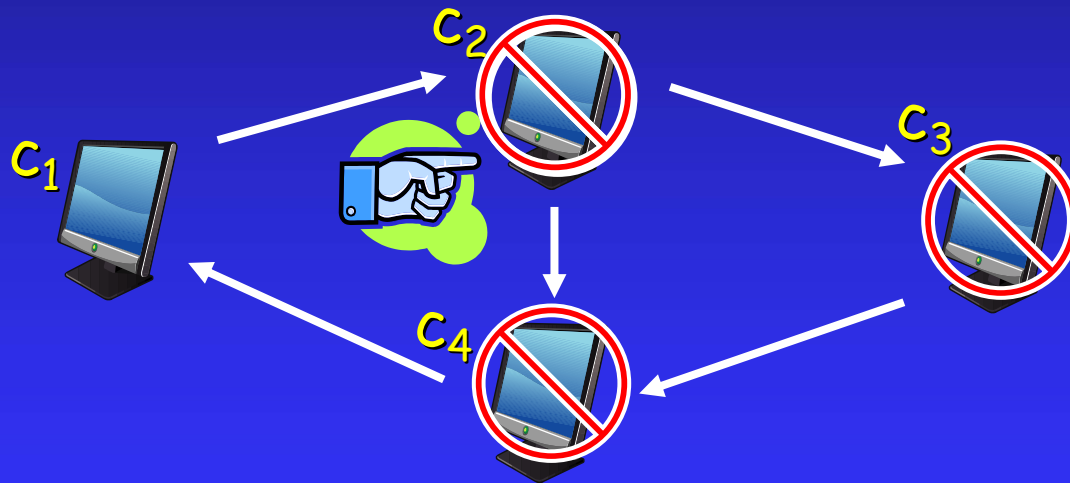


Approximate Solution Techniques for *Factored* FOMDPs

**Scott Sanner & Craig Boutilier
University of Toronto**

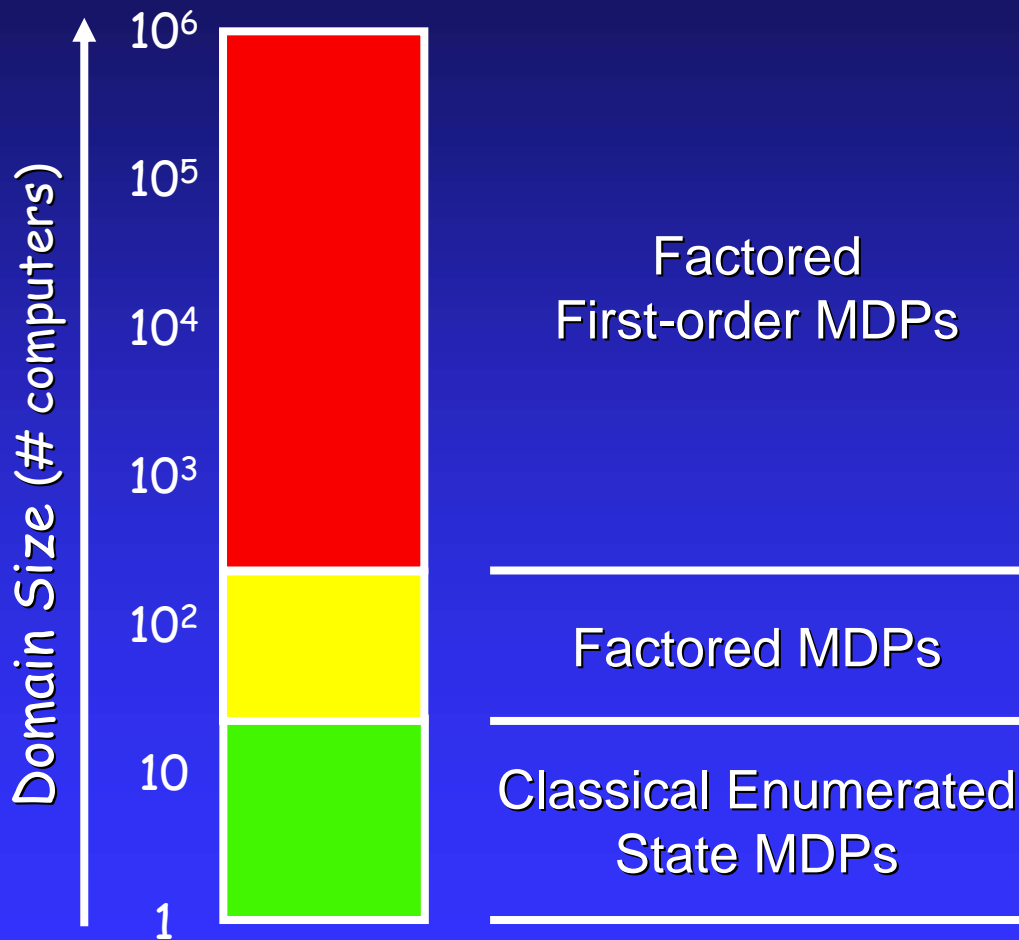
Motivating Example: SysAdmin MDP

- Have n computers $\mathcal{C} = \{c_1, \dots, c_n\}$ in a network
- **State:** each computer c_i is either “up” or “down”
- **Transition:** computer is “up” proportional to its state and # upstream connections that are “up”



- **Action:** manually reboot one computer
- **Reward:** +1 for every “up” computer

How to Solve SysAdmin?



- First-order MDPs cannot represent SysAdmin domain independently
⇒ need factored first-order MDP!
- State-of-the-art (factored) MDP solutions can scale only to ~140 computers

Background: Factored MDPs

Classical MDP Review

■ Representation:

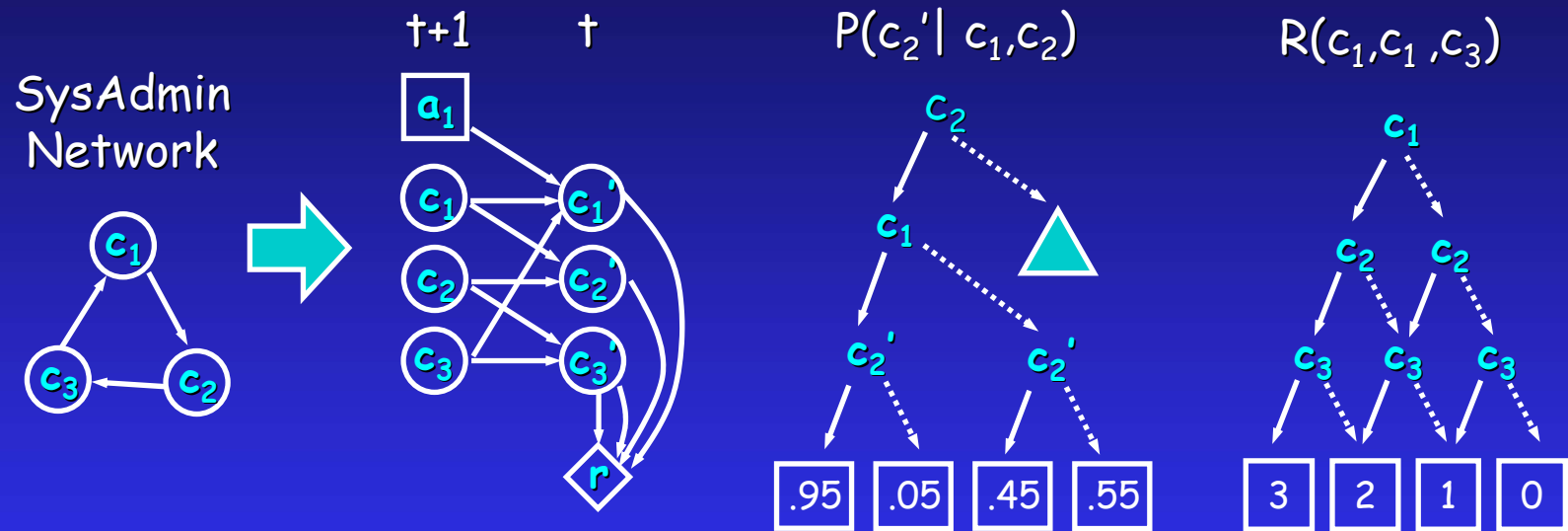
- ◆ $\langle S, A, T, R \rangle$
 - ◆ S : finite set of states
 - ◆ A : finite set of actions
 - ◆ $T: S \times A \times S \rightarrow [0,1]$ transition function
 - ◆ $R: S \times A \rightarrow \mathbb{R}$ reward function
- ◆ Policy $\pi: S \rightarrow A$
- ◆ Value function: $V(s) = E_{\pi} [\sum_{t=0}^{\infty} \gamma^t r^t | s]$

■ Some Solution Methods:

- ◆ Value Iteration
- ◆ Policy Iteration
- ◆ (Approximate) Linear Programming

Factored MDPs and Value Iteration

- Can use DBNs and decision diagrams:



- SPUDD (HSHB, 1999): Value Iteration w/ DDs**

$$V^{t+1}(\vec{x}) = \max_{a \in A} \left\{ R(\vec{x}, a) + \gamma \sum_{\vec{x}'} \left[\prod_{i=1}^n P(x'_i | Par(x'_i), a) V^t(\vec{x}') \right] \right\}$$

Backup Operator

- Define backup operator $B^a[.]$:

$$B^a[V(\vec{x})] = \gamma \sum_{\vec{x}'} \left[\prod_{i=1}^n P(x'_i | Par(x'_i), a) V(\vec{x}') \right]$$

- E.g., rewrite value iteration:

$$V^{t+1}(\vec{x}) = \max_{a \in A} \{ R(\vec{x}, a) + B^a[V^t(\vec{x})] \}$$

- $B^a[.]$ is a *linear* operator:

$$B^a[V_1(\vec{x}) + V_2(\vec{x})] = B^a[V_1(\vec{x})] + B^a[V_2(\vec{x})]$$

Approximate LP for Factored MDPs

- Linear-value function approximation:

$$V(\vec{x}) = \sum_{j=1}^k w_j b_j(\vec{x})$$

- Solve for weights using LP (GKP,01; SP,01):

Variables: w_1, \dots, w_k

Minimize: $\sum_{\vec{x}} V(\vec{x})$

Subject to: $0 \geq R(\vec{x}, a) + B^a[V(\vec{x})] - V(\vec{x}) ; \forall a, \vec{x}$



Exploit
linearity!


Constraint Generation

- Simple, general constraint format:

$$\begin{aligned} 0 &\geq \{F_1(\vec{x}_1) + \dots + F_m(\vec{x}_m)\} ; \forall a, \vec{x} \\ &\geq \max_{\vec{x}} \{F_1(\vec{x}_1) + \dots + F_m(\vec{x}_m)\} ; \forall a \end{aligned}$$

- Efficiently find max in cost network
 - ◆ i.e., variable elimination
- Iteratively solve LP with constraint generation
 - 1) Start with $w_i = 0$
 - 2) Find max violated constraint for each a
 - 3) If violations: add to LP and re-solve, goto (2)

Why Factored FOMDPs?

- Back to SysAdmin...  $c_1 \longrightarrow \dots c_i \dots \longrightarrow c_n$
- As a factored MDP
 - ◆ Instantiate MDP for $n = \#computers$
 - ◆ Specific solutions aside...
 - ◆ MDP repr. is $\Omega(n) \Rightarrow$ any solution is $\Omega(n)$
- As a factored first-order MDP
 - ◆ “Lift” MDP specification
 - ◆ “Lift” (approximate) solution
 - ◆ Solution $O(\text{sub-linear}(n))$ in struct. cases!

Contribution: Factored FOMDP Representation

First-order MDPs (FOMDPs)

- $\langle S, A, T, R \rangle$ for FOMDPs defined in terms of cases
 - ◆ E.g., possible *reward* in SysAdmin ...

$$r\text{Case}(s) = \begin{array}{|c|c|} \hline \exists c \text{ Run}(c,s) & 1 \\ \hline \neg \exists c \text{ Run}(c,s) & 0 \\ \hline \end{array}$$

- **Operators:** Define unary, binary case operations
 - ◆ E.g., can take “cross-sum” \oplus (or \otimes , \ominus) of cases...

$$\begin{array}{|c|c|} \hline \phi & 10 \\ \hline \neg\phi & 20 \\ \hline \end{array} \oplus \begin{array}{|c|c|} \hline \phi & 3 \\ \hline \neg\phi & 4 \\ \hline \end{array} = \begin{array}{|c|c|} \hline \phi \wedge \phi & 13 \\ \hline \phi \wedge \neg\phi & 14 \\ \hline \neg\phi \wedge \phi & 23 \\ \hline \neg\phi \wedge \neg\phi & 24 \\ \hline \end{array}$$

Factored FOMDPs: Additive Reward

- SysAdmin reward scales with domain size:

$$r_{Case}(s) = \begin{array}{|c|c|} \hline Run(c_1, s) & 1 \\ \hline \neg Run(c_1, s) & 0 \\ \hline \end{array} \oplus \dots \oplus \begin{array}{|c|c|} \hline Run(c_n, s) & 1 \\ \hline \neg Run(c_n, s) & 0 \\ \hline \end{array}$$

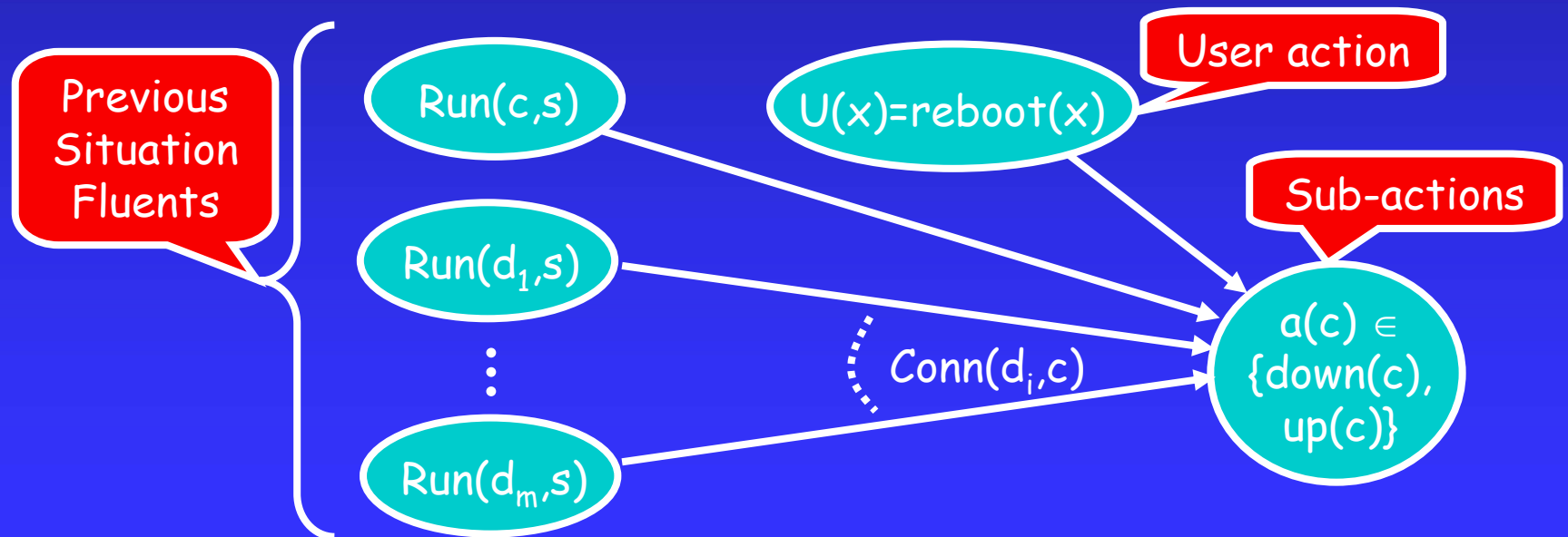
- Beyond expressive power of current FOMDP
- Need language extension for Σ aggregator:

$$r_{Case}(s) = \sum_{c \in \mathcal{C}} \begin{array}{|c|c|} \hline Run(c, s) & 1 \\ \hline \neg Run(c, s) & 0 \\ \hline \end{array}$$

- ◆ Semantics is just the expanded \oplus

Factored FOMDPs: Sub-actions

- Usual FOMDP Frame Assumption:
 - ◆ Def: *Anything not changed by action remains same*
 - ◆ Not true for factored FOMDPs (e.g., SysAdmin)
- Remedy: Make “local” frame assumption (i.e., DBN)
 - ◆ Specify first-order DBN for deterministic *sub-actions*



Factored FOMDP: Transition

- SysAdmin *sub-action* transition probability:

$$P(a(c) = up(c) | U(x) = reboot(x) \wedge x = c, s) = \boxed{\top : 1}$$

$$P(a(c) = up(c) | U(x) = reboot(x) \wedge x \neq c, s) =$$

$$\begin{array}{|l|} \hline Run(c, s) : 0.95 \\ \hline \neg Run(c, s) : 0.05 \\ \hline \end{array} \otimes \frac{1 + \sum_d \left(\begin{array}{|l|} \hline Conn(d, c) \wedge Run(d, s) : 1 \\ \hline \neg Conn(d, c) \vee \neg Run(d, s) : 0 \\ \hline \end{array} \right)}{1 + \sum_d \left(\begin{array}{|l|} \hline Conn(d, c) : 1 \\ \hline \neg Conn(d, c) : 0 \\ \hline \end{array} \right)}$$

- Need \prod aggregator for *joint action* distribution:

$$P(a = a(c_1) \circ \dots \circ a(c_n) | U(\vec{x})) = \prod_{c \in C} P(a(c) | U(\vec{x}), s)$$

SitCalc Extensions and Irrelevance

- Upgrade SSAs to handle factored actions:

$$Run(c, do(a, s)) \equiv a \sqsupseteq up(c) \vee Run(c, s) \wedge \neg a \sqsupseteq down(c)$$

- Example regression in SysAdmin:

$$Regr[Run(c_i, s) ; a = \dots \circ up(c_i) \circ \dots] \equiv \top$$

- **Observation: some sub-actions irrelevant in $Regr[.]$**

- ◆ Def: *all sub-action outcomes \mathbf{B} irrelevant to $\phi(s)$ iff*

$$\forall b \in \mathbf{B}. Regr[\phi(s), b] \equiv \phi(s)$$

- ◆ Need independence of sub-actions (not restrictive)
- ◆ Allows irr. sub-actions to be dropped from $Regr[.]$

Recap: Factored FOMDPs

■ Reward:

- ◆ Can be expressed with Σ aggregator

■ Transition Distribution:

- ◆ Decompose joint action into sub-actions
 - ◆ SSAs expressed in terms of sub-actions
 - ◆ Local distribution for each sub-action
- ◆ Joint distribution uses Π aggregator

■ So far, just syntax...

- ◆ Real problem is exploiting syntactic structure
- ◆ Key ideas: *irrelevance* and *operator linearity*

Contribution: Factored FOMDP Solutions

Factored FOMDP Backup Operator

- Define backup operator:

$$B^{U(\vec{x})}[vCase(s)] = \gamma \bigoplus_{a \in A} \left[P(a|U(\vec{x})) \otimes Repr[vCase(s), a] \right]$$

- Ex: SysAdmin $vCase^0(s) = rCase(s) = \sum_c \begin{array}{|l} Run(c, s) : 1 \\ \hline \neg Run(c, s) : 0 \end{array}$

$$B^{reboot(x)}[vCase^0(s)] = \gamma \sum_c \left[\bigoplus_{a_1 \in A(c_1), \dots, a_n \in A(c_n)} \left(\prod_{i=1}^n P(a_i|U) \right) \otimes \begin{array}{|l} Repr[Run(c, s), a_1 \circ \dots \circ a_n] : 1 \\ \hline Repr[\neg Run(c, s), a_1 \circ \dots \circ a_n] : 0 \end{array} \right]$$

- Result after simplification:

$$B^{reboot(x)}[vCase(s)] = \gamma \sum_c P(up(c)|reboot(x))$$

SDP for Factored FOMDPs

- Complete symbolic dynamic programming step:

Recall Classical MDP: $V^1(s) = R(s) + \gamma \max_a B^a[V^0(s)]$

Upgrade to Factored FOMDP:

$$\begin{aligned} vCase^1(s) &= rCase(s) \oplus \gamma \max \exists x. B^{reboot(x)}[vCase^0(s)] \\ &= \sum_c \left(\frac{Up(c, s) : 1}{\neg Up(c, s) : 0} \right) \oplus \gamma \max \exists x. \sum_c P(up(c)|reboot(x)) \end{aligned}$$

- Caveat: What to do with $\max \exists x$?
- Workaround is to derive a policy (see GKP) to axiomatize optimal x
 - ◆ Need to add policy axioms on every SDP step
 - ◆ Exact representation blows up ☹

Linear-value Approximation

- Approximate value w/ basis fn *classes*:

$$v_{\text{Case}}(s) = w_1 \cdot \sum_c \begin{array}{|c|c|} \hline \phi(c) & 1 \\ \hline \neg\phi(c) & 0 \\ \hline \end{array} \oplus w_2 \cdot \sum_c \begin{array}{|c|c|} \hline \varphi(c) & 1 \\ \hline \neg\varphi(c) & 0 \\ \hline \end{array}$$

- Reduces solution to finding good weights
 - ◆ Weight projection \Rightarrow policy axioms don't accumulate
 - ◆ Only need to do consistency checking!
- Where do basis functions come from?
 - ◆ Use variant of ideas proposed in (GreThi, UAI-04)
- How to find weights?
 - ◆ We provide factored FOALP algorithm

Approximate Linear Programming

■ Factored First-order ALP:

Variables: $w_i ; \forall i \leq k$

Minimize: $\sum_s vCase(s)$

Subject to: $0 \geq rCase(s) \oplus B^{U(\vec{x})}[vCase(s)]$
 $\ominus vCase(s) ; \forall U(\vec{x}), s$

■ Constraint generation solution:

- ◆ Constraints are always of the general form:

$$0 \geq \max_s \exists \vec{x} \left[\sum_c case_1(c, \vec{x}, s) \oplus \dots \oplus \sum_c case_p(c, \vec{x}, s) \right]$$

- ◆ Make a *domain size* assumption
- ◆ Extend var. elimination techniques from FOPI (Poole, 2003; Braz, Amir, & Roth, 2005, 2006) to relation elimination

[Ch. 6] Some FOPI Techniques for Cases

■ Known Eliminations (Poole, 2003; Braz, Amir, & Roth, 2005, 2006)

◆ (Partial) Inversion Elimination:

$$\max \sum_c [\text{case}_1(c) \oplus \text{case}_2(c)] \text{ where } \text{Rel}(\text{case}_1) \cap \text{Rel}(\text{case}_2) = \emptyset$$

◆ Counting Elimination:

$$\max \sum_c \sum_{d \neq c} \text{case}(c,d)$$

■ New Eliminations (SanBout, 2007)

◆ Existential elimination (**max** elim. only):

$$\max \exists x \sum_c [\text{case}_1(c,x) \oplus \dots \oplus \text{case}_p(c,x)]$$

◆ Linear elimination:

$$\max \sum_{c_i} \text{case}(c_i, c_{i+1})$$

Existential Elimination

- Need to compute: $\max \exists x \sum_c [\text{case}(c, x)]$

◆ where $\text{case}(c, x) =$

$x = c$	10
$x \neq c \wedge \dots$	9
$x \neq c \wedge \dots$	0

- Introduce:

◆ $\sum_c e\text{Case}(c, s) = \sum_c \begin{array}{|l} b(c) \supset b(\text{next}(c)) : 0 \\ b(c) \wedge \neg b(\text{next}(c)) : -\infty \end{array}$

◆ $b(c_1), b(c_2), b(c_3), \dots, b(c_{n-1}), b(c_n)$

◆ $\perp \quad \perp \quad \top \quad \top \quad \top \quad \top$

- Replace: $(x = c) \equiv \neg b(c) \wedge b(\text{next}(c))$

- Final constraint:

$$0 \geq \max_s \sum_c [\text{case}_1(c, s) \oplus \dots \oplus \text{case}_p(c, s) \oplus e\text{Case}(c, s)]$$

Linear Elimination $c_1 \longrightarrow \dots \longrightarrow c_i \dots \longrightarrow c_n$

- Need to compute: $r(n) = \max_{c_2 \dots c_n} \sum_{i=1 \dots n} \text{case}(c_i, c_{i+1})$

- where $\text{case}(c_i, c_{i+1}, s) =$

c_i	c_{i+1}	
\perp	\perp	1
\perp	\top	-5
\top	\perp	-5
\top	\top	0

- $r(2) = \max c_2$

c_1	c_2	
\perp	\perp	1
\perp	\top	-5
\top	\perp	-5
\top	\top	0

 $+$

c_2	c_3	
\perp	\perp	1
\perp	\top	-5
\top	\perp	-5
\top	\top	0

 $=$

c_1	c_3	
\perp	\perp	2
\perp	\top	-4
\top	\perp	-4
\top	\top	0

- $r(4) = \max c_2, c_3, c_4$

c_1	c_3	
\perp	\perp	2
\perp	\top	-4
\top	\perp	-4
\top	\top	0

 $+$

c_3	c_5	
\perp	\perp	2
\perp	\top	-4
\top	\perp	-4
\top	\top	0

 $=$

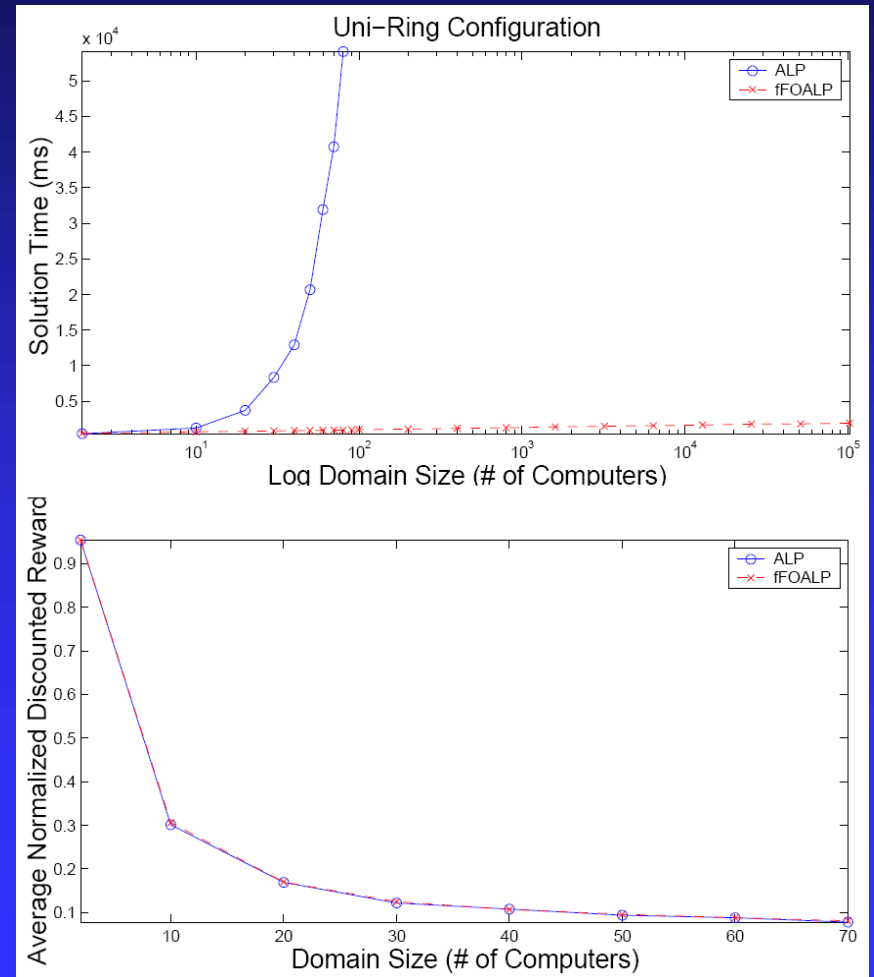
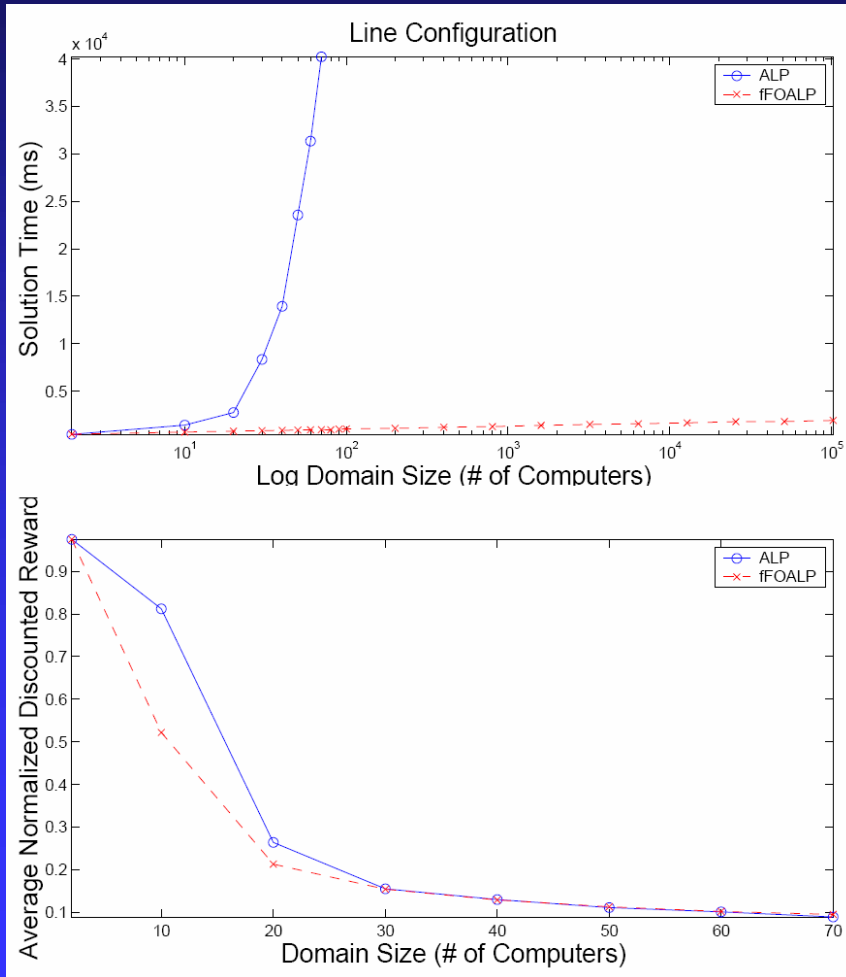
c_1	c_5	
\perp	\perp	4
\perp	\top	-2
\top	\perp	-2
\top	\top	0

- Computation of $r(n)$ takes $O(\log(n))$!

Some Results for fFOALP vs. ALP

■ Line Configuration $c_1 \rightarrow c_2 \rightarrow c_3$

■ Unidirectional Ring



Conclusions

- Introduced factored FOMDP, “lifted” fFOALP solution
 - ◆ Exploited:
 - ◆ backup linearity
 - ◆ first-order irrelevance
 - ◆ FOPI for constraint generation
 - ◆ Solutions $O(\text{sub-linear}(n))$ in structured cases!
- What about *general* factored FOMDPs?
 - ◆ Key to efficient solution is FOPI techniques
 - ◆ FOPI specific to constraint structure induced by
 - ◆ basis function structure ✕
 - ◆ FO-DBN dynamics
 - ◆ Need to catalog “efficient” structures
 - ◆ Identify new structures exploitable by FOPI!