

Learner-aware Teaching: Inverse Reinforcement Learning with Preferences and Constraints

Sebastian Tschiatschek^{1*}, Ahana Ghosh^{2*}, Luis Haug^{3*}, Rati Devidze², Adish Singla²

¹Microsoft Research, ²ETH Zurich, ³Max Planck Institute for Software Systems, *equal contribution

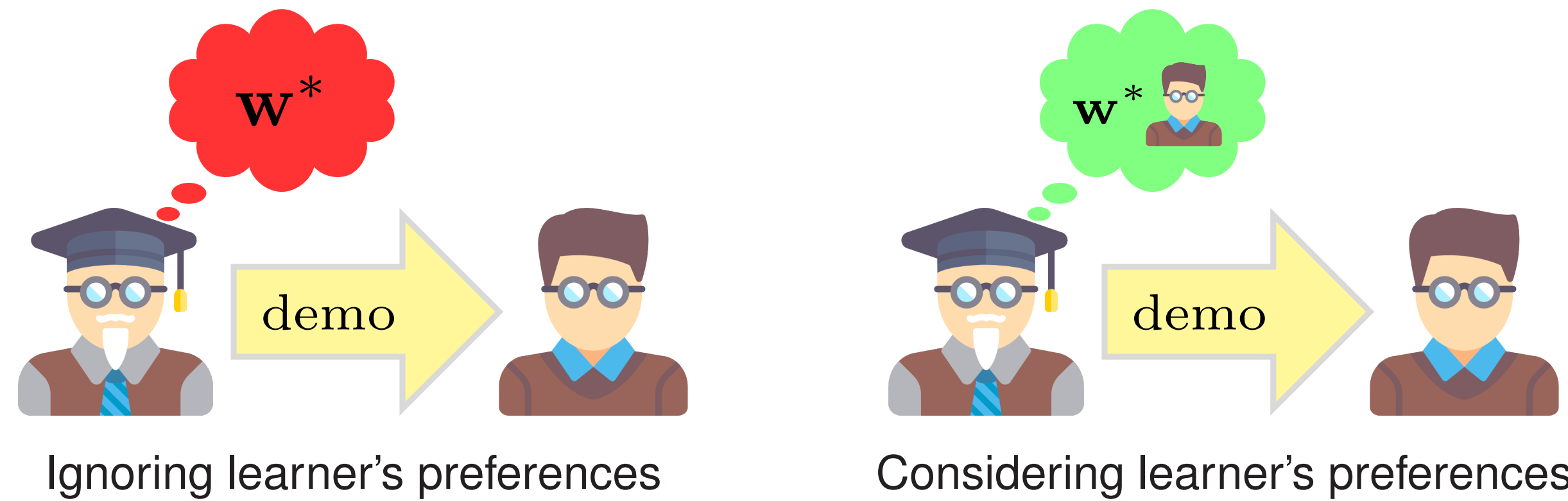
Highlights

Problem setting

- Teaching a *learner with preferences* via demonstrations
- Studied two teaching strategies

Learner-agnostic teaching

Learner-aware teaching



Main results

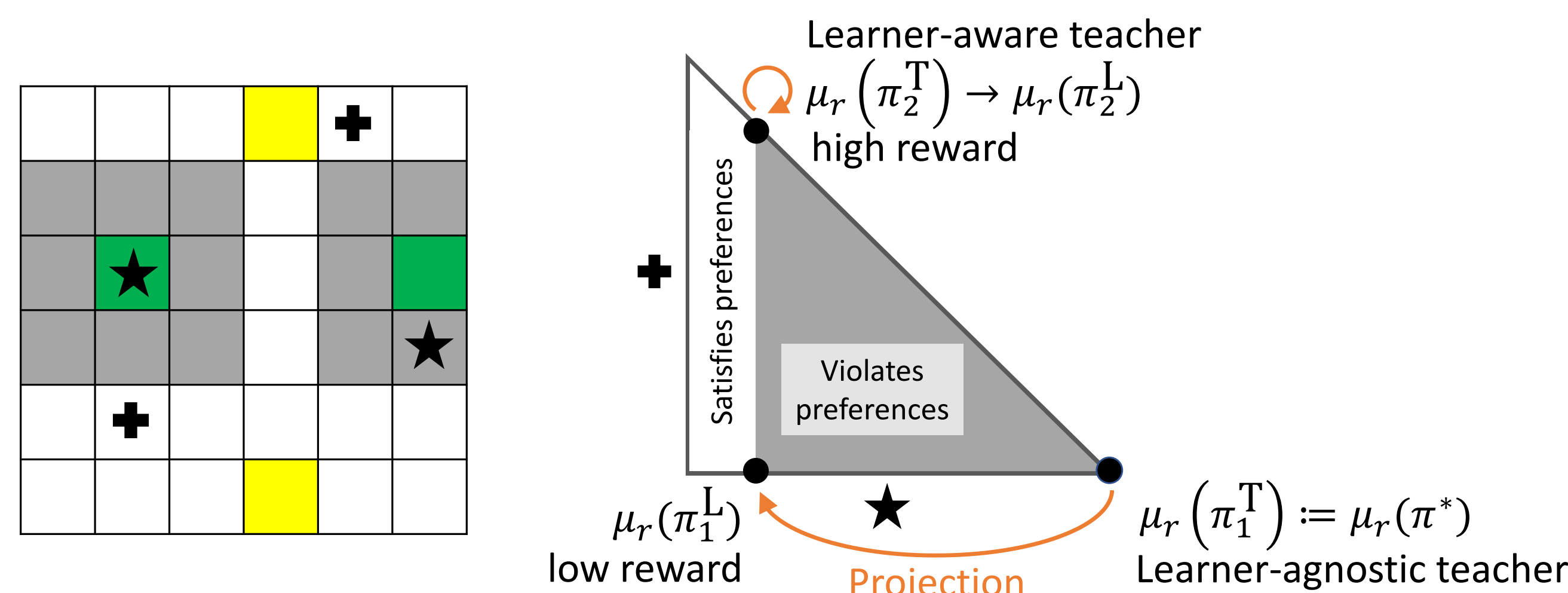
- Learner-agnostic teaching can be arbitrarily bad
- New algorithms for learner-aware teaching achieving high performance

A Teacher and an IRL Learner Without Preferences

- MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, \gamma, R)$ with rewards $R(s) = \langle w^*, \phi_r(s) \rangle$
- Teacher T provides demonstrations using policy π^T in \mathcal{M}
- Policy π has reward $R(\pi) = \langle w^*, \mu_r(\pi) \rangle$, where $\mu_r(\pi) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t \phi_r(s_t) \mid \pi]$
- Learner receives demonstrations and outputs π^L s.t. $\|\mu_r(\pi^L) - \mu_r(\pi^T)\| \leq \epsilon$
- This ensures that $R(\pi^L) \geq R(\pi^T) - \epsilon$

Challenges in Teaching Learners With Preferences

- Object-world gathering game
- \star yields reward 1.0, $+$ yields reward 0.9
- Learner's preference: Avoid frequent proximity of green cells (≤ 1 -cell distance)



- Providing demonstrations from optimal behavioral policy π^* can lead to arbitrarily bad learner's performance!

Learner Models

- Learner's preferences are captured by features $\phi_c(s)$
- Formalized as constraints on $\mu_c(\pi)$, where $\mu_c(\pi) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t \phi_c(s_t) \mid \pi]$

Standard maximum causal entropy IRL learner

$$\begin{aligned} \max_{\pi} \quad & H(\pi) \\ \text{s.t.} \quad & \|\mu_r(\pi) - \mu_r(\pi^T)\| = 0 \end{aligned} \quad \begin{array}{l} \text{causal entropy} \\ \text{feature matching} \end{array}$$

Learner with hard preferences

$$\begin{aligned} \min_{\pi} \quad & \|\mu_r(\pi) - \mu_r(\pi^T)\| \\ \text{s.t.} \quad & g(\mu_c(\pi)) \leq 0 \end{aligned} \quad \begin{array}{l} \text{feature matching} \\ \text{hard preferences} \end{array}$$

Learner with soft preferences

$$\begin{aligned} \max_{\pi, \delta_r^{\text{soft}}, \delta_c^{\text{soft}}} \quad & H(\pi) - C_r \|\delta_r^{\text{soft}}\|_p - C_c \|\delta_c^{\text{soft}}\|_p \\ \text{s.t.} \quad & \|\mu_r(\pi) - \mu_r(\pi^T)\| \leq \delta_r^{\text{soft}} \\ & g(\mu_c(\pi)) \leq \delta_c^{\text{hard}} + \delta_c^{\text{soft}} \end{aligned} \quad \begin{array}{l} \text{feature matching} \\ \text{hard+soft preferences} \end{array}$$

- Learner trades-off reward-feature matching and its own preferences

Learner-Aware Teaching for Known Constraints

Learner-aware teaching for hard preferences: AWARE-CMDP

- Define a set of feasible reward feature expectations $\Omega_r^L = \{\mu_r(\pi) \mid g(\mu_c(\pi)) \leq 0\}$
- Optimal teaching policy = solution of constrained MDP:

$$\max_{\pi^T} \langle w^*, \mu_r(\pi^T) \rangle \quad \text{s.t.} \quad \mu_r(\pi^T) \in \Omega_r^L$$

- Theorem.* The value of learner-aware teaching can be arbitrarily high, given by

$$\max_{\pi \text{ s.t. } \mu_r(\pi) \in \Omega_r^L} \langle w^*, \mu_r(\pi) \rangle - \langle w^*, \text{Proj}_{\Omega_r^L}(\mu_r(\pi^*)) \rangle$$

- For linear $g(\cdot)$, the above problem can be solved via linear programming

Learner-aware teaching for soft preferences: AWARE-BIL

- Optimal teaching problem can be formulated as a bi-level optimization:

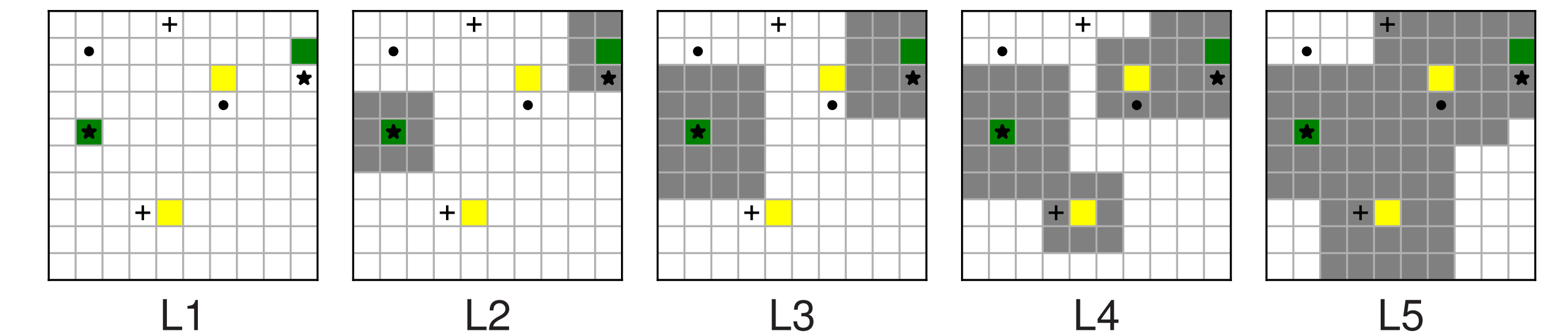
$$\max_{\pi^T} \langle w^*, \mu_r(\pi^T) \rangle \quad \text{s.t.} \quad \pi^L \in \arg \max_{\pi} \text{IRL}(\pi, \mu(\pi^T))$$

- Here $\text{IRL}(\pi, \mu(\pi^T))$ stands for the IRL problem solved by the learner
- Optimal teaching policy is a softmax policy satisfying the learner's constraints
- A challenging non-convex optimization problem
- Proposed a gradient-based optimization approach

Experimental Results

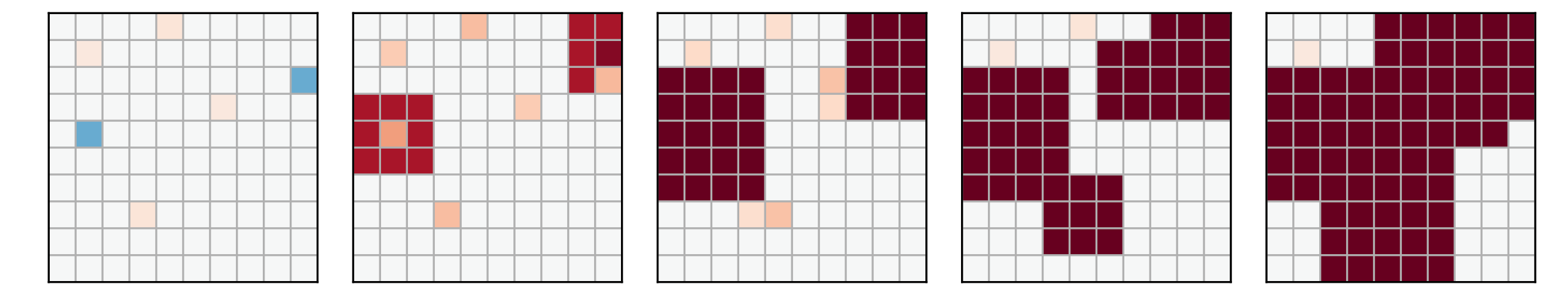
Experimental setup

- Object-world gathering environment:
 - Rewards: \star yields 1.0, $+$ yields 0.9, \bullet yields 0.2
 - Two “green” distractors at 0-cell and 1-cell distance to the \star objects
 - Two “yellow” distractors at 1-cell and 2-cell distance to the $+$ objects
 - Discount factor $\gamma = 0.99$
- Learners with soft preferences ($C_r = 5$, $C_c = 10$) and $\delta_c^{\text{hard}} = 0$
- Environment and learners' preferences for 5 different learners $L1, \dots, L5$

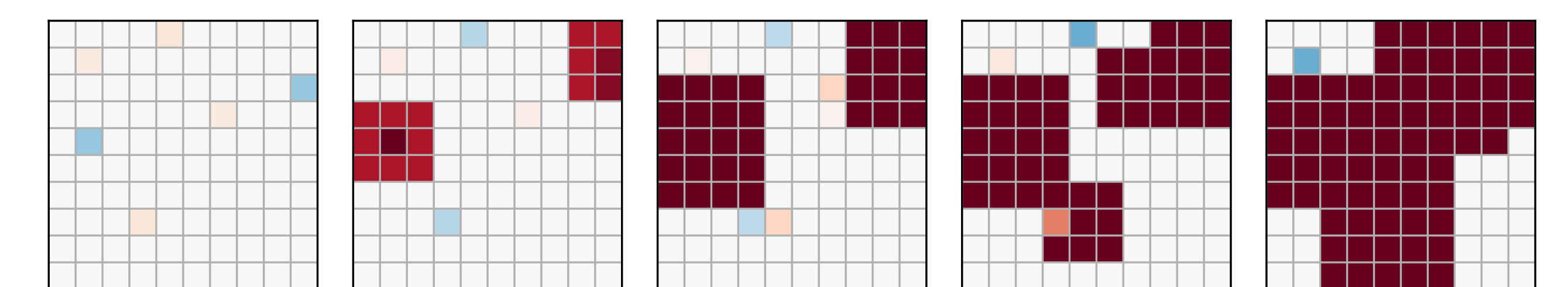


- For instance, L2 has two preference features indicating whether there is a green cell at a distance of 0-cells or 1-cell, respectively

Learner-aware teaching for known constraints



Learners' rewards inferred from learner-agnostic teacher (AGNOSTIC)



Learners' rewards inferred from learner-aware teacher (AWARE-BIL)

Teacher	L1	L2	L3	L4	L5
AGNOSTIC	7.99 ± 0.02	0.01 ± 0.00	0.01 ± 0.00	0.01 ± 0.00	0.00 ± 0.00
AWARE-BIL	8.00 ± 0.02	7.20 ± 0.01	4.86 ± 0.30	3.15 ± 0.27	1.30 ± 0.07

Further Results

- Algorithms for learner-aware teaching with unknown constraints
- Additional experimental results
- Formal statements, proofs, and derivations

