

# Heuristic Search for Multi-Objective Probabilistic Planning

Dillon Chen,<sup>1</sup> Felipe Trevizan,<sup>1</sup> Sylvie Thiébaux<sup>1,2</sup>

<sup>1</sup>School of Computing, The Australian National University

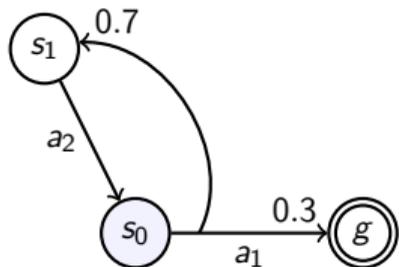
<sup>2</sup>LAAS-CNRS, ANITI, Université de Toulouse

{Dillon.Chen, Felipe.Trevizan, Sylvie.Thieboux}@anu.edu.au

AAAI 2023



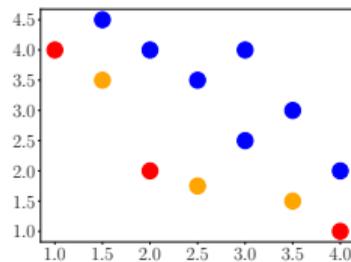
- ▶ stochastic shortest path problems (SSPs) are the defacto model for planning under uncertainty
  - ▶ compute a policy which maps states to actions
  - ▶ want optimal policy minimising the expected cost to reach the goal from an initial state



An SSP with a stochastic action  $a_1$ .

- ▶ heuristic search powerful for implicitly represented SSPs with large state spaces

- ▶ multi-objective stochastic shortest path problems (MOSSPs) generalise SSPs by exhibiting multiple objectives
  - ▶ aim to compute a set of non-dominated vector-valued policies with tradeoffs between objectives



orange+red=non-dominated vectors

- ▶ no existing heuristic search algorithms for MOSSPs
- ▶ algorithms that exist for MOMDPs do not carry over to MOSSPs: MOSSPs require additional assumptions

# Contributions

defining MOSSPs: a generalisation of MOMDPs

assumptions for convergence of MOVI for MOSSPs

heuristics for MOSSPs

heuristic search algorithms for solving MOSSPs: iMOLAO\*, MOLRTDP

experimental evaluation

# Contributions

defining MOSSPs: a generalisation of MOMDPs

assumptions for convergence of MOVI for MOSSPs

heuristics for MOSSPs

heuristic search algorithms for solving MOSSPs: iMOLAO\*, MOLRTDP

experimental evaluation

## Definition (MOSSPs)

An MOSSP is a tuple  $P = (S, A, s_0, G)$

- ▶  $S$  is a set of states
- ▶  $A$  is a set of actions  $a$  which
  - ▶ maps certain states to a probability distribution of states, and
  - ▶ has an associated cost vector  $\vec{C}(a) = [c_1, \dots, c_n] \in \mathbb{R}_{\geq 0}^n$
- ▶  $s_0 \in S$  is an initial state,  $G \subseteq S$  is a set of goal states

Informally,

*A solution is a compact representation of a set of non-dominated vector-valued policies.*

$\vec{u}$  dominates  $\vec{v}$  denoted  $\vec{u} \preceq \vec{v}$  iff  $\vec{u}[i] \leq \vec{v}[i], i = 1, \dots, n$

\*Full definition: too many technical details and sub-definitions to fit into the talk

# MO Bellman Backup

- ▶ fundamental equations for solving (MO)SSPs and (MO)MDPs
- ▶ states assigned  $\mathbf{Q}$  and  $\mathbf{V}$  values; each a finite set of vectors  $\in \mathcal{P}_{<\omega}(\mathbb{R}_{\geq 0}^n)$
- ▶  $\mathbf{Q} : S \times A \rightarrow \mathcal{P}_{<\omega}(\mathbb{R}_{\geq 0}^n) \simeq$  MO expected cost to goal when executing  $a$  at  $s$
- ▶  $\mathbf{V} : S \rightarrow \mathcal{P}_{<\omega}(\mathbb{R}_{\geq 0}^n) \simeq$  MO expected cost to goal from  $s \rightarrow$  *minimum* of  $\mathbf{Q}$  values

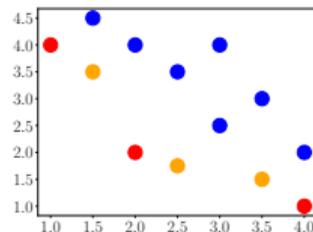
▶ For  $s \in S \setminus G$ :

$$\mathbf{V}^{t+1}(s) = \text{CCS} \left( \bigcup_{a \in A} \mathbf{Q}^{t+1}(s, a) \right)$$

$$\mathbf{Q}^{t+1}(s, a) = \{ \vec{C}(a) \} \oplus \left( \bigoplus_{s' \in S} P(s'|s, a) \mathbf{V}^t(s') \right)$$

▶ For  $g \in G$ :

$$\mathbf{V}^{t+1}(g) = \{ \vec{0} \}$$



red: CCS;  
orange+red: PCS

- ▶ CCS is the convex hull of a Pareto coverage set  $\simeq$  MO generalisation of *min*
- ▶  $\mathbf{V} \oplus \mathbf{U} = \{ \vec{u} + \vec{v} \mid \vec{u} \in \mathbf{U}, \vec{v} \in \mathbf{V} \}$  is setwise sum of vectors
- ▶  $\bigoplus$  generalises  $\oplus$  to several sets
- ▶  $\lim_{t \rightarrow \infty} \mathbf{V}^t = \mathbf{V}^*$  [White, 1982; Barrett and Narayanan, 2008]

# MO Value Iteration

---

## Algorithm: MOVI

---

**Data:** MOSSP problem  $P = (S, A, s_0, G)$ ,  
initial values  $\mathbf{V}(s)$  for each state  $s$   
(default to  $\mathbf{V}(s) = \{\vec{0}\}$ ), and  
consistency threshold  $\epsilon$ .

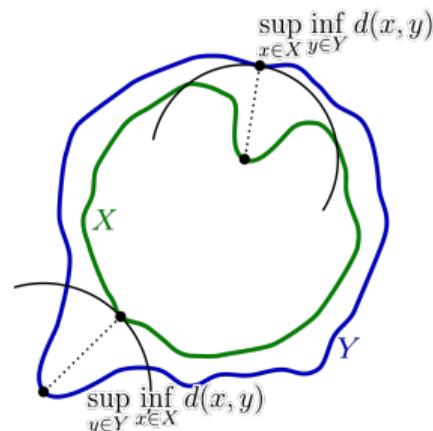
```
1 while  $\max_{s \in S} \text{res}(s) < \epsilon$  do
2   for  $s \in S$  do
3     if  $s \in G$  then  $\mathbf{V}_{\text{new}}(s) \leftarrow \{\vec{0}\}$ ;
4     else  $\mathbf{V}_{\text{new}}(s) \leftarrow$ 
        MOBellmanBackup( $s$ );
5      $\text{res}(s) \leftarrow D(\mathbf{V}, \mathbf{V}_{\text{new}})$ 
6    $\mathbf{V} \leftarrow \mathbf{V}_{\text{new}}$ 
7 return  $\mathbf{V}$ 
```

---

same as (single-objective) Value Iteration  
except

1. MO backups instead of SO backups
2. MO residual defined using Hausdorff metric  
(for finite sets of points):

$$D(\mathbf{U}, \mathbf{V}) = \max \left\{ \begin{array}{l} \max_{\vec{u} \in \mathbf{U}} \min_{\vec{v} \in \mathbf{V}} d(\vec{u}, \vec{v}), \\ \max_{\vec{v} \in \mathbf{V}} \min_{\vec{u} \in \mathbf{U}} d(\vec{u}, \vec{v}) \end{array} \right\}$$



# Contributions

defining MOSSPs: a generalisation of MOMDPs

assumptions for convergence of MOVI for MOSSPs

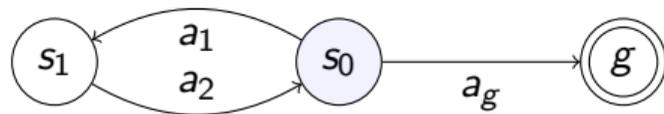
heuristics for MOSSPs

heuristic search algorithms for solving MOSSPs: iMOLAO\*, MOLRTDP

experimental evaluation

# Assumptions for convergence of MOVI for MOSSPs

Consider the MOSSP



MOSSP with action costs  $\vec{C}(a_1) = [1, 0]$ ,  $\vec{C}(a_2) = [1, 0]$ ,  $\vec{C}(a_g) = [0, 1]$ .

- ▶  $\pi_1 : S \rightarrow A, s_0 \mapsto a_g$  is a proper policy
- ▶  $\pi_2 : S \rightarrow A, s_0 \mapsto a_1, s_1 \mapsto a_2$  is an improper policy
- ▶ MOVI does not detect this;  $\mathbf{V}(s_0) = \{[\infty, 0], [0, 1]\}$  in the limit
- ▶ VI able to detect and prune SSP policies which have  $\infty$  cost

# Assumptions for convergence of MOVI for MOSSPs

## Strong improper policy assumption

- ▶ reachability assumption
  - ▶  $\forall s \in S, \exists$  proper policy at  $s$
- ▶ all improper policies cost  $\infty$
- ▶ MOVI + strong assumption is sound and complete
- ▶ implies cycle costs have nonzero components (not easy to guarantee)
- ▶ one method: all action costs have nonzero components
  - ▶ cannot model independent costs
  - ▶ e.g. navigation domain: load, unload consumes 0 fuel, 1 time

---

## Weak improper policy assumption

- ▶ reachability assumption
- ▶ exists a bound  $\vec{b}$  s.t.  $\forall s \in S$ 
  - ▶ for all proper  $\pi$ ,  $V^\pi(s) \preceq \vec{b}$
  - ▶ for all improper  $\pi$ ,  $V^\pi(s) \not\preceq \vec{b}$
- ▶  $\vec{u} \preceq \vec{v}$  iff  $\vec{u}[i] \leq \vec{v}[i], i = 1, \dots, n$
- ▶ *modified MOVI* + weak assumption is sound and complete
- ▶ can derive or estimate upper bound  $\vec{b}$

# Contributions

defining MOSSPs: a generalisation of MOMDPs

assumptions for convergence of MOVI for MOSSPs

heuristics for MOSSPs

heuristic search algorithms for solving MOSSPs: iMOLAO\*, MOLRTDP

experimental evaluation

# Heuristic search

- ▶ heuristic search powerful for (optimal) planning
- ▶ does not require enumerating the whole potentially exponential search space
- ▶ heuristic search algorithms:

	SO	MO
deterministic	A*, BiA* etc.	NAMOA*
stochastic	(i)LAO*, LRTDP etc.	our contributions

# MOSSP Heuristic

- ▶ admissible heuristics in single-objective (SO) deterministic case guarantee optimality with  $A^*$
- ▶ near optimality for SSPs with (i)LAO\* and LRTDP in finitely many iterations
  - ▶ optimality with possibly infinitely many iterations
- ▶ in SO case, heuristic for a state is a scalar  $h(s) \in \mathbb{R}_{\geq 0}$
- ▶ in MO case, heuristic for a state is a set of vectors  $\mathbf{H}(s) \subset \mathbb{R}_{\geq 0}^n$

## Definition (admissible MO heuristic)

$\mathbf{H}$  is *admissible* if  $\forall s \in S \setminus G$ , for all  $\vec{v} \in \mathbf{V}^*(s)$  there exists  $\vec{u} \in \mathbf{H}(s)$  such that  $\vec{u} \preceq \vec{v}$  where  $\mathbf{V}^*$  is the optimal value function, and  $\forall g \in G, \mathbf{H}(g) = \{\vec{0}\}$ .

- ▶ same definition as in deterministic MO case [Mandow and Pérez-de-la-Cruz, 2010]

## Domain independent MOSSP Heuristics

- ▶ no heuristic: the zero heuristic defined by  $\mathbf{H}_{zero}(s) = \{\vec{0}\}$
- ▶ can use SO heuristics
  - *ideal point heuristic*; apply an SO heuristic  $h_i$  to each objective in isolation:

$$\mathbf{H}_{ideal}(s) = \{[h_1(s), \dots, h_n(s)]\}$$

- e.g.  $\mathbf{H}_{ideal}^{max}$ ,  $\mathbf{H}_{ideal}^{pdb2}$ ,  $\mathbf{H}_{ideal}^{pdb3}$  [Bonet and Geffner, 2001; Klößner and Hoffmann, 2021]
- ▶ can use deterministic MO heuristics
  - construct determinised problem by replacing each probabilistic effect with a deterministic action
  - e.g.  $\mathbf{H}_{mo}^{comax}$ ,  $\mathbf{H}_{mo}^{pdb2}$ ,  $\mathbf{H}_{mo}^{pdb3}$  [Geißer et al., 2022]
- ▶ *new* abstraction heuristics
  - combine values from solving smaller projections of the problem
  - considers both MO and stoch. features of MOSSPs
  - e.g.  $\mathbf{H}_{mosspr}^{pdb2}$ ,  $\mathbf{H}_{mosspr}^{pdb3}$
- ▶ other existing SO and MO heuristics can be extended in above ways

# Contributions

defining MOSSPs: a generalisation of MOMDPs

assumptions for convergence of MOVI for MOSSPs

heuristics for MOSSPs

heuristic search algorithms for solving MOSSPs: iMOLAO\*, MOLRTDP

experimental evaluation

- gradually build partial solution via heuristic search

---

**Algorithm: iMOLAO\***


---

**Data:** MOSSP problem

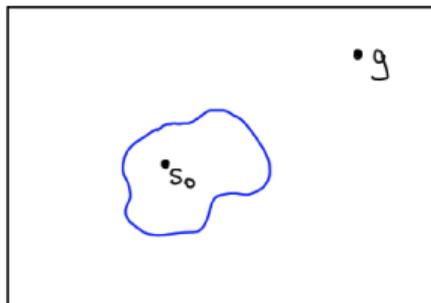
$P = (S, A, s_0, G)$ , heuristic  $\mathbf{H}$ ,  
and consistency threshold  $\epsilon$

```

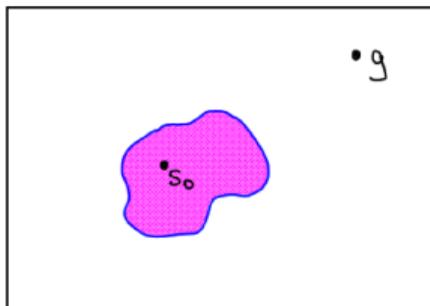
1  $\mathbf{V} \leftarrow \mathbf{H}$ ;  $\Pi \leftarrow \emptyset$ ;  $F \leftarrow \{s_0\}$ ;  $I \leftarrow \emptyset$ ;  $N \leftarrow \{s_0\}$ 
2 while
   $((F \cap N) \setminus G \neq \emptyset) \wedge (\max_{s \in N} \text{res}(s) < \epsilon)$ 
  do
3    $F = \emptyset$ 
4    $N \leftarrow \text{postorderTraversalDFS}(s_0, \Pi)$ 
5   for  $s \in N$  in the computed order do
6      $\mathbf{V}(s) \leftarrow \text{MOBellmanBackup}(s)$ 
7      $\Pi(s) = \text{getActions}(s, \mathbf{V})$ 
8     if  $s \notin I$  then  $F = F \cup \{s\}$ ;
9      $I = I \cup \{s\}$ 
10 return  $\mathbf{V}$ 

```

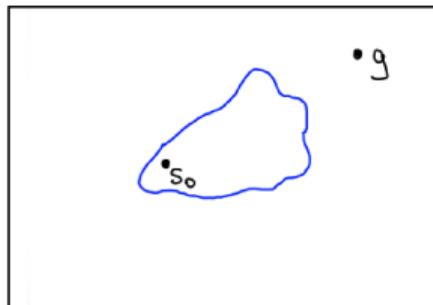
---



1. collect states  $N$  of the current best partial solution  $\Pi$  in postorder traversal DFS order



2. run MO Bellman Backups once on each state in the collected order



3. update new current partial solution  $\Pi$ ; then return to step 1 and repeat until convergence criteria is met

---

## Algorithm: MOLRTDP

---

**Data:** MOSSP problem  $P = (S, A, s_0, G)$ , heuristic  $H$ , and consistency threshold  $\varepsilon$

**procedure** MOLRTDP( $P, \varepsilon, H$ )

```

1  V  $\leftarrow$  H
2  while  $\neg s_0.solved$  do
3     $visited \leftarrow \emptyset$ 
4     $s \leftarrow s_0$ 
5    while  $\neg s.solved$  do
6       $visited.push(s)$ 
7      if  $s \in G$  then break ;
8      V( $s$ )  $\leftarrow$  MOBellmanBackup( $s$ )
9       $a \leftarrow$  sampleUnsolvedGreedyAction( $s$ )
10      $s \leftarrow$  sampleUnsolvedNextState( $s, a$ )
11     while  $\neg visited.empty()$  do
12        $s \leftarrow visited.pop()$ 
13       if  $\neg checkSolved(s)$  then break ;
14  return V
  
```

---

**routine** checkSolved( $s$ )

```

1   $rv \leftarrow true$ ;  $open \leftarrow \emptyset$ ;  $closed \leftarrow \emptyset$ 
2  if  $\neg s.solved$  then  $open.push(s)$  ;
3  while  $\neg open.empty()$  do
4     $s \leftarrow open.pop()$ 
5    if  $res(s) > \varepsilon$  then
6       $rv \leftarrow false$ 
7      continue
8    for  $a \in getActions(s, V)$  do
9      for  $s' \in successors(s, a)$  do
10     if  $\neg s'.solved \wedge s' \notin open \cup closed$ 
11     then  $open.push(s')$  ;
12  if  $rv$  then for  $s \in closed$  do  $s.solved = true$  ;
13  else
14    while  $closed \neq \emptyset$  do
15      $s \leftarrow closed.pop()$ 
16     V( $s$ )  $\leftarrow$  MOBellmanBackup( $s$ )
17  return  $rv$ 
  
```

---

# MOLRTDP

---

## Algorithm: MOLRTDP

---

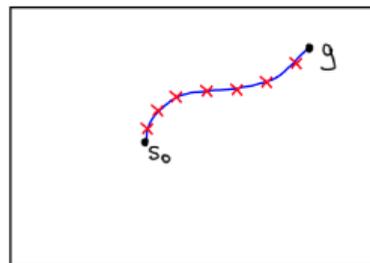
**Data:** MOSSP problem  $P = (S, A, s_0, G)$ ,  
heuristic  $H$ , consistency threshold  $\varepsilon$

**procedure**  $MOLRTDP(P, \varepsilon, H)$

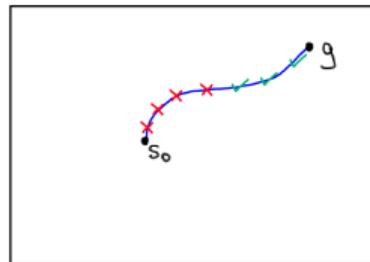
```
1  V  $\leftarrow H$ 
2  while  $\neg s_0.solved$  do
3     $visited \leftarrow \emptyset$ 
4     $s \leftarrow s_0$ 
5    while  $\neg s.solved$  do
6       $visited.push(s)$ 
7      if  $s \in G$  then break ;
8       $V(s) \leftarrow MOBellmanBackup(s)$ 
9       $a \leftarrow sampleUnsolvedGreedyAction(s)$ 
10      $s \leftarrow sampleUnsolvedNextState(s, a)$ 
11   while  $\neg visited.empty()$  do
12      $s \leftarrow visited.pop()$ 
13     if  $\neg checkSolved(s)$  then break ;
14 return V
```

---

- a state is 'solved' if its own and all its ancestors' values under the current partial solution have converged



1. random greedy walk to the goal on **unsolved states** if possible, while performing backups along the way



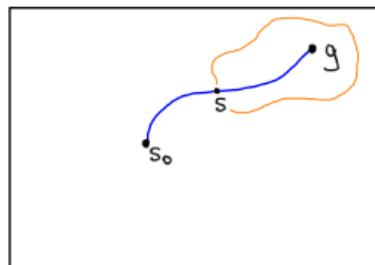
2. update whether states are **solved** in reverse order of the trialed walk using  $checkSolved$ ; return to step 1 and repeat until  $s_0$  is solved

# MOLRTDP - checkSolved

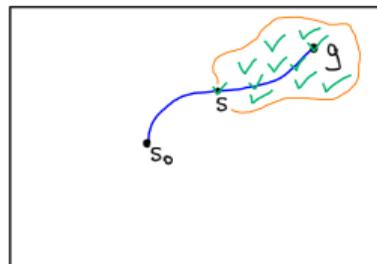
**routine** checkSolved( $s$ )

```
1   $rv \leftarrow true$ ;  $open \leftarrow \emptyset$ ;  $closed \leftarrow \emptyset$ 
2  if  $\neg s.solved$  then  $open.push(s)$ ;
3  while  $\neg open.empty()$  do
4     $s \leftarrow open.pop()$ 
5    if  $res(s) > \epsilon$  then
6       $rv \leftarrow false$ 
7      continue
8    for  $a \in getActions(s, \mathbf{V})$  do
9      for  $s' \in successors(s, a)$  do
10       if  $\neg s'.solved \wedge s' \notin open \cup closed$ 
11        then  $open.push(s')$ ;
12  if  $rv$  then for  $s \in closed$  do
13     $s.solved = true$ ;
14  else
15    while  $closed \neq \emptyset$  do
16       $s \leftarrow closed.pop()$ 
17       $\mathbf{V}(s) \leftarrow MOBellmanBackup(s)$ 
18  return  $rv$ 
```

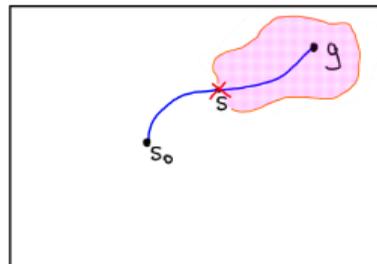
- a state is 'solved' if its own and all its ancestors' values under the current partial solution have converged



2.1. collect all greedy ancestor states of  $s$  under the current partial solution



2.2.a. if values for all ancestor states have converged, set the state and its ancestors to be solved



2.2.b. else run backups on the state and all its ancestors

# Contributions

defining MOSSPs: a generalisation of MOMDPs

assumptions for convergence of MOVI for MOSSPs

heuristics for MOSSPs

heuristic search algorithms for solving MOSSPs: iMOLAO\*, MOLRTDP

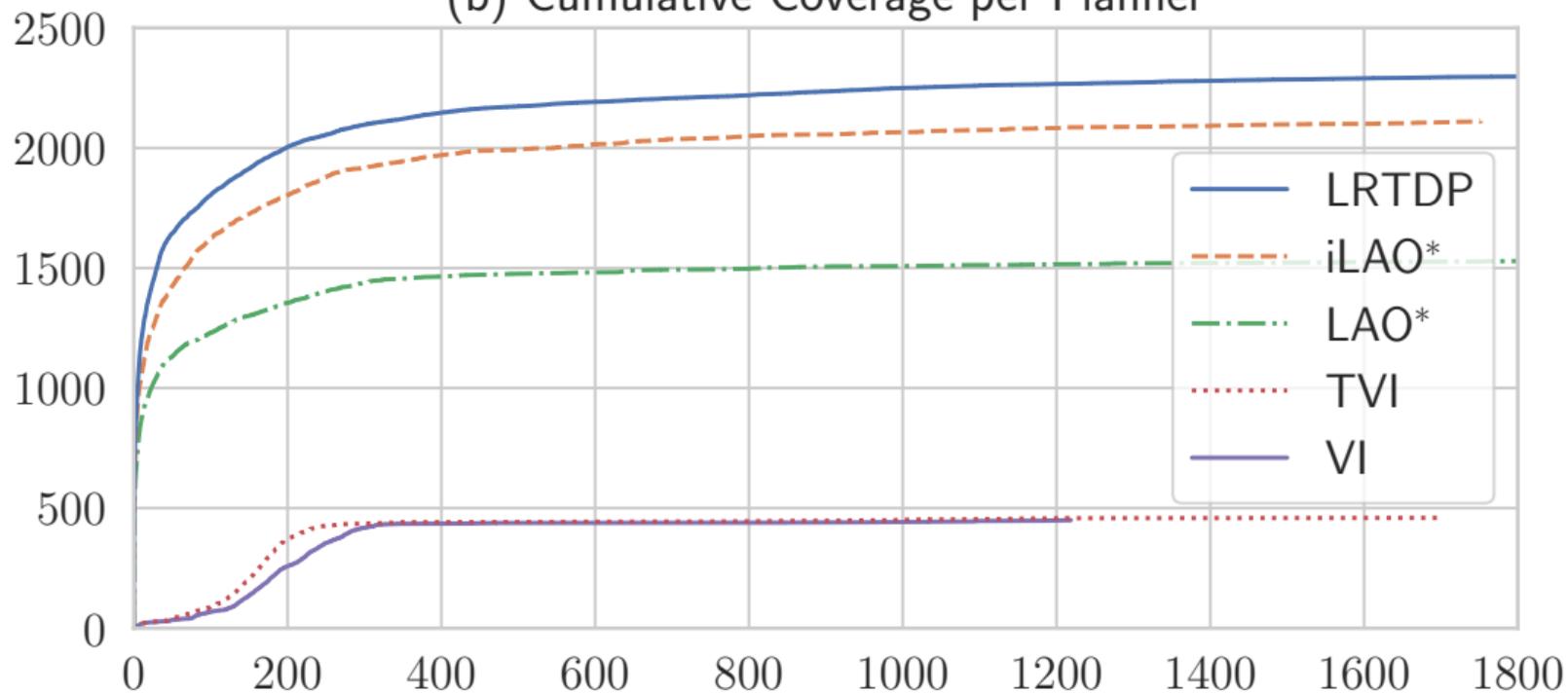
**experimental evaluation**

## Experimental setup

- ▶ 5 solvers (MO prefix omitted): VI, TVI, LAO\*, iLAO\*, LRTDP
- ▶ 9 heuristics: zero,  $H_{ideal}^{max}$ ,  $H_{mo}^{comax}$ ,  $H_{ideal}^{pdb2}$ ,  $H_{ideal}^{pdb3}$ ,  $H_{mo}^{pdb2}$ ,  $H_{mo}^{pdb3}$ ,  $H_{mossp}^{pdb2}$ ,  $H_{mossp}^{pdb3}$
- ▶ 30 minute timeout, single CPU core, 4GB memory
- ▶ 7 domains, mix of MO and probabilistic interesting domains
- ▶ 610 total problems
- ▶  $5 \times 9 \times 610 = 27450$  possible experimental configurations

## Best planner

(b) Cumulative Coverage per Planner



Cumulative coverage per planner, marginalised by heuristics.

## Best heuristic

Heuristic	Normalised coverage
$H_{\text{mossp}}^{\text{pdb3}}$	19.1
$H_{\text{mossp}}^{\text{pdb2}}$	17.9
$H_{\text{mo}}^{\text{pdb3}}$	17.4
$H_{\text{mo}}^{\text{pdb2}}$	17.1
$H_{\text{ideal}}^{\text{max}}$	14.9
$H_{\text{ideal}}^{\text{pdb3}}$	14.7
$H_{\text{ideal}}^{\text{pdb2}}$	14.4
blind	12.1
$H_{\text{mo}}^{\text{comax}}$	10.7

Heuristic	Unnormalised coverage
$H_{\text{mossp}}^{\text{pdb3}}$	893.5
$H_{\text{mossp}}^{\text{pdb2}}$	893.3
$H_{\text{mo}}^{\text{pdb3}}$	871.2
$H_{\text{mo}}^{\text{pdb2}}$	851.8
$H_{\text{ideal}}^{\text{pdb3}}$	768.7
$H_{\text{ideal}}^{\text{max}}$	755.2
$H_{\text{ideal}}^{\text{pdb2}}$	737.2
blind	572.9
$H_{\text{mo}}^{\text{comax}}$	504.5

- ▶ marginalise by planner
- ▶ normalisation is done over domain due to uneven number of problems
- ▶ ranking of top 3 and bottom 3 heuristics same

## Heuristic accuracy

	Critical Path				Abstractions				
	blind	$H_{ideal}^{max}$	$H_{mo}^{comax}$	$H_{ideal}^{pdb2}$	$H_{ideal}^{pdb3}$	$H_{mo}^{pdb2}$	$H_{mo}^{pdb3}$	$H_{moss}^{pdb2}$	$H_{moss}^{pdb3}$
SAR-4	100	97	44	93	92	44	44	38	<b>26</b>
SAR-5	100	97	45	93	92	45	45	39	<b>28</b>
ExBw-2d	100	59	<b>24</b>	52	<b>45</b>	52	<b>45</b>	52	45
ExBw-3d	100	58	<b>22</b>	52	<b>44</b>	52	<b>44</b>	52	44
Tireworld	100	100	68	100	100	68	68	57	<b>12</b>
VisitAll	100	100	54	100	100	61	53	22	<b>12</b>
VisitAllTire	100	100	50	100	100	66	<b>45</b>	66	<b>45</b>

- ▶ mean relative error (%) of heuristic value at the initial state relative to the optimal value for solved instances
- ▶ error calculated by directed Hausdorff metric difference divided by the norm of the largest vector of the optimal value:

$$\max_{\vec{v} \in \mathbf{V}} \min_{\vec{u} \in \mathbf{H}} d(\vec{v}, \vec{u}) / \max_{\vec{v} \in \mathbf{V}} \|\vec{v}\|$$

⇒ MO features more important than stoch. features to capture in a heuristic

Thank you for your attention!