# SMoSE: Sparse Mixture of Shallow Experts
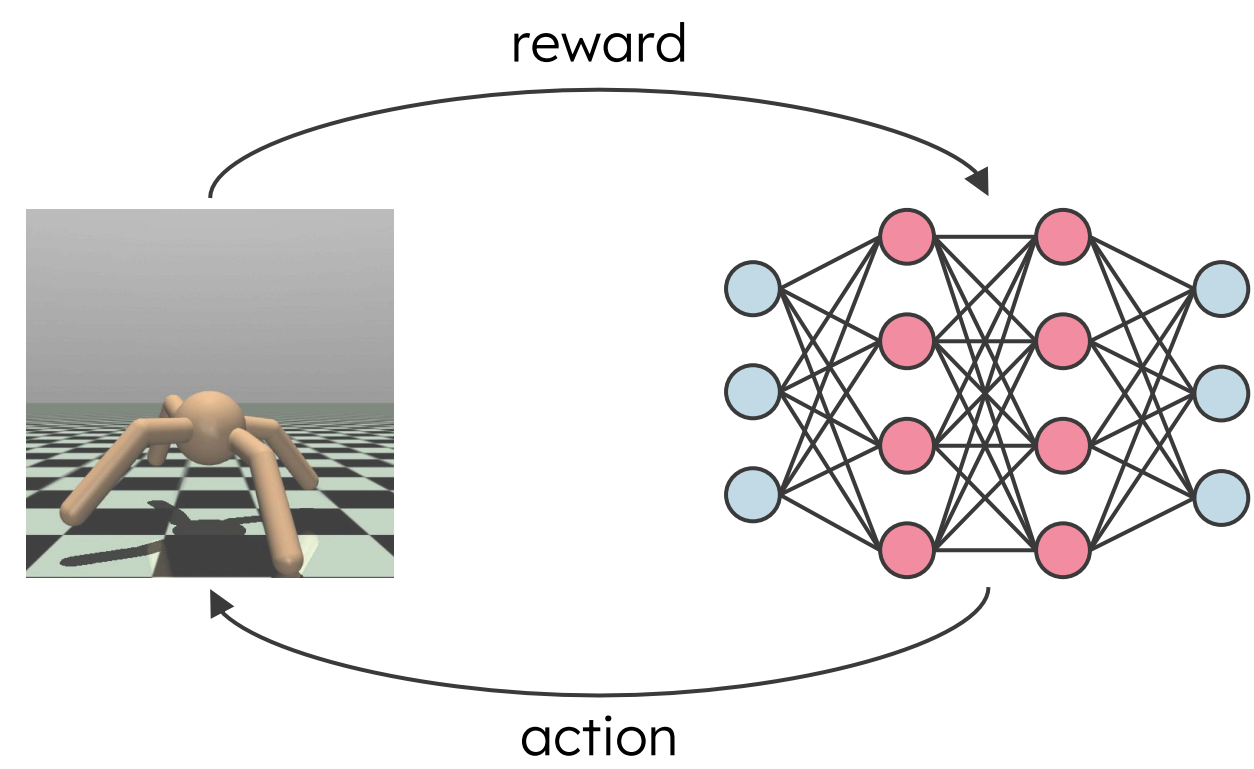## for Interpretable Reinforcement Learning in Continuous Control Tasks

Mátyás Vincze[1,2]   Laura Ferrarotti[2]   Leonardo Lucio Custode[1]

Bruno Lepri[2]   Giovanni Iacca[1]

AAAI-25 / IAAI-25 / EAAI-25

## Motivation

### Unlock safe and efficient RL

#### SOTA approaches are not interpretable



reward

action

- Scaling limits interpretability
- Closed-box models only allow explainability

### Interpretable approaches do not work in continuous control

- Evolutionary solutions are sample-inefficient (10x environment interactions)
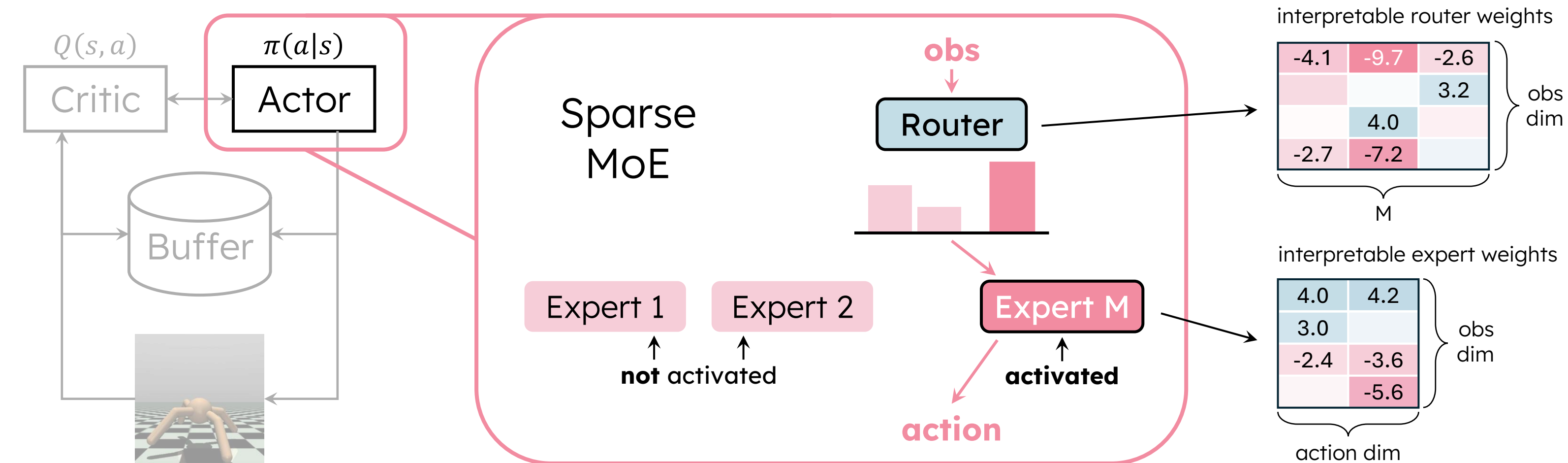- Huge performance gap compared to SOTA



Paper

## Method

### Sparse MoE actor, Linear experts, Post-training distillation

#### Architecture : Linear router, linear experts

Router partitions the state space while the experts specialize on simple skills



#### Training stabilization

- Load balancing with auxiliary loss

$$L_{aux} = 0.1 * \begin{bmatrix} f_{imp}(S) = \frac{1}{2}\left(\frac{std(Imp(S))}{mean(Imp(S))}\right)^2 \\ + \\ f_{load}(S) = \frac{1}{2}\left(\frac{std(Load(S))}{mean(Load(S))}\right)^2 \end{bmatrix}$$
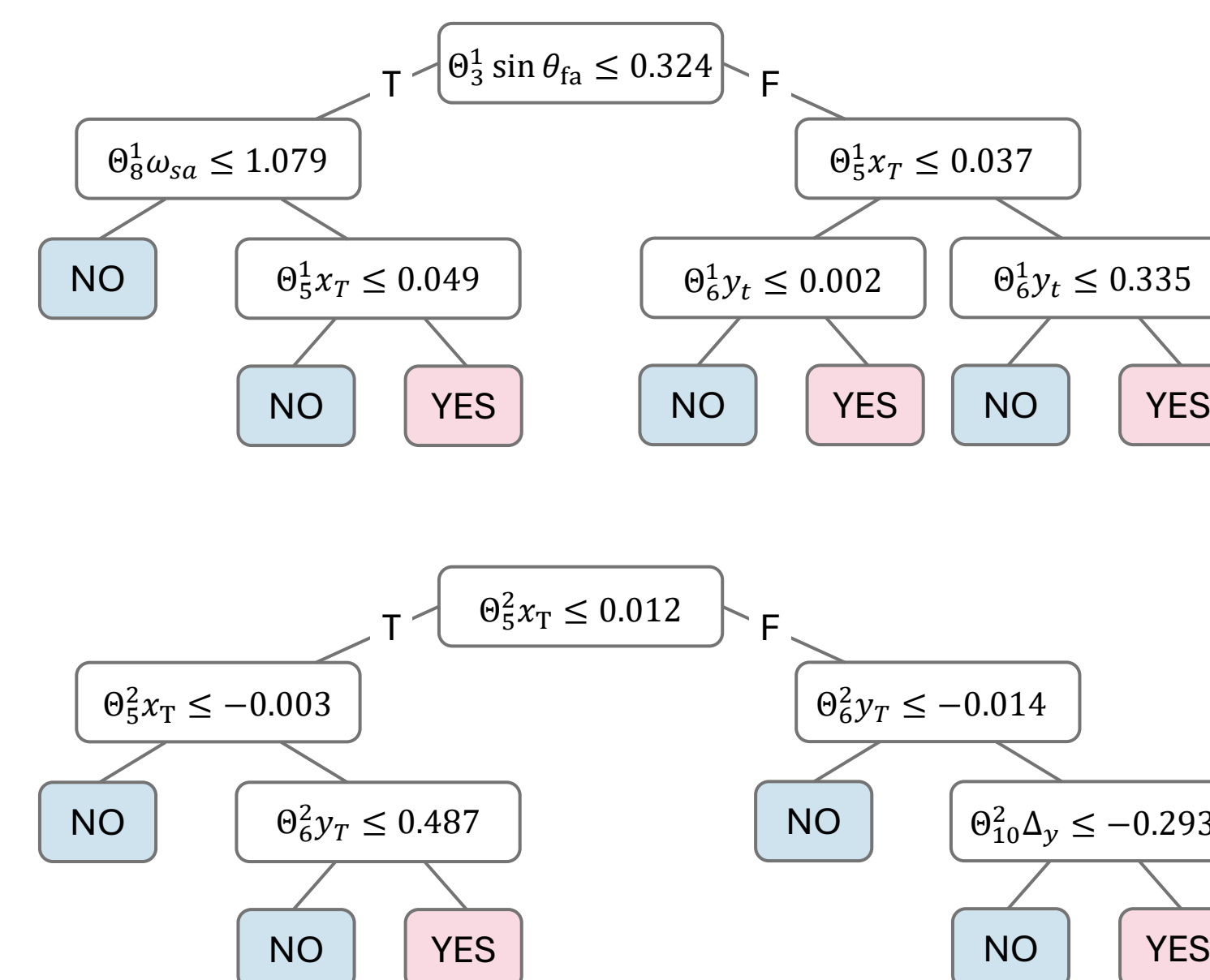
- Forced expert-space exploration

$$\varepsilon \sim \mathcal{N}\left(0, 1/M^2\right)$$

$$Load_m(S) = \sum_{s_k \in S} \mathbb{P}(\varepsilon_{new} \geq \tau(s_k) - \pi_m(s_k))$$

$$Imp_m(S) = \sum_{s_k \in S} softmax(\pi_m(s_k|\theta_m, \sigma_m))$$

#### Router distillation

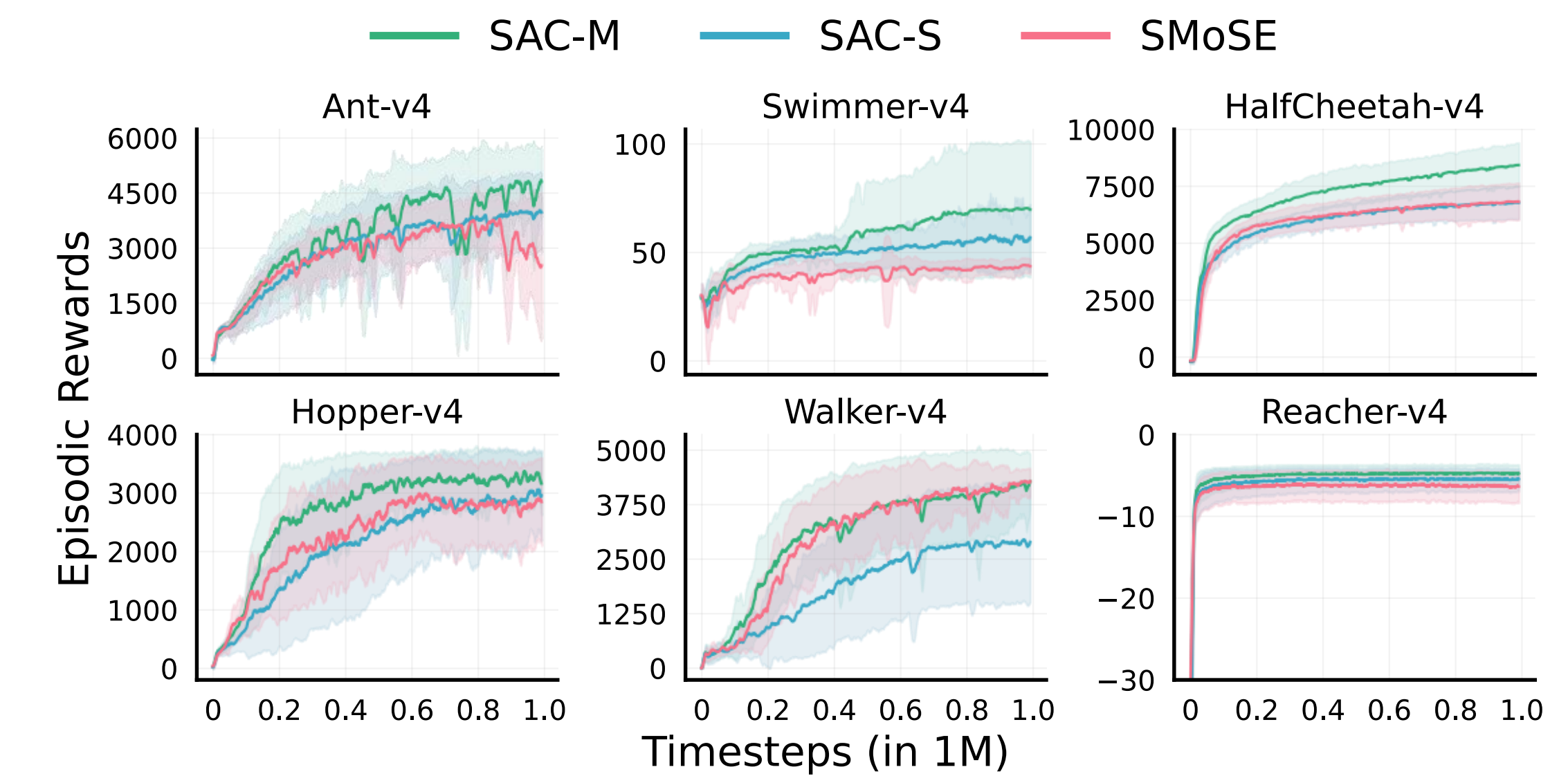- Per-expert binary decision tree for "free"



## Results

### Strong performance on Mujoco tasks

- Fully interpretable
- **2x performance** compared to SOTA interpretable approaches
- **99% less parameters** than SAC-L
- **Closes the gap** with close-box solutions

#### Training speed

SAC-M   SAC-S   SMoSE



Episodic Rewards

Timesteps (in 1M)

#### Episodic Reward on Mujoco

| | Walker2d | Hopper | Ant | HalfCheetah | Reacher | Swimmer |
|---|---|---|---|---|---|---|
| SAC-L | **4358.06** | 2636.49 | **5255.46** | **11809.87** | **-3.75** | 68.59 |
| SAC-M | 4020.51 | **3224.25** | 4894.18 | 8992.22 | -4.02 | 71.94 |
| SAC-S | 2967.14 | 3076.09 | 4162.97 | 7214.3 | -4.82 | 59.42 |
| PPO | 3362.16 | 2311.9 | 2327.12 | 2308.29 | -6.57 | 93.26 |
| CGP | 1090.00 | 1150.00 | 1130.00 | 6375.00 | -68.50 | **280.00** |
| LGP | 1080.00 | 1120.00 | 1210.00 | 6388.50 | -58.50 | 278.50 |
| Metric-40 | 775.00 | 2005.00 | 2210.50 | 2210.50 | x | x |
| Ours | 4224.29 | 2816.08 | 3245.43 | 7310.17 | -5.49 | 45.4 |