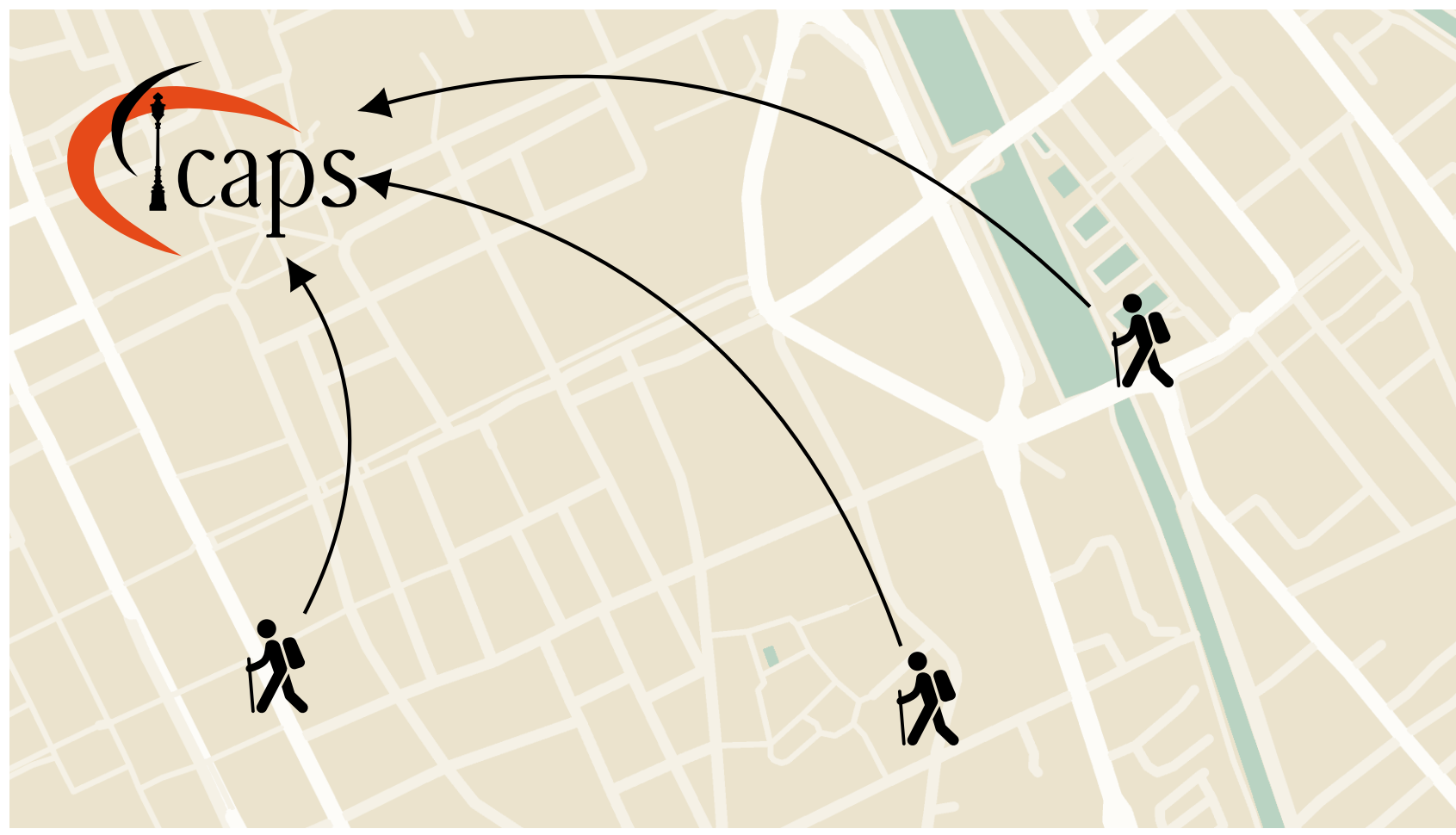


## Motivation



Frequently, we have to **solve** tasks from the **same state space** with the **same goal**.

⇒ Once learn a good heuristic, speed up future searches.

## FDR Task

An FDR planning task is a tuple

$\Pi = (V, A, I, G)$  with:

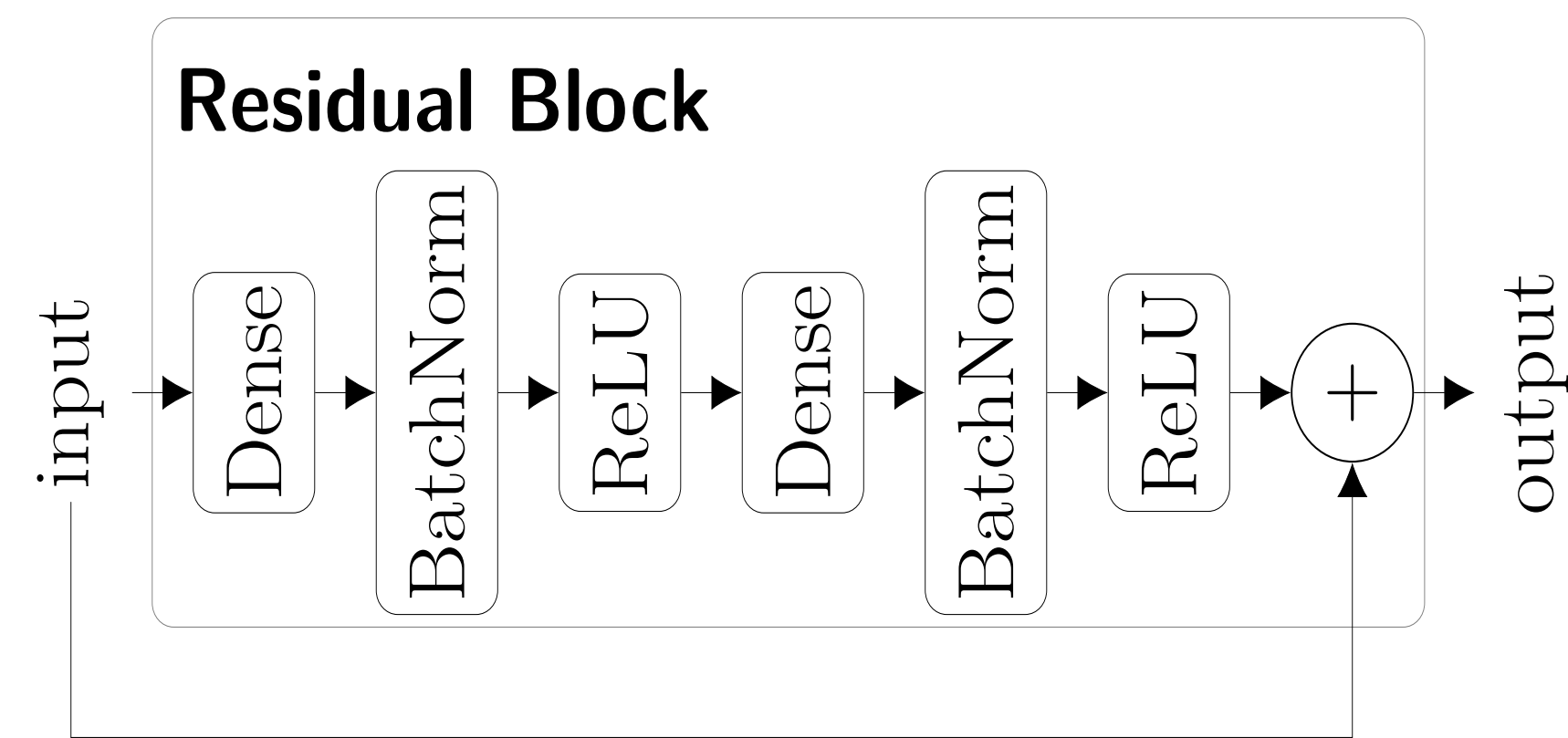
- ▶  $V$ : a set of multi-valued variables
- ▶  $A$ : a set of actions
- ▶  $I$ : an initial state
- ▶  $G$ : a partial variable assignment

States assign every variable a value.

A state  $s$  can be represented as vector:

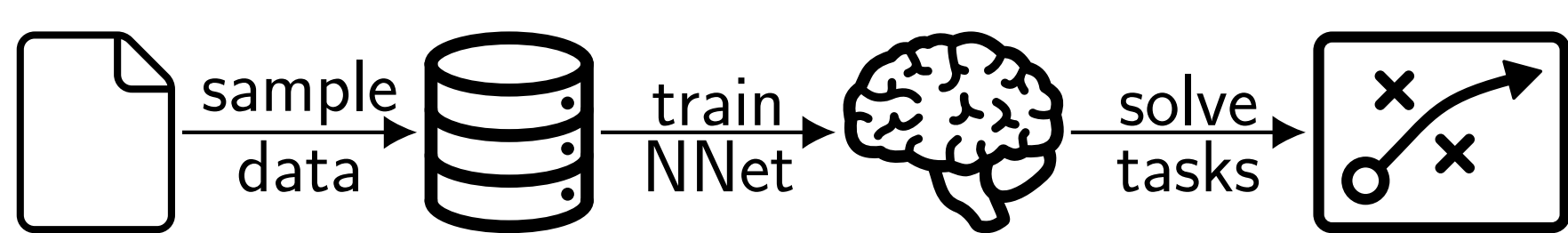
$[s(v) = d \text{ for } v \in V \text{ for } d \in \text{dom}(V)]$

## Residual Network



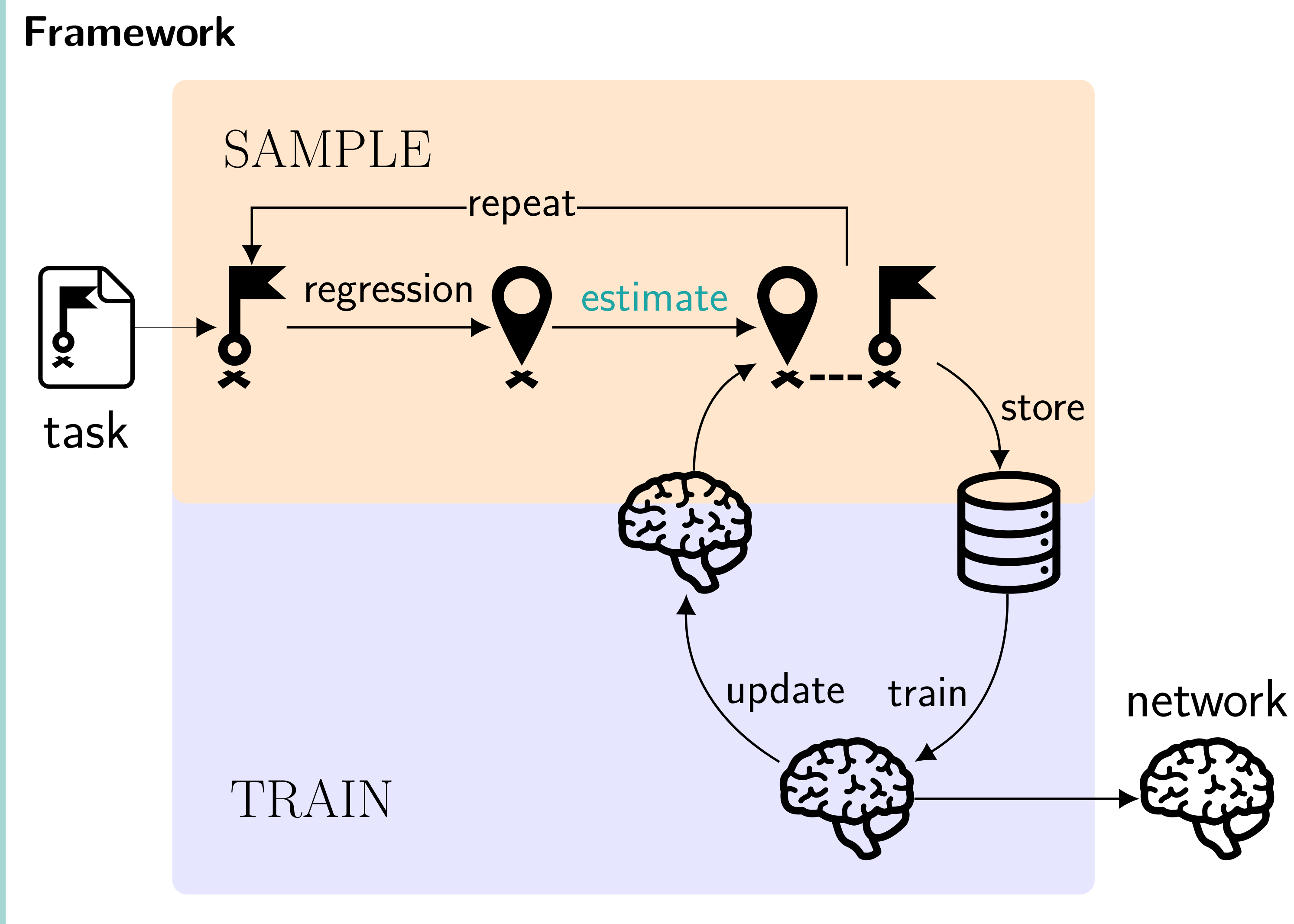
We use one residual block and two dense layers. Each dense layer has 250 neurons.

## Our Previous Work (SL)



1. Sample for 400 hours data via progression and solve with GBFS(FF).
2. Use supervised learning to train NNet.

# Reinforcement Learning can be superior to Supervised Learning for learning heuristics.

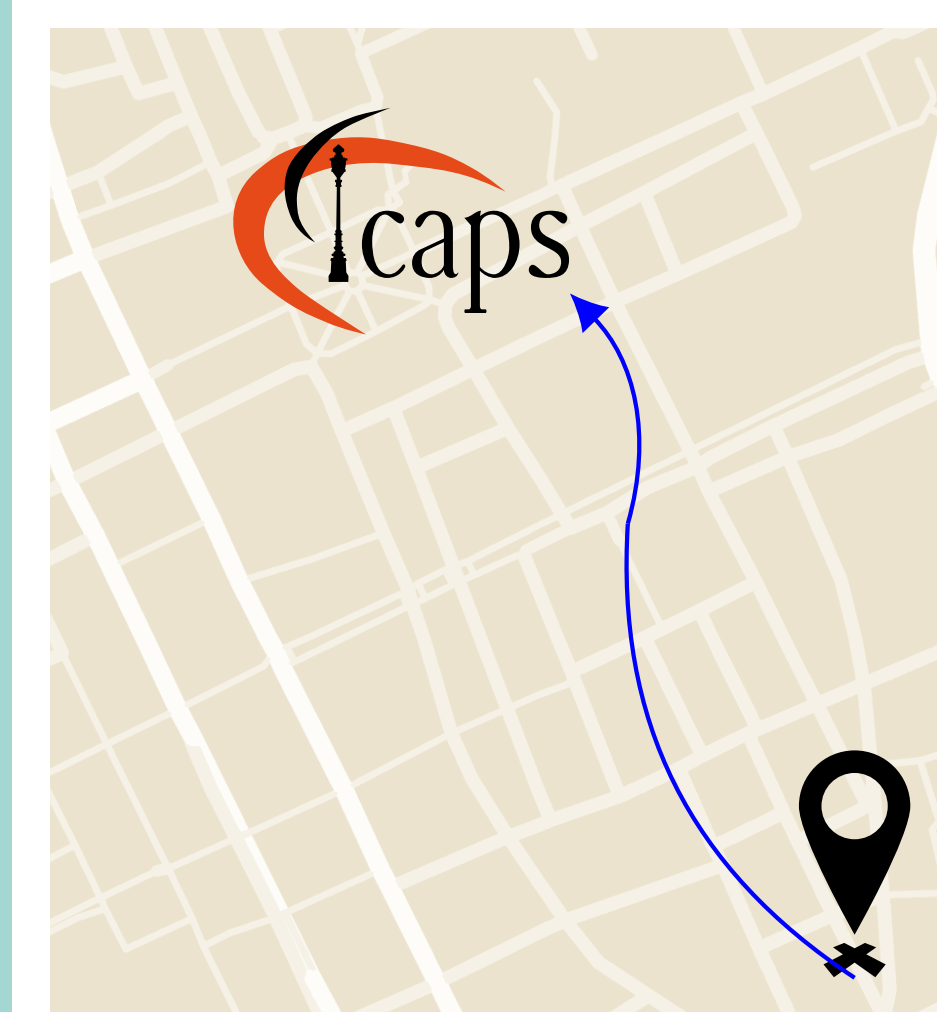


## Estimate: Approx. Value Iteration (AVI)



$$\text{estimate}(\varrho) = \min_{x \in \text{neighbors}} (\text{network}(x) + \text{distance}(\varrho, x))$$

## Estimate: Sampling Search (SaSe)



- ▶ Use the network (brain icon) as heuristic<sup>1</sup> to find a **plan**.
- ▶  $\text{estimate}(\varrho) = |\text{plan}|$
- ▶ Increase the number of regression steps over time.

<sup>1</sup>GBFS with 10s timeout

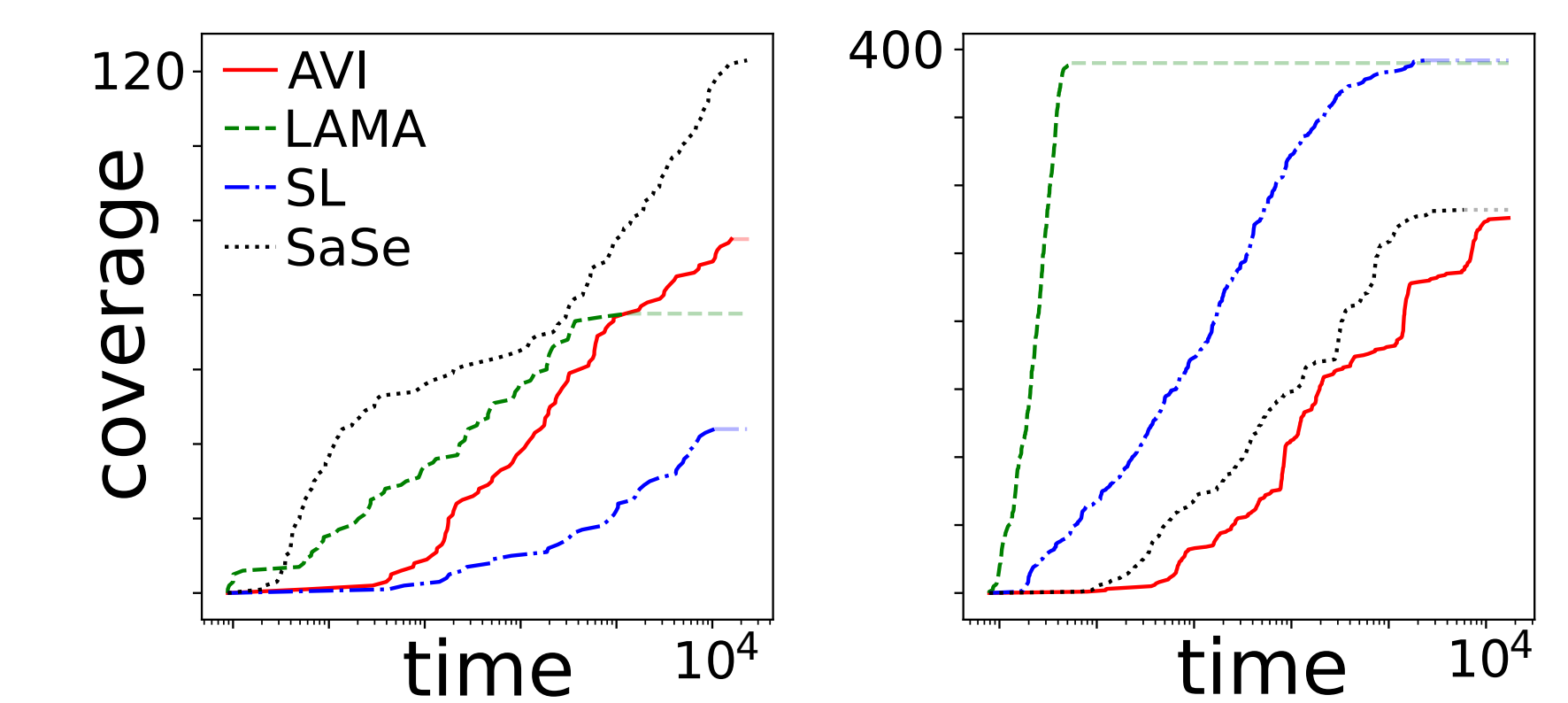
## Results

Domain	AVI	SaSe	SL	Lama
blocks	0.0	0.0	<b>98.0</b>	96.8
depots	17.7	39.7	64.3	<b>98.7</b>
grid	51.0	86.0	74.0	<b>97.0</b>
npuzzle	1.0	1.5	0.0	<b>97.8</b>
pipesworld-nt	29.8	50.4	92.8	<b>97.2</b>
rovers	25.8	35.8	12.5	<b>98.0</b>
scanalyzer	83.3	33.3	77.7	<b>97.7</b>
storage	47.5	<b>71.5</b>	22.0	37.5
transport	69.0	70.5	<b>98.0</b>	97.5
visittall	13.0	30.7	0.7	<b>95.0</b>
<b>Average</b>	<b>33.8</b>	<b>41.9</b>	<b>54.0</b>	<b>91.3</b>

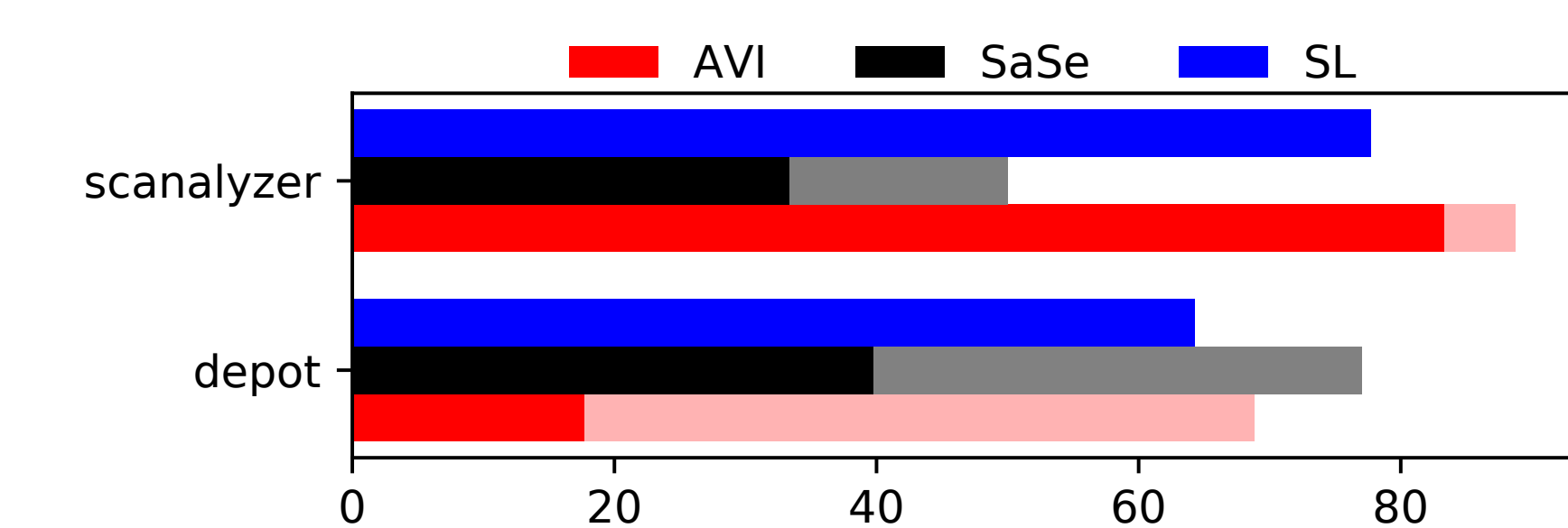
Coverage (in %) on the tasks used by Ferber et al. 2020.

Domain	AVI	SaSe	Lama
depots	15.1	6.9	<b>80.6</b>
grid	0.0	0.0	<b>90.0</b>
npuzzle	0.0	0.0	<b>84.0</b>
pipesworld-nt	1.4	25.1	<b>68.7</b>
rovers	0.1	0.8	<b>97.7</b>
scanalyzer	34.0	3.3	<b>98.7</b>
storage	18.8	<b>26.5</b>	11.0
visittall	0.0	36.0	<b>98.0</b>
<b>Average</b>	<b>8.7</b>	<b>12.3</b>	<b>78.6</b>

Coverage (in %) on the tasks too hard for training data generation for Ferber et al. 2020.



Cumulative coverage on the tasks used by Ferber et al. 2020 for the storage and transport domains.



Coverage (in %) improvements after identifying good models for the tasks by Ferber et al. 2020.

## Conclusions & Future Work

- ▶ We have no single best NN technique.
- ▶ Identifying successful models leads to significant improvements.
- ▶ Future Work: Predicting expansions and incorporating unsolved samples of SaSe.

# Reinforcement Learning for Planning Heuristics

Patrick Ferber, Malte Helmert, and Jörg Hoffmann

University of Basel, Switzerland, and Saarland University, Germany

