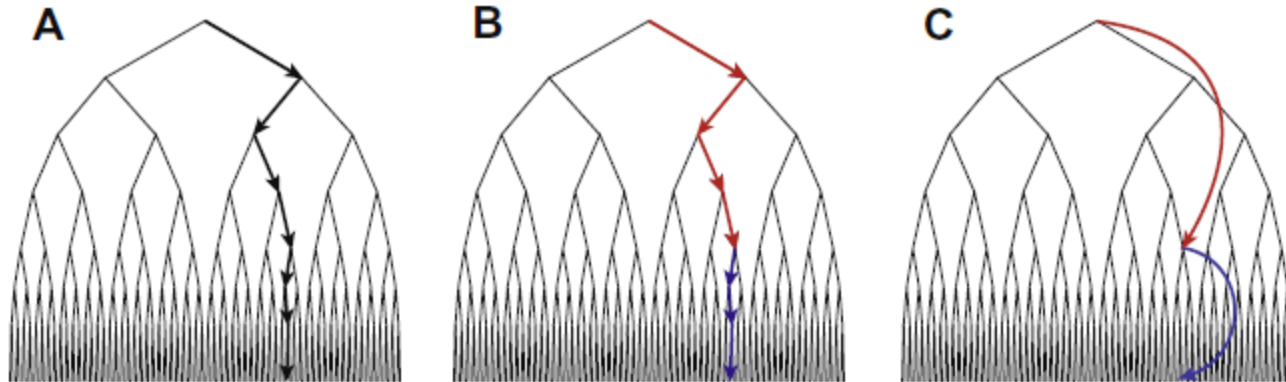
- How would you tell someone how to perform this?
- How do you actually get the job done? Using what?

An old idea:

sequential behavior cannot be understood as a chain of stimulus–response associations... behavior displays hierarchical structure, comprising nested subroutines.

- Karl Lashley (1951)

On the Complexity of RL



- RL agents must learn from exploring the environment, trying out different actions in different states
- Complexity grows quickly; how to get around?
 - Balance exploration/exploitation
 - Abstractions

Example: *Instrumental* Hierarchy

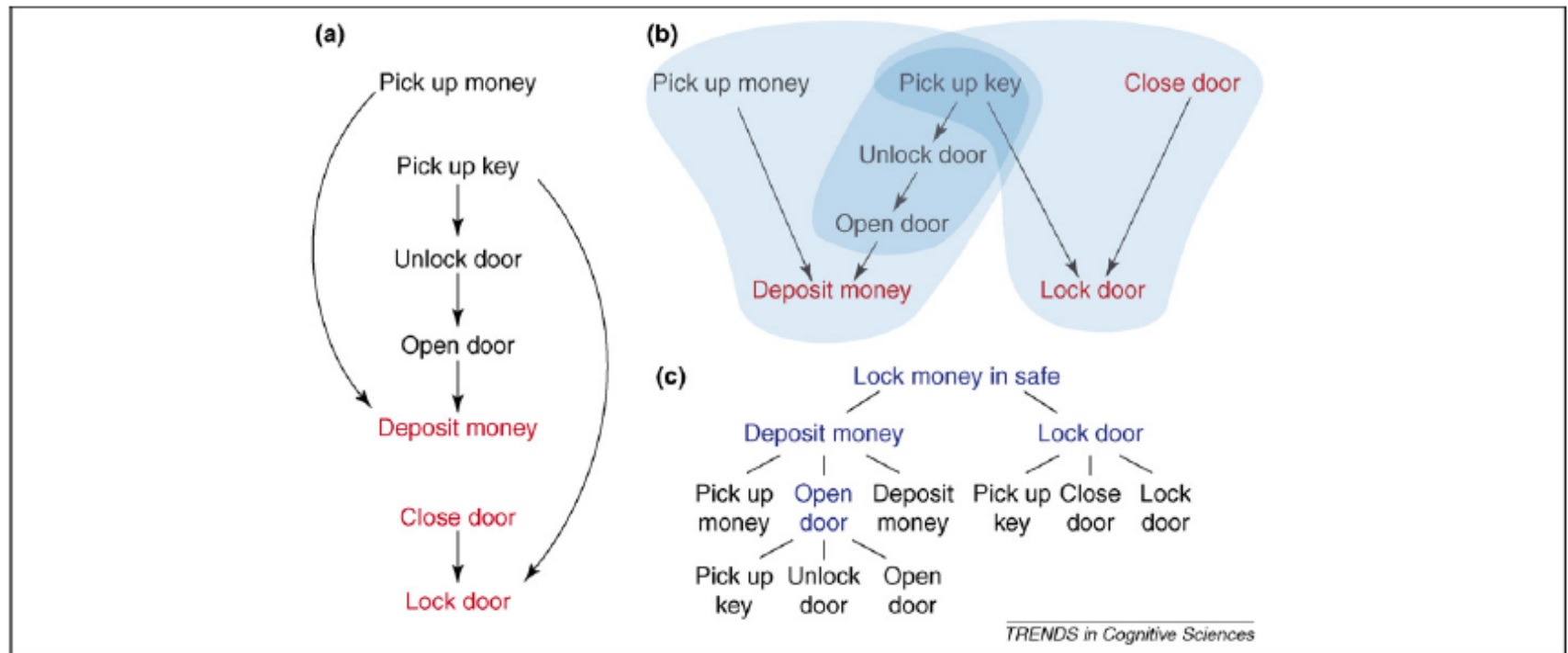


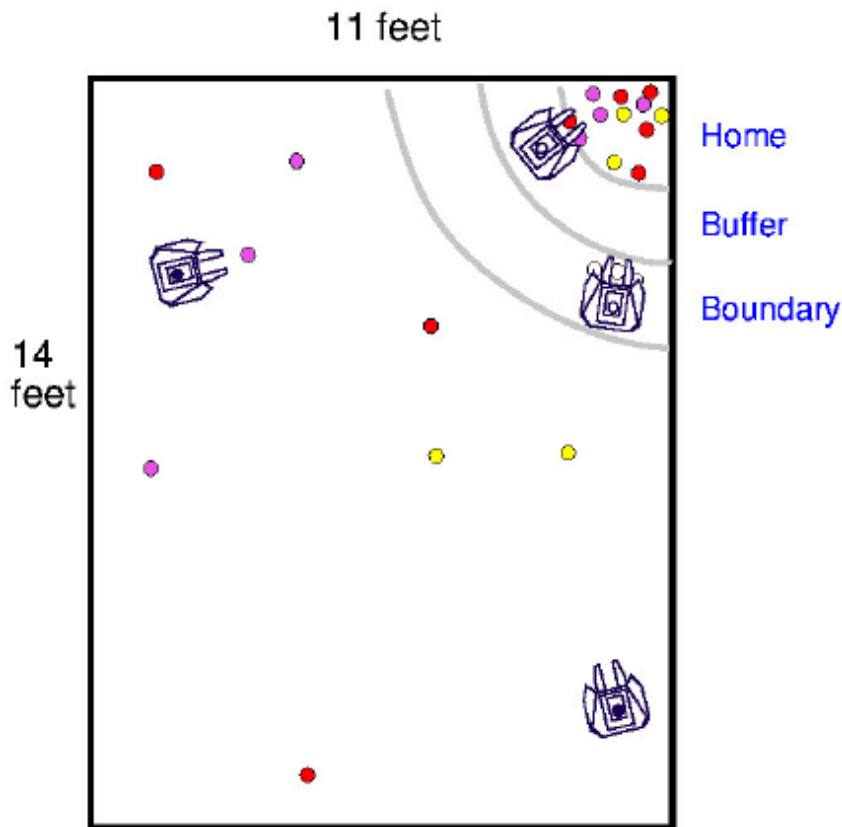
Figure 1. An illustration of hierarchical instrumental structure. (a) An action sequence for locking money in a safe. Arrows denote means–ends relationships. Red indicates that the action accomplishes one component of the goal (money in the safe with the door closed and locked). (b) The sequence in part (a) redrawn to highlight the presence of a part-whole structure. Blue fields indicate coherent parts and sub-parts of the action sequence. At the coarsest level, the sequence breaks down into two parts, one organized around the sub-goal of depositing the money, the other around the sub-goal of locking the safe door. The action ‘pick up key’ subserves both goals. Also indicated is a subordinate or nested sequence, which is organized around opening the safe door. (c) One way of representing the sequence as a schema-,sub-task- or sub-goal hierarchy. Temporarily abstract actions are in blue.

[Source: M. Botvinick, Hierarchical models of behavior and prefrontal function, Trends in Cognitive Science 12(5), 2008]

In Computational Terms, Why Hierarchy?

- Knowledge **transfer**/injection
 - Concise encoding of temporally extended actions enables transfer of that knowledge to a new task needing the same
- Biases **exploration**
 - Rather than search only at the finest resolution, algorithm can compute value of alternatives at a larger scale
- **Faster** solutions (even if model known)
 - Following well understood principles of computer science, hierarchy enables modularity of policies and models

An Early Idea: Reward Shaping



- The robots' objective is to collectively find pucks and bring them home.
- Represent the 12-dim environment by state variables (features?):
 - have-puck?
 - at-home?
 - near-intruder?
- What should the immediate reward function be?

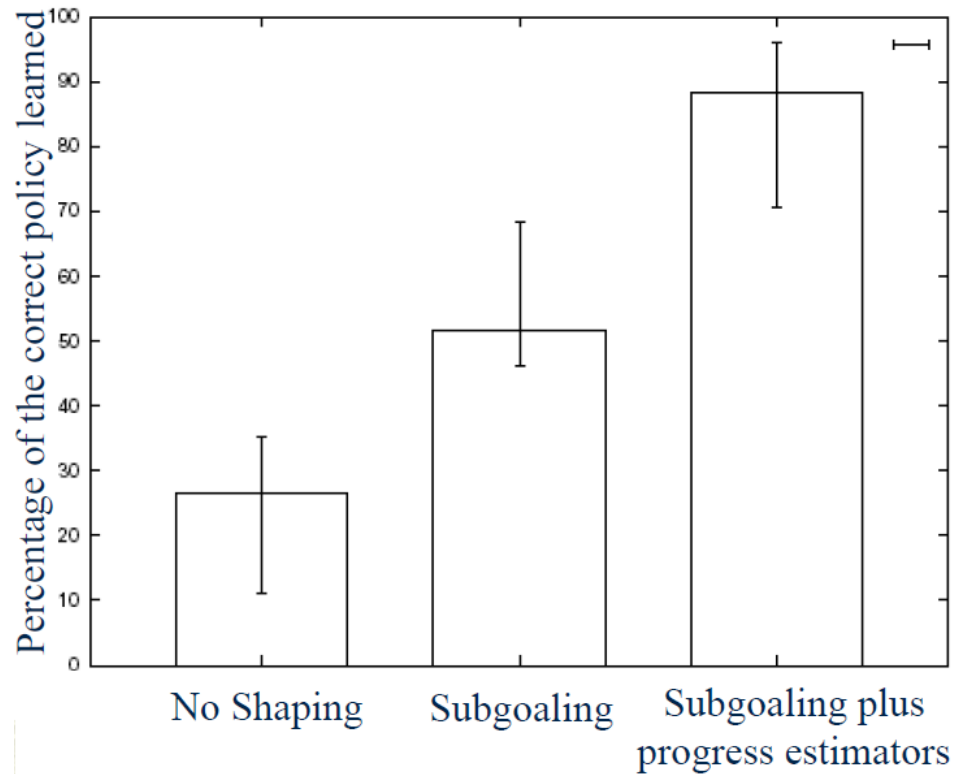
[Source: M. Mataric, Reward Functions for Accelerated Learning, ICML 1994]

Reward Shaping

- If a reward is given only when a robot drops a puck at home, learning will be extremely difficult.
 - The delay between the action and the reward is large.
- Solution: Reward shaping (intermediate rewards).
 - Add rewards/penalties for achieving sub-goals/errors:
 - subgoal: grasped-puck
 - subgoal: dropped-puck-at-home
 - error: dropped-puck-away-from-home
- Add progress estimators:
 - Intruder-avoiding progress function
 - Homing progress function
- Adding intermediate rewards will potentially allow RL to handle more complex problems.

Reward Shaping Results

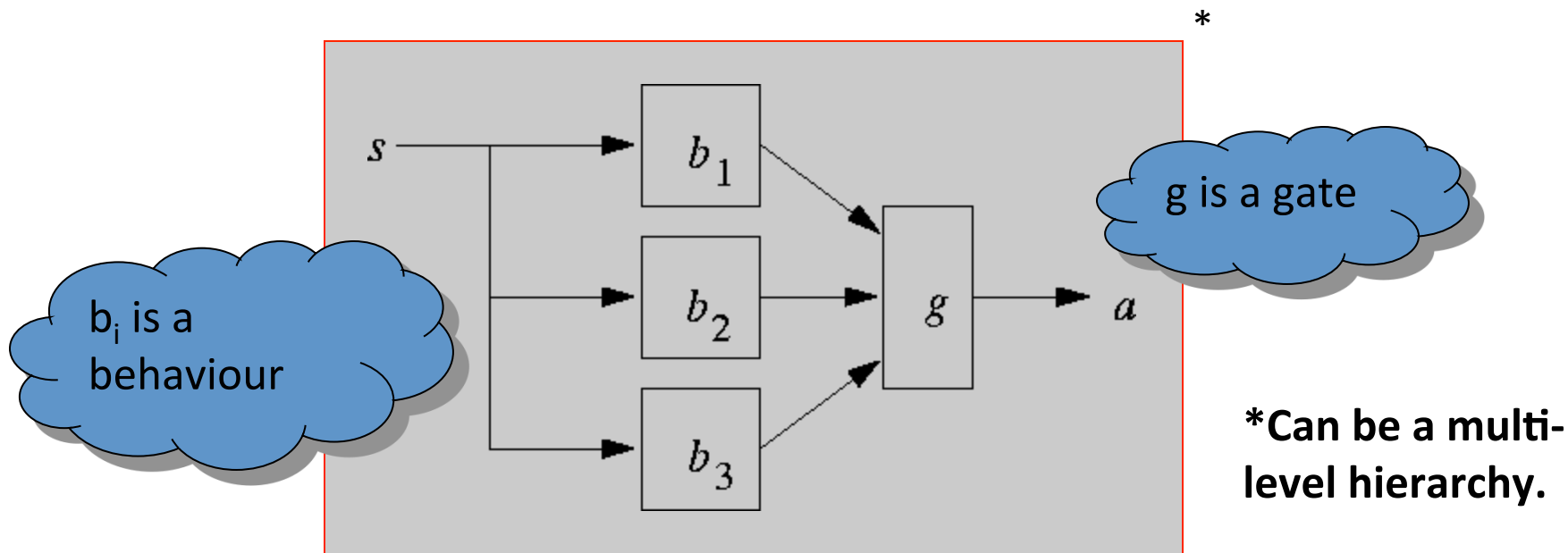
Percentage of policy learnt after 15 minutes:



Another Early Idea: Hierarchical Algorithms - Gating Mechanisms

Hierarchical Learning (especially popular with neural network models)

- *Learn the gating function*
- *Learn the individual behaviours*
- *Learn both simultaneously*



Temporal Abstraction

- What's the issue?
 - Want “macro” actions (multiple time steps)
 - Advantages:
 - Avoid dealing with (exploring/computing values for) less desirable states
 - Reuse experience across problems/regions
- What's not obvious
 - Dealing with the Markov assumption
 - Getting the calculations right (e.g., stability and convergence)

Semi-Markov Decision Processes

Semi-Markov Decision Processes

- A generalization of MDPs:
The amount of **time** between one decision and the next is a **random variable** (either real or integer valued)
- Treat the system as remaining in each state for a random waiting time
 - after which, transition to next state is instantaneous
- Real valued case: continuous time, discrete events
- Discrete case: Decisions only made an integer multiple of an underlying time step

Semi-Markov Decision Processes

- SMDP is defined in terms of
 - $P(s', \tau | s, a)$: Transition probability (τ is the waiting time)
 - $R(s, a)$ or just r : Reward, amount expected to accumulate during waiting time, τ , in particular state and action
- Bellman equation can then be written down as, for all s :

$$V^*(s) = \max_{a \in \mathcal{A}_s} \left[r + \sum_{s', \tau} \gamma^\tau P(s', \tau | s, a) V^*(s') \right]$$

Note the need to sum over waiting time, as well.

Semi-Markov Decision Processes

- Likewise, we can write down the Bellman equation for the state-action value function as,

$$Q^*(s, a) = r + \sum_{s', \tau} \gamma^\tau P(s', \tau | s, a) \max_{a' \in \mathcal{A}_s} Q^*(s', a')$$

$$\forall s \in \mathcal{S}, a \in \mathcal{A}_s$$

- So, Dynamic Programming algorithms can be naturally extended to the case of SMDPs as well

Q-Learning with SMDPs

- Can we also modify sampling based algorithms accordingly?
- Consider the standard Q-learning algorithm, rewritten slightly in the following form,

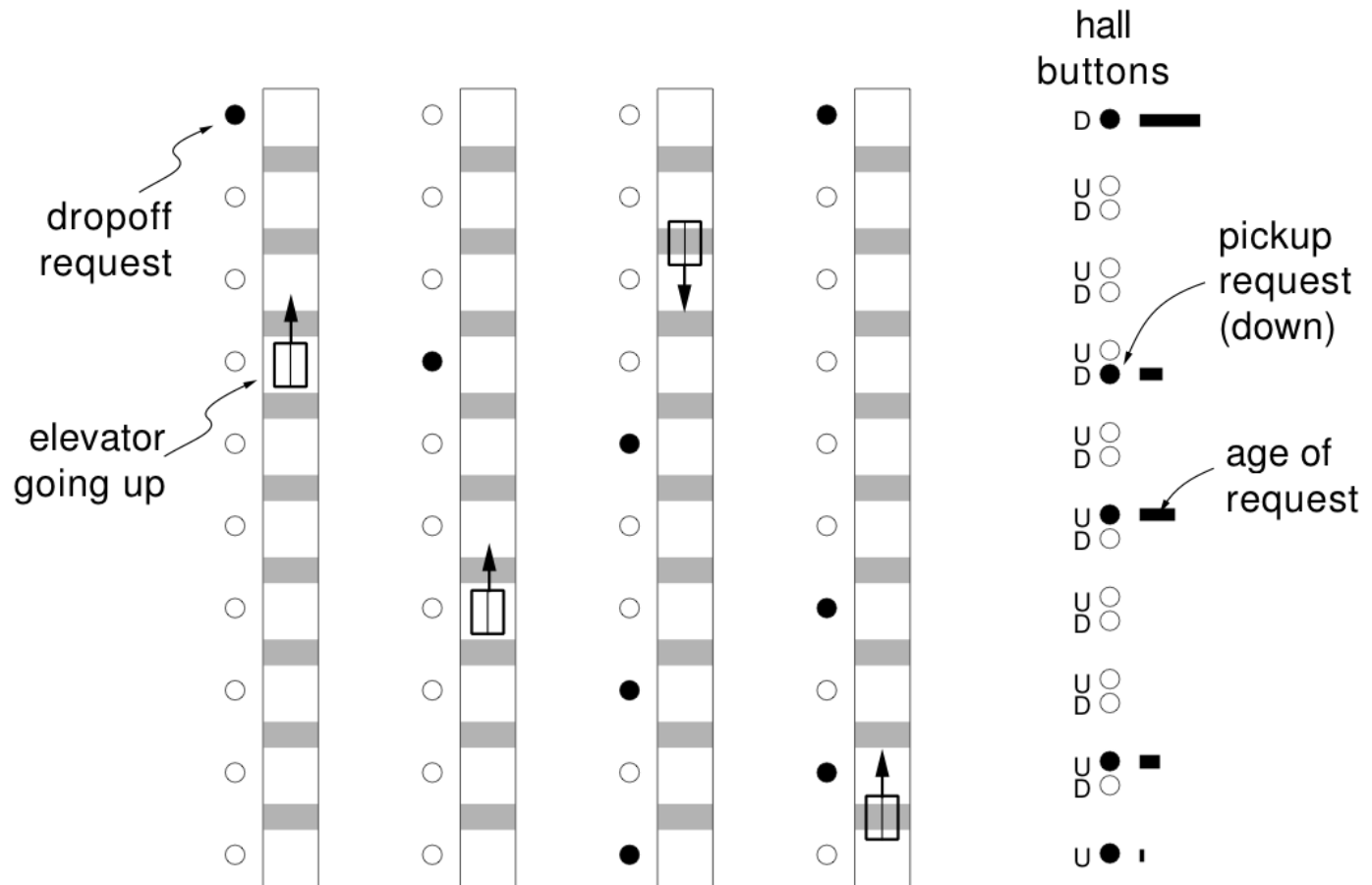
$$Q_{k+1}(s, a) = (1 - \alpha_k)Q_k(s, a) + \alpha_k[r + \gamma \max_{a' \in \mathcal{A}_s} Q_k(s', a')]$$

- If we write down the reward sum, in brackets, for the entire waiting time duration, then we will have

$$Q_{k+1}(s, a) = (1 - \alpha_k)Q_k(s, a) + \alpha_k[r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{\tau-1} r_{t+\tau} + \gamma^\tau \max_{a' \in \mathcal{A}_s} Q_k(s', a')]$$

Case Study: Elevator Dispatching

[Crites and Barto, 1996]



Semi-Markov Q-Learning

Continuous-time problem but decisions in discrete jumps.
For this SMDP, the expression for returns can be written as,

$$R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \quad \text{or} \quad R_t = \int_0^{\infty} e^{-\beta\tau} r_{t+\tau} d\tau$$

Note that the meaning of quantity r differs in the two expressions:

- reward at a discrete time step in discrete case
- reward “rate” in continuous case

The negative exponential has a similar role as the discount factor as we have been using it so far.

Semi-Markov Q-Learning

Suppose system takes action a from state s at time t_1 ,
and next decision is needed at time t_2 in state s' :

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \left[\int_{t_1}^{t_2} e^{-\beta(\tau-t_1)} r_\tau d\tau + e^{-\beta(t_2-t_1)} \max_{a'} Q(s', a') \right]$$

Problem Setup: Passenger Arrival Patterns

Up-peak and Down-peak traffic

- Not equivalent: down-peak handling capacity is much greater than up-peak handling capacity; so up-peak capacity is limiting factor.
- Up-peak easiest to analyse: once everyone is onboard at lobby, rest of trip is determined. The only decision is when to open and close doors at lobby. Optimal policy for pure case is: close doors when threshold number on; threshold depends on traffic intensity.
- More policies to consider for two-way and down-peak traffic.
- We focus on down-peak traffic pattern.

Various Extant Control Strategies

- Zoning: divide building into zones; park in zone when idle. Robust in heavy traffic.
- Search-based methods: greedy or non-greedy. Receding Horizon control.
- Rule-based methods: expert systems/fuzzy logic; from human “experts”
- Other heuristic methods: Longest Queue First (LQF), Highest Unanswered Floor First (HUFF), Dynamic Load Balancing (DLB)
- Adaptive/Learning methods: NNs for prediction, parameter space search using simulation, DP on simplified model, non-sequential RL

The Elevator Model (Lewis, 1991)

Discrete Event System: continuous time,
asynchronous elevator operation

Parameters:

- Floor Time (time to move one floor at max speed): 1.45 secs.
- Stop Time (time to decelerate, open and close doors, and accelerate again): 7.19 secs.
- TurnTime (time needed by a stopped car to change directions): 1 sec.
- Load Time (the time for one passenger to enter or exit a car): a random variable with range from 0.6 to 6.0 secs, mean of 1 sec.
- Car Capacity: 20 passengers

Traffic Profile:

- Poisson arrivals with rates changing every 5 minutes; down-peak

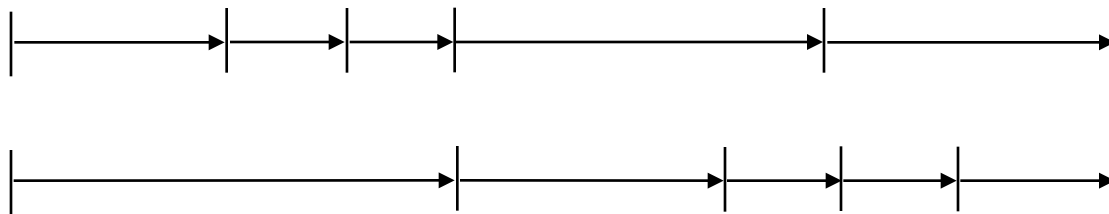
State Space

- 18 hall call buttons: 2^{18} combinations
- positions and directions of cars: 18^4 (rounding to nearest floor)
- motion states of cars (accelerating, moving, decelerating, stopped, loading, turning): 6
- 40 car buttons: 2^{40}
- Set of passengers waiting at each floor, each passenger's arrival time and destination: unobservable. However, 18 real numbers are available giving elapsed time since hall buttons pushed; we discretize these.
- Set of passengers riding each car and their destinations: observable only through the car buttons

Conservatively about 10^{22} states

Actions

- When moving (halfway between floors):
 - stop at next floor
 - continue past next floor
- When stopped at a floor:
 - go up
 - go down
- Asynchronous



Constraints

standard

- A car cannot pass a floor if a passenger wants to get off there
- A car cannot change direction until it has serviced all onboard passengers traveling in the current direction
- Don't stop at a floor if another car is already stopping, or is stopped, there

special heuristic

- Don't stop at a floor unless someone wants to get off there
- Given a choice, always move up



Stop and Continue

Performance Criteria

Minimize:

- Average wait time
- Average system time (wait + travel time)
- % waiting > T seconds (e.g., T = 60)
- Average squared wait time (to encourage fast and fair service)



Average Squared Wait Time

Instantaneous cost, p individuals:

$$r_\tau = \sum_p (\text{wait}_p(\tau))^2$$

Define return as an integral rather than a sum (Bradtke and Duff, 1994):

$$\int_0^{\infty} e^{-\beta\tau} r_\tau d\tau$$

Computing Rewards

Must calculate

$$\int_0^{\infty} e^{-\beta(\tau-t_s)} r_{\tau} d\tau$$

- “Omniscient Rewards”: the simulator knows how long each passenger has been waiting.
- “On-Line Rewards”: Assumes only arrival time of first passenger in each queue is known (elapsed hall button time); estimate arrival times

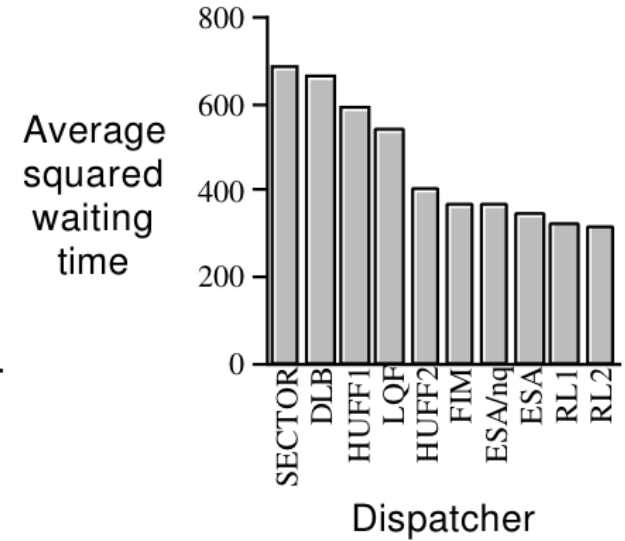
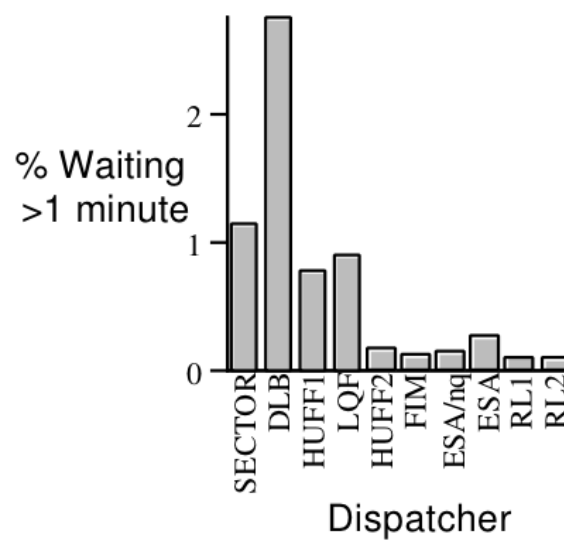
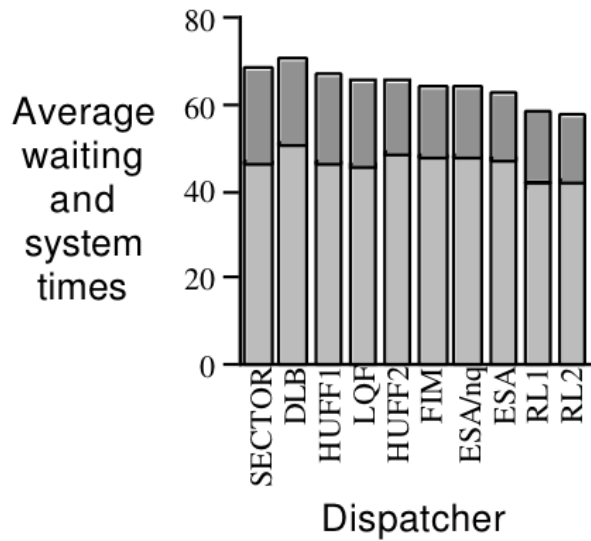
Neural Networks

47 inputs, 20 sigmoid hidden units, 1 or 2 output units

Inputs:

- 9 binary: state of each hall down button
- 9 real: elapsed time of hall down button if pushed
- 16 binary: one on at a time: position and direction of car making decision
- 10 real: location/direction of other cars: “footprint”
- 1 binary: at highest floor with waiting passenger?
- 1 binary: at floor with longest waiting passenger?
- 1 bias unit $\equiv 1$

Elevator Results



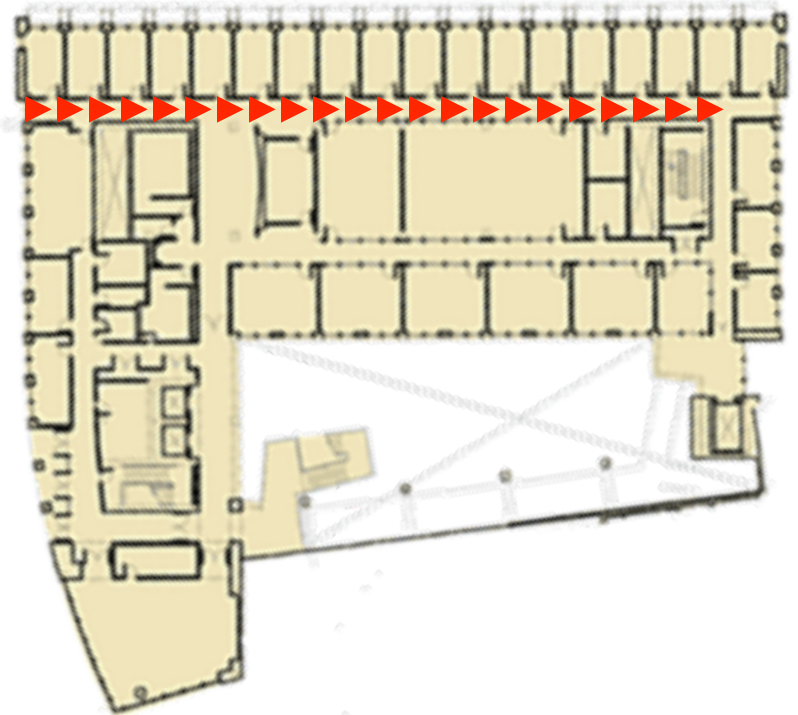
Options Framework

Options example: Move until end of hallway

Start : Any state in the hallway.

Execute : policy π as shown.

Terminate : when state s is the end of hallway.



Options can take **variable number of steps**

[Reference: R.S. Sutton, D. Precup, S. Singh, Between MDPs and Semi-MDPs: A framework for temporal Abstraction in reinforcement learning, Artificial Intelligence Journal 112:181-211, 1999.]

Options [Sutton, Precup, Singh '99]

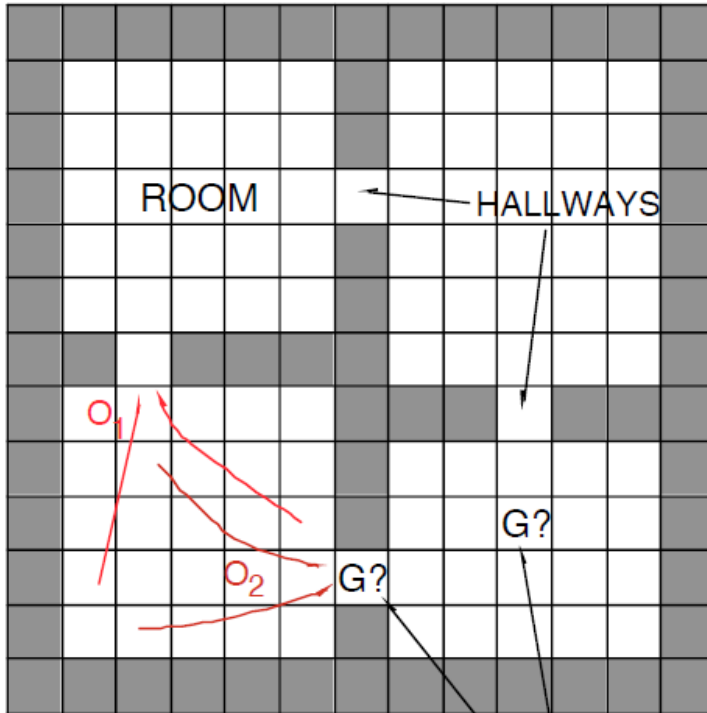
- An option is a *behaviour* defined in terms of:

$$o = \{ I_o, \pi_o, \beta_o \}$$

- I_o : Set of states in which o can be initiated.
- $\pi_o(s)$: Policy (mapping S to A)[§] when o is executing.
- $\beta_o(s)$: Probability that o terminates in s .

§Can be a policy over lower level options.

Rooms Example

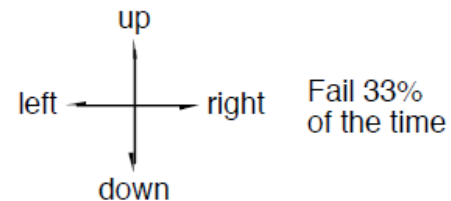


Goal states are given a terminal value of 1

4 rooms

4 hallways

4 unreliable primitive actions



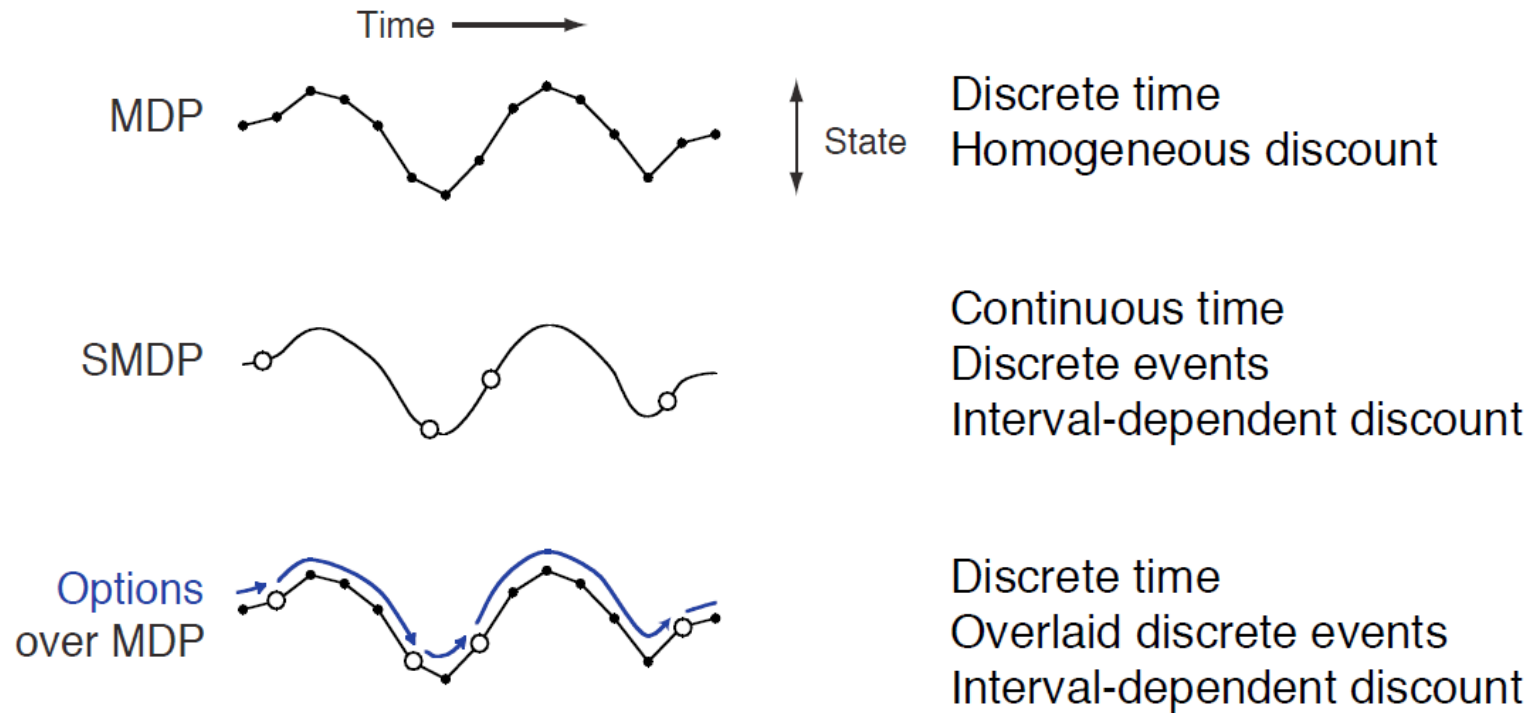
8 multi-step options

(to each room's 2 hallways)

Given goal location, quickly plan shortest route

All rewards zero
 $\gamma = .9$

Options Define a Semi-MDP



A discrete-time SMDP overlaid on an MDP
Can be analyzed at either level

MDP + Options = SMDP

Theorem:

*For any MDP,
and any set of options,
the decision process that chooses among the options,
executing each to termination,
is an SMDP.*

Thus all Bellman equations and DP results extend for value functions over options and models of options (cf. SMDP theory).

Why is this Useful?

- We can now define policy over options as well:

$$\mu : S \times O \rightarrow [0, 1]$$

- And redefine all value functions appropriately:

$$V^\mu(s), Q^\mu(s, o), V_O^*(s), Q_O^*(s, o)$$

- All policy learning methods discussed so far, e.g., Value and Policy Iteration, can be defined over S and O
- Coherent theory of learning and planning, with **courses of action at variable time scales**, yet at the same level

Value Functions Over Options

We can write the expression for optimal value as,

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

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

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

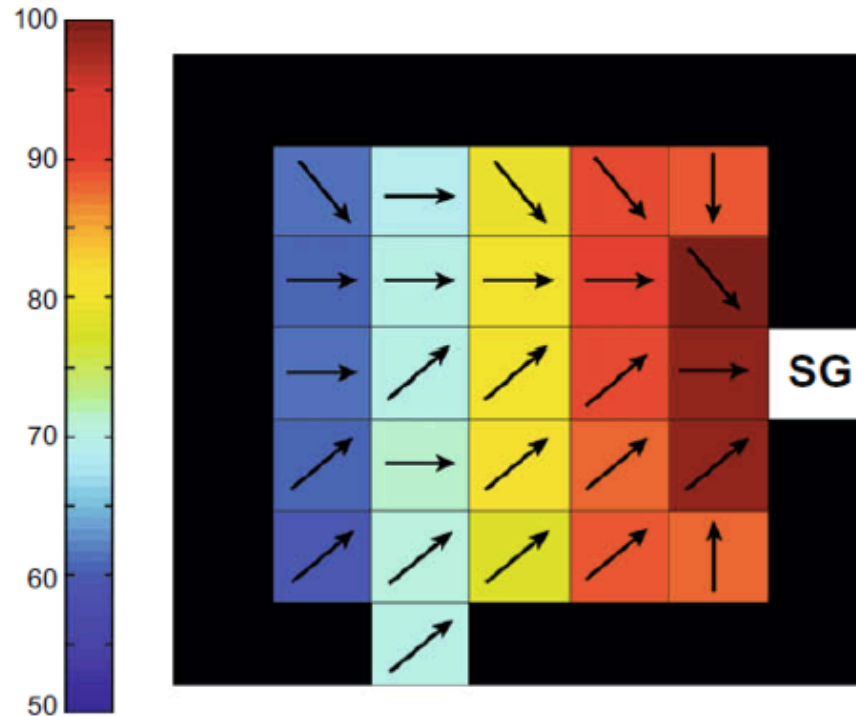
k being the duration of o when taken in s ; conditioning is over the event that the option is initiated at that state and time.

Motivations for Options Framework

- Add temporally extended activities to choices available to RL agent, without precluding planning and learning at finer grained MDP level
- Optimal policies over primitives are not compromised due to addition of options
- However, if an option is useful, learning will quickly find this out – prevent prolonged and useless ‘flailing about’

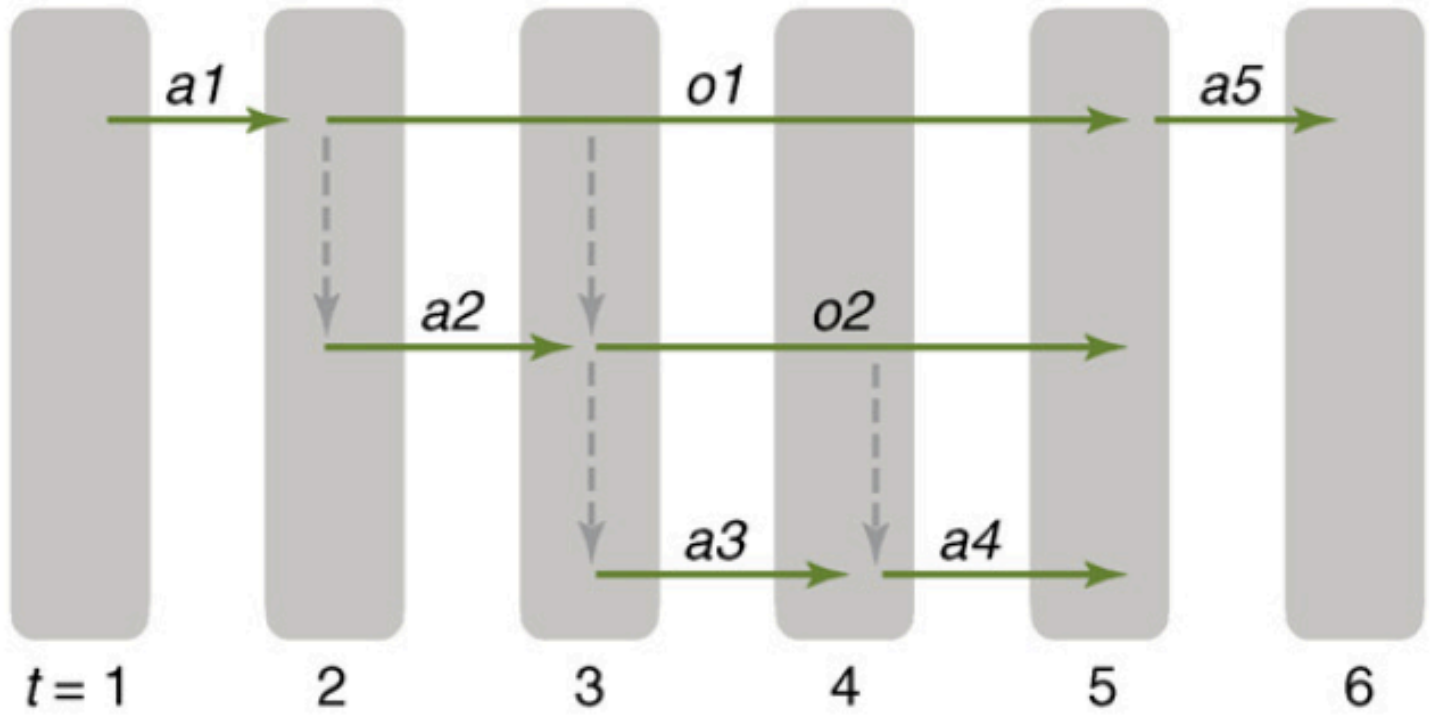
PS: If all options are 1-step, you recover the core MDP

Rooms Example – Policy within One Room



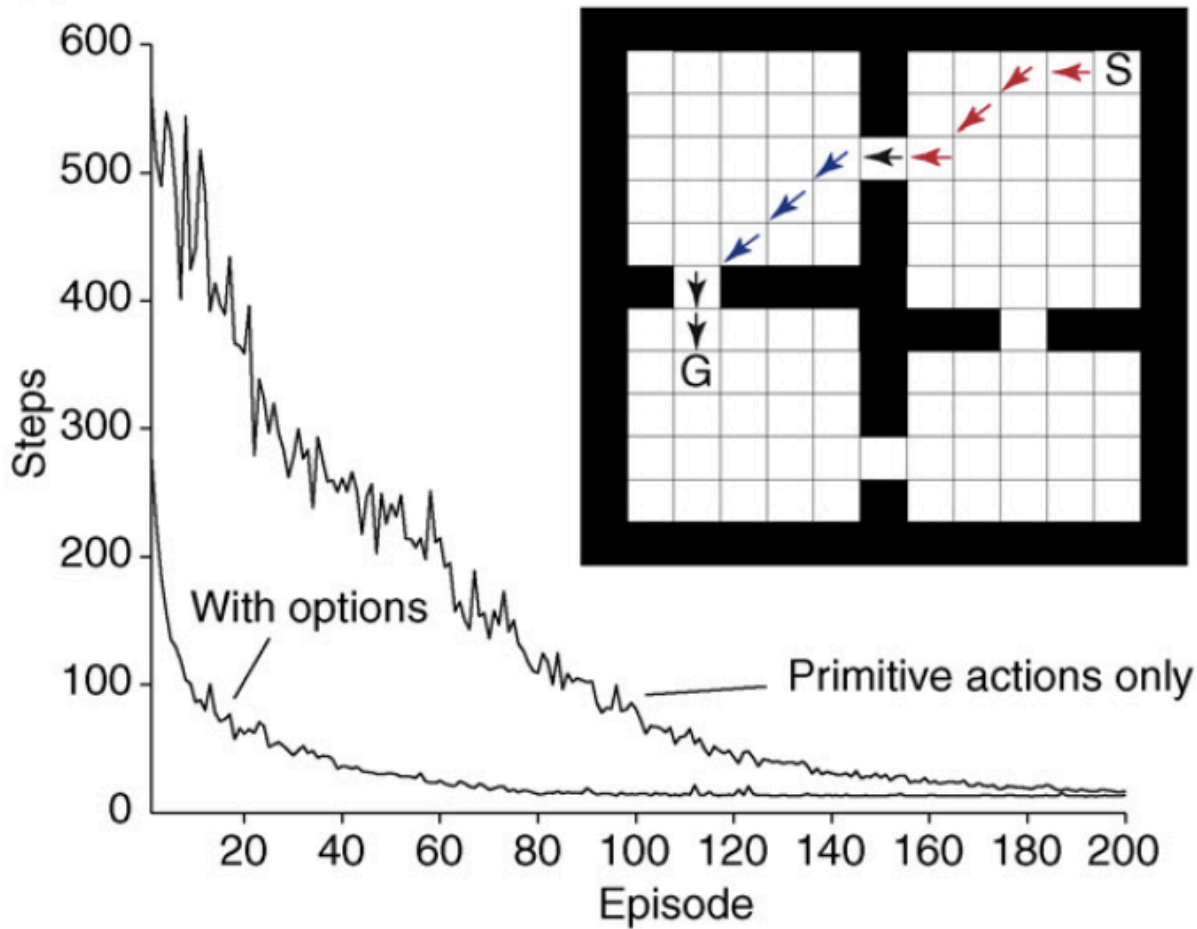
Colour = option specific value

Time Course of Use of Action/Option



[Source: M. Botvinick, Hierarchical models of behavior and prefrontal function, Trends in Cognitive Science 12(5), 2008]

Performance Improvement with Options



[Source: M. Botvinick, Hierarchical models of behavior and prefrontal function, Trends in Cognitive Science 12(5), 2008]

Learning Options

Goal/State Abstraction

- Why are these together?
 - Abstract goals typically imply abstract states
- Makes sense for classical planning
 - Classical planning uses state sets
 - Implicit in use of state variables
- Does this make sense for RL? Things to think about:
 - No goals
 - Markov property issues

Automatic Hierarchy Discovery

- Hard in contexts such as classical planning
- Within a single problem:
 - Battle is lost if all states considered (polynomial speedup at best)
 - If fewer states considered, when to stop?
- Across problems
 - Considering all states OK for few problems?
 - Generalize to other problems in class (**transfer learning**)
- How best to measure progress? Remains an open problem...

Bottlenecks: Principle for Option Learning

- **Subgoals** are created based on commonalities across multiple paths to a solution
- Cast the task of finding these commonalities as a multiple-instance learning problem
- System attempts to identify a target concept on the basis of “bags” of instances: positive bags have at least one positive instance, while negative bags consist of all negative instances.
- A successful trajectory corresponds to a positive bag, where the instances are the agent’s observations along that trajectory.
- A negative bag consists of observations made over an unsuccessful trajectory.

[Reference: A. McGovern, A.G. Barto, Accelerating reinforcement learning through discovery of useful subgoals, Proc. i-SAIRAS 2001.

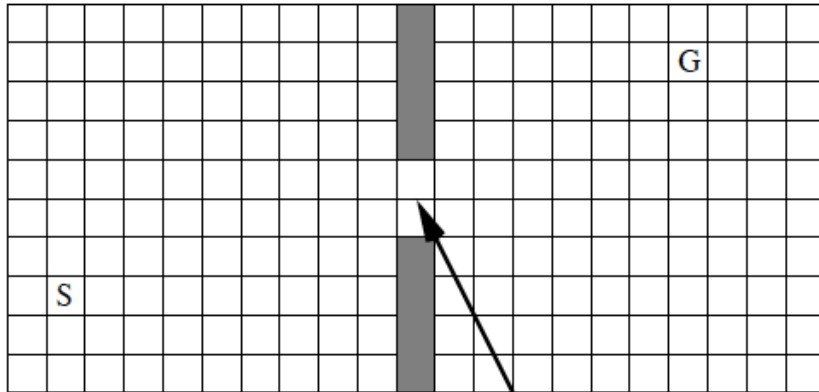
Motivation for Bottlenecks

- Agent searches for bottleneck regions in observation space.
- The idea of looking for bottleneck regions was motivated by studying room-to-room navigation tasks where the agent should quickly discover the utility of doorways as subgoals.
- IF the agent can recognize that a doorway is a kind of bottleneck by detecting that the sensation of being in the doorway always occurred somewhere on successful trajectories
 - BUT not always on unsuccessful ones,
- THEN it can create an option to reach the doorway.

Bottlenecks - Computationally

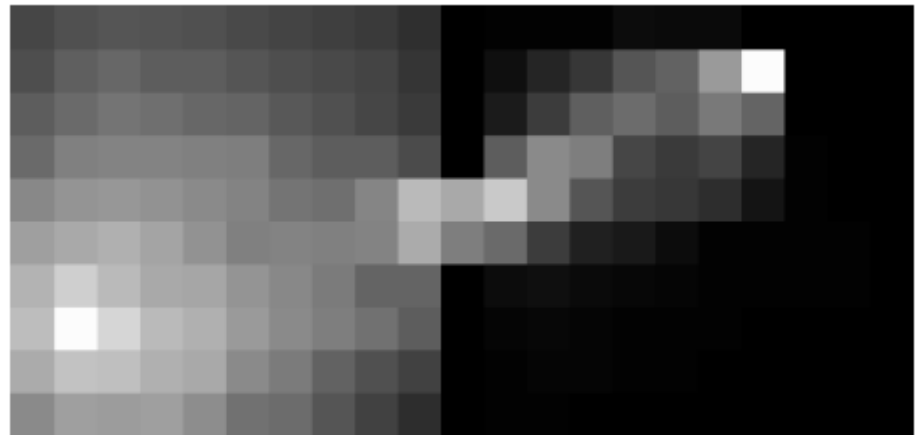
- If the agent uses some form of randomness to select exploratory primitive actions, it is likely to remain within the more strongly connected regions of the state space.
- An *option* for achieving a bottleneck region, on the other hand, will tend to connect separate strongly connected areas.
- For example, in a room-to-room navigation task, navigation using primitive movement commands produces relatively strongly connected dynamics within each room but not between rooms.
- A doorway links two strongly connected regions. By adding an option to reach a doorway sub-goal, the rooms become more closely connected.
- This allows the agent to more uniformly explore its environment.

Example Data in a Simulated Environment



Bottleneck area
and useful subgoal location

First-visit histogram:



Concept of Diverse Density

- Due to Maron and Lozano Perez
- Compute the probability $Pr(t)$ that the t^{th} concept from a set of concepts $\{c_i\}$ is the correct one
- If B_i^{\pm} denote the positive and negative bags of experienced traces,

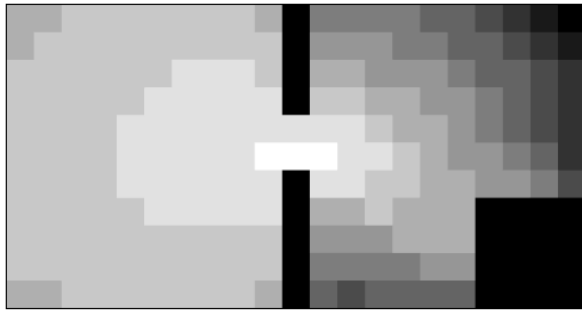
$$DD(t) = Pr(t|B_1^+, \dots, B_n^+, B_1^-, \dots, B_m^-)$$

$$DD(t) = \prod_{1 \leq i \leq n} Pr(t|B_i^+) \prod_{1 \leq j \leq m} Pr(t|B_j^-)$$

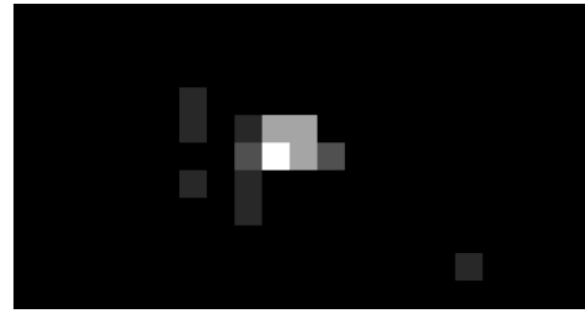
Options Discovery

- For the simulated environment being considered, when DD is calculated over all traces:

A: Average Diverse Density



B: Subgoals Discovered



- From this, options can be constructed by learning local policies to the subgoals