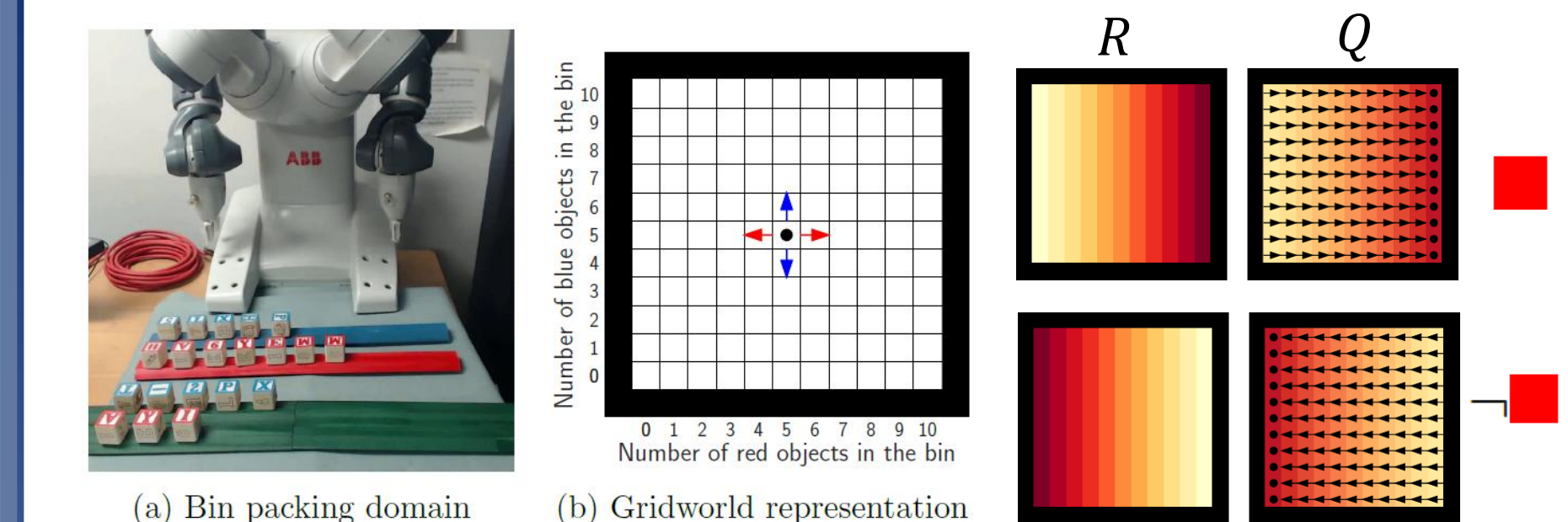


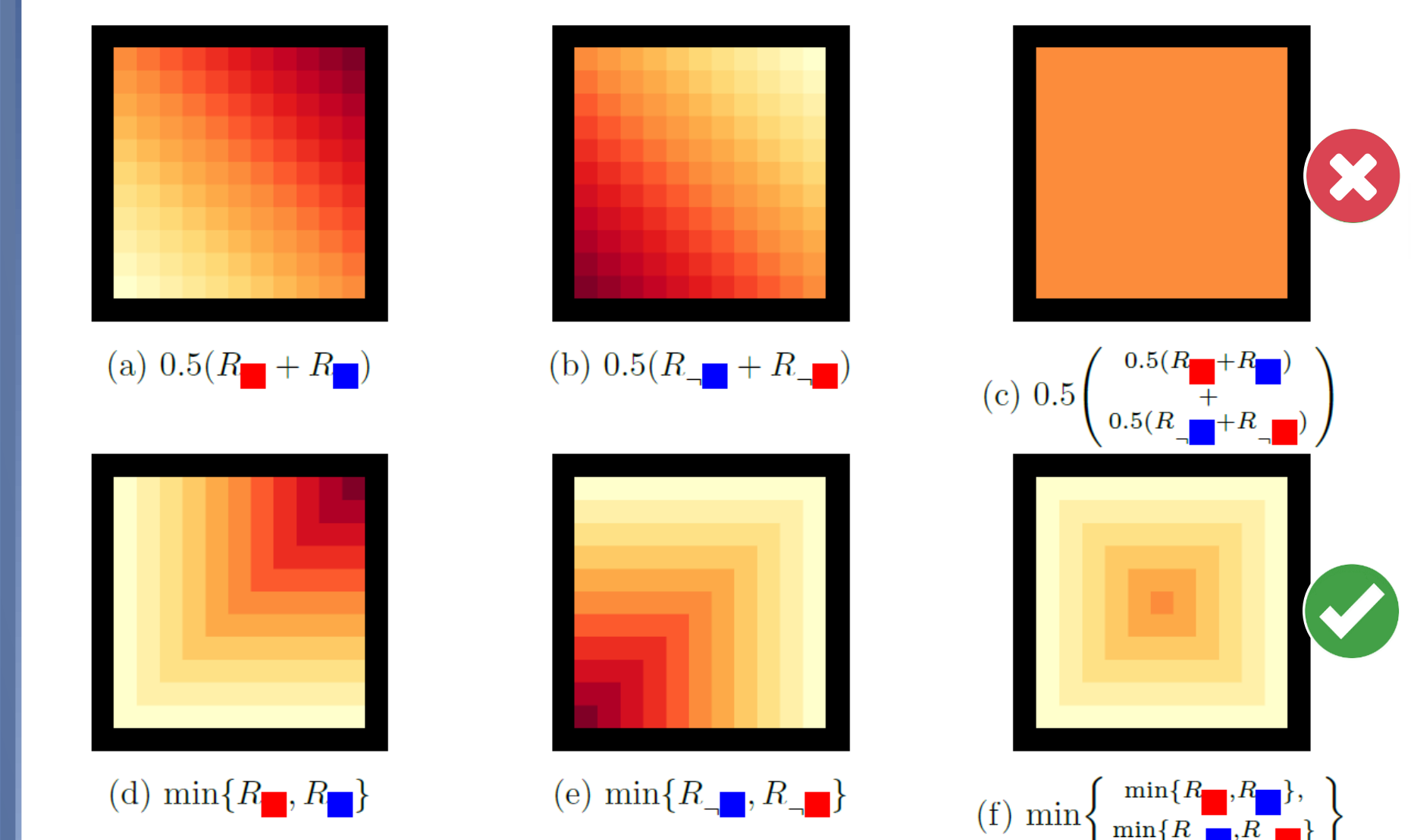
We can provably: specify tasks with arbitrary rewards using logics, construct a basis for arbitrary rewards bounds, and solve tasks zero-shot by leveraging World Value Functions.

Motivation

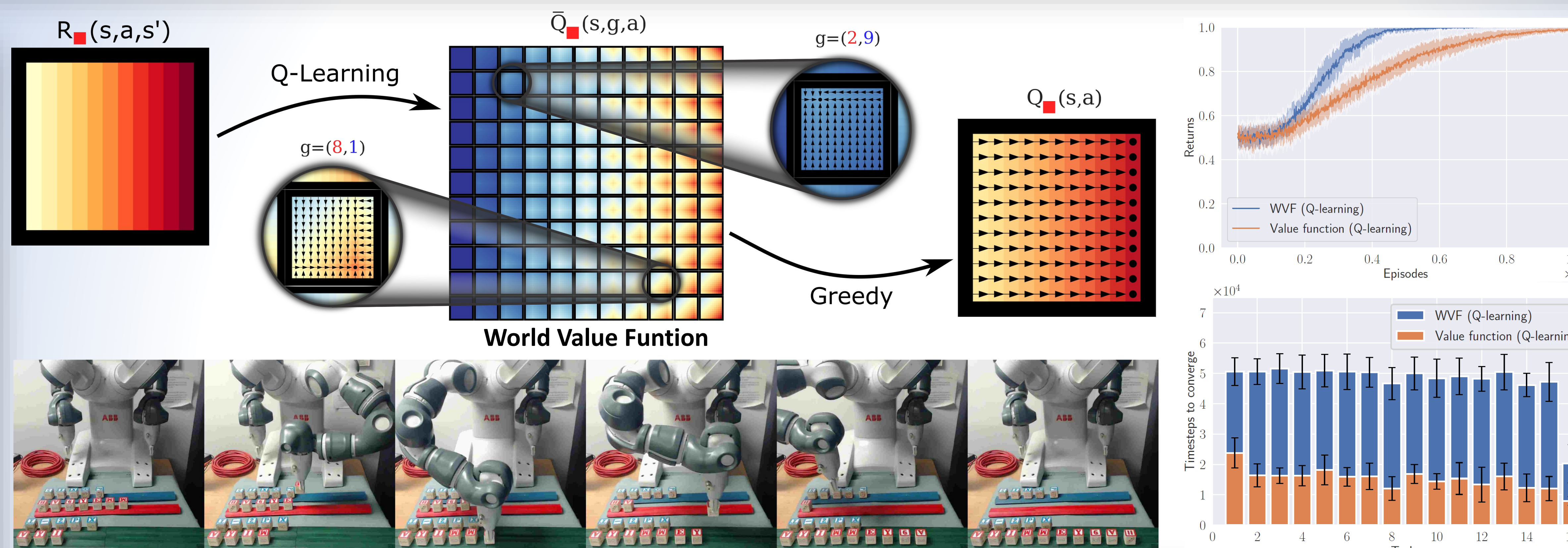
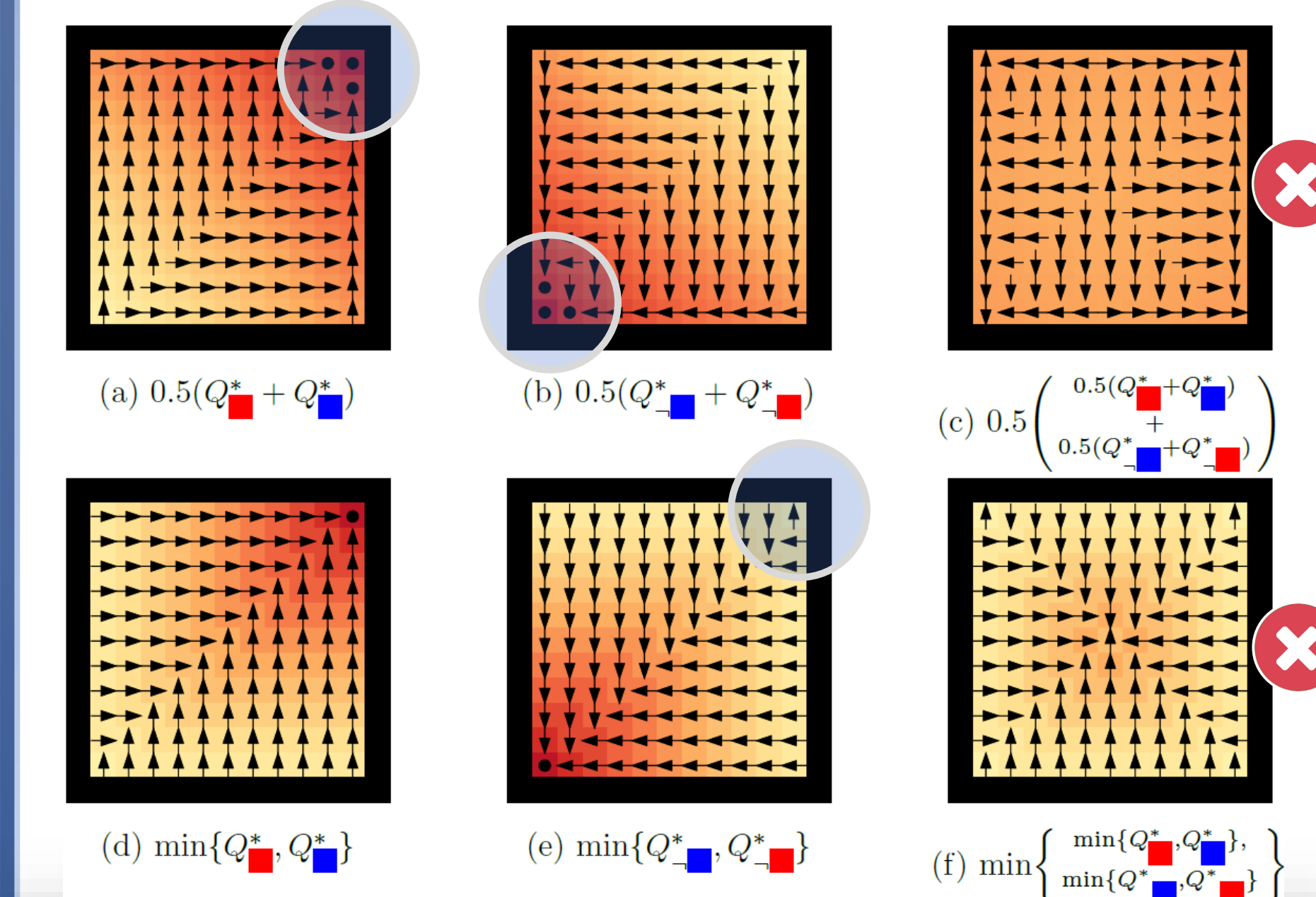
We want **instructable agents** (e.g. via language) that can **solve tasks beyond Boolean rewards** [1,2], and **generalise compositionally** new tasks.



Task/Reward specification is hard!
Arbitrary reward compositions maybe not result in the desired tasks



Skill/Value function/Policy learning is hard!
Arbitrary skill compositions maybe not result in the desired solutions



Task (Rewards) composition

Logic Operators

- Disjunction (OR):** $A \vee B := \max\{R_A(s, a, s'), R_B(s, a, s')\}$
- Conjunction (AND):** $A \wedge B := \min\{R_A(s, a, s'), R_B(s, a, s')\}$
- Negation (NOT):** $\neg A := (R_{MAX}(s, a, s') + R_{MIN}(s, a, s')) - R_A(s, a, s')$

Task space bounds $\{R_{MAX}, R_{MIN}\}$ change semantics:

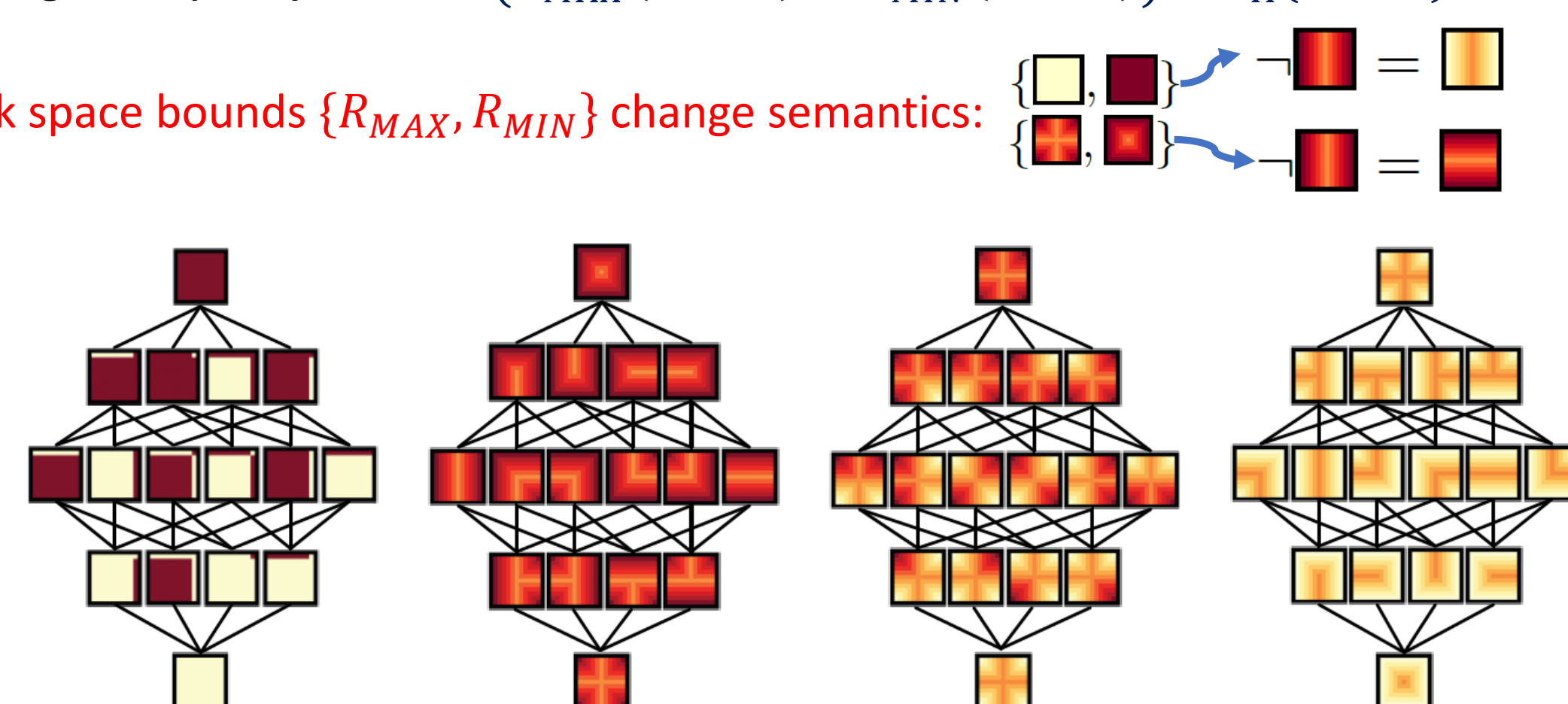


Figure 9: Examples of Boolean task sub-algebra with basis $\{0, 1\}$, $\{0.5, 0.5\}$ and task space bounds $\{0, 1\}$, $\{0.5, 0.5\}$ respectively.

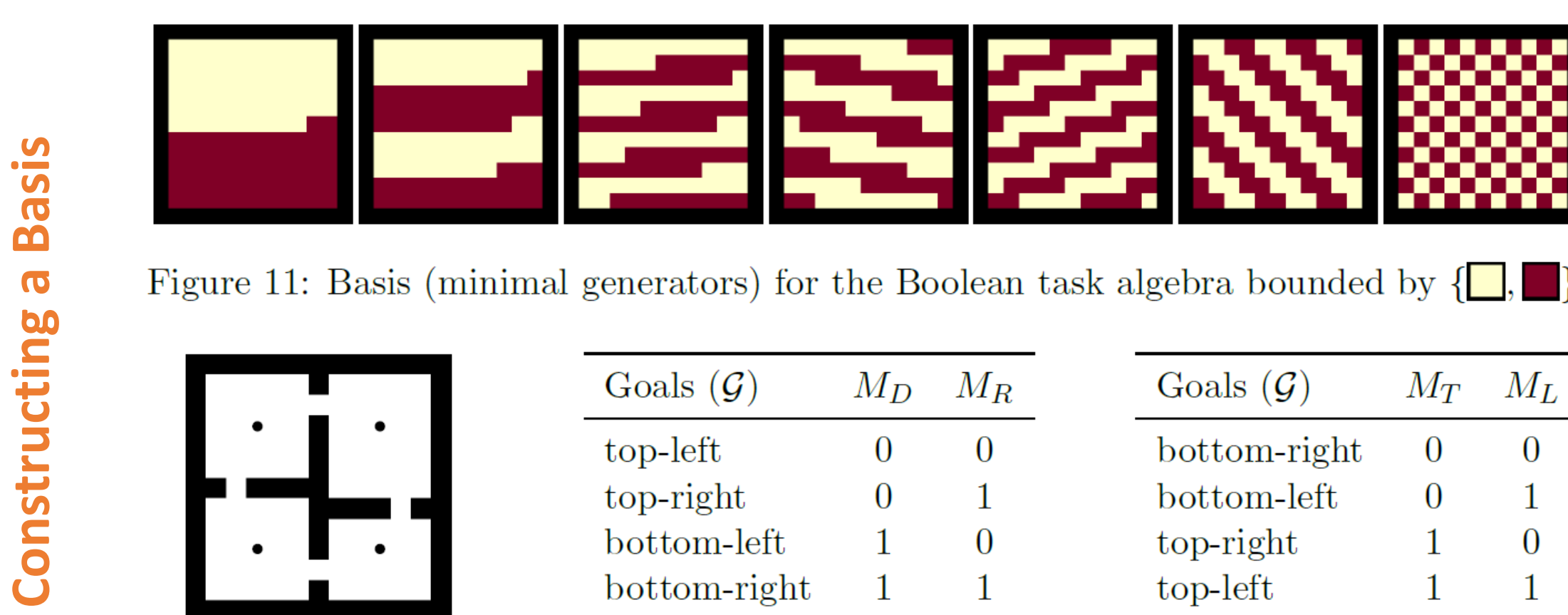


Figure 11: Basis (minimal generators) for the Boolean task algebra bounded by $\{0, 1\}$. (a) Four Rooms domain. The circles indicate goals the agent must reach (G). (b) Goals labeled by the well order \leq given by: top-left \leq top-right \leq bottom-left \leq bottom-right. (c) Goals labeled by the well order \leq given by: bottom-right \leq bottom-left \leq top-right \leq top-left.

Skill (WVF) composition

Logic Operators

- Disjunction (OR):** $Q_A \vee Q_B := \max\{Q_A(s, a, g), Q_B(s, a, g)\}$
- Conjunction (AND):** $Q_A \wedge Q_B := \min\{Q_A(s, a, g), Q_B(s, a, g)\}$
- Negation (NOT):** $\neg Q_A := (Q_{MAX}(s, a, g) + Q_{MIN}(s, a, g)) - Q_A(s, a, g)$

Skill space bounds $\{Q_{MAX}, Q_{MIN}\}$ similarly change semantics.

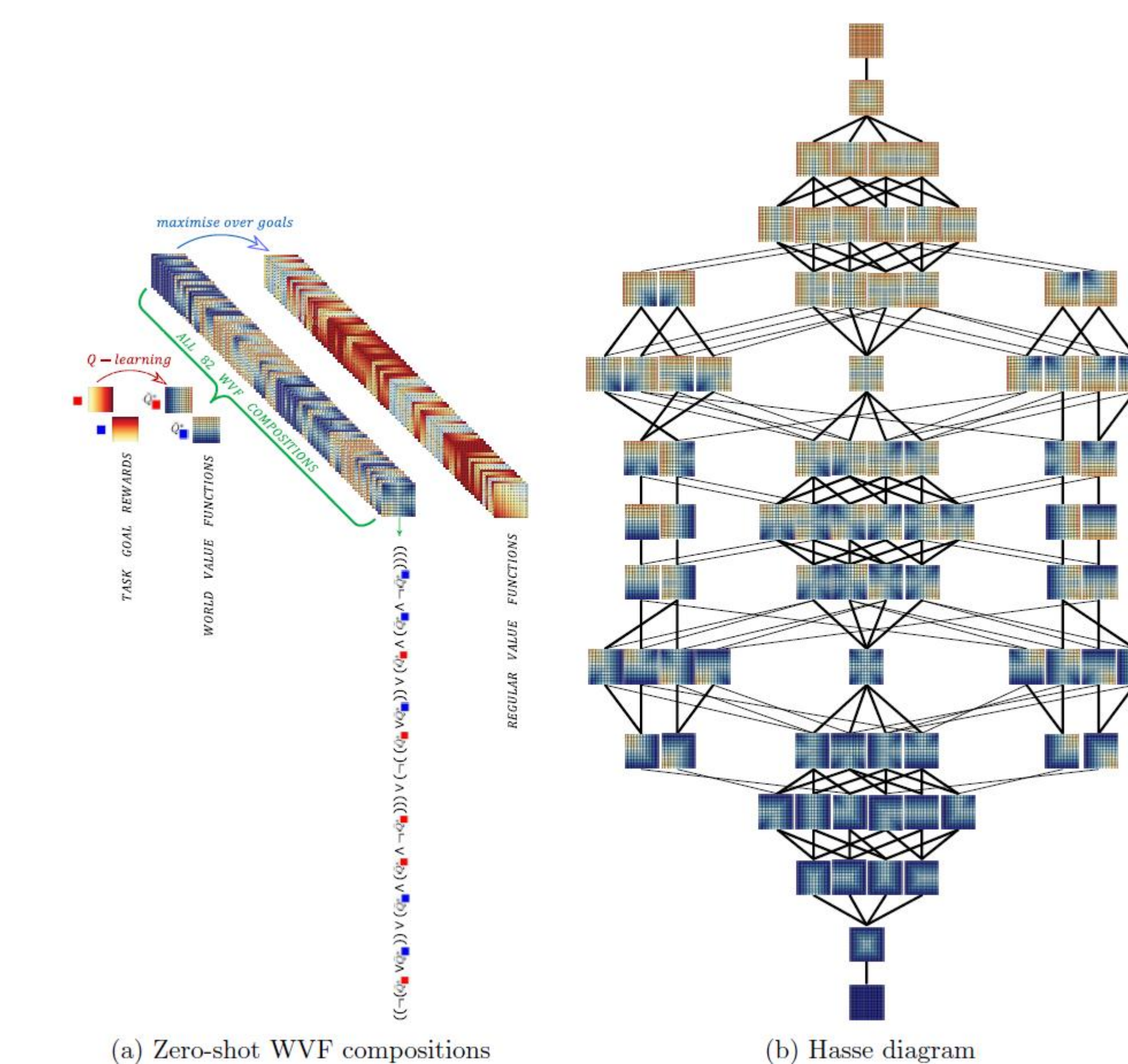
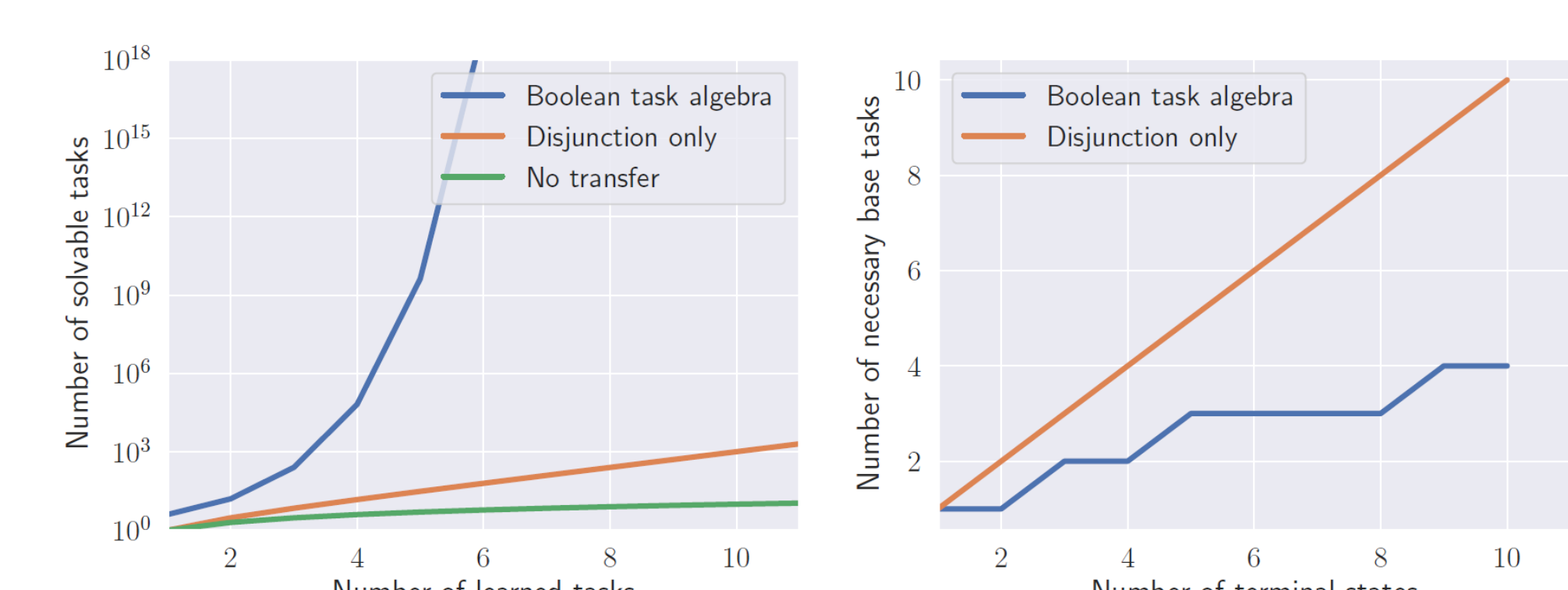


Figure 12: Skill (WVF) composition. (a) Zero-shot WVF compositions. (b) Hasse diagram. Illustration of the De Morgan sub-lattice generated by composing \bar{Q}_L^* and \bar{Q}_T^* .

Blessing of Dimensionality

Exponential explosion of skills from a logarithmic amount of learning



Empirically works even when theoretical assumptions do not hold

Assumptions:

- Deterministic dynamics,
- Same dynamics (i.e. even termination must be the same)
- Same non-terminal rewards (i.e. goal-reaching tasks)

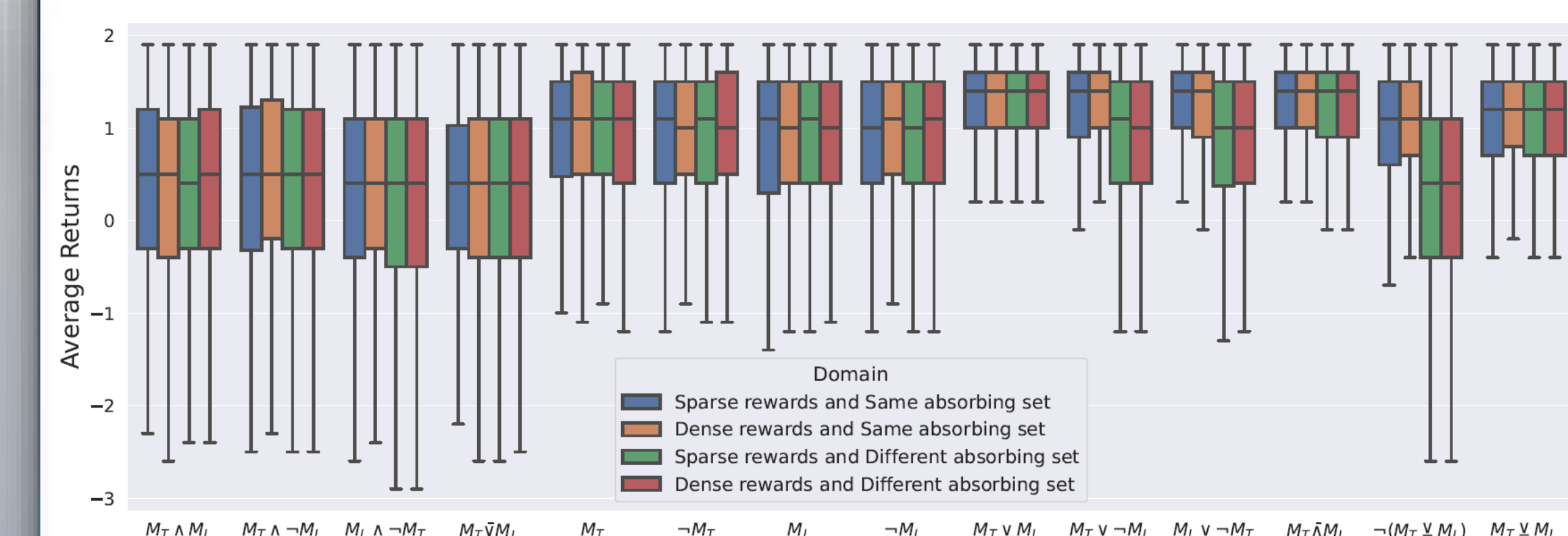


Figure 14: Average Returns. Trained and evaluated with noisy transitions (0.3% noise).

Empirically works even in continuous control with function approximation

