

Hierarchy Discovery in Partially Observable Domains

April 9th, 2010

Reinforcement Learning Workshop, Barbados

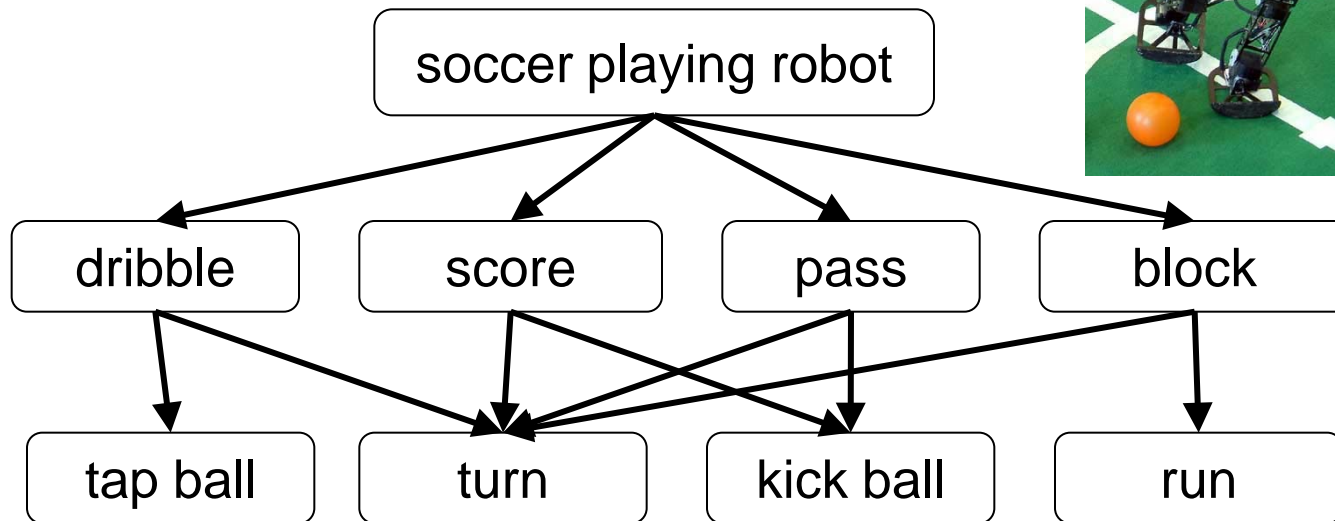


Presented by Pascal Poupart
University of Waterloo, Canada

Joint work with Marc Toussaint and Laurent Charlin₁

Hierarchies in Planning

- Idea: task decomposed into subtasks arranged hierarchically



- Benefits: temporal/spatial abstraction, sub-policy reuse, intuitive policy representation

Where do hierarchies come from?

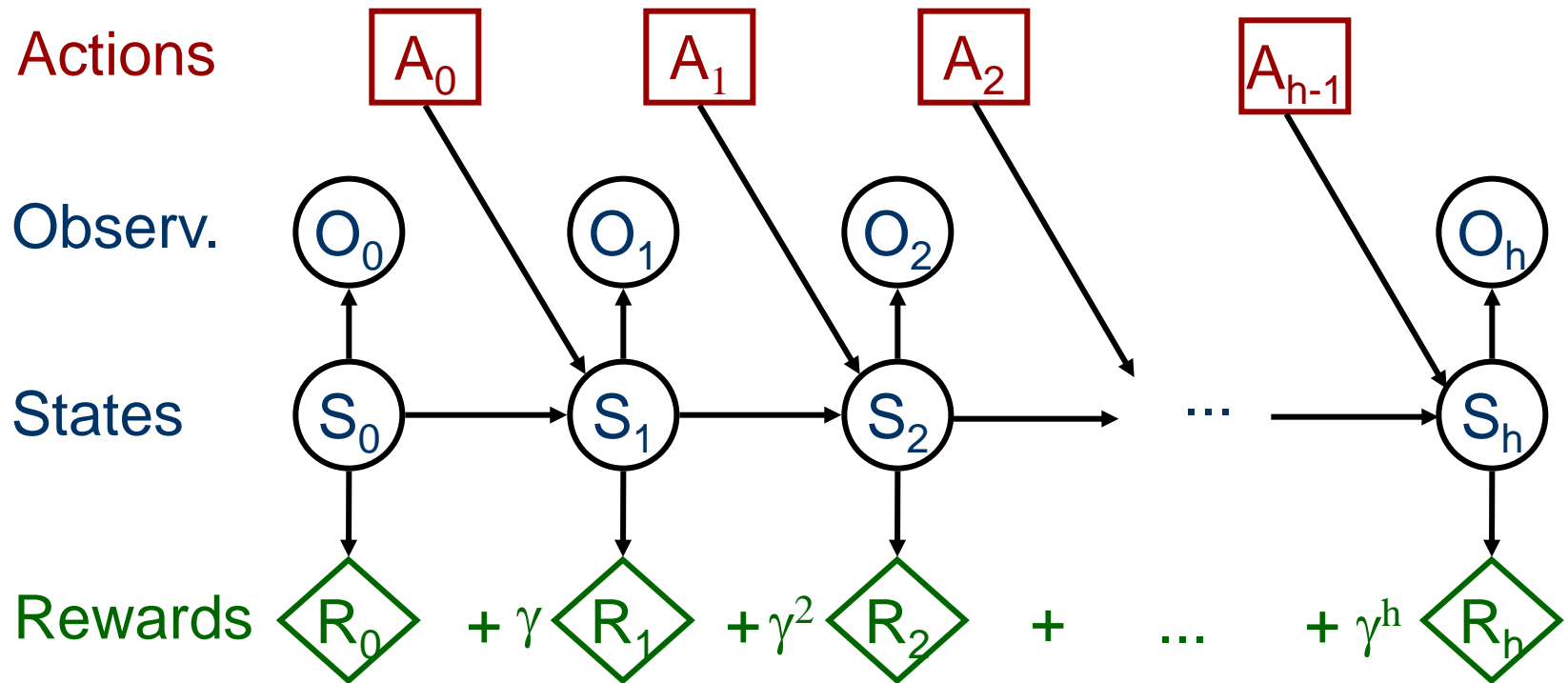
- Most of the time: expert knowledge
- Can be discovered automatically:
 - Local search: traverse or
 - Complete search: states or satisfy necessary conditions
- How can we discover hierarchies in partially observable domains?

Common assumption:
observable states/conditions

Outline

- Background:
 - Partially observable Markov decision processes
 - Finite state controllers
 - Hierarchical controllers
- Hierarchy discovery by
 - Non-convex optimization
 - Likelihood maximization
- Experiments
- Conclusion

POMDP Graphical Representation



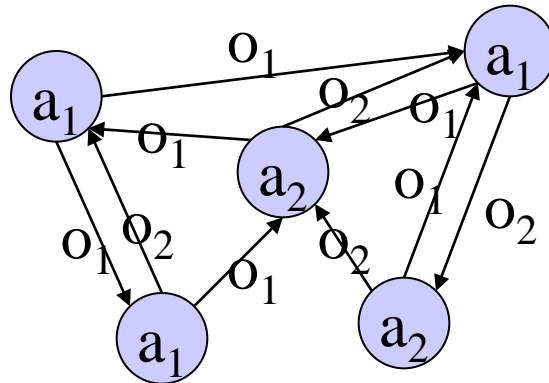
Solution: **policy** π maximizes **expected total rewards**

Policy Optimization

- Beliefs b : distribution over states
 - Bayes filtering: $b_{ao'}(s') = k \sum_s b(s) \Pr(s'|s,a) \Pr(o|s',a)$
 - Sufficient statistic of past observations and actions
- Policy $\pi : B \rightarrow A$
 - mapping from beliefs to actions
- Value function $V^\pi(b) = \sum_t \gamma^t \mathbf{E}_{b_t|\pi} [R]$
- Optimal policy $\pi^*: V^*(b) \geq V^\pi(b)$ for all π, b
 - $V^*(b) = \max_a \mathbf{E}_b[R] + \gamma \sum_{o'} \Pr(o'|s,a) V^*(b_{ao'})$

Finite State Controllers

- Alternative policy representation: controllers
 - Action mapping: $\alpha: N \rightarrow A$ or $\Pr(a|n)$
 - Next node mapping: $\sigma: N \times O \rightarrow N$ or $\Pr(n'|a,n)$



- Policy optimization: select best α and σ

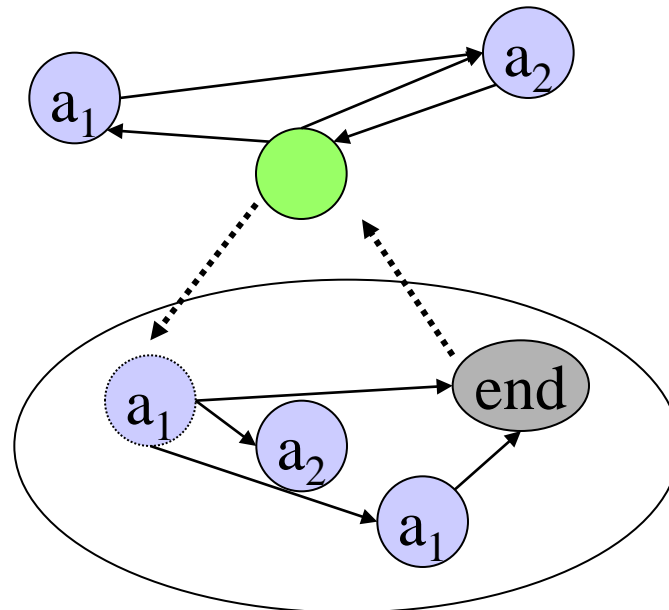
Controller Optimization

- Policy search: alternate between
 - Policy evaluation
 - Policy improvement
- Optimization problem (Amato et al., 2007)

$$\begin{aligned}
 & \max_{x,y} \sum_s b_o(s) \underbrace{V_{n_o}(s)}_y \\
 \text{s.t. } & \underbrace{V_n(s)}_y = \sum_{a,n'} \left[\underbrace{\Pr(a, n' | n, o_k)}_x R(s, a) + \sum_{s',o} \Pr_\gamma(s' | s, a) \Pr(o | s', a) \underbrace{\Pr(n', a | n, o)}_x \underbrace{V_{n'}(s')}_y \right] \quad \forall n, s \\
 & \underbrace{\Pr(n', a | n, o)}_x \geq 0 \quad \forall n', a, n, o \quad \sum_{n',a} \underbrace{\Pr(n', a | n, o)}_x = 1 \quad \forall n, o \\
 & \sum_{n'} \underbrace{\Pr(n', a | n, o)}_x = \sum_{n'} \underbrace{\Pr(n', a | n, o_k)}_x \quad \forall a, n, o
 \end{aligned}$$

Hierarchical Controllers

- Hansen & Zhou, 2005: let some nodes be sub-controllers
 - Action mapping: $\alpha: N \rightarrow A$ or $\Pr(a|n)$
 - Next node mapping: $\sigma: N \times O \rightarrow N$ or $\Pr(n'|n,o)$
 - Child mapping (abstract nodes): $\phi: N \rightarrow N$ or $\Pr(n'|n^{\text{par}})$
 - Special exit nodes



Hierarchical Controllers

- Local reward function: $R(s, \hat{a}) = \sum_n \Pr(n|n^{\text{par}}) V_n(s)$
 - $V_n(s) = \sum_a \Pr(a|n) [R(s, a) + \gamma \sum_{s', n', o'} \Pr(s'|s, a) \Pr(o'|a, s') \Pr(n'|n, o') V_{n'}(s')]$
- Local transition prob: $\Pr(s'|\hat{a}, s) = \text{oc}(s', n^{\text{end}})$
 - Discounted occupancy frequency
 - $\text{oc}(s', n') = \Pr(n'|n^{\text{par}}) + \gamma \sum_{s, n, a, o'} \text{oc}(s, n) \Pr(a|n) \Pr(s'|s, a) \Pr(o'|a, s') \Pr(n'|n, o')$

Optimization Problem

- Charlin, Poupart & Shioda, 2006

$$\max_{w,x,y,z} \sum_{s \in S} b_0(s) \underbrace{V_{n_0}(s)}_y \quad (3)$$

$$\text{s.t. } \underbrace{V_n(s)}_y = \sum_{a,n'} \left[\underbrace{\Pr(n', a|n, o_k)}_x R(s, a) + \sum_{s',o} \Pr_\gamma(s'|s, a) \Pr(o|s', a) \underbrace{\Pr(n', a|n, o)}_x \underbrace{V_{n'}(s')}_y \right] \quad \forall s, n \quad (4)$$

$$\underbrace{V_{\bar{n}}(s)}_y = \sum_{n_{beg}} \underbrace{\Pr(n_{beg}|\bar{n})}_z \left[\underbrace{V_{n_{beg}}(s)}_y + \sum_{s_{end}, a, n'} \underbrace{oc(s_{end}, n_{end}|s, n_{beg})}_w \left[\underbrace{\Pr(n', a|\bar{n}, o_k)}_x R(s_{end}, a) \right. \right. \\ \left. \left. + \sum_{s',o} \Pr_\gamma(s'|s_{end}, a) \Pr(o|s', a) \underbrace{\Pr(n', a|\bar{n}, o)}_x \underbrace{V_{n'}(s')}_y \right] \right] \quad \forall s, \bar{n} \quad (5)$$

$$\underbrace{oc(s', n'|s_0, n_0)}_w = \delta(s', n', s_0, n_0) + \sum_{s,o,a} \left[\right. \quad (6)$$

$$\left. \sum_n \underbrace{oc(s, n|s_0, n_0)}_w \Pr_\gamma(s'|s, a) \Pr(o|s', a) \underbrace{\Pr(n', a|n, o)}_x \right] \quad \left. \vphantom{\sum_n} \right\} n \text{ concrete (6a)}$$

$$\left. + \sum_{s_{end}, n_{beg}, \bar{n}} \underbrace{oc(s, \bar{n}|s_0, n_0)}_w \Pr_\gamma(s'|s_{end}, a) \Pr(o|s', a) \right. \\ \left. \underbrace{oc(s_{end}, n_{end}|s, n_{beg})}_w \underbrace{\Pr(n', a|\bar{n}, o)}_x \underbrace{\Pr(n_{beg}|\bar{n})}_z \right] \quad \forall s_0, s', n_0, n' \quad \left. \vphantom{\sum_{s_{end}, n_{beg}, \bar{n}}} \right\} \bar{n} \text{ abstract (6b)}$$

$$\Pr(\bar{n}'|\bar{n}) = 0 \text{ if } label(\bar{n}') \leq label(\bar{n}), \forall \bar{n}, \bar{n}' \quad (7)$$

Optimization Techniques

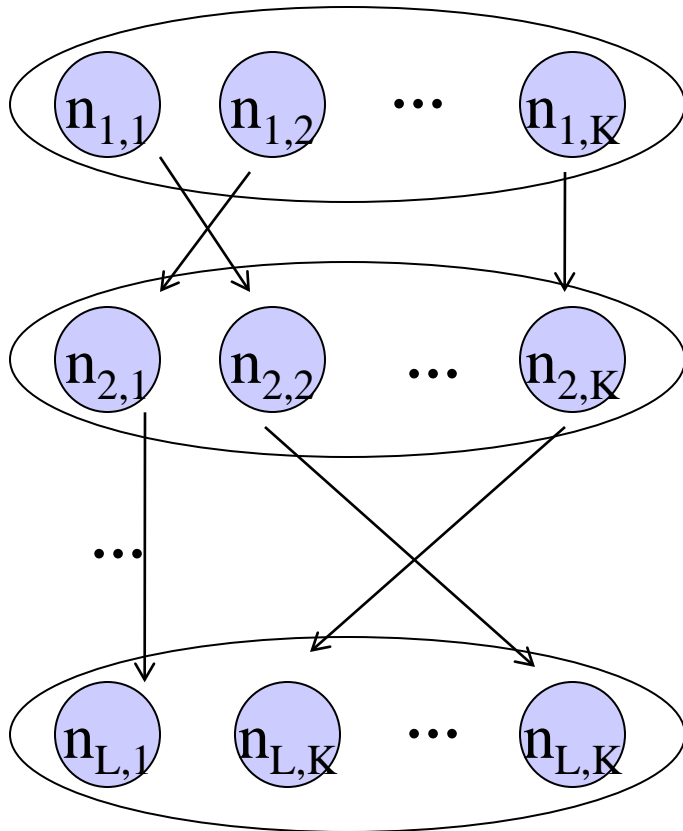
- Non-convex quartic optimization problem
 - Non-convex solvers do not scale
- Approximations:
 - Mixed-integer non-linear programming
 - Bounded hierarchical policy iteration
 - Do not scale either

Hierarchy Discovery

- Why discover the hierarchy since this seems to make the problem harder?
- Problem complexity remains the same
 - Search in a different policy space
 - Bias the search towards hierarchical policies
 - May find exponentially more compact policies
- Reveal interesting structure
 - Facilitate policy explanation for humans

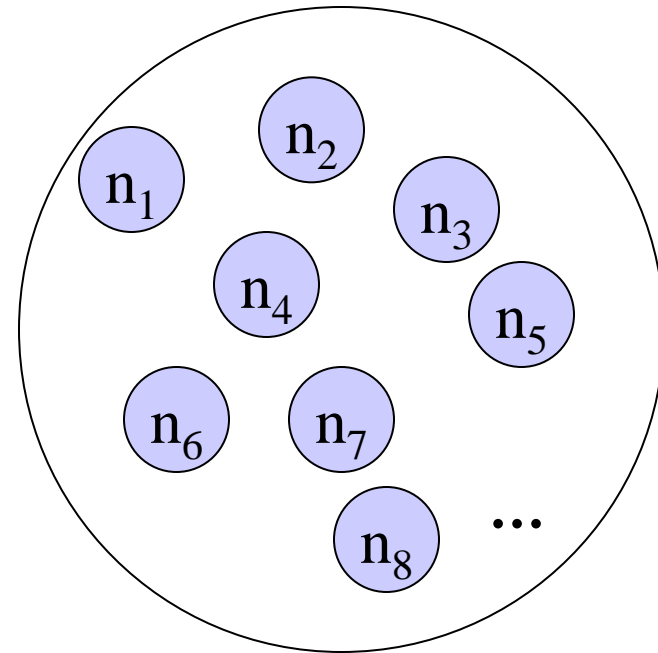
Hierarchical vs Flat Controllers

Hierarchical



$K * L$ nodes

Flat



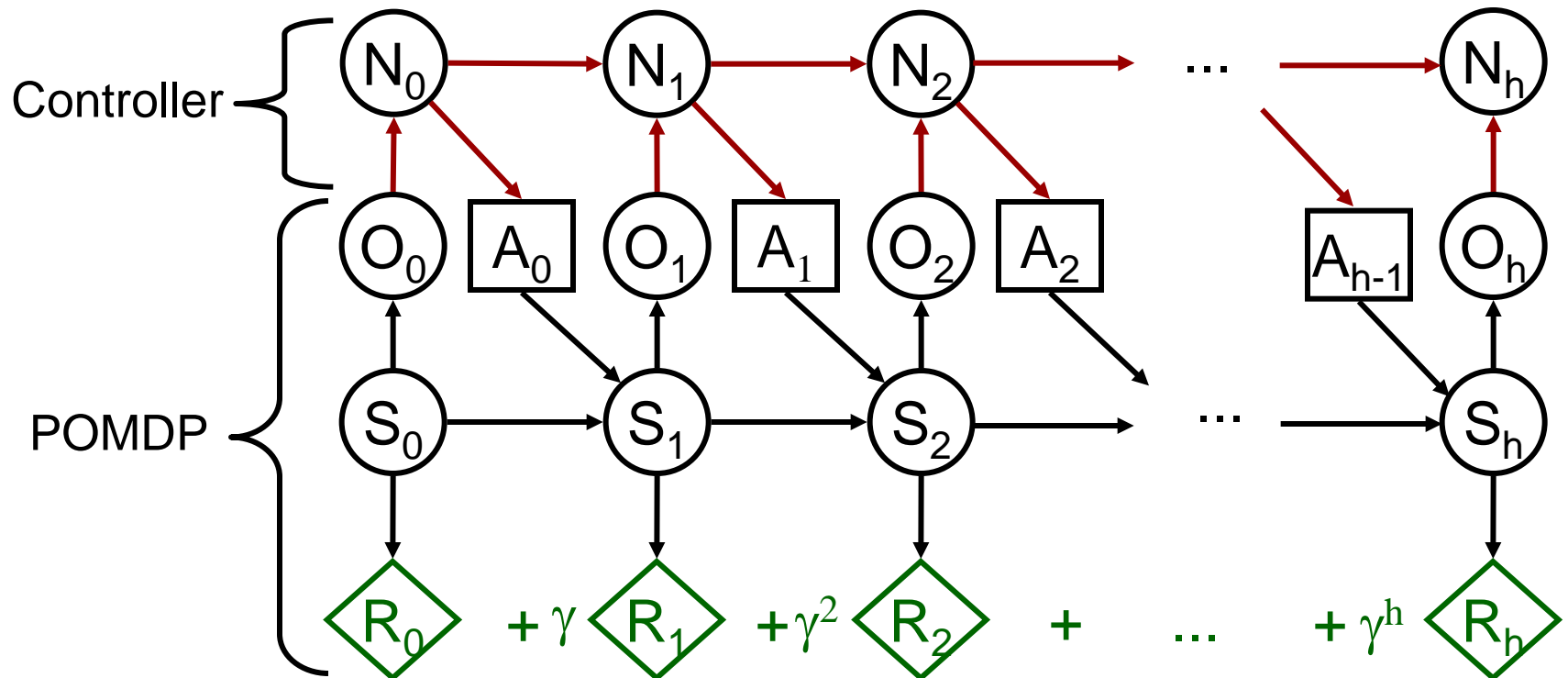
K^L nodes

Controllers as DBNs

- Alternative: planning as inference
 - Convert POMDP with hierarchical controller into a mixture of DBNs
 - Reduce policy optimization to a maximum likelihood estimation problem
 - Use EM
 - Scalable!

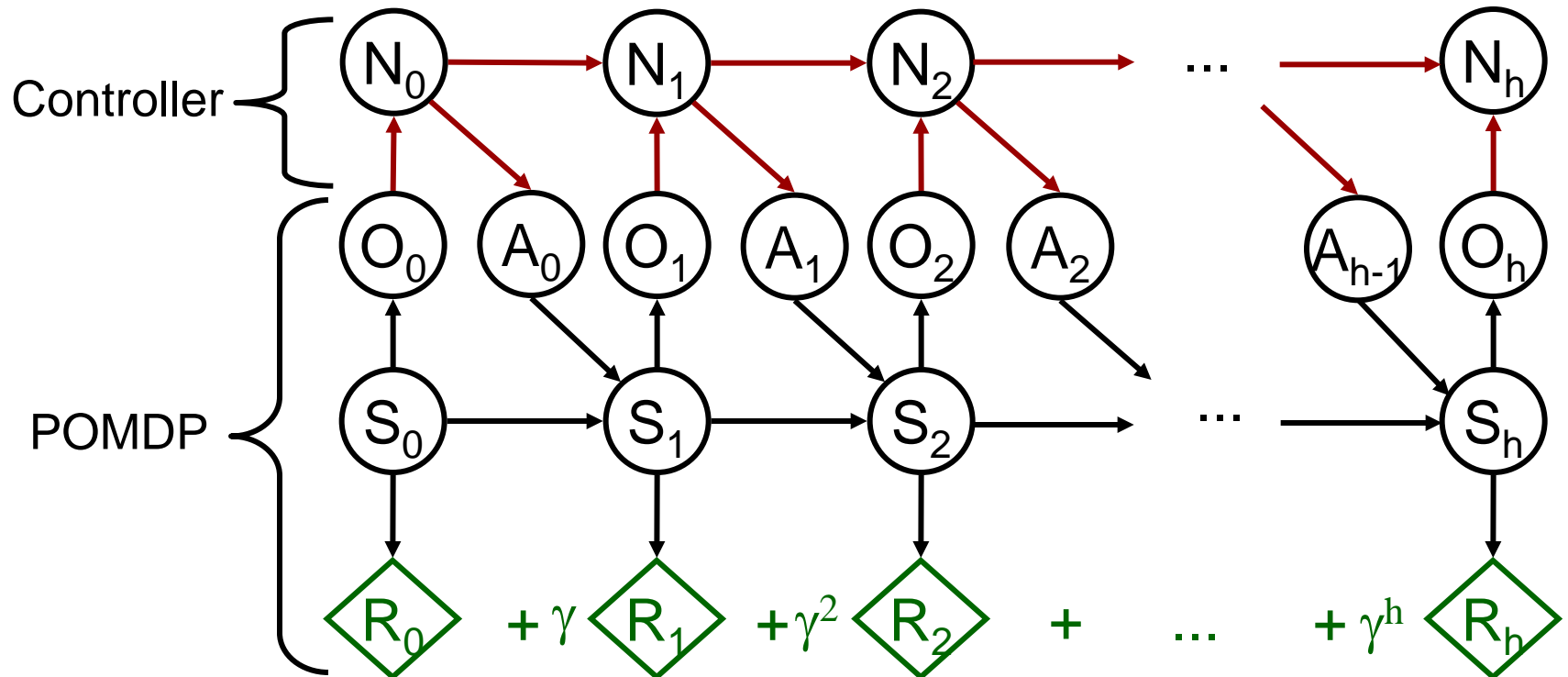
Graphical Model

- Meuleau et al., 1999



Graphical Model

- Toussaint et al., 2006:
 - Optimize $\Pr(A_t|N_t)$ & $\Pr(N_{t+1}|N_t, O_{t+1})$ by EM

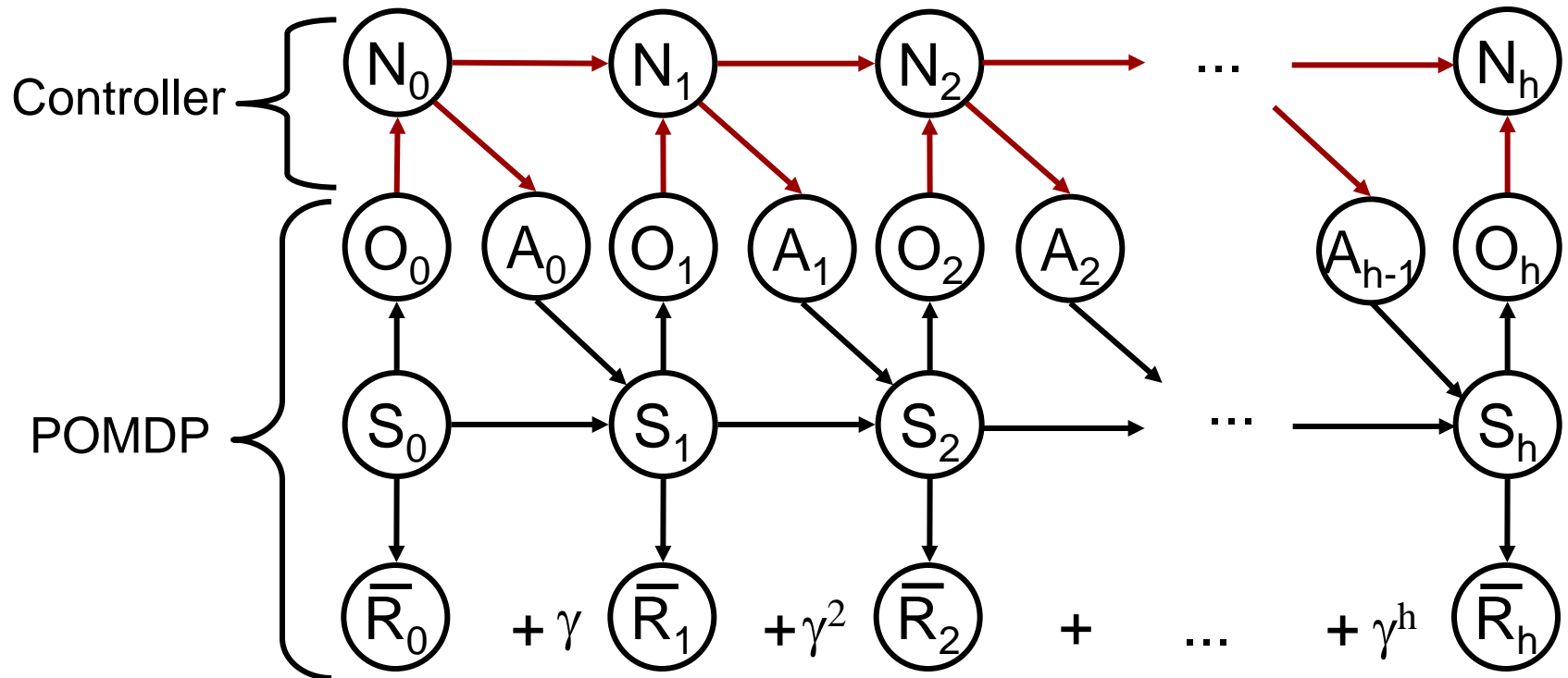


Graphical Model

- Normalize rewards

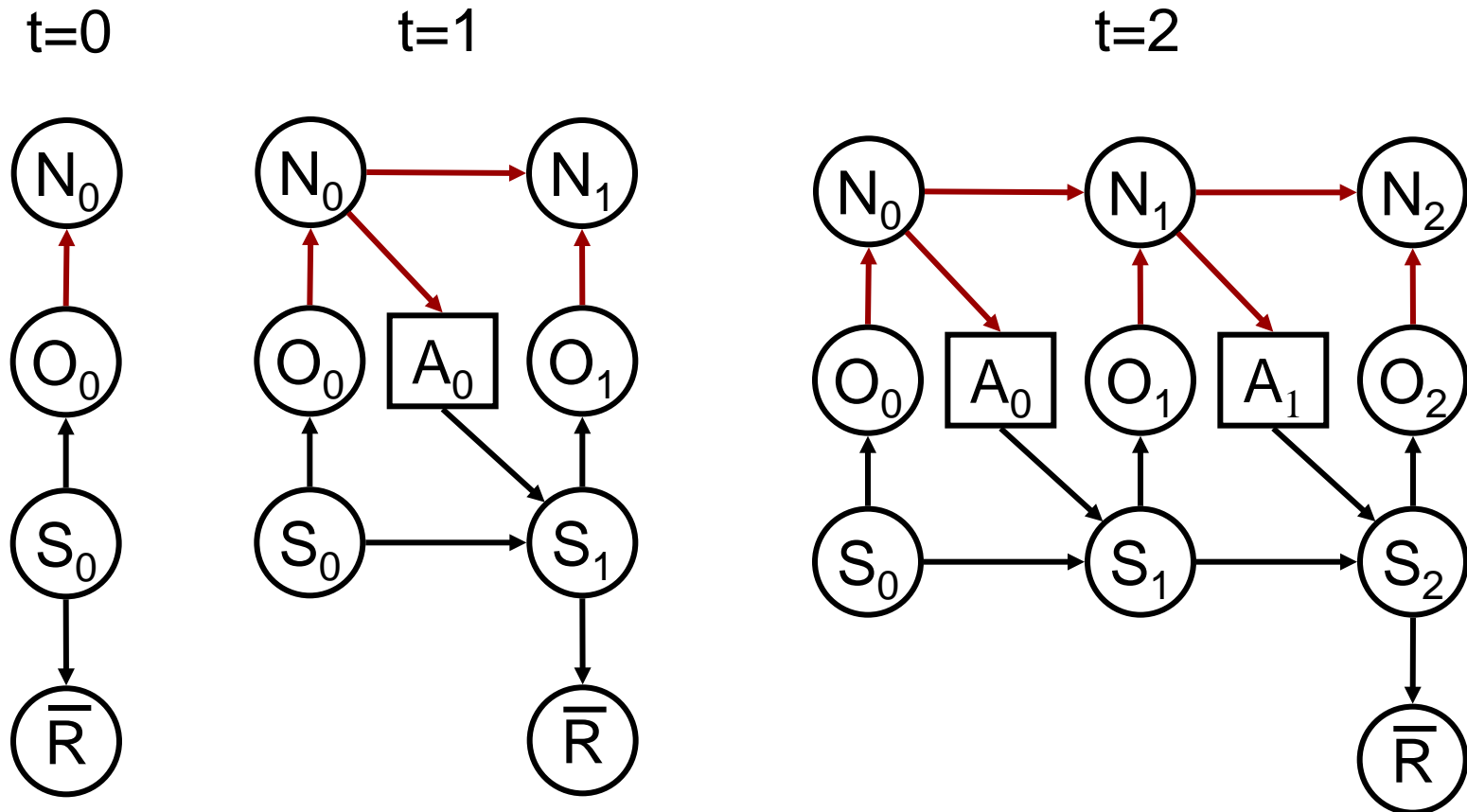
$$- \bar{R}_t \in \{0,1\}$$

$$\Pr(\bar{R}_t=1) = \frac{R(s_t) - r_{\min}}{r_{\max} - r_{\min}}$$



Graphical Model

- Mixture of DBNs: $\Pr(t) = k \gamma^t$

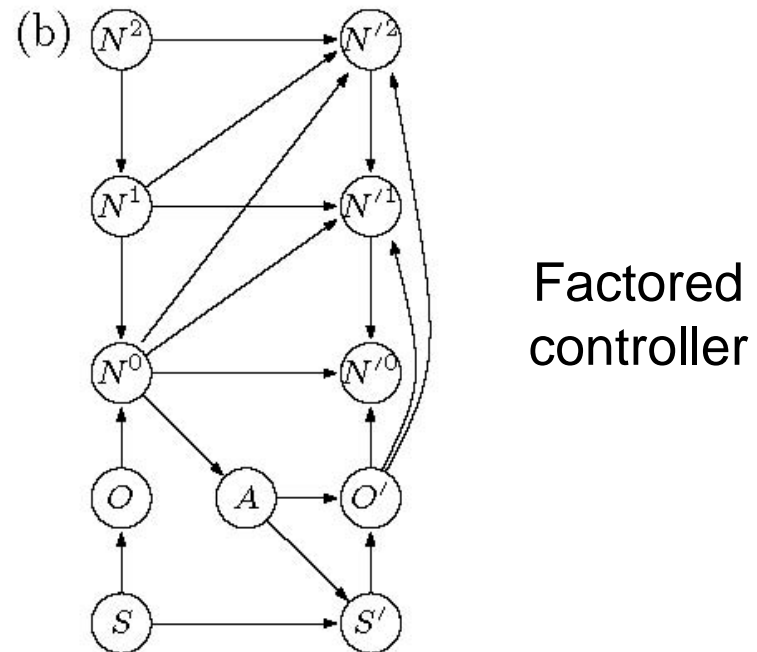
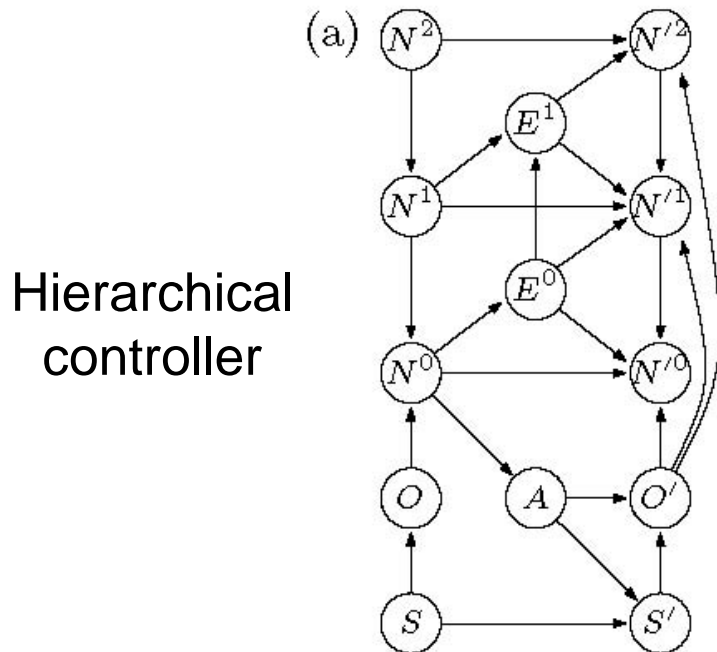


Controllers as DBNs

- Policy optimization
 - Maximize likelihood of $\bar{R}=1$:
 - $\operatorname{argmax}_{\Pr(a|n), \Pr(n'|n, o')} \Pr(\bar{R}=1)$
 - Expectation maximization
- **Advantage: any inference algorithm can be used**

Factored Controllers

- Toussaint, Charlin & Poupart, 2008
 - Hierarchical controllers \Leftrightarrow DBN mixture
 - More generally: factored controllers



Factored Controllers

- Hierarchy discovery & policy optimization
 - Maximize likelihood of $\bar{R}=1$
 - Expectation Maximization
- Same problem as non-convex quartic opt. prob.
 - EM much faster than optimization-based techniques
 - But EM gets stuck in local optima more easily

Results

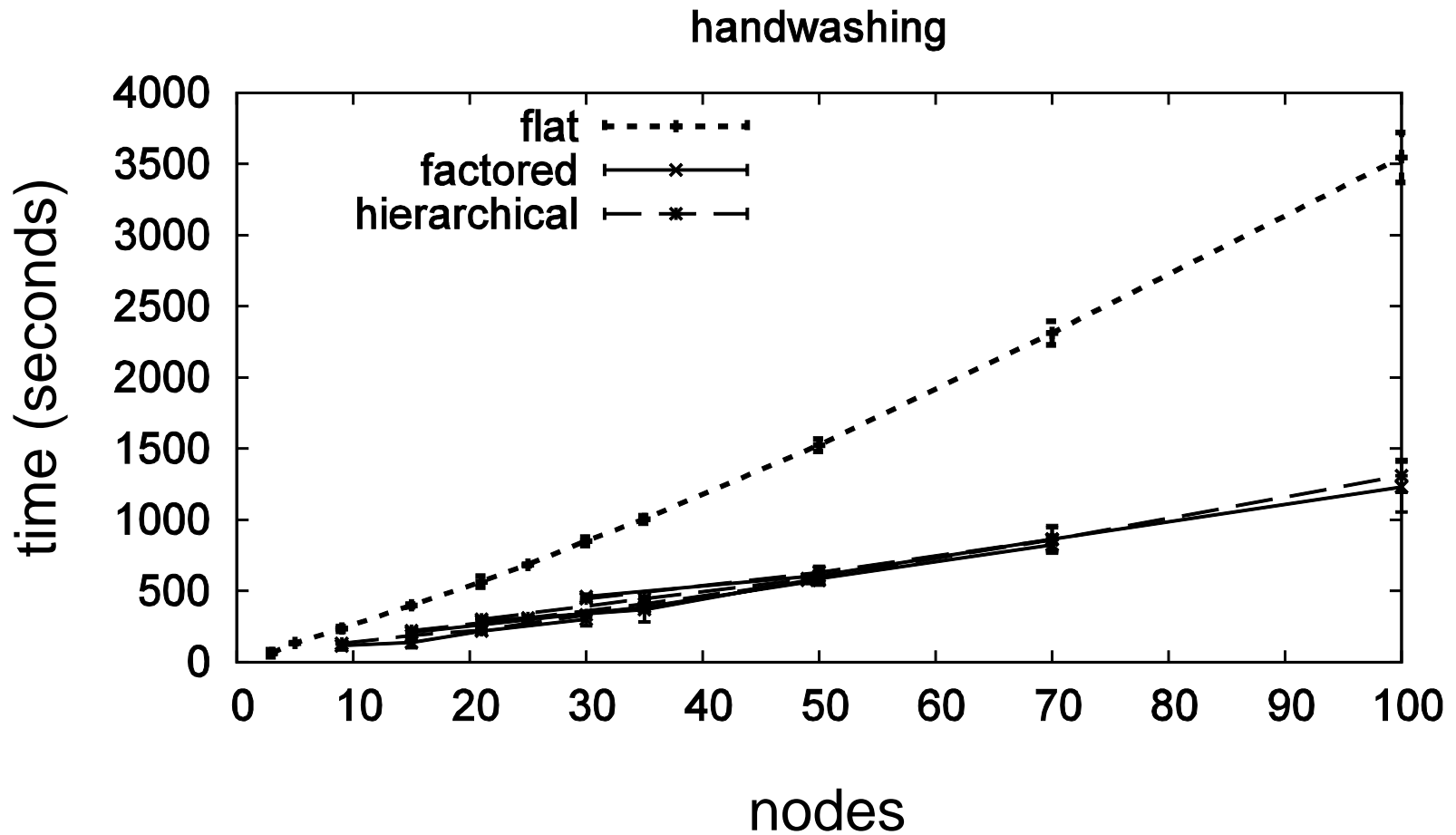
- Benchmark problems

Non-convex
optimization
(NEOS server)

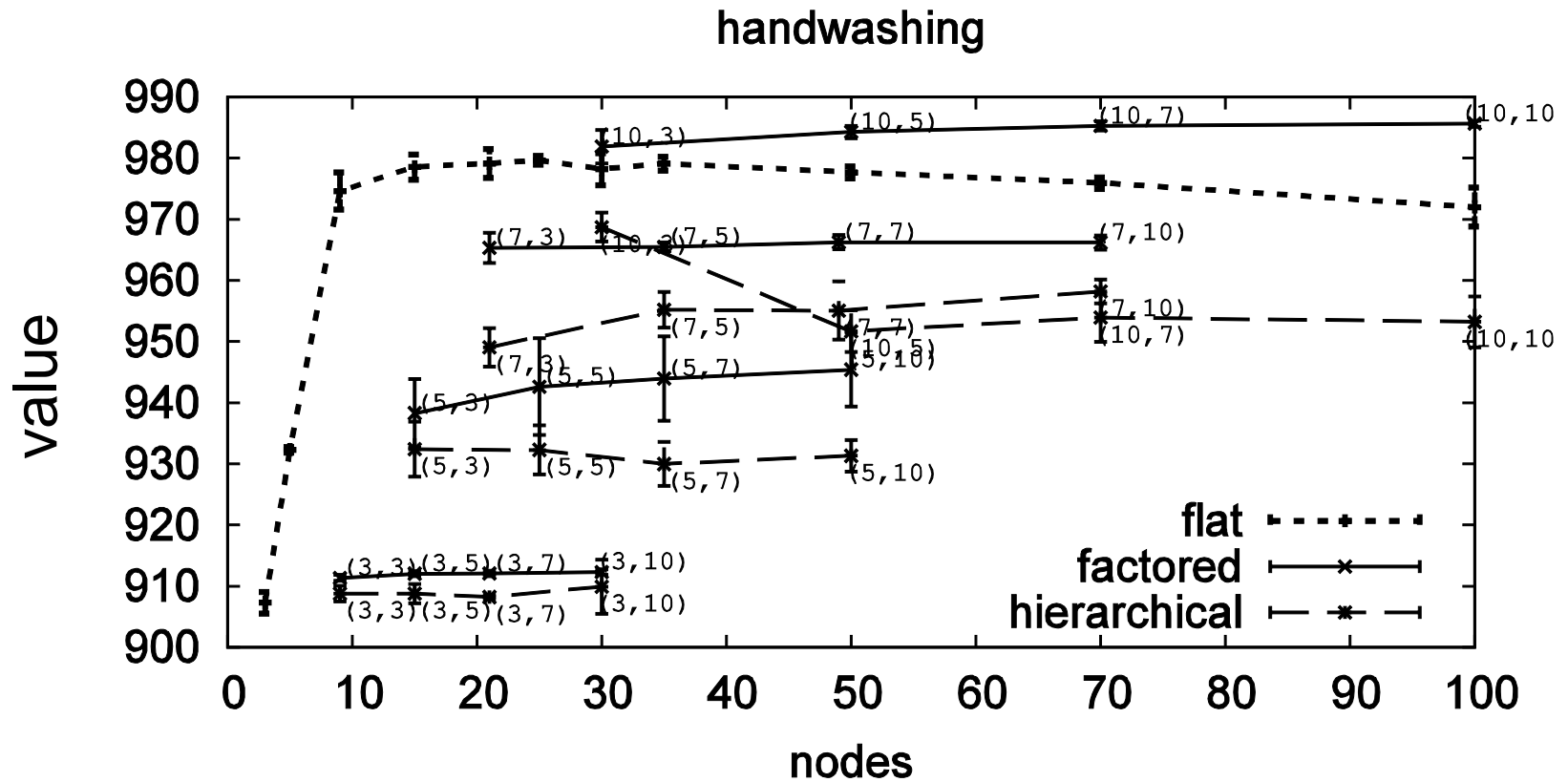
EM

Problem	$ \mathcal{S} , \mathcal{A} , \mathcal{O} $	V^*	HSVI2 V	Best results from [4]			ML approach (avg. over 10 runs)		
				nodes	t(s)	V	nodes	t(s)	V
paint	4, 4, 2	3.28	3.29 ± 0.04	(1,3)	< 1	3.29	(5,3)	0.96 ± 0.3	3.26 ± 0.004
shuttle	8, 3, 5	32.7	32.9 ± 0.8	(1,3)	2	31.87	(5,3)	2.81 ± 0.2	31.6 ± 0.5
4x4 maze	16, 4, 2	3.7	3.75 ± 0.1	(1,2)	30	3.73	(3,3)	2.8 ± 0.8	$3.72 \pm 8e-5$
chain-of-chains	10, 4, 1	157.1	157.1 ± 0	(3,3)	10	0.0	(10,3)	6.4 ± 0.2	151.6 ± 2.6
handwashing	84, 7, 12	≤ 1052	N/A			N/A	(10,5)	655 ± 2	984 ± 1
cheese-taxi	33, 7, 10	≤ 5.3	2.53 ± 0.3			N/A	(10,3)	311 ± 14	$-9 \pm 11 (2.25^*)$

Results (EM)



Results

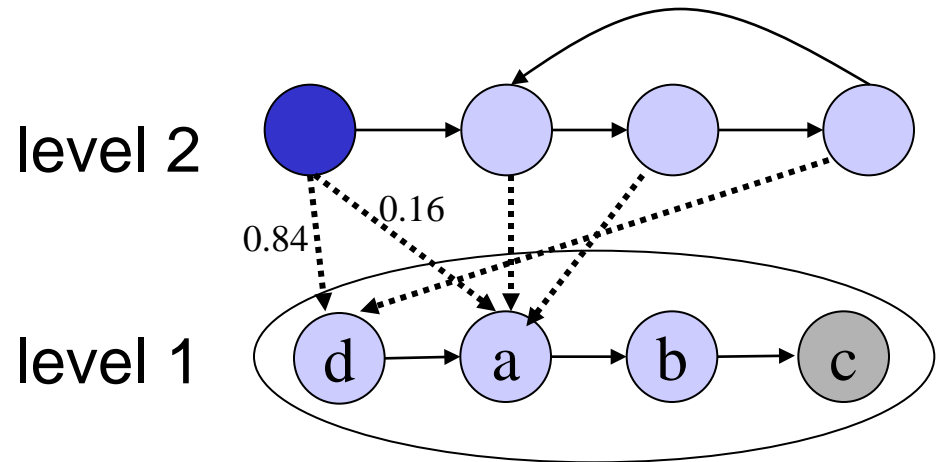
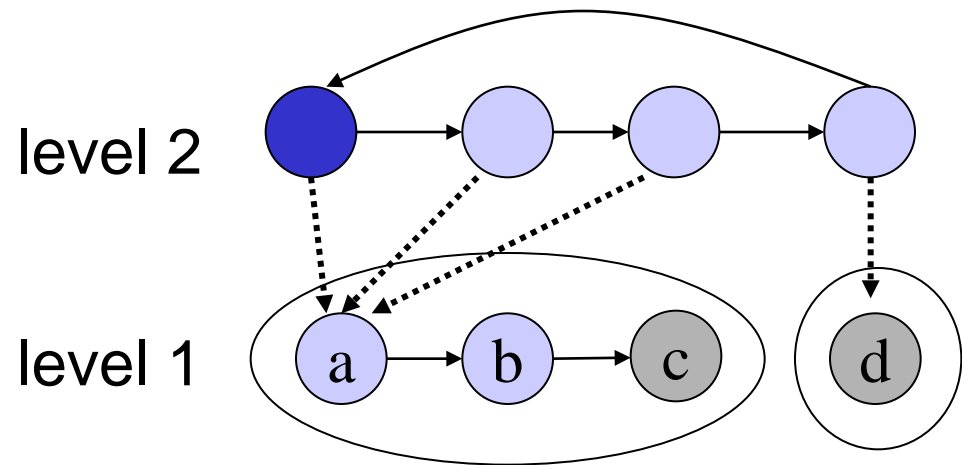


Non-intuitive hierarchy

- Simple unobservable POMDP:

- Reward when executing d after doing a-b-c three times

- Controller found:



Conclusion

- Simultaneous hierarchy discovery and planning
 - Direct optimization: **not scalable**
 - Planning as inference with EM: **scalable but subject to local optima**
 - Hierarchies discovered are **unintuitive**
- Future work:
 - Technique to escape local optima
 - Reinforcement learning (Vlassis & Toussaint, 2009)

Selected work relevant to workshop

- Feature extraction for Walker User Behaviour Recognition
 - UAI submission
- Bayesian RL
 - An Analytic Solution to Discrete Bayesian Reinforcement Learning, ICML 2006
- Value-directed information bottleneck
 - Value-directed Compression of POMDPs, NIPS 2002
- Policy explanation
 - Minimal Sufficient Explanations for Factored Markov Decision Processes, ICAPS 2009