## An Algebraic Approach to Abstraction in Reinforcement Learning

Doctoral Dissertation Defense
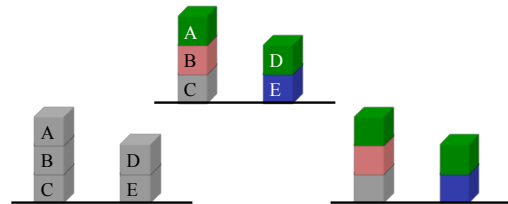
Balaraman Ravindran

Advisor: Andrew G. Barto

Committee: Roderic A. Grupen

Sridhar Mahadevan

Neil E. Berthier (Dept. of Psychology)

## Abstraction



- Ignore information irrelevant for the task at hand.
- Form simpler representation.

## Abstraction

- A key reason that humans are effective problem solvers
  - Learn and plan at a higher level
  - Knowledge transfer
  - c.f. macros, chunks, skills, behaviors, . . .

- **Temporal abstraction** or plan abstraction
- **Spatial abstraction**
- **Combination of the two**

## Motivation

- Well studied problem in AI
- Focus of thesis:
  - Decision theoretic setting
    - Markov decision processes
  - General framework
    - Accommodate different notions of abstraction
      - Aggregation, symmetry (Zinkevich and Balch ’01, Popplestone and Grupen ’00), projections, structured abstractions (Boutilier et al. ’94, ’95, ’01)
  - Formal algebraic framework
    - Group theory, model minimization, operations research
  - Combination of temporal and spatial abstraction
    - Behaviors in a relative frame of reference
      - Efficient knowledge transfer

## Outline of Thesis

- Abstraction in decision making
  - Algebraic framework
  - Exploiting symmetry and structure
  - Approximate equivalence
- Abstraction in hierarchical reinforcement learning
  - Hierarchical task decomposition
  - Relativized options
  - Algorithms for dynamic abstraction
    - Choosing transformations
    - Deictic representation

## Outline of Thesis

- Abstraction in decision making
  - ➡ Algebraic framework
  - Exploiting symmetry and structure
  - ➡ Approximate equivalence
- Abstraction in hierarchical reinforcement learning
  - Hierarchical task decomposition
  - ➡ Relativized options
  - Algorithms for dynamic abstraction
    - ➡ Choosing transformations
    - Deictic representation

## Outline of Talk

- Abstraction in decision making
  - → Algebraic framework

  - → Approximate equivalence
- Abstraction in hierarchical reinforcement learning

  - → Relativized options
  - Algorithms for dynamic abstraction
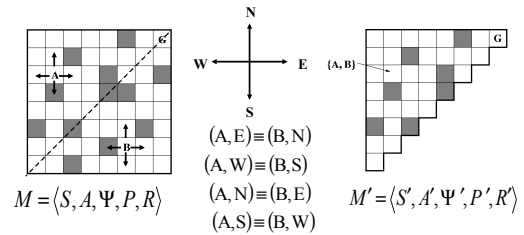    - → Choosing transformations

  - Summary

---

## Outline of Talk

- Abstraction in decision making
  - → Algebraic framework
    - Markov Decision Processes
    - MDP homomorphisms
    - Some theoretical results
  - Approximate equivalence
- Abstraction in hierarchical reinforcement learning

  - Relativized options
  - Algorithms for dynamic abstraction
    - Choosing transformations

  - Summary

---

## Markov Decision Processes

- MDP, $M$, is the tuple: $M = \langle S, A, \Psi, P, R \rangle$
  - $S$ : set of states.
  - $A$ : set of actions.
  - $\Psi \subseteq S \times A$ : set of admissible state-action pairs.
  - $P : \Psi \times S \to [0,1]$ : probability of transition.
  - $R : \Psi \to \Re$ : expected reward.
- Policy $\pi : S \to A$  (can be stochastic)
- Maximize total expected reward.

---

## Example



$M = \langle S, A, \Psi, P, R \rangle$

$(A, E) \equiv (B, N)$
$(A, W) \equiv (B, S)$
$(A, N) \equiv (B, E)$
$(A, S) \equiv (B, W)$

$M' = \langle S', A', \Psi', P', R' \rangle$

---

## Homomorphisms

Group homomorphism

Let $G$ and $G'$ be groups with operations $+$ and $+'$ respectively

$h : G \to G'$ is a group homomorphism iff

$h(x + y) = h(x) +' h(y) \quad \forall x, y \in G$

$$
\begin{array}{ccc}
G \times G & \xrightarrow{\;+\;} & G \\
{\scriptstyle h \times h}\downarrow & & \downarrow{\scriptstyle h} \\
G' \times G' & \xrightarrow{\;+'\;} & G'
\end{array}
$$

---

## Homomorphisms (cont.)

Automaton homomorphism

in the autonomous case:

$M = \langle S, \delta \rangle, \quad M' = \langle S', \delta' \rangle$

state set — transition function

$$
\begin{array}{ccc}
S & \xrightarrow{\;\delta\;} & S \\
{\scriptstyle h}\downarrow & & \downarrow{\scriptstyle h} \\
S & \xrightarrow{\;\delta'\;} & S
\end{array}
$$

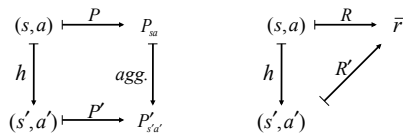$h(\delta(s)) = \delta'(h(s))$

induces equivalence classes in $S$

2

## MDP Homomorphism
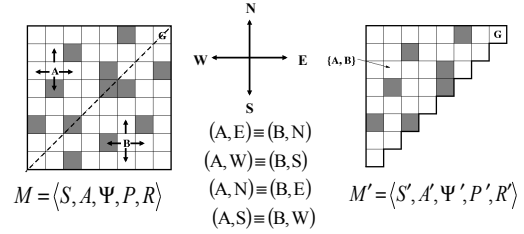
MDPs $M = \langle S, A, \Psi, P, R \rangle$, $M' = \langle S', A', \Psi', P', R' \rangle$

surjection $h : \Psi \to \Psi'$ defined by $h((s,a)) = (f(s), g_s(a))$ where :
$f : S \to S'$, $g_s : A_s \to A'_{f(s)}$, for all $s \in S$, are surjections such that
for all $s, \bar{s} \in S$, and $a \in A_s$ :

(1) $\quad P'(f(s), g_s(a), f(\bar{s})) = \sum_{t \in [\bar{s}]_f} P(s, a, t)$

(2) $\quad R'(f(s), g_s(a)) = R(s,a)$

$$(s,a) \xmapsto{\;P\;} P_{sa} \qquad (s,a) \xmapsto{\;R\;} \bar{r}$$

$$h \downarrow \quad agg. \downarrow \qquad\qquad h \downarrow \quad R' \nearrow$$

$$(s',a') \xmapsto{\;P'\;} P'_{s'a'} \qquad (s',a')$$

## Example



$M = \langle S, A, \Psi, P, R \rangle$

$(A,E) \equiv (B,N)$
$(A,W) \equiv (B,S)$
$(A,N) \equiv (B,E)$
$(A,S) \equiv (B,W)$

$M' = \langle S', A', \Psi', P', R' \rangle$

$h(A,E) = h(B,N) = (\{A,B\}, E)$

State dependent action recoding

## Some Theoretical Results

[generalizing those of Dean and Givan, 1997]

- Optimal Value equivalence:
  If $h(s,a) = (s',a')$ then $Q^*(s,a) = Q^*(s',a')$.
- Corollary:
  If $h(s_1, a_1) = h(s_2, a_2)$ then $Q^*(s_1, a_1) = Q^*(s_2, a_2)$.

**Theorem:** If $M'$ is a homomorphic image of $M$, then a policy optimal in $M'$ induces an optimal policy in $M$.

- Solve homomorphic image and *lift* the policy to the original MDP. ▪

## Model Minimization

- Finding reduced models that preserve some aspects of the original model
- Various modeling paradigms
  - Finite State Automata (Hartmanis and Stearns '66)
    - Transition Behavior
  - Model Checking (Emerson and Sistla '96, Lee and Yannakakis '92)
    - Correctness of system models
  - Markov Chains (Kemeny and Snell '60)
    - Steady state distribution
  - MDPs (Dean and Givan '97, Ravindran and Barto '02)
    - Optimal solutions
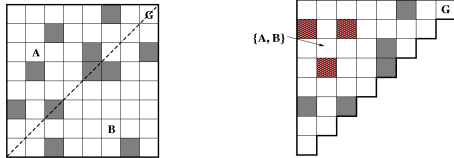
## MDP Minimization

- In general, NP-hard
  - Polynomial time algorithm for computing homomorphic image, under certain assumptions
    - Extends Dean and Givan '97, Lee and Yannakakis '92

- State dependent action recoding
  - Greater reduction in problem size
  - Model symmetries
    - Reflections, rotations, permutations ▪

## Outline of Talk

- Abstraction in decision making
  - Algebraic framework
    - Approximate homomorphisms
    - Error bounds
  - Approximate equivalence
    - Bounded parameter approximations
- Abstraction in hierarchical reinforcement learning

  - Relativized options
  - Algorithms for dynamic abstraction
    - Choosing transformations
  - Summary

## Approximate Notions of Equivalence

- Complete and exact equivalence often do not exist.
- Approximate equivalence.
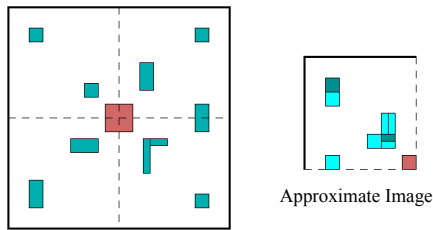  - "Equivalent" state-action pairs have nearly same behavior.



---

## Approximate Homomorphisms

- Use averages
- Relax homomorphism criteria:
  - $P'\big(f(s), g_s(a), f(\bar{s})\big) = \sum_{t \in [\bar{s}]_f} P(s, a, t)$
  - Compute $\sum_{t \in [\bar{s}]_f} P(s, a, t)$ for all $(s,\ a)$

$$P'\big(f(s), g_s(a), f(\bar{s})\big) = \frac{1}{\big|[(s,a)]_h\big|} \sum_{(q,b) \in [(s,a)]_h} \sum_{t \in [\bar{s}]_f} P(q, b, t)$$

  - Similar computation for the reward function.

---

## Example



Approximate Image

Task is to reach red goal area.

---

## Error Bound

- Approximate homomorphism between arbitrarily different MDPs!
- Useful when loss in performance is acceptable.
- Bound the maximum difference in optimal value function in $M$ and the value of the lifted optimal policy.
  - Specializes Whitt '78.
  - Function of maximum difference in the probabilities and rewards that are averaged.
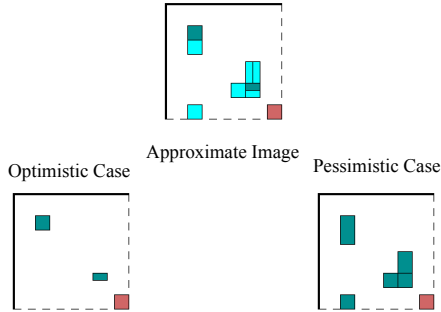
---

## Error bound (cont.)

- $K_p$ – maximum difference between $P'\big(f(s), g_s(a), f(\bar{s})\big)$ and $\sum_{t \in [\bar{s}]_f} P(s, a, t)$
- $K_r$ – corresponding difference in reward
- $\Delta$ – the range of the reward function
- $\gamma$ – the discount factor, $0 \le \gamma < 1$

$$\big\| V^* - V'^* \big\|_{\max} \le \frac{2}{1-\gamma}\left( K_r + \frac{\gamma}{1-\gamma} \Delta \frac{K_p}{2} \right)$$

---

## Bounded Parameter Approximation

- Model as a map onto a *Bounded-parameter MDP* (Givan, Leach and Dean '00)
  - Transition probabilities and rewards given by bounded intervals
  - Upper and lower bounds on optimal values of states
  - Loose bounds

## Example Revisited

Approximate Image

Optimistic Case

Pessimistic Case

---

## Outline of Talk

- Abstraction in decision making
  - Algebraic framework

  - Approximate equivalence
- Abstraction in hierarchical reinforcement learning
  - • Semi-Markov decision processes
  - • Options framework
  - Relativized options  • Relativized options
  - Algorithms for dynamic abstraction
    - • Choosing transformations
  - Summary

---

## (discrete-time) semi-Markov Decision Process

- SMDP, $M$, is the tuple: $M = \langle S, A, \Psi, P, R \rangle$
  - $S$ : set of states.
  - $A$ : set of actions.
  - $\Psi \subseteq S \times A$ : set of admissible state-action pairs.
  - $P : \Psi \times S \times N \to [0,1]$ : transition probabilities.
  - $R : \Psi \times N \to \Re$ : expected reward.
- Policy (stationary, stochastic): $\pi : \Psi \to [0,1]$
- Maximize expected return.
- Generalize MDP homomorphism.

---

## Hierarchical Reinforcement Learning

Options (Sutton, Precup, & Singh, 1999): **A generalization of actions to include temporally-extended courses of action**
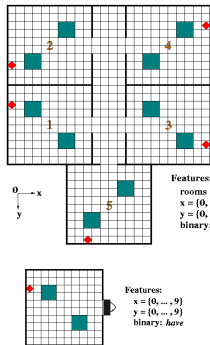
An option is a triple $o = < I, \pi_o, \beta >$

- $I \subseteq S$ is the set of states in which $o$ may be started
- $\pi_o : \Psi \to [0,1]$ is the (stochastic) policy followed during $o$
- $\beta : S \to [0,1]$ is the probability of terminating in each state
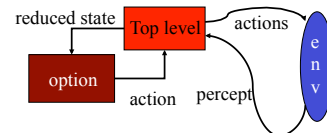
Example: robot docking

$I$ : all states in which charger is in sight
[?][?][?] : pre-defined controller
[?] : terminate when docked or charger not visible

---

## Sub-goal Options

Features:
rooms = {0, 1, 2, 3, 4, 5}
x = {0, ... , 9}
y = {0, ... , 19}
binary: $have_i$, $i = 1, ... , 5$

0 — x
y

Features:
x = {0, ... , 9}
y = {0, ... , 9}
binary: $have$

- Task is to collect all objects in the world
- 5 options – one for each room.
- Markov, subgoal options
- Implicitly define option policy
- Employ option specific abstraction

---

## Relativized Options

reduced state    Top level    actions

option    action    percept    env
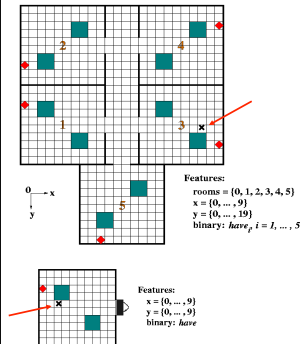
Relativized option:
$$O = \langle h, M_O, I, \beta \rangle$$
$h$ : Option homomorphism.
$M_O$ : Option SMDP. (Image of $h$.)
$I \subseteq S$ : Initiation set.
$\beta : S_O \to [0,1]$ : Termination criterion.

## Rooms world task



**Features:**
rooms = {0, 1, 2, 3, 4, 5}
x = {0, ... , 9}
y = {0, ... , 19}
binary: $have_i$, i = 1, ... , 5

**Features:**
x = {0, ... , 9}
y = {0, ... , 9}
binary: *have*

- Task is to collect all objects in the world
- 5 options – one for each room
- Single relativized option – *get-object-exit-room*
- Partial homomorphism
- Especially useful when learning option policy
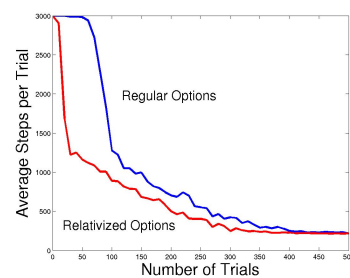  - Speed up.
  - Knowledge transfer.

## Experimental Setup

- Regular Agent
  - 5 options, one for each room
  - Option reward of +1 on exiting room with object
- Relativized Agent
  - 1 relativized option, known homomorphism
  - Same option reward
- Global reward of +1 on completing task
- Actions fail with probability 0.1
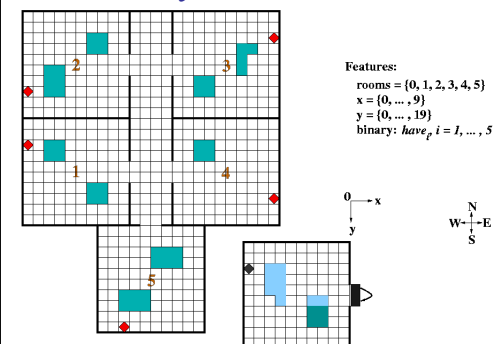
## Learning Algorithm

- Hierarchical SMDP Q-learning (Dietterich ' 00b)
  - Q-learning at the lowest level (Watkins ' 89)
  - SMDP Q-learning at the higher levels (Bradtke and Duff ' 95)
- Simultaneous learning at all levels
  - Converges to recursively optimal policy
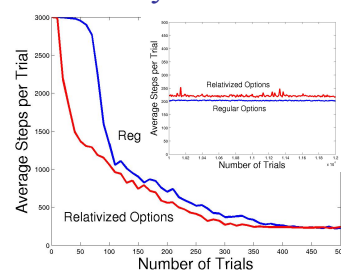    - Using results from Dietterich ' 00a

## Results



- Average over 100 runs

## Asymmetric Testbed



**Features:**
rooms = {0, 1, 2, 3, 4, 5}
x = {0, ... , 9}
y = {0, ... , 19}
binary: $have_i$, i = 1, ... , 5

## Results – Asymmetric Testbed



- Still significant speed up in initial learning
- Asymptotic performance slightly worse

6

## Outline of Talk

- Abstraction in decision making
  - Algebraic framework
  - Approximate equivalence
- Abstraction in hierarchical reinforcement learning
  - Relativized options
  - Algorithms for dynamic abstraction
    - Choosing transformations
  - Summary

## Choosing Transformations
### Motivation

- Relax prior knowledge requirement
  - Unknown homomorphism
- Option SMDP and policy can be viewed as a *policy schema* (Schmidt '75, Arbib '95)
  - Template of a policy
  - Acquire schema in a prototypical setting
  - Learn bindings of sensory inputs and actions to schema
- Assume set of possible bindings available

## Choosing Transformations
### Problem Formulation

- Given:
  - $M_O, I, \beta$ of a relativized option
  - $H$, a family of transformations
- Identify the option homomorphism $h$
- Formulate as a parameter estimation problem
  - One parameter, takes values from $H$
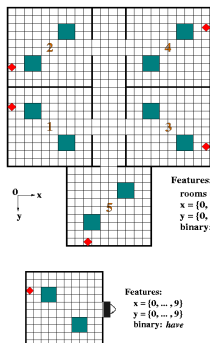  - Samples: $\langle s_1, a_1, s_2, a_2, \cdots \rangle$
  - Bayesian learning

## Choosing Transformations
### Algorithm

- Assume uniform prior: $p_0(h, \bar{s})$
- Experience: $\langle s_n, a_n, s_{n+1} \rangle$

$$P(\langle s_n, a_n, s_{n+1} \rangle | h, \bar{s}) = P_O(f(s_n), g_{s_n}(a_n), f(s_{n+1}))$$
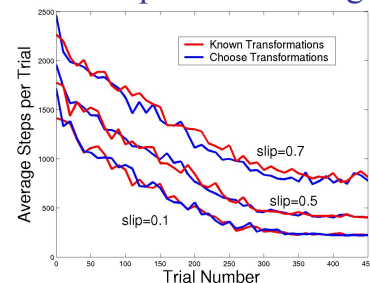
- Update Posteriors:

$$p_n(h, \bar{s}) = \frac{P_O(f(s_n), g_{s_n}(a_n), f(s_{n+1})) \cdot p_{n-1}(h, \bar{s})}{\text{Normalizing Factor}}$$

## Rooms world task



Features:
rooms = {0, 1, 2, 3, 4, 5}
x = {0, ... , 9}
y = {0, ... , 19}
binary: $have_i$, i = 1, ... , 5

Features:
x = {0, ... , 9}
y = {0, ... , 9}
binary: have

- Train in room 1
- 8 candidate transformations
  - Reflections about x and y axes and the x=y and x=-y lines
  - Rotations by integer multiples of 90 degrees

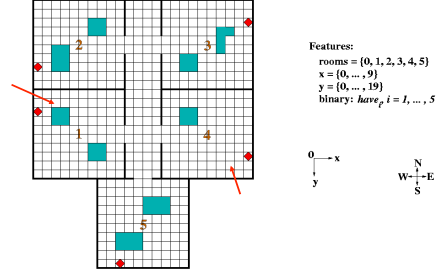## Results – Speed of convergence



- Not much of a difference since the task is too easy
- Correct transformation identified in 15 iterations

## Choosing Transformations
### Approximate Equivalence

- More complex domains
- Problem with Bayesian update
  - Use prototypical room as option schema
  - Susceptible to incorrect samples
- Use a heuristic lower bound

---

## Example



Features:
rooms = {0, 1, 2, 3, 4, 5}
x = {0, ... , 9}
y = {0, ... , 19}
binary: $have_i$, $i = 1, ... , 5$

$$p_n(h, \bar{s}) = \frac{P_O(f(s_n), g_{s_n}(a_n), f(s_{n+1})) \cdot p_{n-1}(h, \bar{s})}{\text{Normalizing Factor}}$$
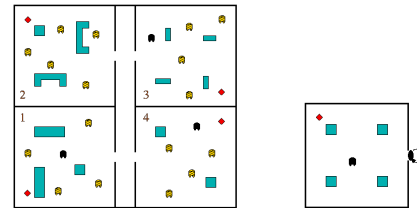
---

## Choosing Transformations
### Heuristic Update Rule

- Use a heuristic update rule:

$$w_n(h, \bar{s}) = \frac{\overline{P}(f(s_n), g_{s_n}(a_n), f(s_{n+1})) \cdot w_{n-1}(h, \bar{s})}{\text{Normalizing Factor}}$$

where, $\overline{P}(s, a, s') = \max(\nu, P_O(s, a, s'))$

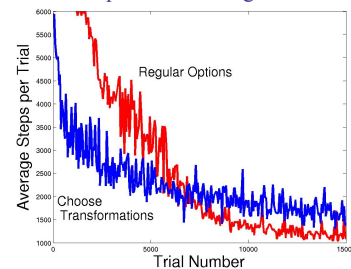and $\nu$ is a small positive constant.

---

## Complex Game World



- Gather all 4 diamonds in the world
- $25 \times 10^{55}$ states
- 40 transformations
  - 8 spatial transformations combined with 5 projections

---

## Experimental Setup

- Regular agent
  - 4 sub-goal options
- Relativized agent
  - Uses option MDP shown earlier
  - Chooses from 40 transformations
- Room 2 has no right transformation
- Hierarchical SMDP Q-learning
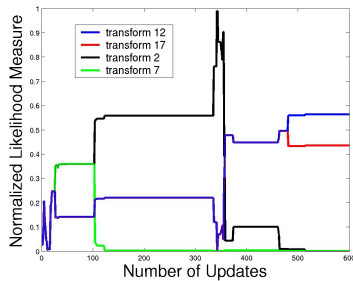
---

## Results
### Speed of Convergence



- Learning the policy is more difficult than learning the correct transformation!
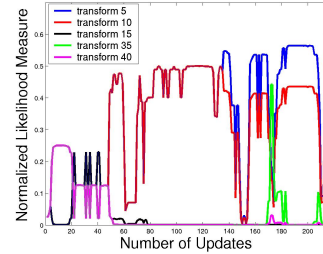
## Results
### Transformation Weights in Room 4



- Transformation 12 eventually converges to 1

## Results
### Transformation Weights in Room 2



- Weights oscillate a lot
- Some transformation dominates eventually
  - Changes from one run to another

## Choosing Transformations

- Related work
  - Multiple forward models (Haruno et al. ' 01, Doya et al. '02)
  - Dynamic control models (Coelho and Grupen ' 98)
  - Variably bound controllers (Huber and Grupen ' 99)

- Representations can be designed to implicitly perform transformations
  - Formalizes such representations
  - E.g. Deictic representations

## Summary of Contributions

- Developed an abstraction framework for MDPs
  - Introduced MDP homomorphisms
    - State dependent action recoding
  - Theoretical results
- Approximate homomorphisms
  - Bound maximum loss
  - Upper and lower bound performance

## Summary of Contributions
### (cont.)

- Abstraction in hierarchical systems
  - Relativized options
    - An option defined in a relative frame of reference
    - Uses partial homomorphisms
  - Policy schema
    - Policy template
  - Bayesian algorithm for choosing the right bindings
    - Heuristic modification for approximate equivalence
    - Complex game domain

## Other Contributions

- Exploiting structure and symmetry
  - Structured morphisms
  - Symmetry groups
    - Reflections, rotations and permutations
  - Polynomial time algorithm
- Hierarchical decomposition framework
  - Based on SMDP homomorphisms
  - Relation to safe state abstraction (Dietterich ' 00a)
- Deictic option schema
  - Representation based on pointers (Agre ' 88)
  - Modification of Bayesian algorithm

## Future Work

- Practical application of framework
  - Humanoid experiments
- Abstraction algorithms
  - Symbolic representations (Feng et al. ' 02,' 03)
- Relation to partial observability
- Relation to other abstract representations
  - Probabilistic relational models (Getoor et al. ' 01