

# EDHuCATING MAZE-NAVIGATING AGENTS: HUMAN-IN-THE-LOOP EXPLORATION AND OPTIMIZATION

Kevin Godin-Dubois<sup>1</sup>, Karine Miras<sup>1</sup>, Anna Kononova<sup>2</sup>

<sup>1</sup> Vrije Universiteit Amsterdam, Amsterdam, Netherlands

<sup>2</sup> Leiden University, Leiden, Netherlands

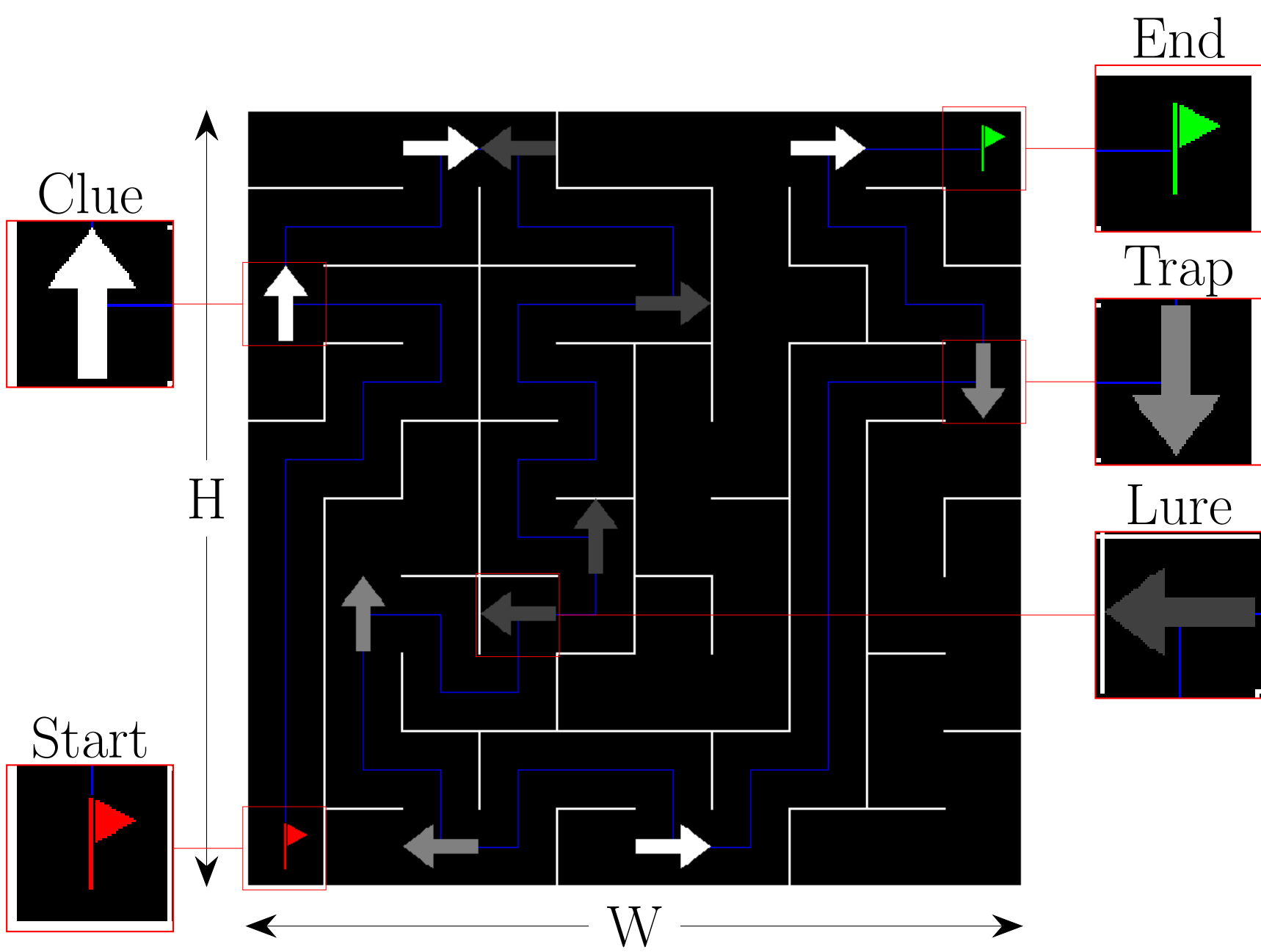
## Context

Light-weight testbed: maze navigation as a visual task with sparse rewards

Leveraging human creativity/reactivity with **Environment-Driven Human-Controlled Automated Training**

Loosely embodied robots with discrete local inputs and outputs

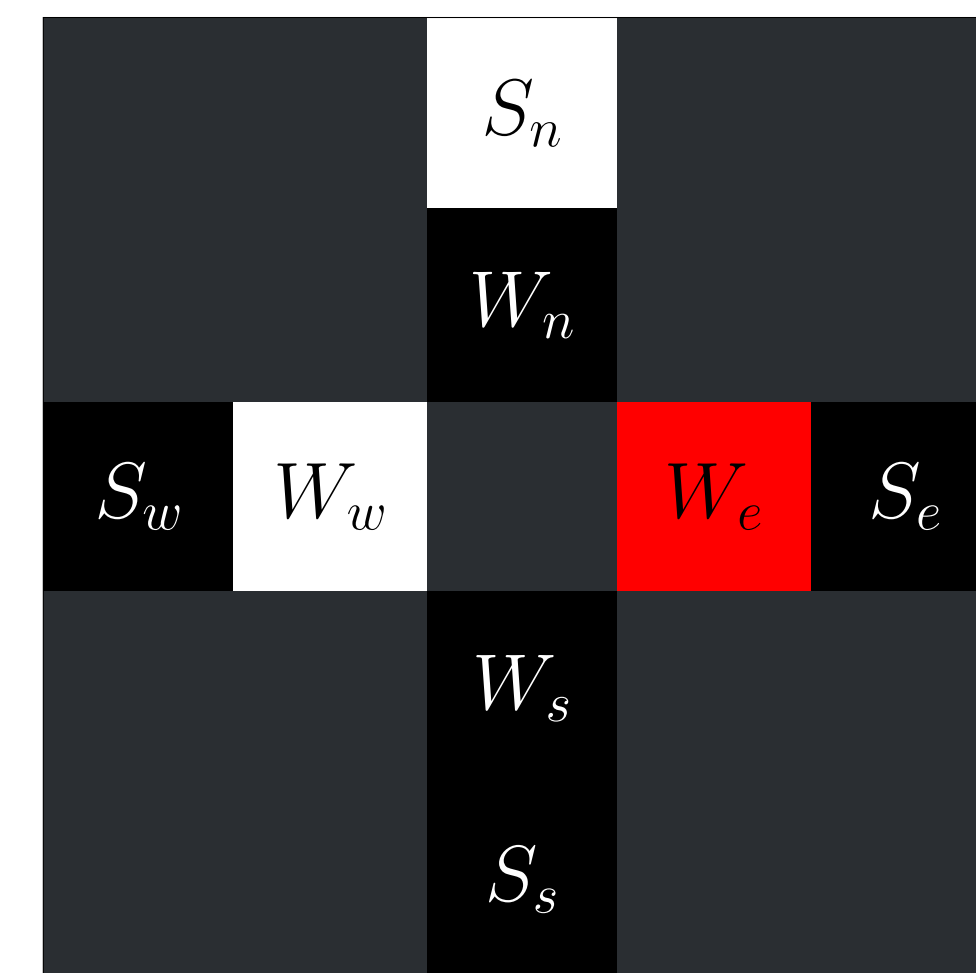
## 2D Mazes with visual information



Additional parameters:

- Random Number Generator seed
- Unicursive

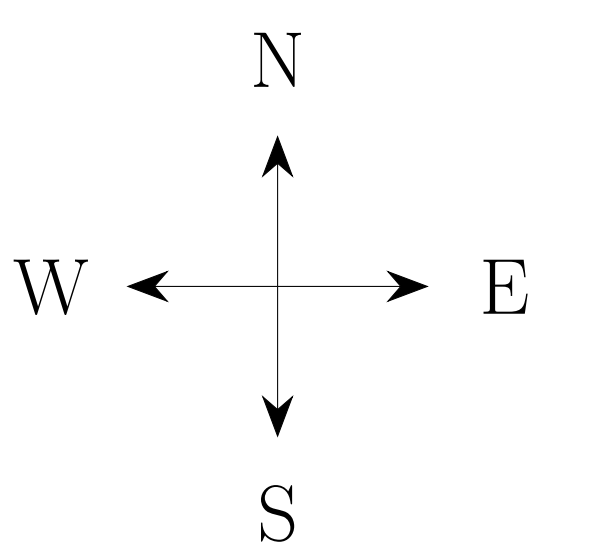
## Agents inputs/outputs



Inputs

$W_*$   $\begin{cases} \text{Wall if 1} \\ \text{Empty if 0} \\ \text{Origin if 0.5} \end{cases}$

$S_*$   $\begin{cases} \text{Sign direction} \\ \text{(if any)} \end{cases}$



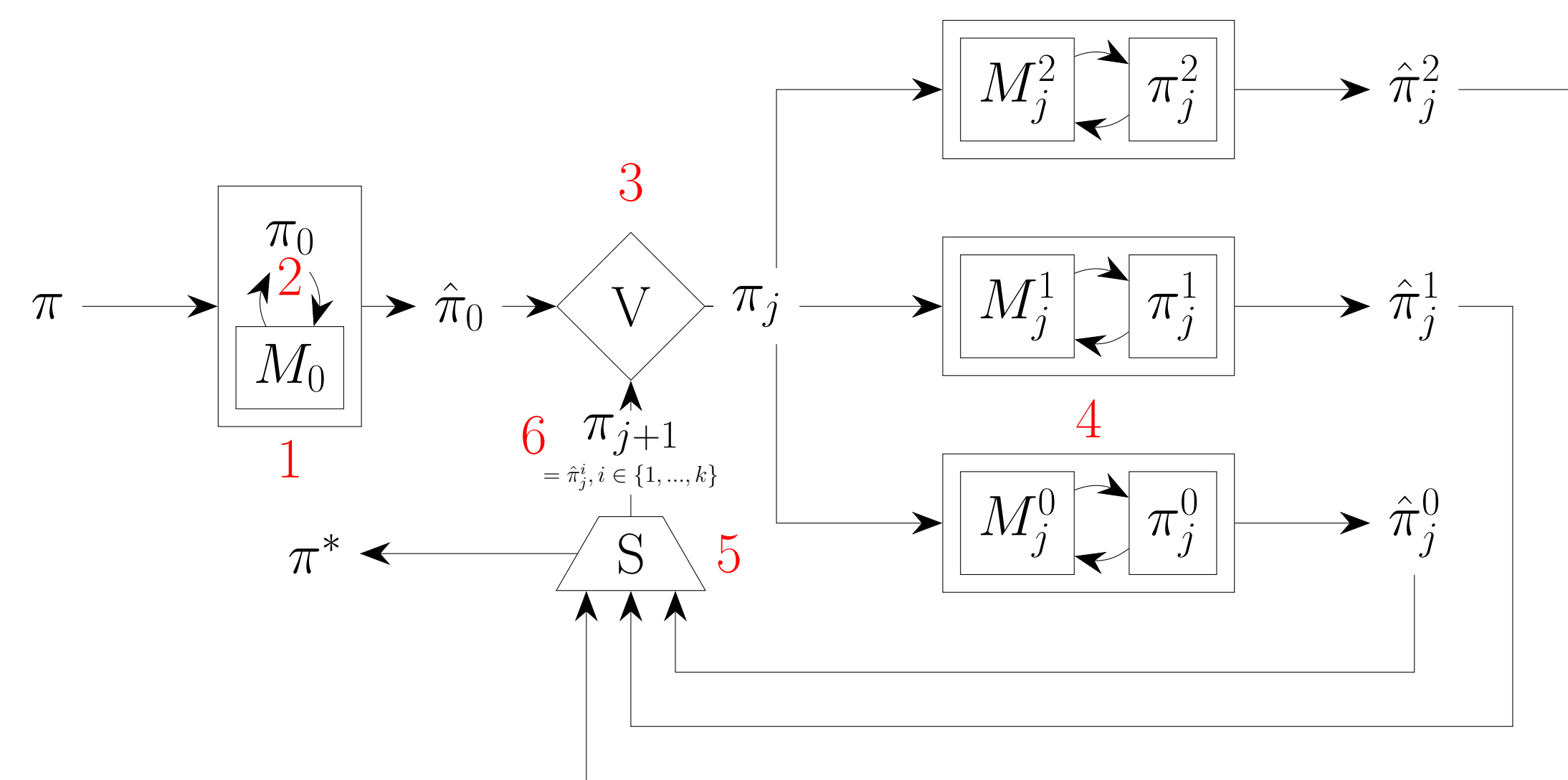
Outputs

## Training with varying levels of human involvement

### Framework

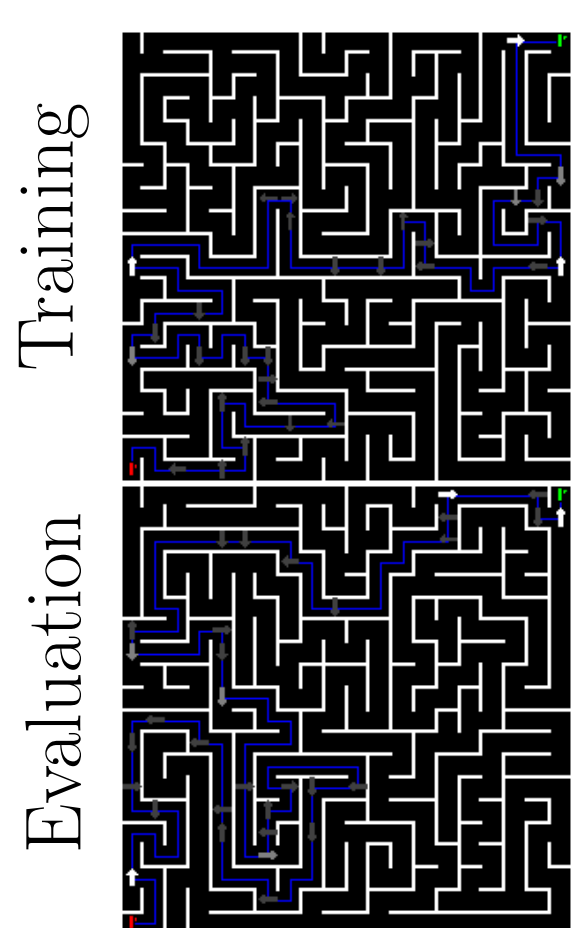
- Off-the-shelf Stable Baseline with maze environment
- Multi-Layer Perceptron control
- Training algorithms: A2C and PPO
- All parameters kept to defaults
- Total budget of 3'000'000 timesteps
- All mazes rotations
- Clues, Lures and Traps set to 1, 0.25 and 0.5, respectively
- Trained and evaluated on different (but similar) mazes

### EDHuCAT

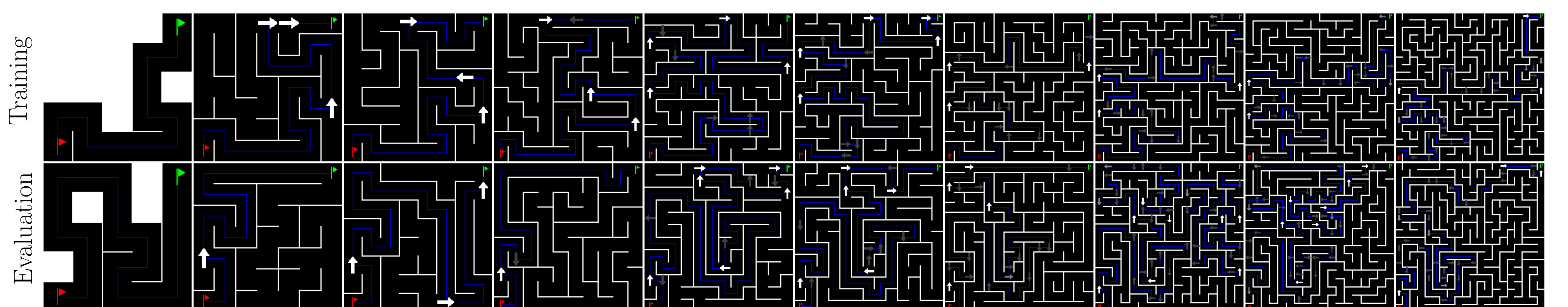


1. Trivial initial training/evaluations mazes
2. Small training period
3. Human-based generation of new mazes
4. Concurrent training
5. Human-based selection of "best" agent
6. Loop back to 3

### Direct

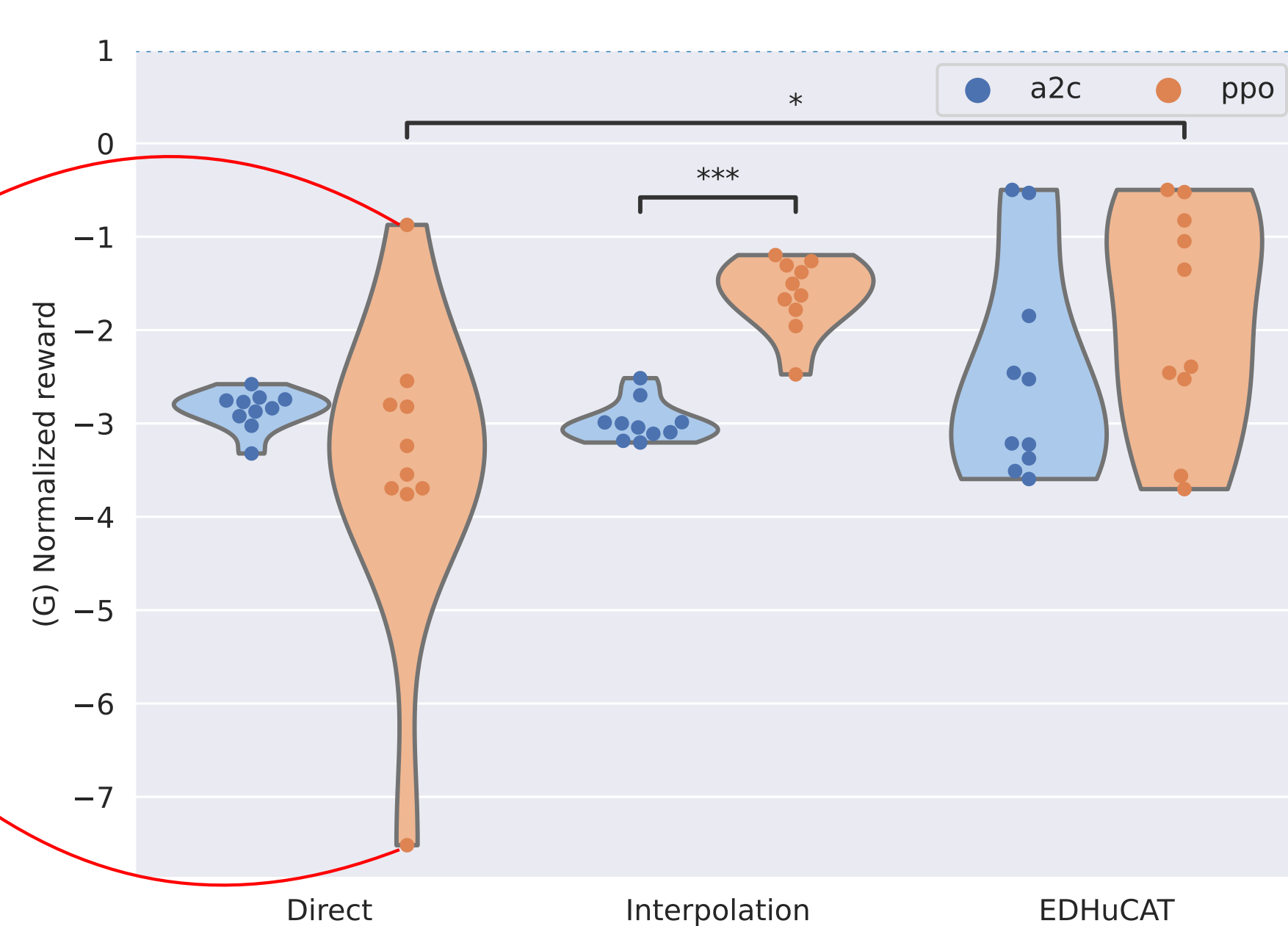
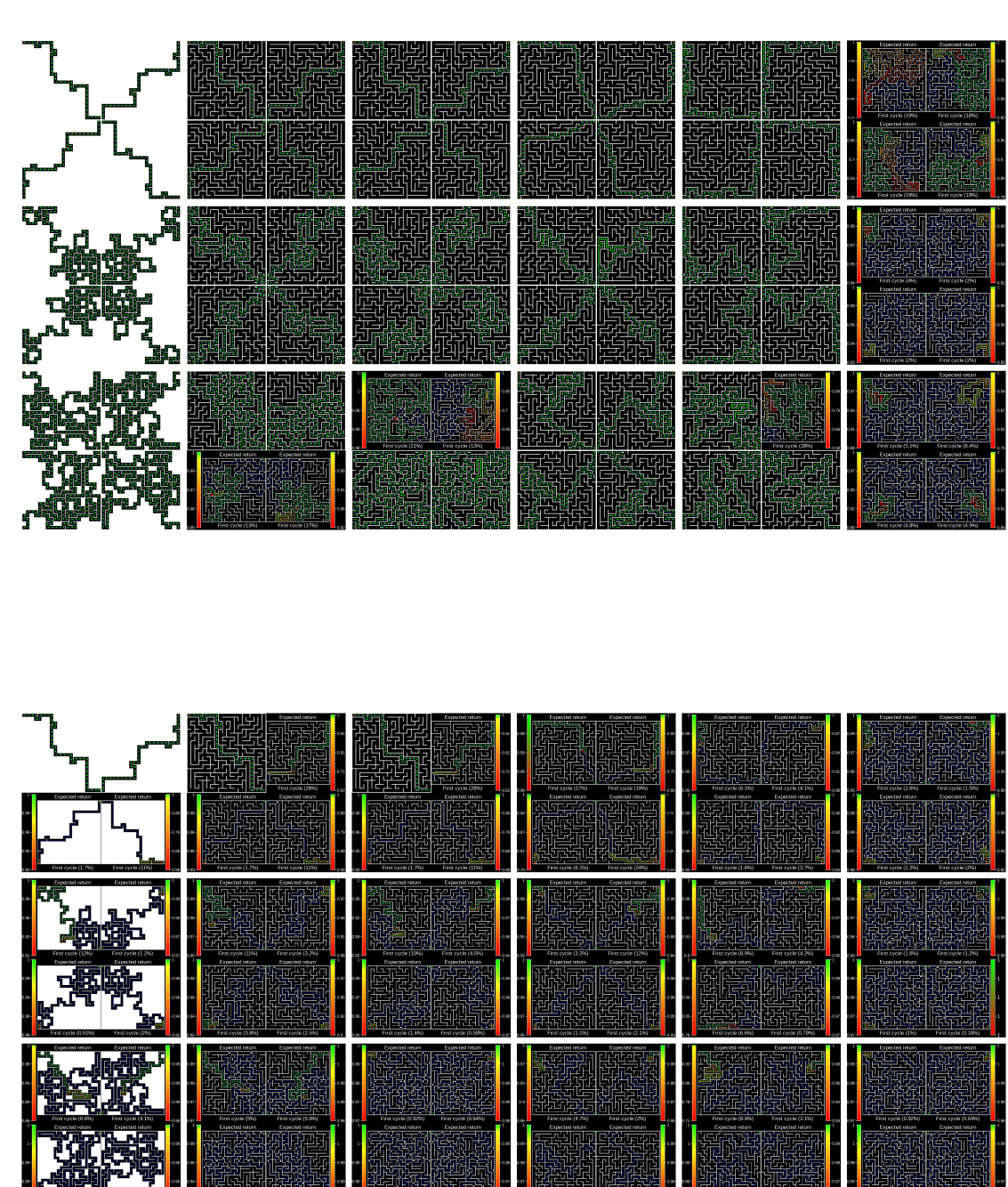


### Interpolation: "smooth" incremental training



## Generalization on unseen mazes and input combinations

### Maze-solving



### Inputs processing

