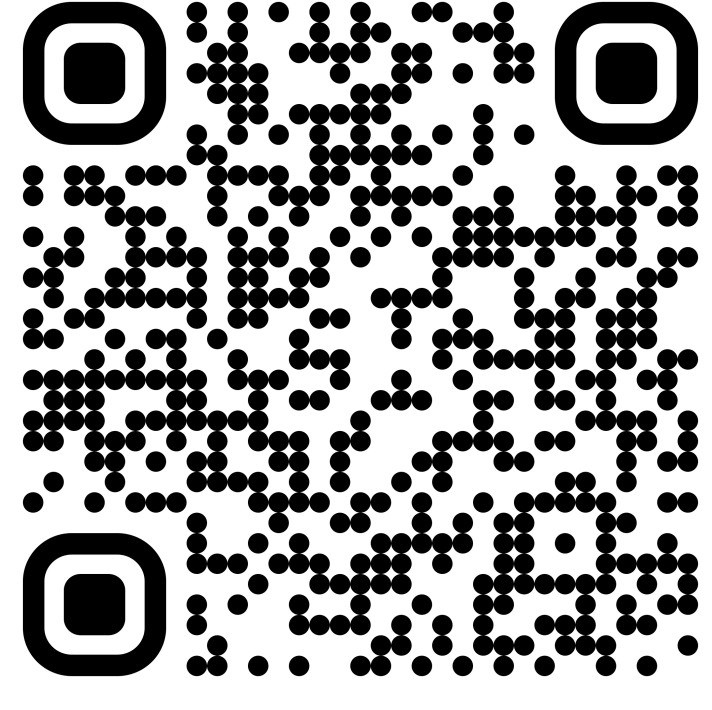


Towards Model Repair by Human Opinion-Guided Reinforcement Learning



Kyanna Dagenais

McMaster University – Hamilton, Canada
 dagenaik@mcmaster.ca | kyannadagenais.ca

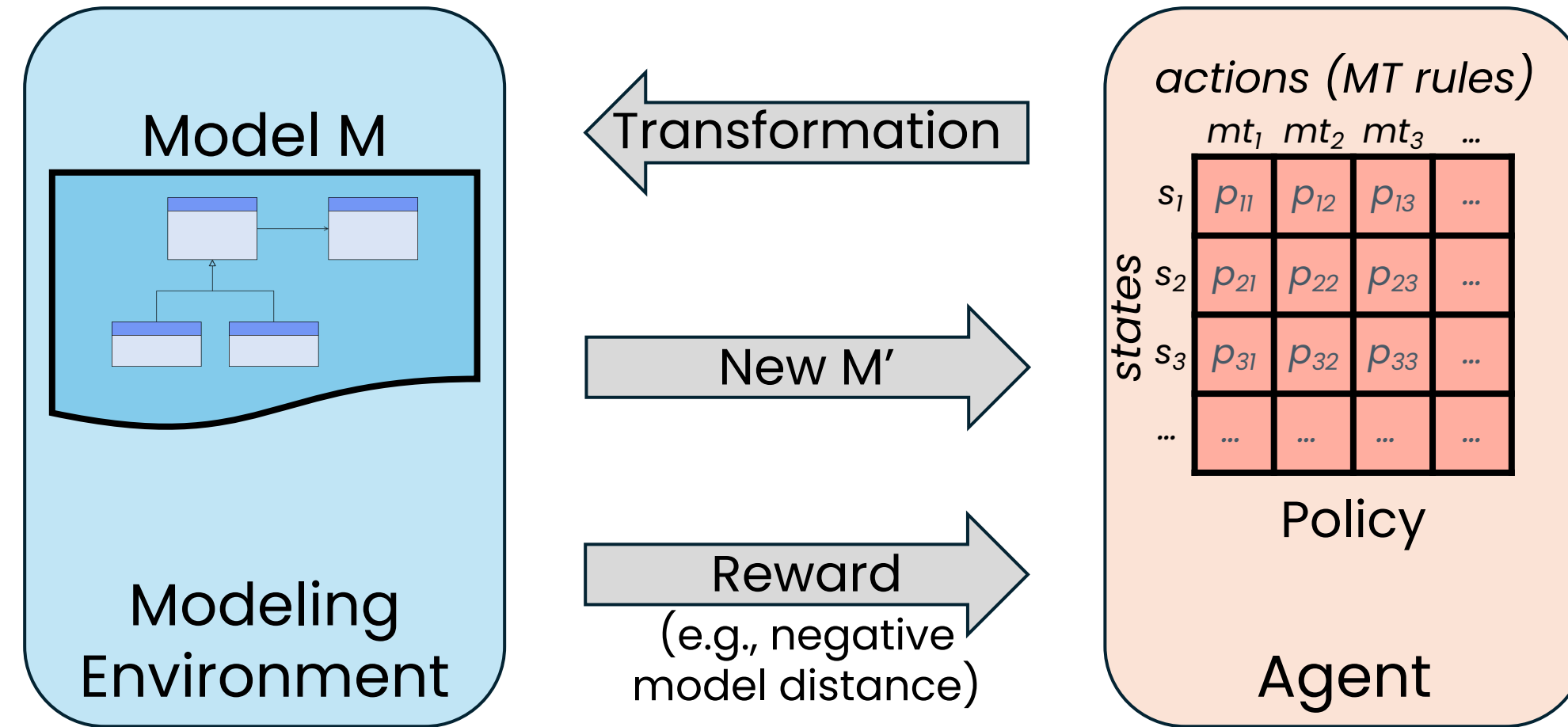


Sustainable Systems and Methods

Context and Research Problem

Model repair

- Goal**
- ▶ Repair invalid models by model transformations
- Problem**
- ▶ Complex models → long repair sequences
 - ▶ Automation is needed
- Opportunity**
- ▶ Modeling projects are longitudinally extensive
 - ▶ Learn repair patterns as we go

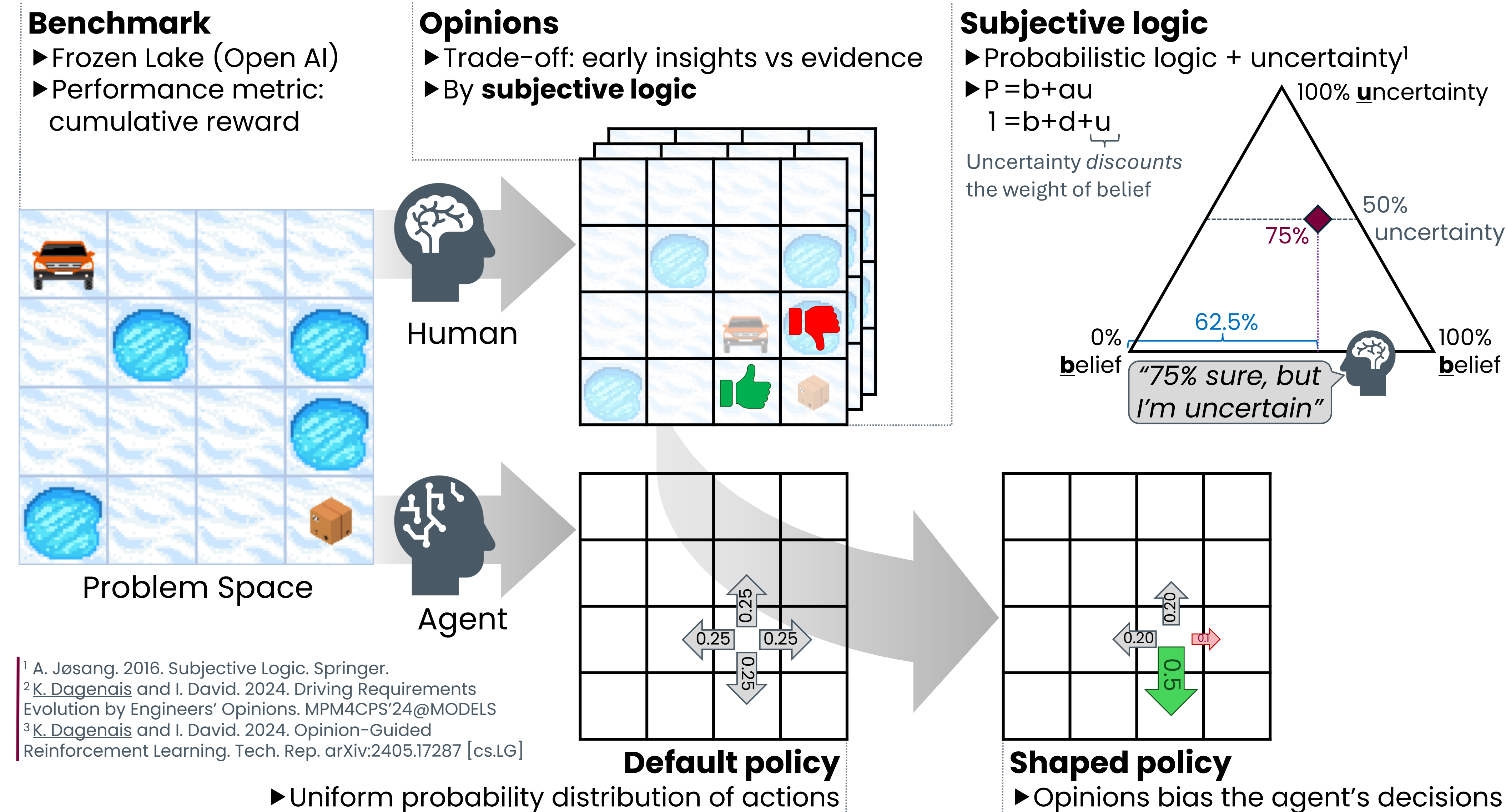


Reinforcement learning (RL)

- Basics**
- ▶ Agent learns by trial and error
 - ▶ **Pro:** does not require historical data
 - ▶ Policy: state to MT-rule mapping (probability)
- Limitation**
- ▶ **Shallow learning curve** (learning takes time)
- Goal**
- ▶ Guide RL by rapidly emerging (uncertain) opinions

Methods

Approach



Mapping between RL and MDE

RL	MDE
Frozen lake	Design space
Agent	Current model state
Action (step)	Model transformation
(0,0)	Invalid model
Goal state	Valid model
Terminal state	Model beyond repair

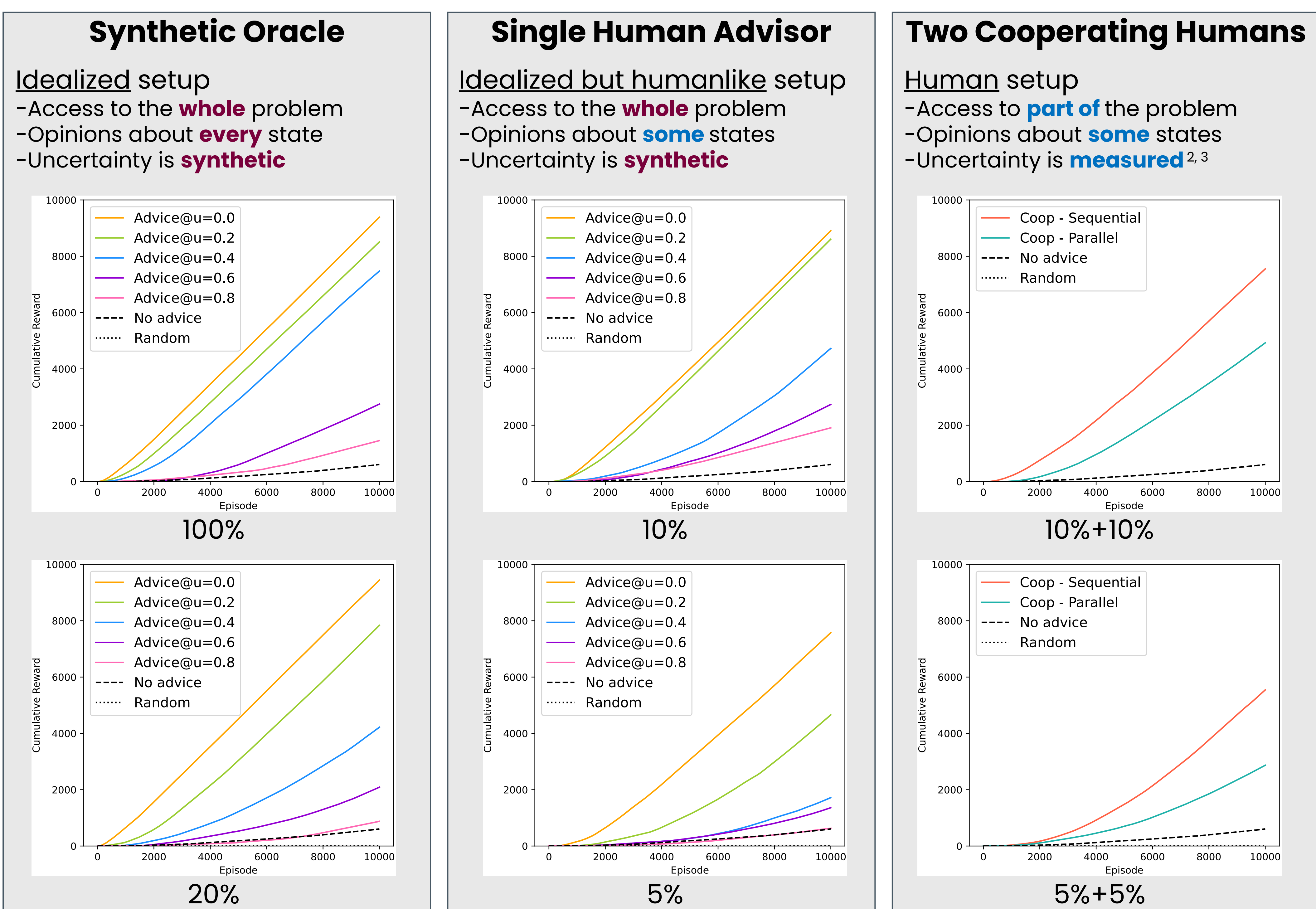
Evaluation

- ▶ 12x12 version of the Frozen Lake
 - ▶ Start → goal: ≥22 steps
 - ▶ 20% of states marked terminal
- ▶ Agent
 - ▶ Algorithm: discrete policy gradient
 - ▶ Learning rate: 0.9
 - ▶ Discount factor: 1.0
- ▶ Learning on 10 000 episodes
- ▶ Reward model
 - ▶ Terminal state: reward = 0
 - ▶ Goal state: reward = +1

Results and takeaways

Results

*% denotes opinion quota, i.e., the number of opinions compared to the number of states



Takeaways

All opinions, even if uncertain, can be of **high utility**.

In the charts, cumulative reward tends to be significantly higher than "No advice" and "Random".

A single human advisor is **as effective** as a synthetic oracle.

In the "Single Human Advisor" charts, cumulative reward tends to be similar to that of the "Synthetic Oracle" charts.

A single human advisor can be **more efficient** than a synthetic oracle.

In the "Single Human Advisor" charts, the same high reward is obtained at lower advice quotas (only 10% and 5% vs 100% and 20%).

Real cooperating humans' performance is **comparable** to that of synthetic advisors.

In the "Two Cooperating Humans" charts, cumulative reward tends to be similar to the other two cases.