# Synthesis of Search Heuristics for Temporal Planning via Reinforcement Learning

Andrea Micheli and Alessandro Valentini

Embedded Systems Unit, Fondazione Bruno Kessler, Italy

## MOTIVATION

Once deployed, a temporal planner will solve several different problems on the same domain
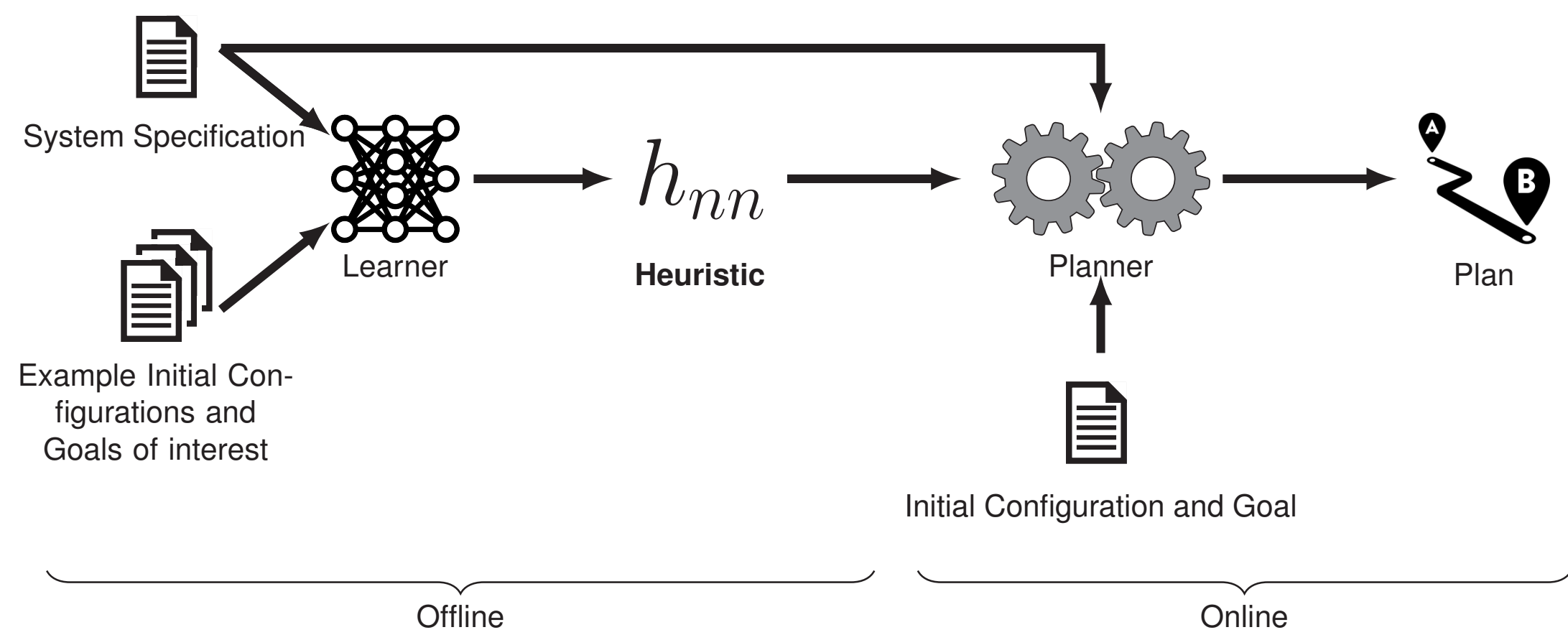
### Example

- ► Organize the logistics of the same factory once a day
- ► Operate the same drone in the same area for different missions from different initial states

### Key Intuition

Instead of resorting to pure reasoning each time, can we learn characteristics of the domain and exploit them for efficiency?
Analogous to a worker that gets accustomed to a certain workplace and gains dexterity

## PIPELINE



System Specification — Learner — $h_{nn}$ Heuristic — Planner — Plan

Example Initial Configurations and Goals of interest

Initial Configuration and Goal

Offline                    Online

We learn a specialized heuristic keeping a fully-functional planner online

## MDP FOR A BOUNDED PLANNING PROBLEM SET

### Bounded Planning Problem Set

A bounded planning problem set with at most $k$ objects for a planning domain $\mathcal{D}$ (written $\mathcal{P}_{\mathcal{D}}^k$) is a finite set of planning problems $P_i$ for $\mathcal{D}$ each having less than $k$ objects.
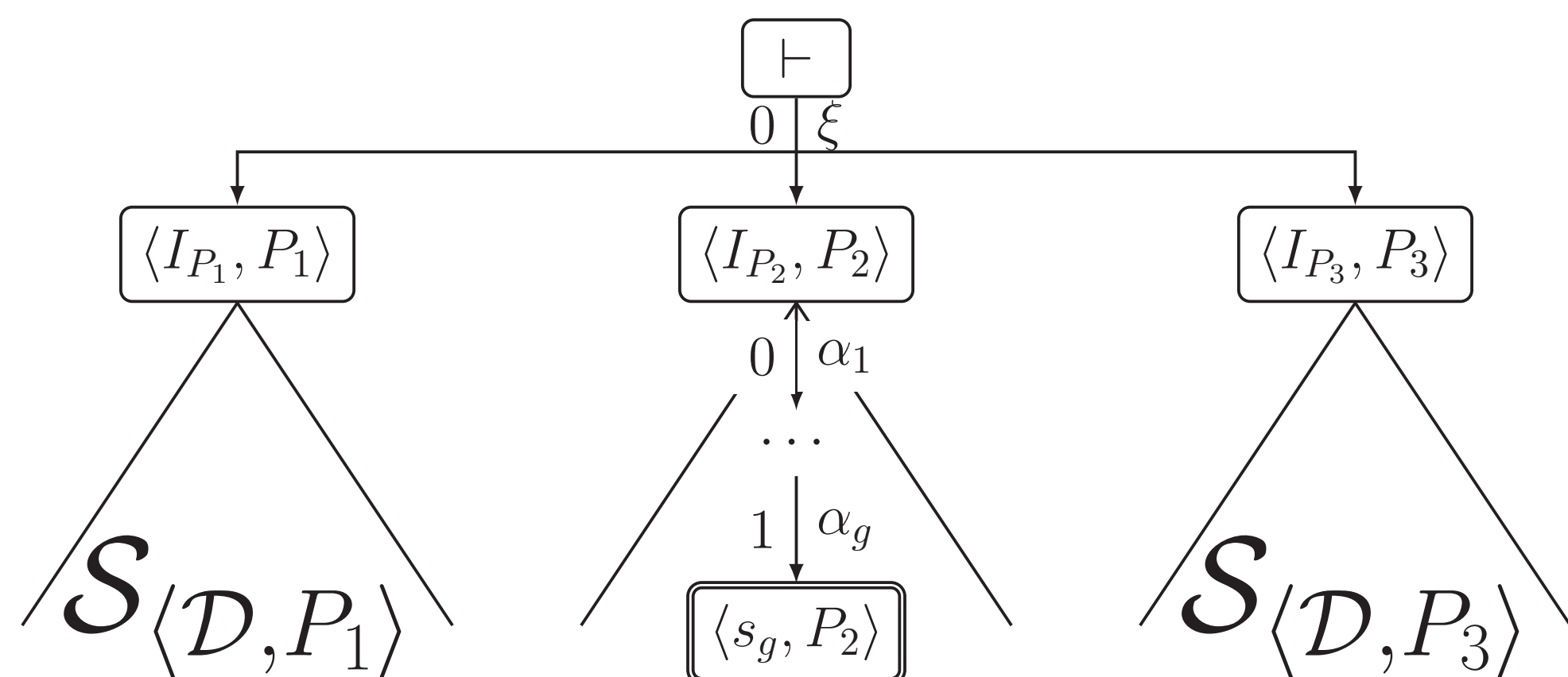
Given a planning domain and a bounded planning problem set we define an MDP $\mathcal{M}_{\mathcal{P}_{\mathcal{D}}^k} \doteq \langle S, A, T, R, \vdash \rangle$ is s.t.

$S \doteq \{\vdash\} \cup$ all planner states for all instances

$A \doteq \{\xi\} \cup$ all actions (events) for all instances;

$$T(s,a) \doteq \begin{cases} \{\langle I_{P_i}, \frac{1}{|\mathcal{P}_{\mathcal{D}}^k|}\rangle \mid P_i \in \mathcal{P}_{\mathcal{D}}^k\} & \text{if } s = \vdash, a = \xi \\ \{\langle a[s], 1\rangle\} & \text{if } s \neq \vdash \end{cases}$$

$$R(s,a,s') \doteq \begin{cases} 1 & \text{if } s' \text{ is a goal state} \\ -1 & \text{if } s' \text{ is a dead-end} \\ 0 & \text{otherwise.} \end{cases}$$
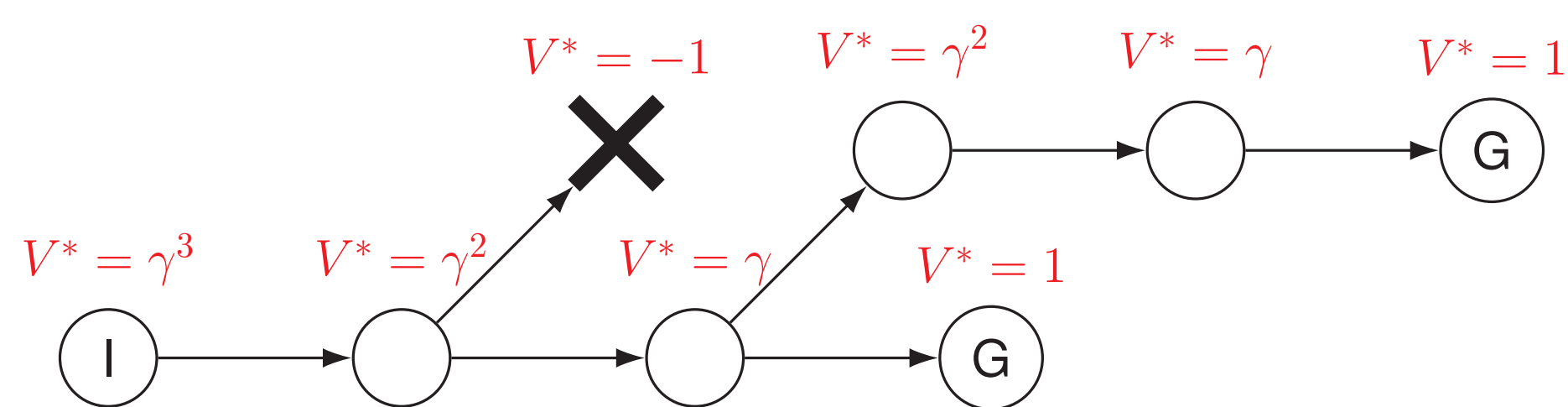


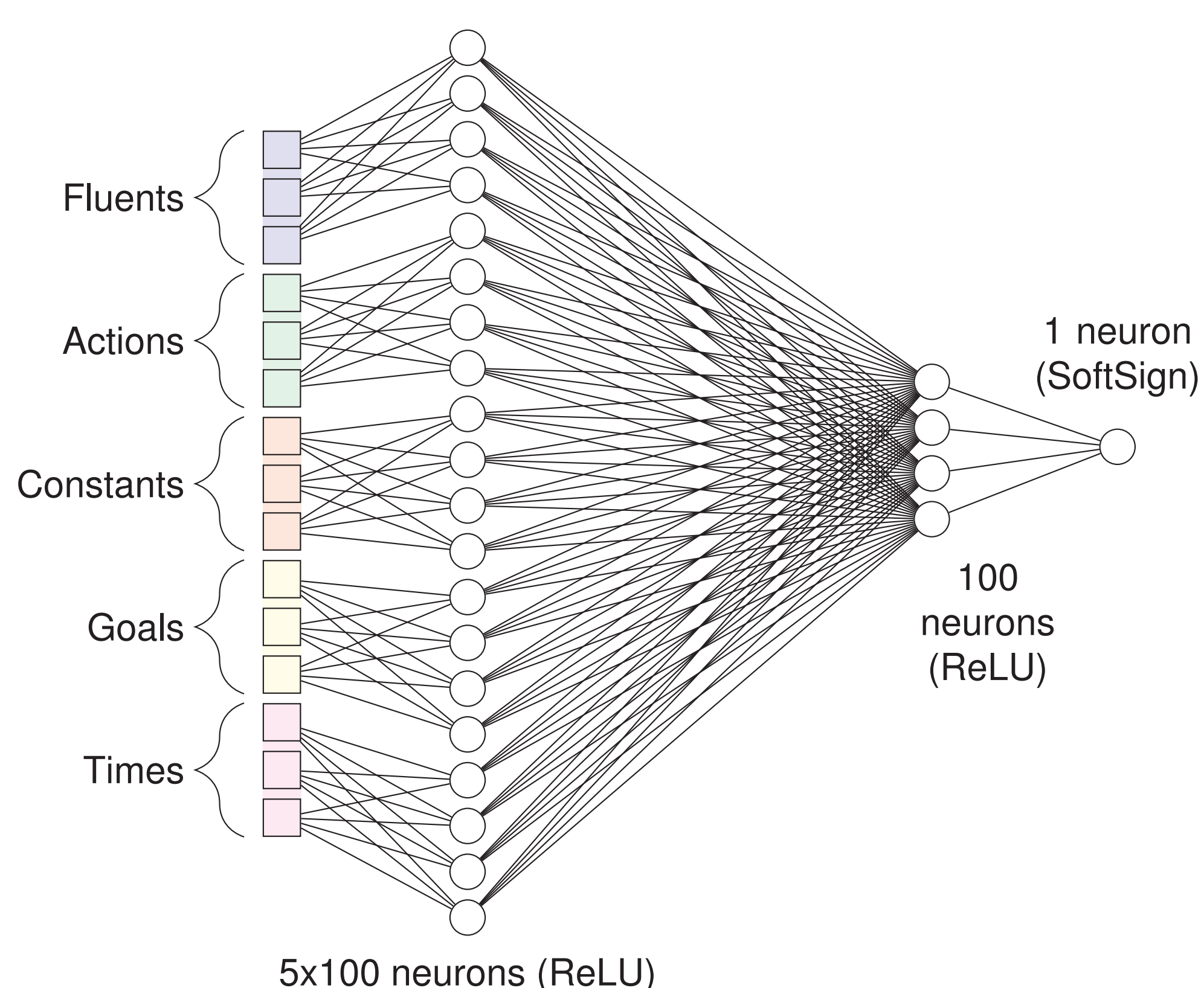## FROM THE OPTIMAL VALUE FUNCTION ($V^*$) TO THE OPTIMAL HEURISTIC ($h^*$)

For a bounded planning problem set $\mathcal{P}_{\mathcal{D}}^k$ the following equation holds.

$$h_{\mathcal{P}_{\mathcal{D}}^k}^*(s) = \begin{cases} \log_\gamma(V_{\mathcal{M}_{\mathcal{P}_{\mathcal{D}}^k}}^*(s)) & \text{if } V_{\mathcal{M}_{\mathcal{P}_{\mathcal{D}}^k}}^*(s) > 0 \\ \infty & \text{otherwise} \end{cases}$$
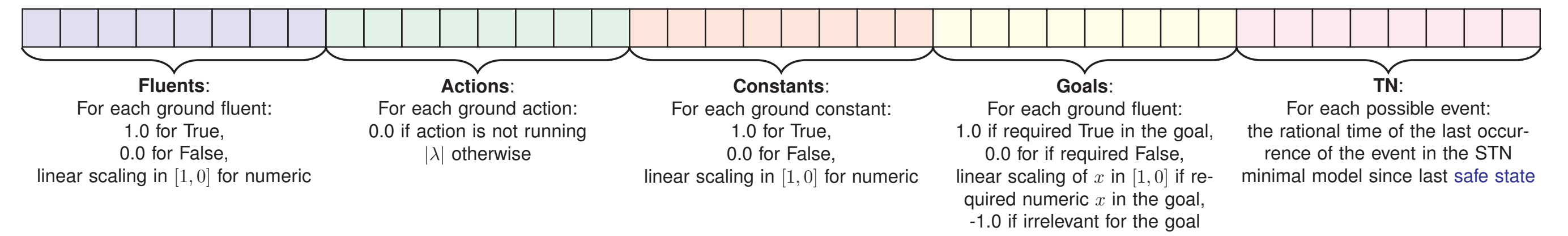
### Intuition



## NEURAL ARCHITECTURE: PREDICTING $V^*$



Fluents
Actions
Constants
Goals
Times

5x100 neurons (ReLU)

100 neurons (ReLU)

1 neuron (SoftSign)

## STATE VECTORIZATION

Given a planning state $s$ we derive a vector $\vec{s}$ in $\mathbb{R}^x$



**Fluents:** For each ground fluent: 1.0 for True, 0.0 for False, linear scaling in $[0,1]$ for numeric

**Actions:** For each ground action: 0.0 if action is not running $|A|$ otherwise

**Constants:** For each ground constant: 1.0 for True, 0.0 for False, linear scaling in $[1,0]$ for numeric

**Goals:** For each ground fluent: 1.0 if required True in the goal, 0.0 if required False, linear scaling of $x$ in $[1,0]$ if required numeric $x$ in the goal, -1.0 if irrelevant for the goal

**TN:** For each possible event: the rational time of the last occurrence of the event in the STN minimal model since last safe state

## RL-BASED HEURISTIC LEARNER

Basically, Deep-Q-Learning on $\mathcal{M}_{\mathcal{P}_{\mathcal{D}}^k}$ with some adjustments:

- ► State value function, single output network instead of DQN
- ► Heuristic-proportional random action selection
- ► Bias in problem selection
- ► Memory replay with positive bias
- ► Fixed max depth of episodes

```
1:  procedure RL2PLANHEURISTIC(tis, N_episodes)
2:     V_nn ← INITNN( )
3:     mem ← LIST( )
4:     i2s ← {i ↦ 0 | i ∈ tis}
5:     for i ∈ 1, ..., N_episodes do
6:        (s, goals) = inst ← PICKKEYINVPROPORTIONALLYTOVALUE(i2s)
7:        ⟨done, solved⟩ ← ⟨False, False⟩
8:        π ← ⟨s⟩
9:        while not done do
10:           ε ← ε_max × ε^(⌊(ln(ε_min/ε_max))/N_episodes⌋ × i)
11:           if RANDOM( ) < ε then
12:              α ←SELECTACTIONUSINGHEURISTIC(s)
13:           else
14:              α ←SELECTACTIONUSINGPOLICY(V_nn, s)
15:           ⟨s', done, ρ⟩ ←DOSTEP(π, s, α, inst)
16:           APPEND(mem, ⟨s, ρ⟩)
17:           if ρ[α] = 1 then
18:              solved ← True
19:              APPEND(π, ⟨s'⟩)
20:           s ← s'
21:        V_nn ←REPLAY(V_nn, mem)
22:        if solved then
23:           i2s[inst] ← i2s[inst] + 1
24:     return V_nn
```

### Learned heuristic

The learned heuristic $h_{nn}$ is an approximation of $h^*$

$$h_{nn}(s) \doteq \begin{cases} \min(\log_\gamma(V_{nn}(\vec{s})), \Delta_h) & \text{if } V_{nn}(\vec{s}) > 0 \\ \Delta_h & \text{if } V_{nn}(\vec{s}) = 0 \\ 2\Delta_h - \min(\log_\gamma(-V_{nn}(\vec{s})), \Delta_h) & \text{otherwise} \end{cases}$$

Where $\Delta_h$ is bigger than the pre-fixed cutoff length of episodes set in learning

**Soundness:** $h_{nn}$ never returns $\infty$ because learning can be imperfect and we do not want unsound pruning.

## EXPERIMENTAL EVALUATION

### Case studies

- ► **MaJSP**: A fleet of AGVs with logistics tasks in a warehouse.
  - ▷ The problems differ for the number of items to be moved and the intermediate steps.
- ► **Kitting**: A single robot serving a continuous production line with kits of components taken from shelves.
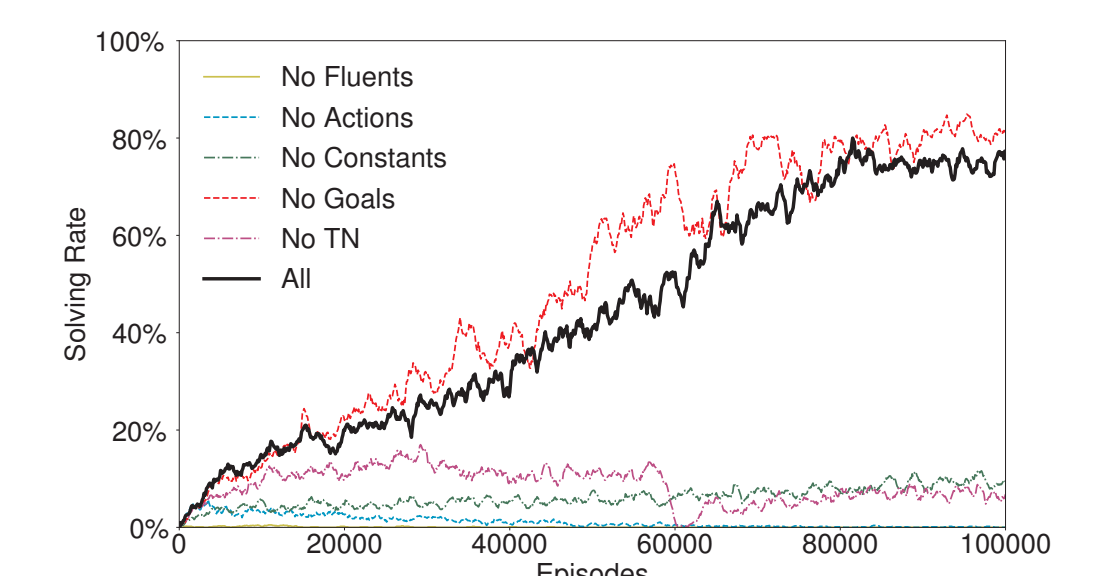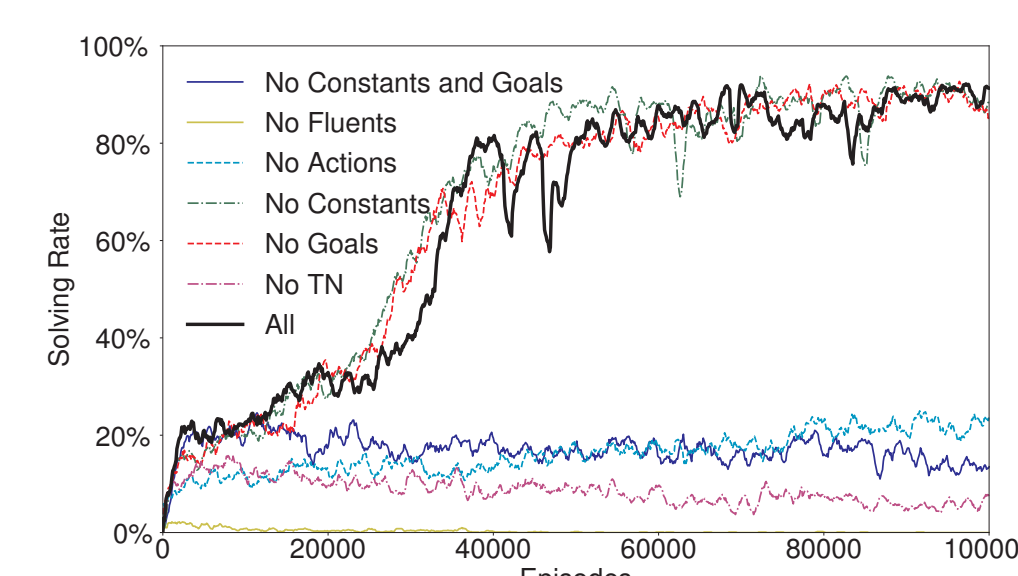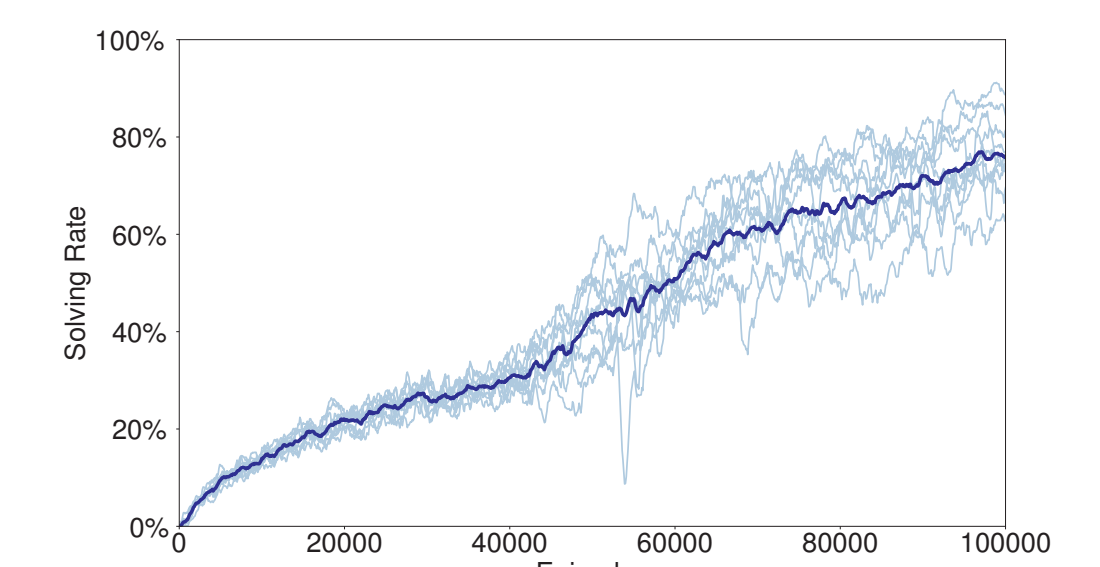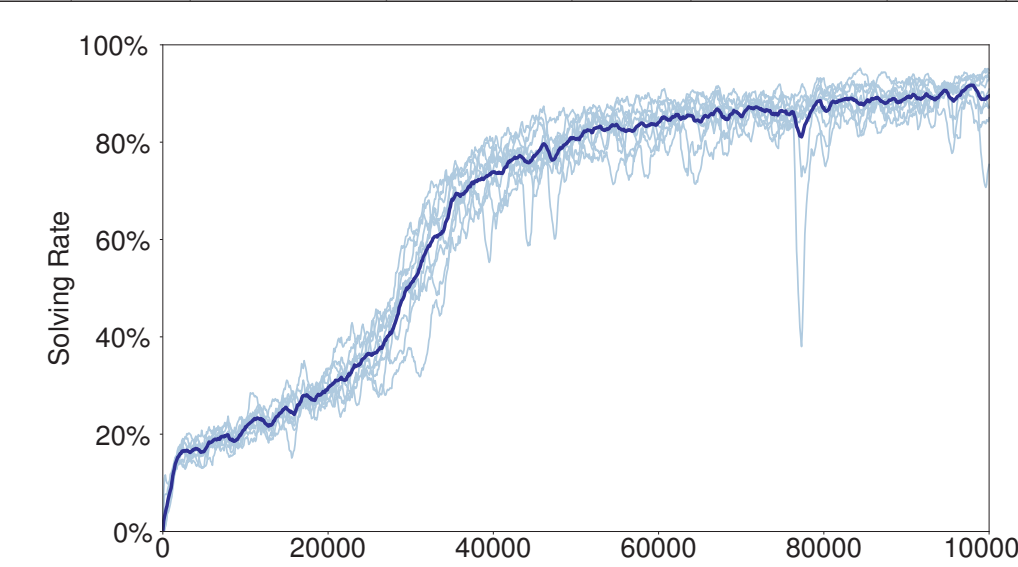  - ▷ The problems require different sequences of kits to be delivered.

### Competitors

- ► TAMER ($h_{add}$): our fully-symbolic state-of-the-art planning
- ► $\pi_{nn}$: The learned RL policy executed without backtracking
- ► TAMER ($h_{nn}$): our planner equipped with the learned heuristic $h_{nn}$

### Results

10-fold cross-validation and sensitivity analysis over 100K RL episodes

**MaJSP**

| fold (size: 77) | TAMER ($h_{add}$) | | $\pi_{nn}$ | | | TAMER ($h_{nn}$) | |
|---|---|---|---|---|---|---|---|
| | solved | avg plan size | # episodes | solved | avg plan size | solved | avg plan size |
| 1 | 52 | 14 | 50k | 66 | 25 | 73 | 18 |
| | | | 100k | 71 | 22 | 73 | 18 |
| 2 | 58 | 14 | 50k | 70 | 22 | 75 | 17 |
| | | | 100k | 70 | 19 | 72 | 17 |
| 3 | 58 | 14 | 50k | 70 | 23 | 73 | 17 |
| | | | 100k | 73 | 19 | 75 | 17 |
| 4 | 57 | 13 | 50k | 66 | 21 | 72 | 17 |
| | | | 100k | 68 | 20 | 75 | 17 |
| 5 | 55 | 15 | 50k | 66 | 25 | 75 | 19 |
| | | | 100k | 66 | 23 | 69 | 19 |
| 6 | 60 | 14 | 50k | 66 | 23 | 76 | 17 |
| | | | 100k | 69 | 17 | 77 | 17 |
| 7 | 54 | 14 | 50k | 66 | 21 | 73 | 18 |
| | | | 100k | 69 | 21 | 73 | 18 |
| 8 | 57 | 14 | 50k | 61 | 23 | 73 | 18 |
| | | | 100k | 73 | 20 | 69 | 18 |
| 9 | 57 | 14 | 50k | 71 | 25 | 74 | 18 |
| | | | 100k | 70 | 21 | 70 | 18 |
| 10 | 52 | 14 | 50k | 72 | 21 | 77 | 19 |
| | | | 100k | 65 | 22 | 54 | 16 |
| all | 560 | 14 | 50k | 676 | 23 | **744** | 18 |
| | | | 100k | 699 | 20 | 708 | 18 |

**Kitting**

| fold (size: 109) | TAMER ($h_{add}$) | | $\pi_{nn}$ | | | TAMER ($h_{nn}$) | |
|---|---|---|---|---|---|---|---|
| | solved | avg plan size | # episodes | solved | avg plan size | solved | avg plan size |
| 1 | 44 | 15 | 50k | 66 | 18 | 99 | 21 |
| | | | 100k | 97 | 21 | 107 | 21 |
| 2 | 35 | 15 | 50k | 66 | 21 | 95 | 22 |
| | | | 100k | 82 | 21 | 97 | 21 |
| 3 | 38 | 15 | 50k | 55 | 18 | 83 | 20 |
| | | | 100k | 97 | 20 | 99 | 20 |
| 4 | 45 | 15 | 50k | 68 | 20 | 98 | 19 |
| | | | 100k | 88 | 22 | 100 | 21 |
| 5 | 47 | 15 | 50k | 85 | 19 | 101 | 19 |
| | | | 100k | 88 | 19 | 101 | 19 |
| 6 | 38 | 15 | 50k | 53 | 20 | 85 | 20 |
| | | | 100k | 78 | 22 | 108 | 23 |
| 7 | 30 | 15 | 50k | 44 | 18 | 75 | 19 |
| | | | 100k | 90 | 24 | 106 | 23 |
| 8 | 42 | 15 | 50k | 65 | 19 | 93 | 20 |
| | | | 100k | 95 | 21 | 104 | 21 |
| 9 | 36 | 15 | 50k | 59 | 18 | 70 | 17 |
| | | | 100k | 89 | 22 | 91 | 20 |
| 10 | 40 | 14 | 50k | 71 | 19 | 95 | 21 |
| | | | 100k | 92 | 21 | 102 | 21 |
| all | 395 | 15 | 50k | 617 | 19 | 896 | 20 |
| | | | 100k | 896 | 21 | **1015** | 21 |





## CONCLUSION

### Take-Away Message

- ► Strict correlation between planning heuristics and state value functions in RL
- ► Use RL to automatically synthesize planning heuristics looks promising

### Future work

- ► Extend the approach to overcome limitations
  - ▷ Fixed state size, Fixed network architecture, Bounded numeric values, Incomplete temporal information
- ► Supervised learning from search spaces