# On Prediction and Planning in Partially Observable Markov Decision Processes with Large Observation Sets

## Pablo Samuel Castro

pcastr@cs.mcgill.ca

### McGill University

Joint work with: Doina Precup and Prakash Pananganden

# Motivation

- Interested in sequential decision making under uncertainty

- Agent must infer its "state" based on observations of environment

- A larger observation space gives more information, but increases complexity of problem

- Hardware is cheap and small => many sensors/observations!

# Our contribution

- Allow subsets of observation space to be specified for planning/learning.

- Provide theoretical foundations when planning/learning using this idea.

- Will address questions such as:

  - How is agent's behaviour affected by using only a subset of all observations?

  - How are agent's predictions affected?
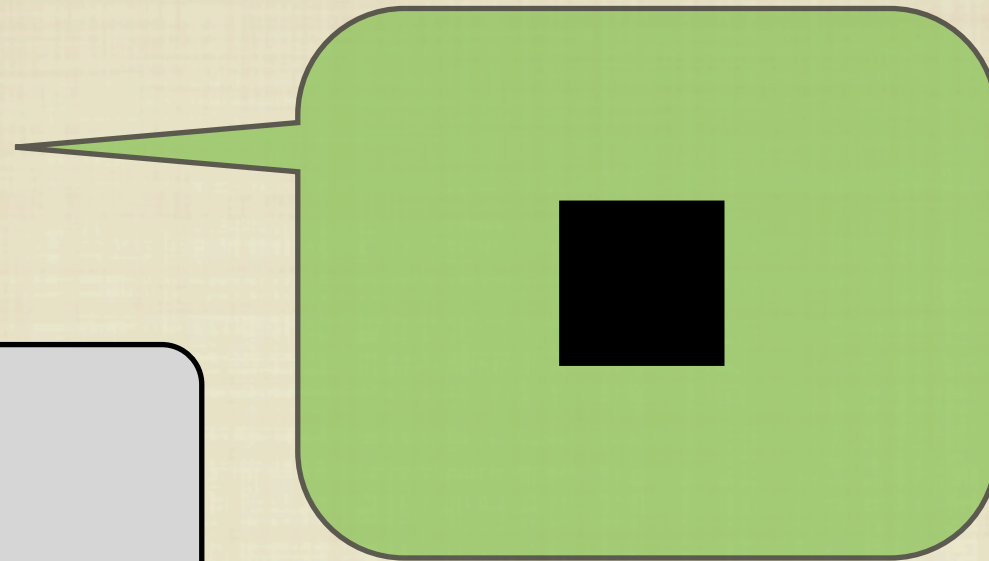
# Outline

- **POMDP** review

- **New POMDP** formulation

- **Equivalence relations**

    - **Value functions**

    - **Trajectory predictions**

    - **Bisimulation**
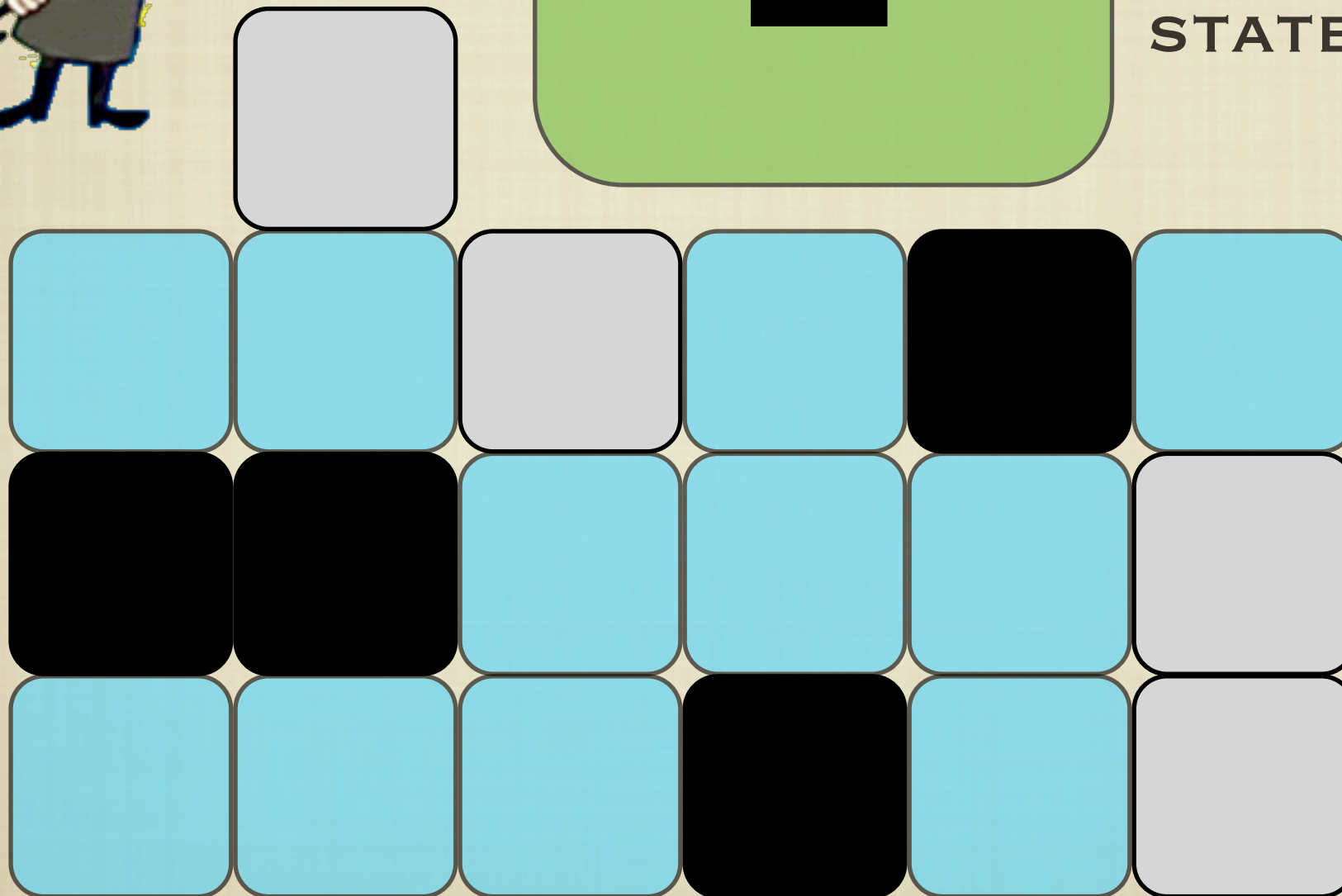
- **Conclusions and future work**

# OUTLINE

- **POMDP REVIEW**

- NEW POMDP FORMULATION

- EQUIVALENCE RELATIONS

    - VALUE FUNCTIONS

    - TRAJECTORY PREDICTIONS

    - BISIMULATION

- CONCLUSIONS AND FUTURE WORK

# Partially observable MDPs (POMDPs)



Maintain a distribution over states based on clues

# Standard POMDPs

- 6-tuple $\langle S, A, P, R, \Omega, O \rangle$ consisting of

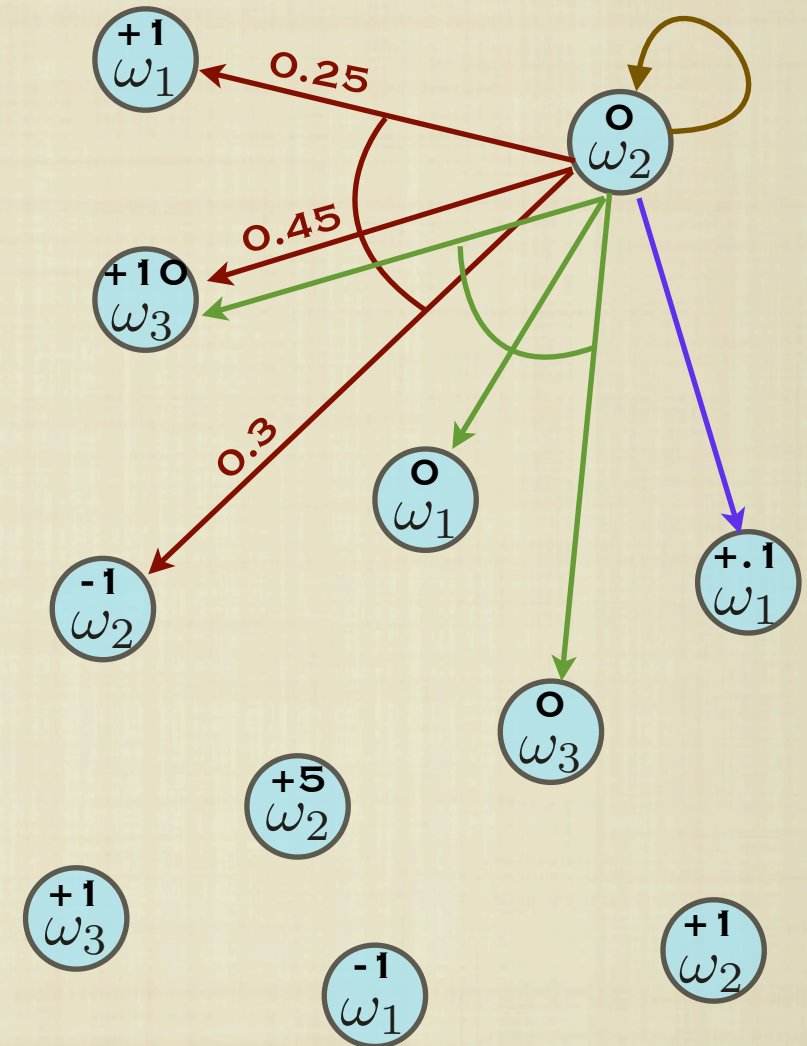  - Set of states $S$, $(s, s', t, \ldots)$

  - Set of actions $A$, $(a, b, \ldots)$

  - Probabilistic transition function $P(s, a)(s')$
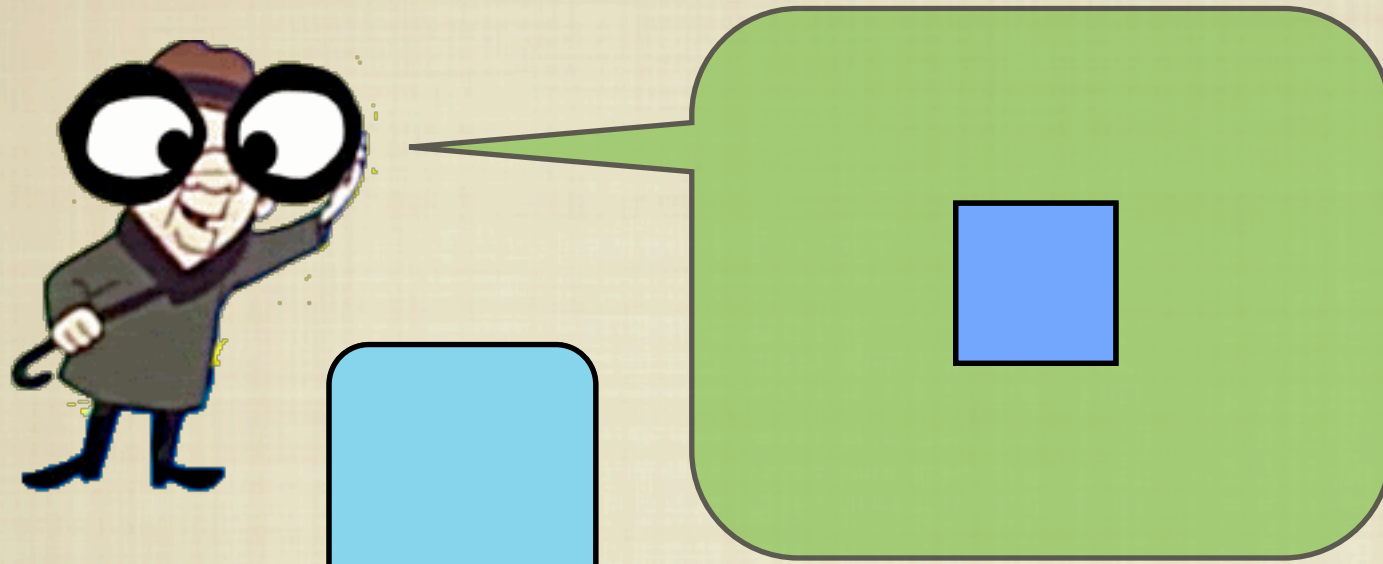
  - Bounded reward function $R(s, a)$

  - Set of observations $\Omega$
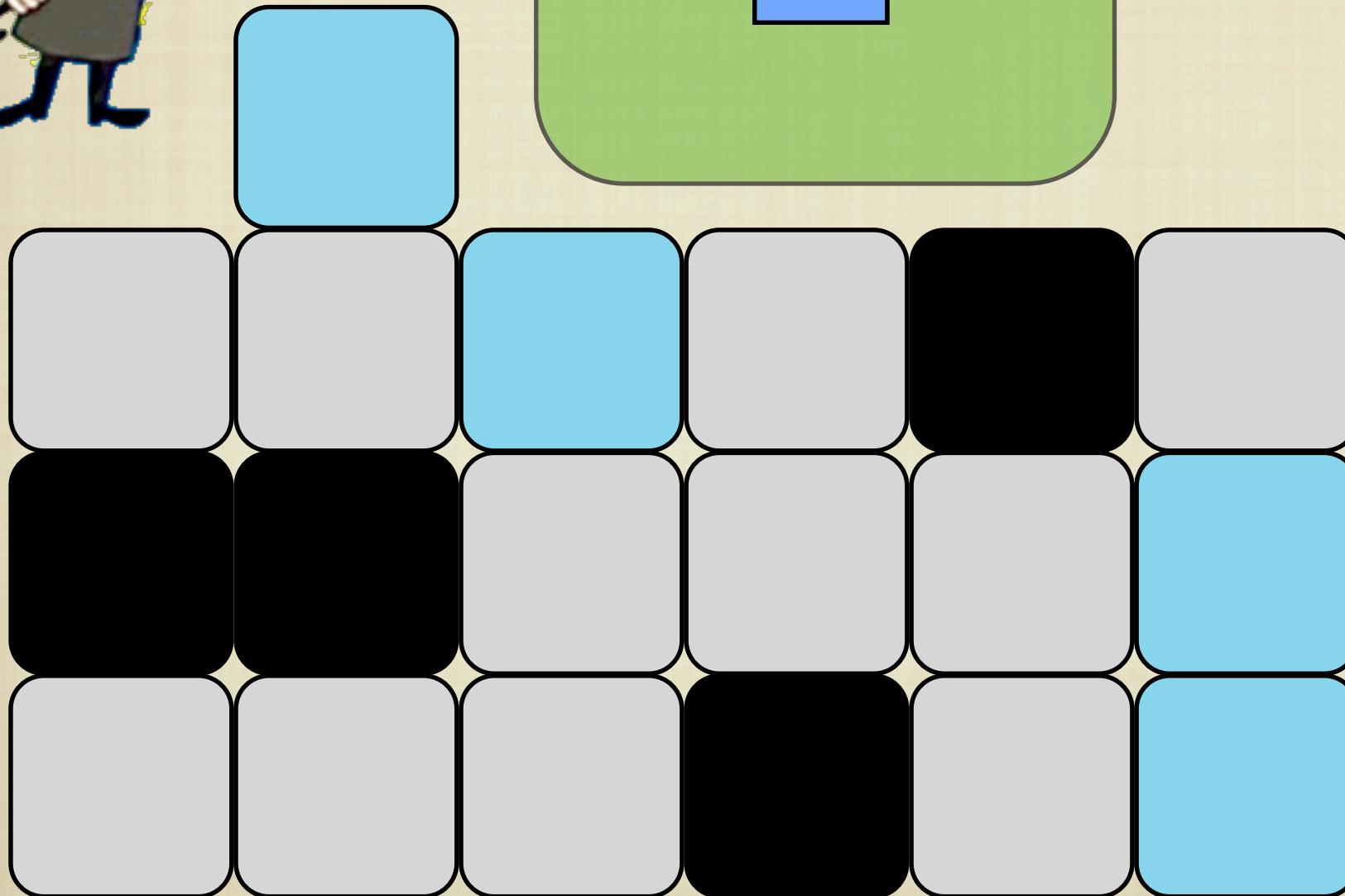
  - Observation function $O(a, s)(\omega)$

  - Discount factor $0 \leq \gamma < 1$

# Belief states



Move forward

# Belief states

- A belief state $\mu$ is a distribution over $S$.

- Given $\mu$, action $a$ and observation $\omega$, there is a unique next belief state $\tau(\mu, a, \omega)$.

- Can also compute probability of next observations.

# Outline

- **POMDP** review

- **New POMDP formulation**

- Equivalence relations

  - Value functions

  - Trajectory predictions

  - Bisimulation

- Conclusions and future work

# POMDPs

- **5-tuple $\langle S, A, P, \Omega, O \rangle$ consisting of**

  - **Set of states $S$, $(s, s', t, \ldots)$**

  - **Set of actions $A$, $(a, b, \ldots)$**
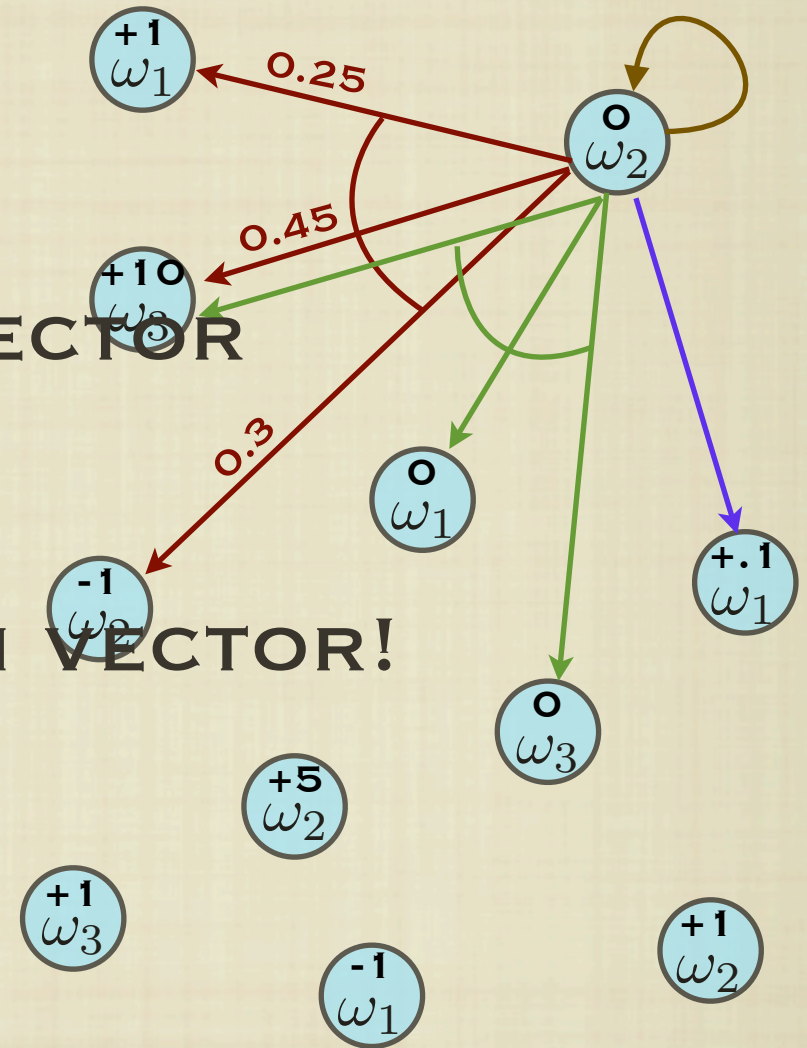
- **k-dimensional observation vector**

  - **Probabilistic transition function** $= \Omega_1 \times \Omega_2 \times \cdots \times \Omega_k$ $P(s, a)(s')$

- **Rewards part of observation vector!**

  - **Set of observations $\Omega$**

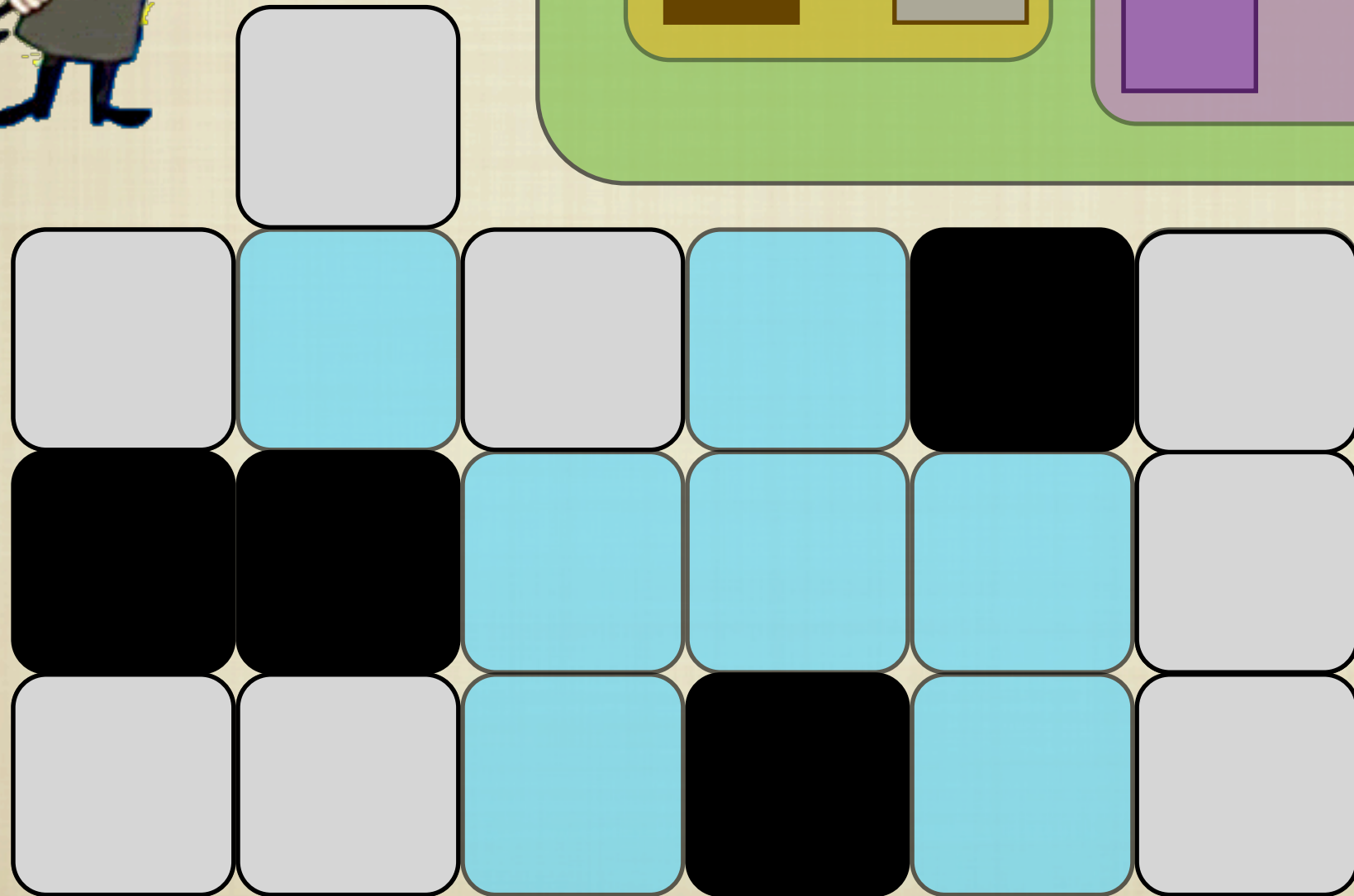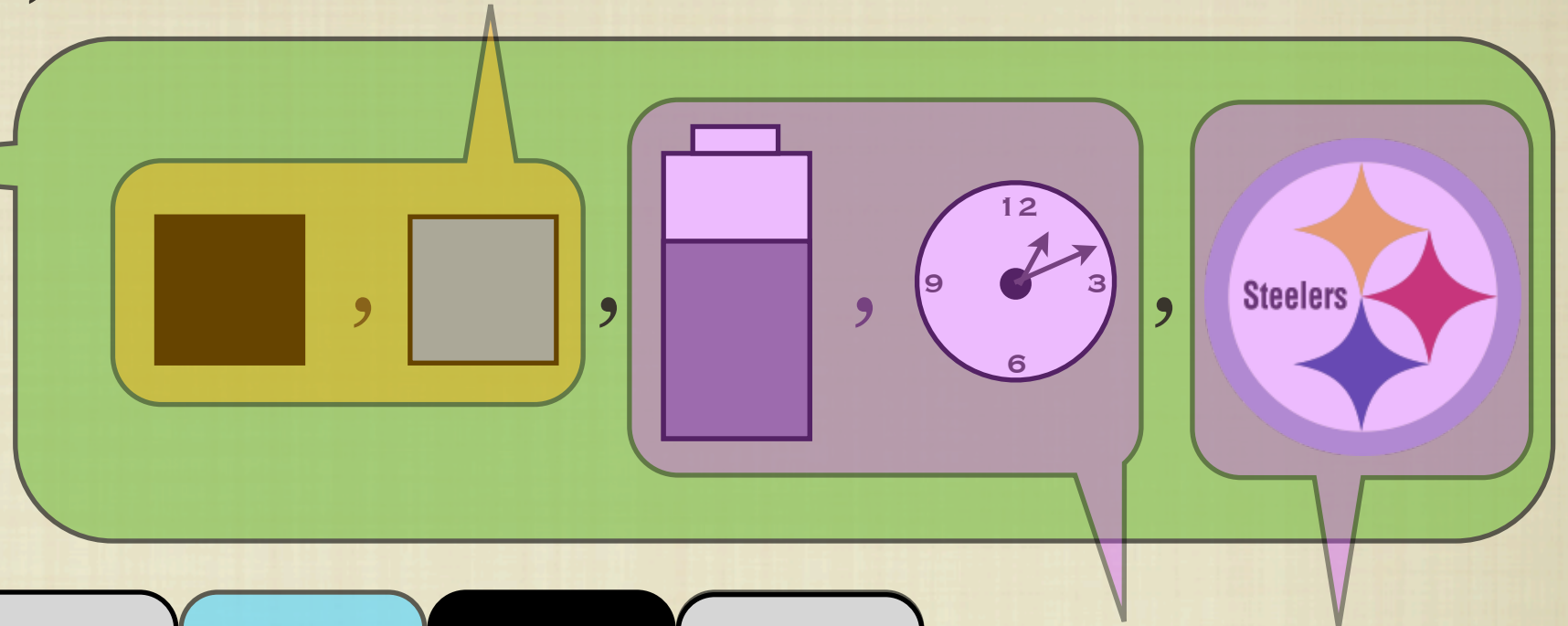  - **Observation function $O(s, a, s)(\omega)$**

  - **Discount factor $0 \leq \gamma < 1$**

# Partially observable MDPs (POMDPs)

State updates

Performance

# Specifying data and interest

- Let $\mathcal{D} \subseteq \{1, 2, \ldots, k\}$ be indices of observation coordinates used for belief updates

- Let $\mathcal{I} \subseteq \{1, 2, \ldots, k\}$ be indices of observation coordinates that are observables of interest for planning/prediction.

- Let $\Omega_{\mathcal{D}}$ be set of observations containing only observations from $\mathcal{D}$. similarly for $\Omega_{\mathcal{I}}$.

# New POMDP dynamics

- **We project observation functions with binary projection matrices** $\Phi_\mathcal{D}$: $O_\mathcal{D} = O\Phi_\mathcal{D}$

- **Unique next beliefs specific to choice of** $\mathcal{D}$:
$$\tau_\mathcal{D}(\mu, a, \omega) = \frac{\mu P^a O_\mathcal{D}^\omega}{\mu P^a O_\mathcal{D}^\omega \mathbf{e}^T}$$

- **Can define a transition fn. between belief states** $T_\mathcal{D}(\mu, a)(\mu')$.

# Measuring performance

- Elements from $\Omega_{\mathcal{I}}$ may be of many different types.

- Need a way to quantify an agent's performance.

- We assume a function $f : \Omega_{\mathcal{I}} \to \mathbb{R}$ that maps observations of interest to a real number.

# Policies and value functions

- **Closed-loop policies: Map belief states to actions** $(\pi \in \Pi)$

- **Value of a belief state** $\mu$ **under** $\pi$: $\mathbb{E}^{\pi}\left[\sum_{i=0}^{H} \gamma^i r_i | \mu\right]$

$$V_{\mathcal{D},\mathcal{I}}^{\pi} = \sum_{\omega_{\mathcal{I}} \in \Omega_{\mathcal{I}}} Pr(\omega_{\mathcal{I}} | \mu, \pi(\mu))) f(\omega_{\mathcal{I}}) + \gamma \sum_{\mu' \in \mathcal{B}} T_{\mathcal{D}}(\mu, \pi(\mu))(\mu') V_{\mathcal{D},\mathcal{I}}^{\pi}(\mu')$$

- **Optimal value function:**

$$V_{\mathcal{D},\mathcal{I}}^{*}(\mu) = \max_{a \in A} \left\{ \sum_{\omega_{\mathcal{I}} \in \Omega_{\mathcal{I}}} Pr(\omega_{\mathcal{I}} | \mu, a) f(\omega_{\mathcal{I}}) + \gamma \sum_{\mu' \in \mathcal{B}} T_{\mathcal{D}}(\mu, a)(\mu') V_{\mathcal{D},\mathcal{I}}^{*}(\mu') \right\}$$

# Is new POMDP definition suitable?

- Are fully observable MDPs still expressible?

- Does the definition properly follow intuition? (e.g. do larger observation subsets yield improved performance)?
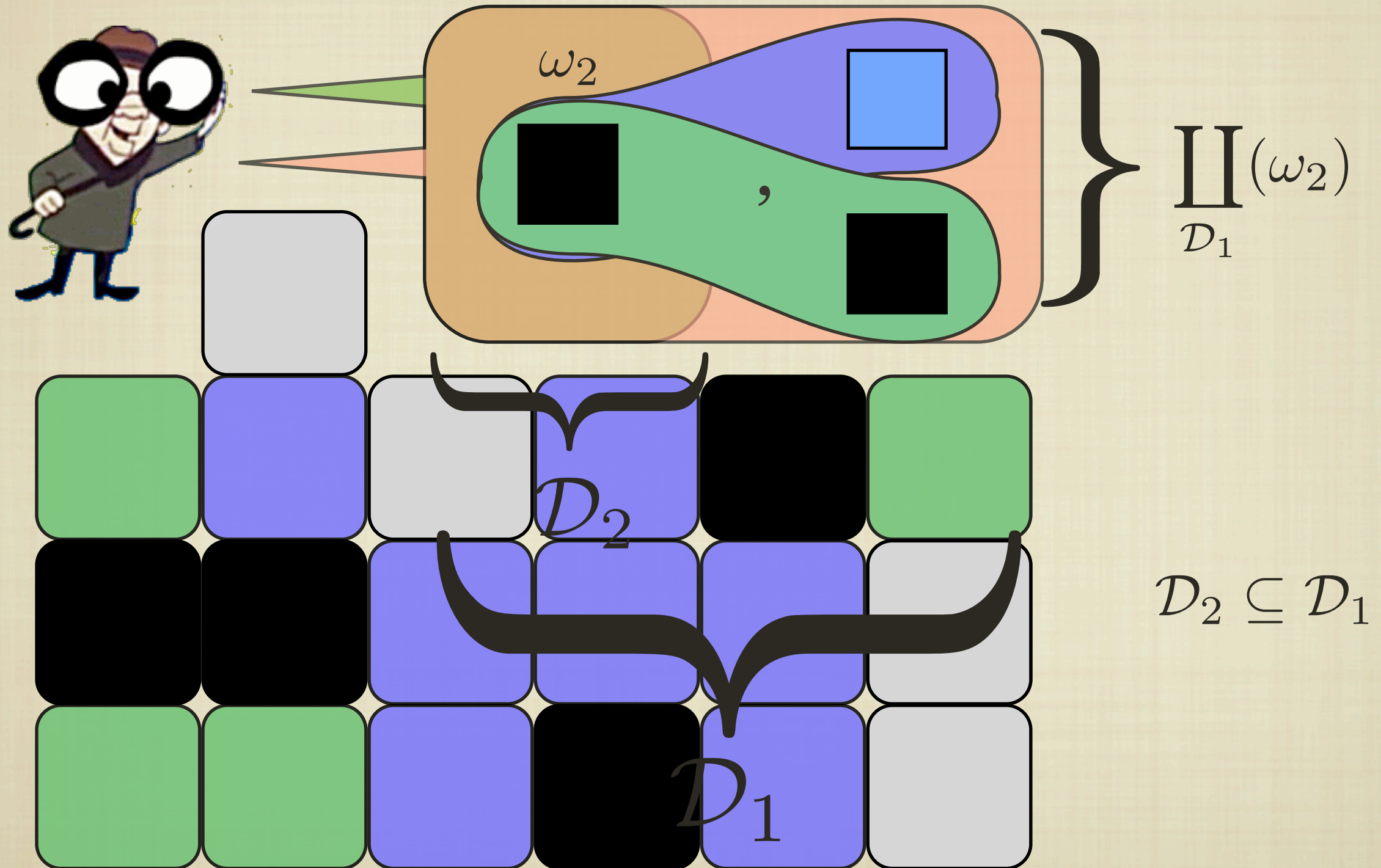
# MDPs

- Assume the observations are just $S \times \mathbb{R}$

- $\mathcal{D}$ points to $S$ and $\mathcal{I}$ points to $\mathbb{R}$.

# Optimal value functions

- **Proposition:** Given indexing sets $\mathcal{D}_2 \subseteq \mathcal{D}_1$ and $\mathcal{I}$, then

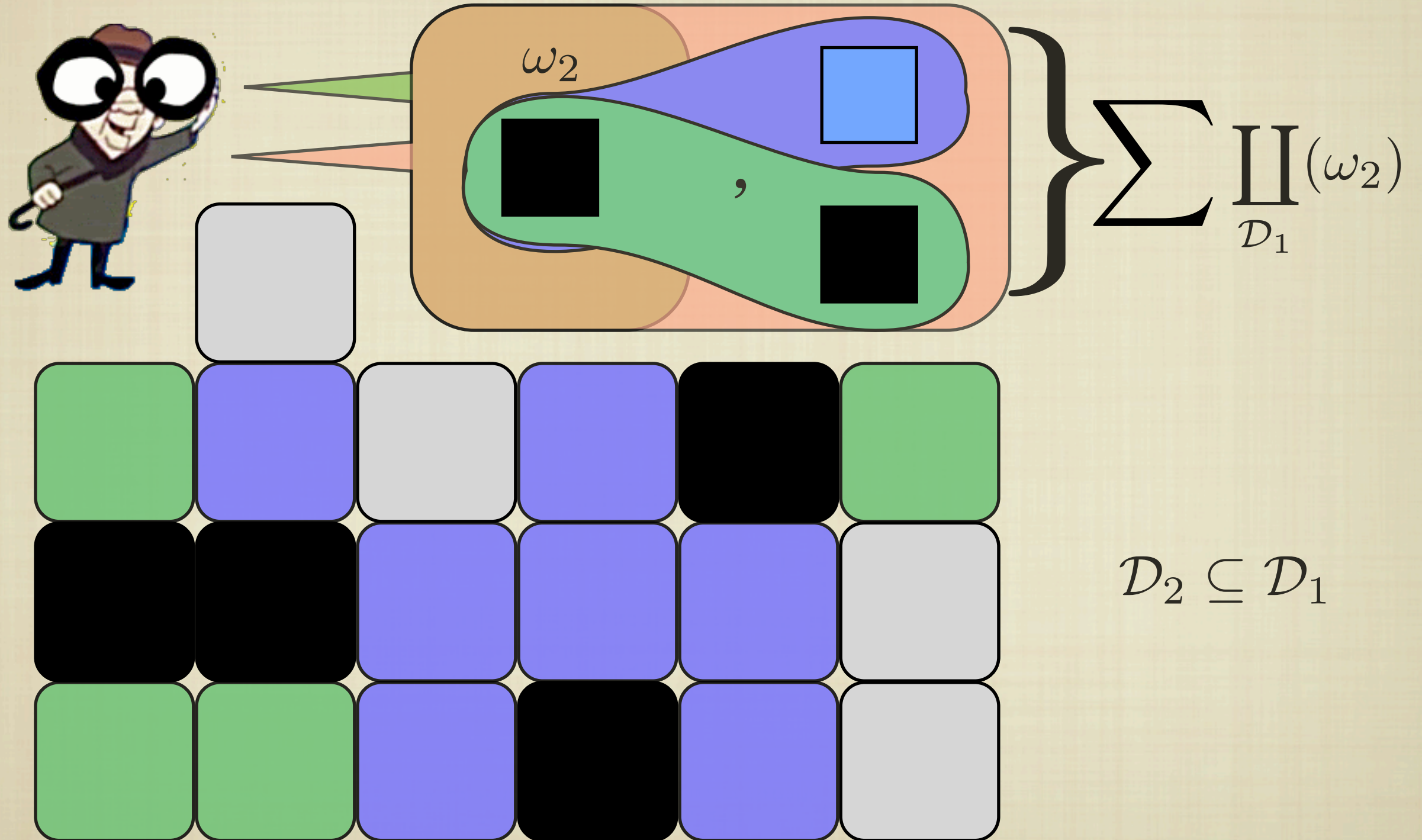$$V^*_{\mathcal{D}_2, \mathcal{I}} \leq V^*_{\mathcal{D}_1, \mathcal{I}}$$

# Convexity of Beliefs



$$\coprod_{\mathcal{D}_1}(\omega_2)$$

$$\mathcal{D}_2 \subseteq \mathcal{D}_1$$

# Convexity of Beliefs

**Theorem:** Given a belief $\mu$, an action $a$, $\mathcal{D}_2 \subseteq \mathcal{D}_1$, and observation $\omega_2 \in \Omega_{\mathcal{D}_2}$, the unique next belief $\tau_{\mathcal{D}_2}(\mu, a, \omega_2)$ can be expressed as a convex combination of the belief states $\left\{ \tau_{\mathcal{D}_1}(\mu, a, \omega_1) \right\}_{\omega_1 \in \coprod_{\mathcal{D}_1}(\omega_2)}$.

$$\left.\right\} \sum_{\mathcal{D}_1} \coprod (\omega_2)$$

$$\mathcal{D}_2 \subseteq \mathcal{D}_1$$

# Outline

- **POMDP** review

- **New POMDP formulation**

- **Equivalence relations**

  - Value functions

  - Trajectory predictions

  - Bisimulation

- Conclusions and future work

# Equivalence relations

- Partition belief space into equivalence classes

- Capture some form of behavioural equivalence

- Two beliefs in same equivalence are behaviourally indistinguishable

# Outline

- **POMDP** review

- **New POMDP formulation**

- **Equivalence relations**

- **Value functions**

- **Trajectory predictions**

- **Bisimulation**

- **Conclusions and future work**

# Value function equivalences

- For all belief states $\mu, \nu$ let $\Pi_{\mu,\nu}$ be the set of all policies $\pi \in \Pi$ where $\pi(\mu) = \pi(\nu)$

- Belief states $\mu, \nu$ are $(\mathcal{D}, \mathcal{I})$-**closed value equivalent** if for all $\pi \in \Pi_{\mu,\nu}$,
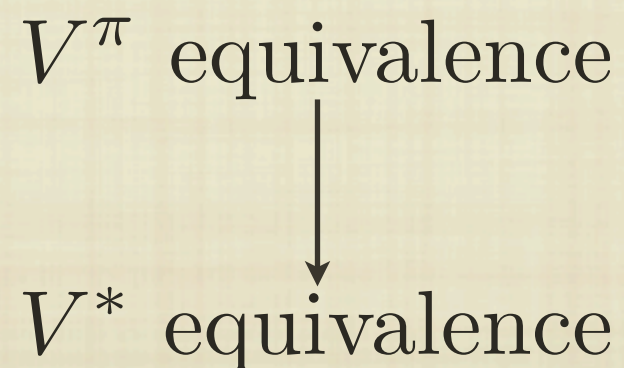
$$V_{\mathcal{D},\mathcal{I}}^{\pi}(\mu) = V_{\mathcal{D},\mathcal{I}}^{\pi}(\nu)$$

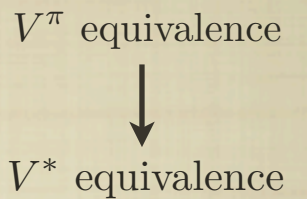- Belief states $\mu, \nu$ are $(\mathcal{D}, \mathcal{I})$-**optimal value equivalent** if

$$V_{\mathcal{D},\mathcal{I}}^{*}(\mu) = V_{\mathcal{D},\mathcal{I}}^{*}(\nu)$$

# Closed and optimal value equivalences

**Theorem: If two states are closed value equivalent, then they are necessarily optimal value equivalent.**

$$V^\pi \text{ equivalence}$$

$$\downarrow$$

$$V^* \text{ equivalence}$$

# CLOSED AND OPTIMAL VALUE EQUIVALENCES

**LEMMA:** IF $s_0, t_0$ ARE $V^\pi$ EQUIVALENT AND $V^*(s_0) > V^*(t_0)$, THEN PROB. OF REACHING $t_0$ FROM $s_0$ UNDER $\pi^*$ IS STRICTLY POSITIVE.

LET $\Pi_{CV}$ BE SET OF ALL POLICIES $\pi$ CONSTRUCTED FROM SOME OPTIMAL POLICY $\pi^*$ AS FOLLOWS:

$$\pi(s') = \pi^*(s_0) \text{ if } s' = t_0$$

$$\pi(s') = \pi^*(s') \text{ otherwise}$$

# Closed and optimal value equivalences

- $B$ is set of bounded functions $V : S \times \Pi_{CV} \to [0, 1]$

- $\mathcal{R} \in B$, $\mathcal{R}(s, \pi) = R(s, \pi(s))$

- $\Upsilon : B \to B,\ \Upsilon(V)(s, \pi) = \gamma \sum_{s' \neq t_0} P(s, \pi(s))(s') V(s', \pi) + P(s, \pi(s))(t_0) V(t_0, \pi)$

- $\tau(e) = \mathcal{R} + \Upsilon(e)$ has least fixed pt $e^*(s, \pi) = V^\pi(s)$
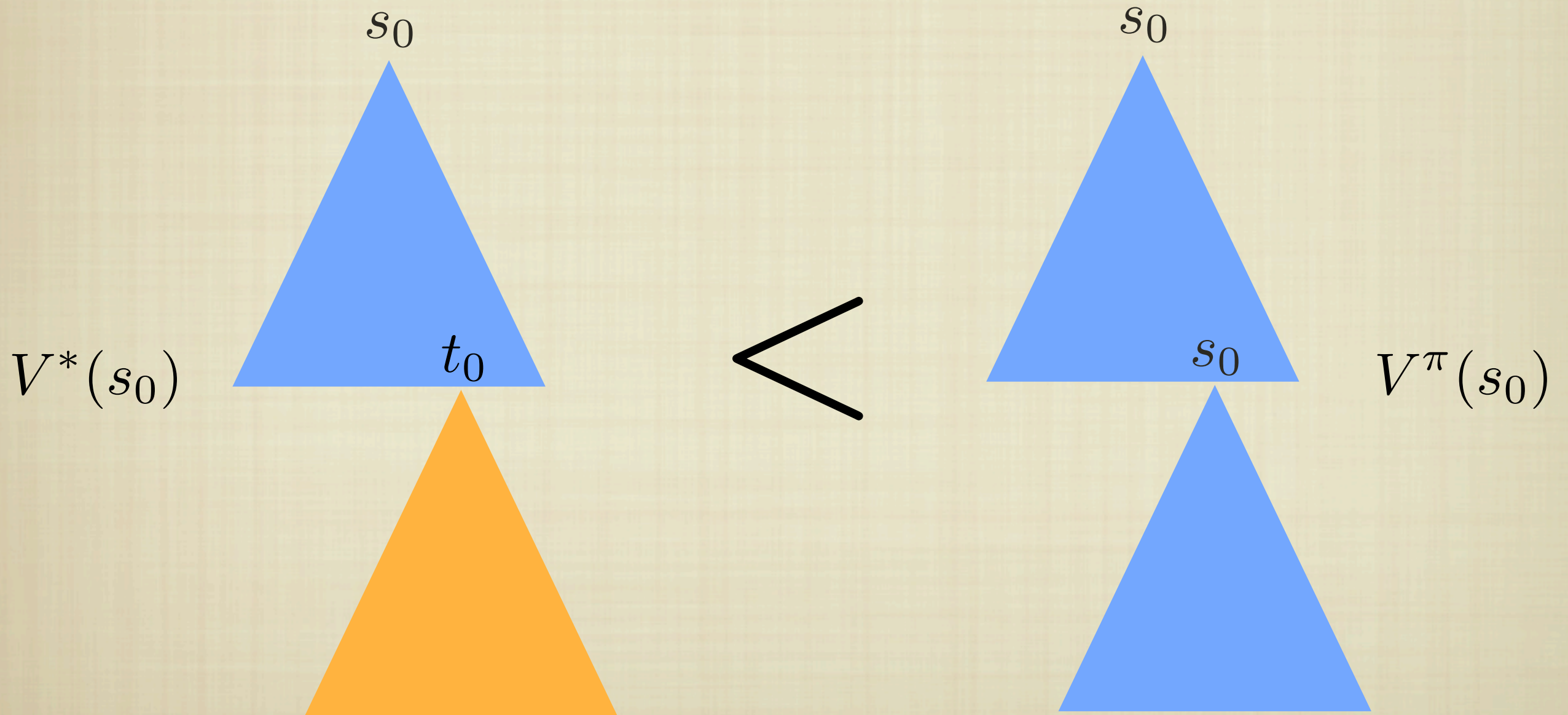
**Theorem (based on (Kozen, 2007))**

Define $\varphi \subseteq B$ as $V \in \varphi \Rightarrow \forall \pi \in \Pi_{CV}.V(s, \pi) \geq V^*(s)$, then
if $\varphi \neq \emptyset$ and $e \in \varphi \Rightarrow \tau(e) \in \varphi$, then $e^* \in \varphi$.

# Closed and optimal value equivalences

**Lemma:** If $s_0$ and $t_0$ are $V^\pi$ equivalent then $V^*(s_0) \leq V^*(t_0)$.

**Proof:** Assume $V^*(s_0) > V^*(t_0)$



$V^*(s_0) \qquad\qquad\qquad < \qquad\qquad\qquad V^\pi(s_0)$

# Closed and optimal value equivalences

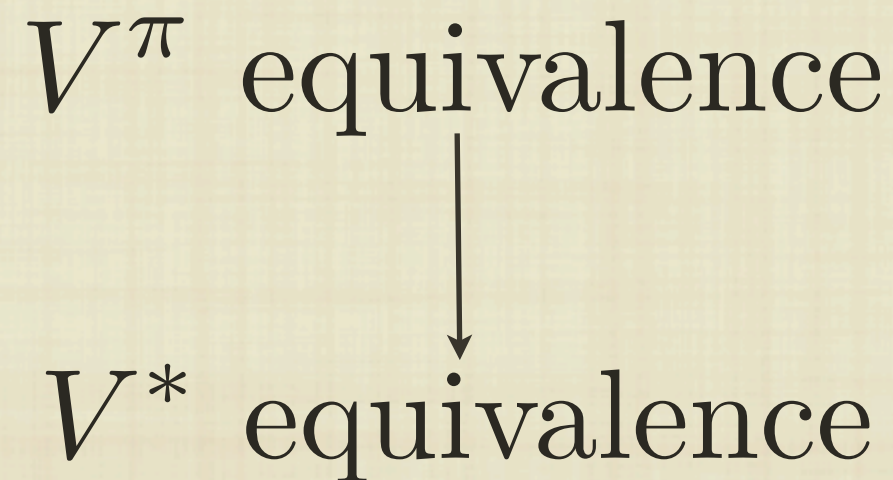**Lemma:** If $s_0$ and $t_0$ are $V^\pi$ equivalent then $V^*(s_0) \leq V^*(t_0)$.

**Lemma:** If $s_0$ and $t_0$ are $V^\pi$ equivalent then $V^*(s_0) \geq V^*(t_0)$.

# CLOSED AND OPTIMAL VALUE EQUIVALENCES

**LEMMA:** If $s_0$ and $t_0$ are $V^\pi$ equivalent then $V^*(s_0) \leq V^*(t_0)$.

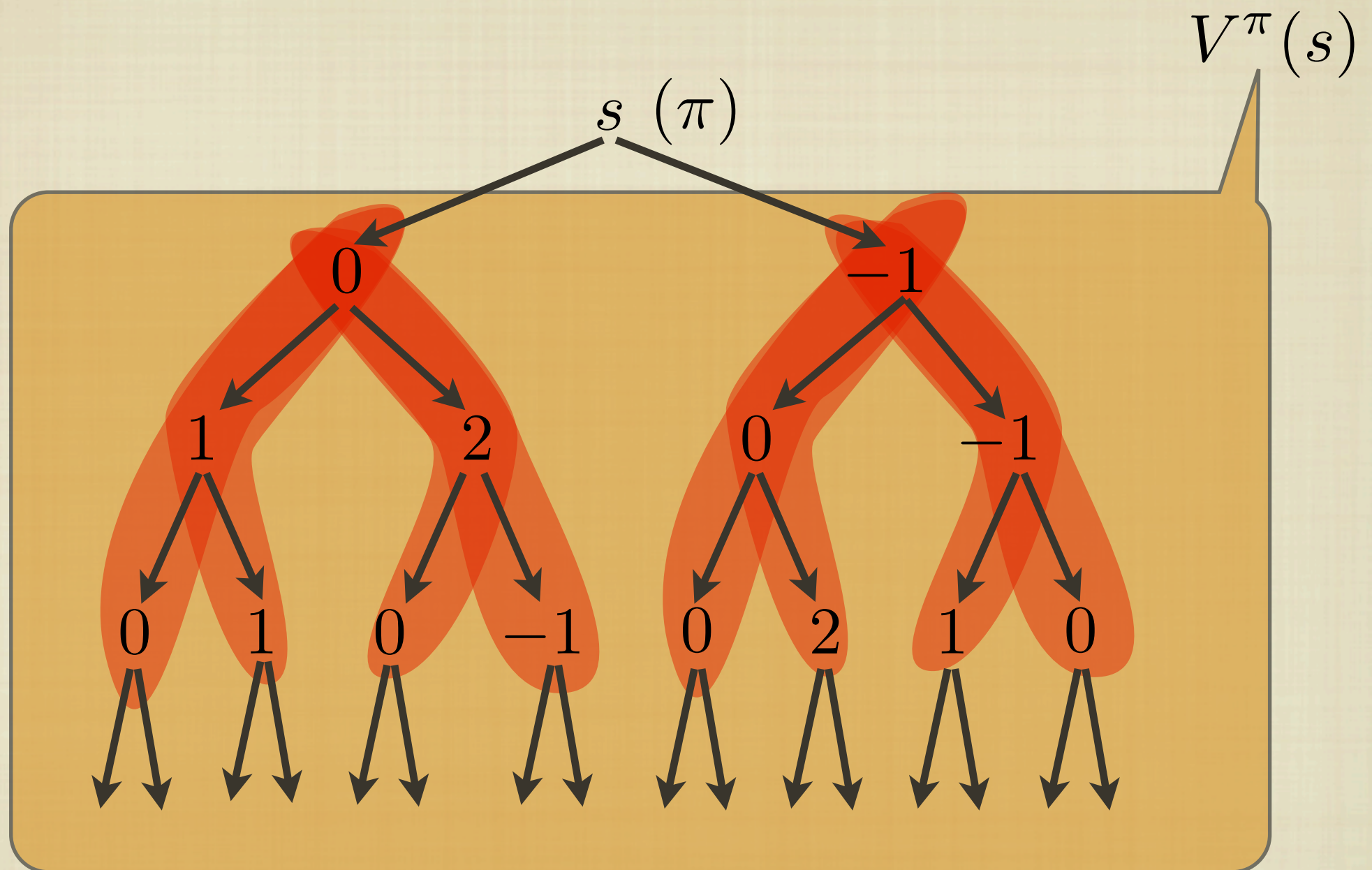**LEMMA:** If $s_0$ and $t_0$ are $V^\pi$ equivalent then $V^*(s_0) \geq V^*(t_0)$.

**THEOREM:** If $s_0$ and $t_0$ are $V^\pi$ equivalent then $V^*(s_0) = V^*(t_0)$.

$$V^\pi \text{ equivalence}$$

$$\downarrow$$

$$V^* \text{ equivalence}$$

# Outline

- POMDP review

- **New POMDP formulation**

- Equivalence relations

  - Value functions

  - **Trajectory predictions**

  - Bisimulation

- Conclusions and future work

# Trajectory equivalence

- Two belief states $\mu, \nu$ are $\mathcal{I}$-closed trajectory equivalent if for all $\pi \in \Pi_{\mu,\nu}$ and all finite observation trajectories, $\alpha = \langle \omega_1, \omega_2, \ldots, \omega_n \rangle \in \Omega_{\mathcal{I}}^*$

$$Pr(\alpha|\mu, \pi) = Pr(\alpha|\nu, \pi)$$

# TRAJECTORY EQUIVALENCE

- **OPEN-LOOP POLICIES** $\theta \in \Theta$ **MAP TIME STEPS TO ACTIONS**

- **TWO BELIEF STATES** $\mu, \nu$ **ARE** $\mathcal{I}$**-OPEN TRAJECTORY EQUIVALENT IF FOR ALL** $\theta \in \Theta$ **AND ALL FINITE OBSERVATION TRAJECTORIES** $\alpha = \langle \omega_1, \omega_2, \ldots, \omega_n \rangle \in \Omega_{\mathcal{I}}^*$,
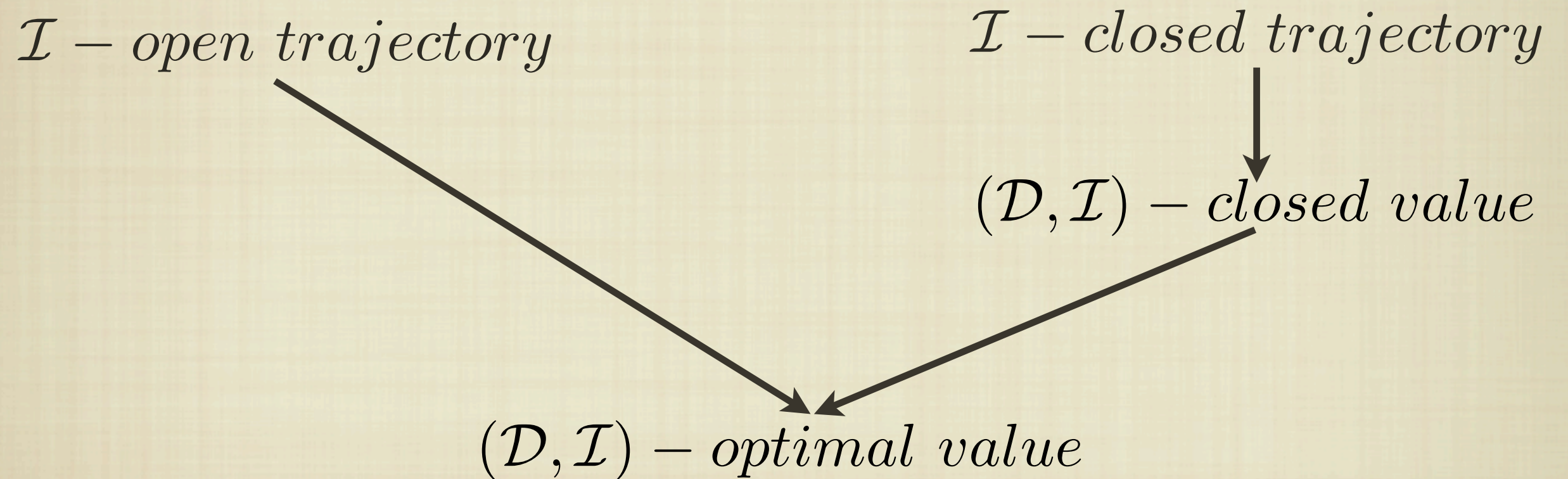
$$Pr(\alpha | \mu, \theta) = Pr(\alpha | \nu, \theta)$$

- **A TRAJECTORY** $\alpha$ **AND OPEN LOOP POLICY** $\theta$ **CONSTITUTE A PSR TEST (LITTMAN ET AL., 2002)!**

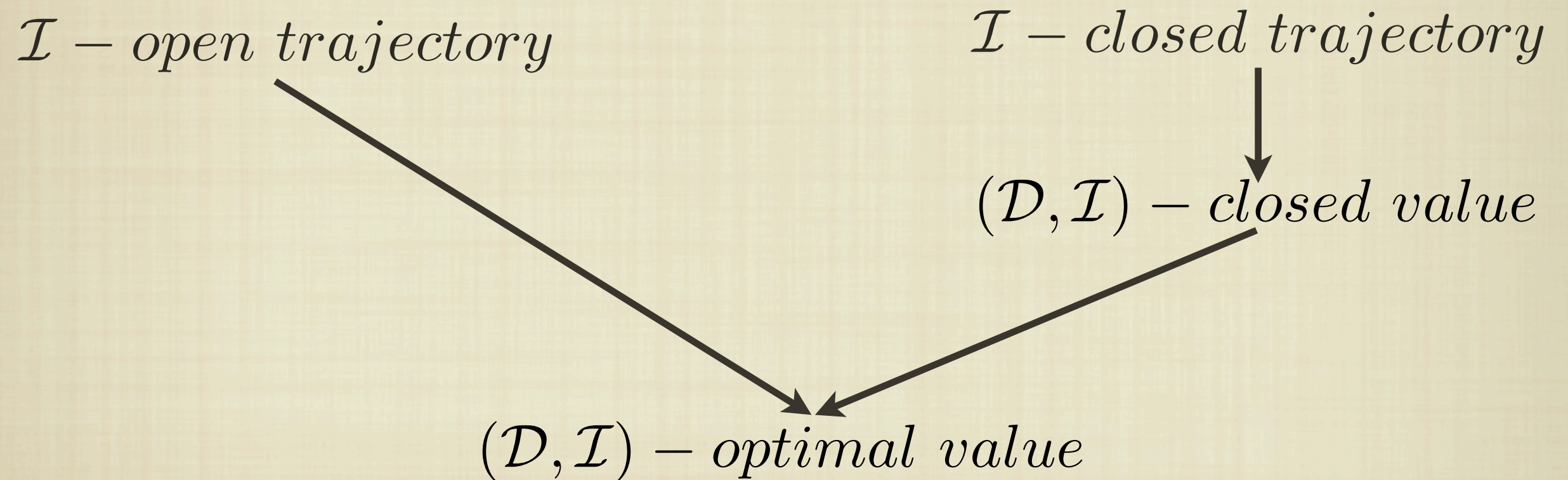$$\langle a_1, \omega_1, a_1, \omega_2, \ldots, a_n, \omega_n \rangle$$

# Hierarchy

- **If $\mathcal{D} \subseteq \mathcal{I}$, then the following hierarchy is obtained**

$\mathcal{I} - open\ trajectory$

$\mathcal{I} - closed\ trajectory$

$(\mathcal{D}, \mathcal{I}) - closed\ value$

$(\mathcal{D}, \mathcal{I}) - optimal\ value$

# Hierarchy

- If $\mathcal{D} \nsubseteq \mathcal{I}$, then the following hierarchy is obtained

$$\mathcal{I} - open\ trajectory$$

$$\mathcal{I} - closed\ trajectory$$

$$(\mathcal{D}, \mathcal{I}) - closed\ value$$

$$(\mathcal{D}, \mathcal{I}) - optimal\ value$$

# HIERARCHY

$$\mathcal{I} - open\ trajectory \longrightarrow\!\!\!\!/\!\!\!\!\longrightarrow (\mathcal{D}, \mathcal{I}) - optimal\ value$$

# HIERARCHY

S AND T ARE OPEN TRAJECTORY EQUIVALENT

$s = 0.6$

$V^*(s) = \dfrac{\langle a, 0, a, 1 \rangle}{1 - \gamma}$

a, b 0.6

a, b 0.4

a

b

a

b

a, b

a, b

a, b
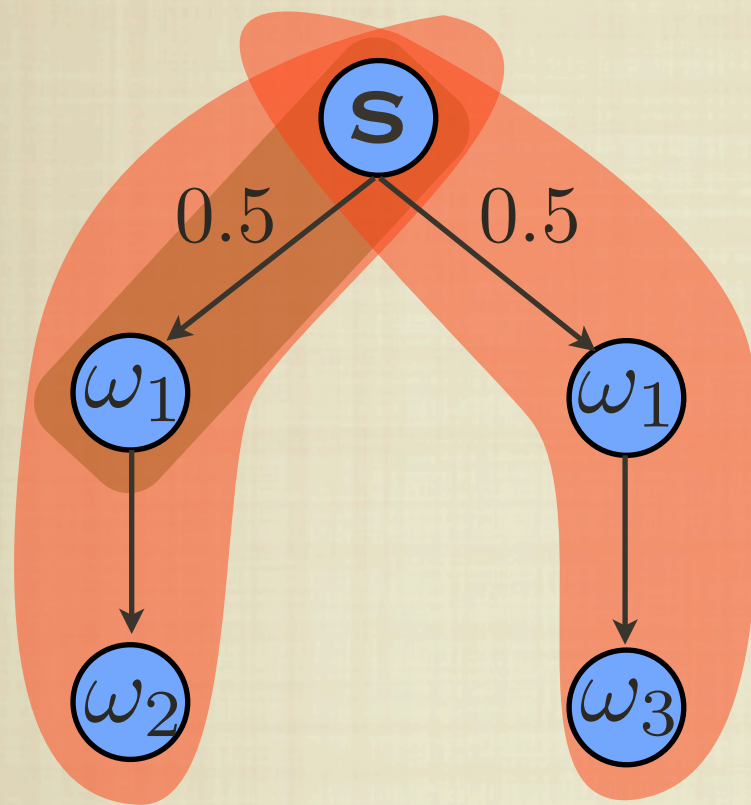
0.6 a

0.4 a

0.6 b

0.4 b

S AND T ARE NOT OPTIMAL VALUE EQUIVALENT!

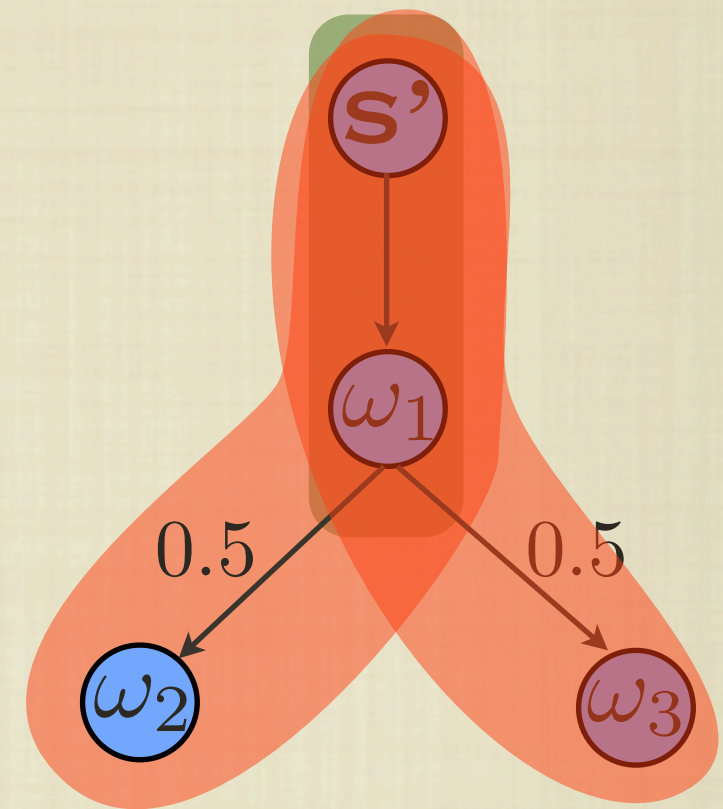$V^*(t) = \dfrac{0.6\gamma}{1 - \gamma}$

a, b

$t = 0.6$

# Outline

- **POMDP** review

- **New POMDP formulation**

- **Equivalence relations**

  - **Value functions**

  - **Trajectory predictions**

  - **Bisimulation**

- **Conclusions and future work**

# Bisimulation

# Bisimulation

- An equivalence relation $E$ is a $(\mathcal{D}, \mathcal{I})$-bisimulation relation if whenever $\mu, \nu$ are $(\mathcal{D}, \mathcal{I})$-bisimilar then

  - For all $\omega \in \Omega_{\mathcal{I}}, a \in A$, $Pr(\omega | \mu, a) = Pr(\omega | \nu, a)$

  - For all $c \in \mathcal{B}/_E, a \in A$,

  $$\sum_{\mu' \in c} T_{\mathcal{D}}(\mu, a)(\mu') = \sum_{\mu' \in c} T_{\mathcal{D}}(\nu, a)(\mu')$$

- If $\mu$ and $\nu$ are $(\mathcal{D}, \mathcal{I})$-bisimilar we will write $\mu \sim \nu$.

# DETERMINISTIC BISIMULATION

- AN EQUIVALENCE RELATION $E$ IS A DETERMINISTIC $(\mathcal{D}, \mathcal{I})$-BISIMULATION RELATION IF WHENEVER $\mu, \nu$ ARE DETERMINISTIC $(\mathcal{D}, \mathcal{I})$-BISIMILAR THEN

  - FOR ALL $\omega \in \Omega_{\mathcal{I}}, a \in A, \; Pr(\omega | \mu, a) = Pr(\omega | \nu, a)$

  - FOR ALL $\omega \in \Omega_{\mathcal{D}}, \; a \in A, \; \tau_{\mathcal{D}}(\mu, a, \omega) E \tau_{\mathcal{D}}(\nu, a, \omega)$

- IF $\mu$ AND $\nu$ ARE DETERMINISTIC $(\mathcal{D}, \mathcal{I})$-BISIMILAR WE WILL WRITE $\mu \simeq \nu$.

# HIERARCHY

$\mathcal{D} \subseteq \mathcal{I}$

$Deterministic$

$(\mathcal{D}, \mathcal{I}) - bisimulation$

$\overline{\overline{=}}$

$\mathcal{I} - open\ trajectory$

$\mathcal{I} - closed\ trajectory$

$(\mathcal{D}, \mathcal{I}) - bisimulation$

$(\mathcal{D}, \mathcal{I}) - closed\ value$

$(\mathcal{D}, \mathcal{I}) - optimal\ value$

$$Deterministic \atop (\mathcal{D}, \mathcal{I}) - bisimulation \quad \longleftrightarrow\!\!\!\!/\!\!\!\!\longrightarrow \quad (\mathcal{D}, \mathcal{I}) - bisimulation$$

# Hierarchy

$$s \sim t$$
$$s \not\approx t$$

$$\mathcal{D} \nsubseteq \mathcal{I}$$

$$Deterministic$$
$$(\mathcal{D}, \mathcal{I}) - bisimulation$$
$$\overline{=}$$
$$\mathcal{I} - open\ trajectory \qquad\qquad \mathcal{I} - closed\ trajectory$$

$$(\mathcal{D}, \mathcal{I}) - bisimulation \qquad\qquad (\mathcal{D}, \mathcal{I}) - closed\ value$$

$$(\mathcal{D}, \mathcal{I}) - optimal\ value$$

# HIERARCHY

$$\begin{array}{c} Deterministic \\ (\mathcal{D}, \mathcal{I}) - bisimulation \end{array} \quad \longrightarrow\!\!\!\!/\!\!\!\longrightarrow \quad (\mathcal{D}, \mathcal{I}) - bisimulation$$

# HIERARCHY

*Deterministic*
$(\mathcal{D}, \mathcal{I}) - bisimulation$
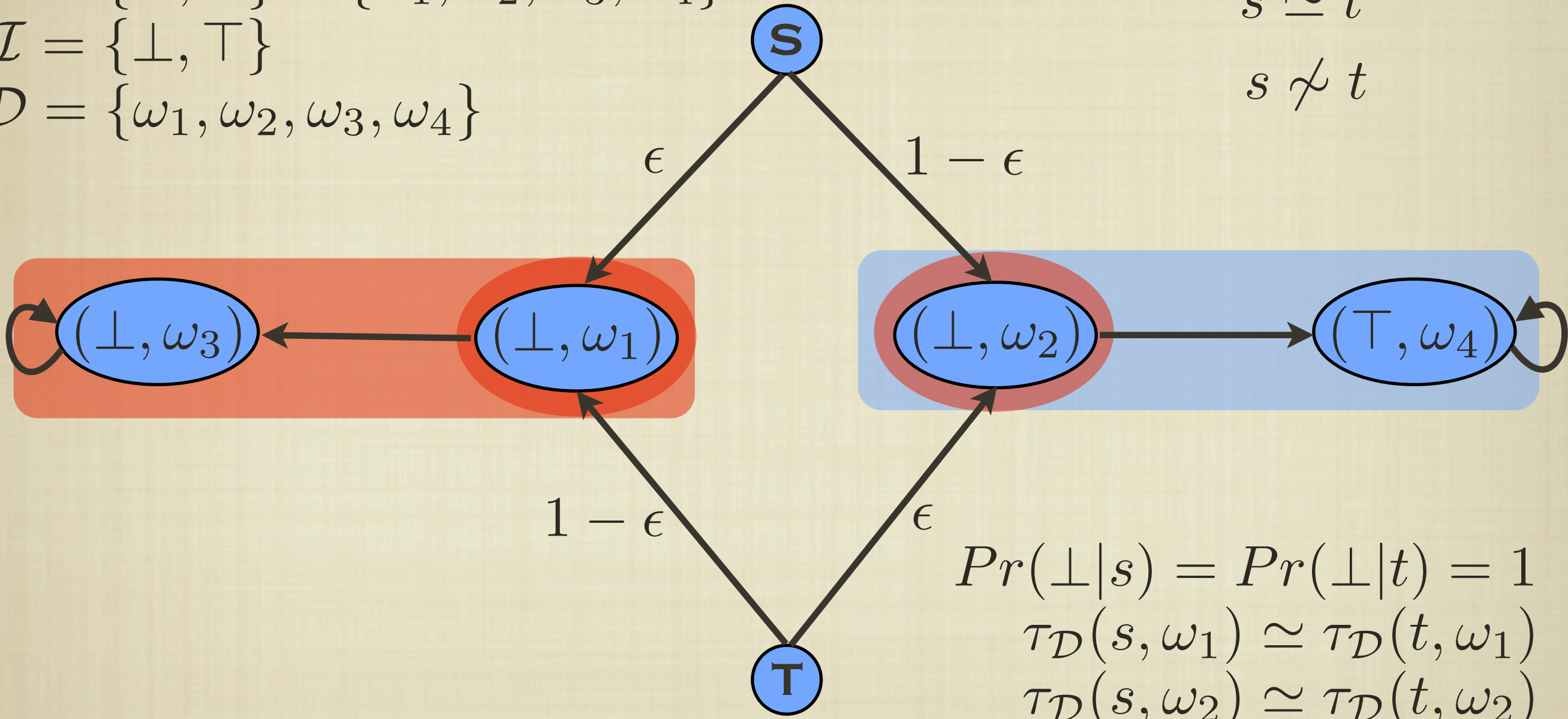$(\mathcal{D}, \mathcal{I}) - bisimulation$

$\Omega = \{\bot, \top\} \times \{\omega_1, \omega_2, \omega_3, \omega_4\}$
$\mathcal{I} = \{\bot, \top\}$
$\mathcal{D} = \{\omega_1, \omega_2, \omega_3, \omega_4\}$

$s \simeq t$
$s \not\sim t$



$\epsilon$     $1 - \epsilon$

$1 - \epsilon$     $\epsilon$

$Pr(\bot|s) = Pr(\bot|t) = 1$
$\tau_{\mathcal{D}}(s, \omega_1) \simeq \tau_{\mathcal{D}}(t, \omega_1)$
$\tau_{\mathcal{D}}(s, \omega_2) \simeq \tau_{\mathcal{D}}(t, \omega_2)$

$\mathcal{D} \nsubseteq \mathcal{I}$

*Deterministic*

$(\mathcal{D}, \mathcal{I}) - bisimulation$
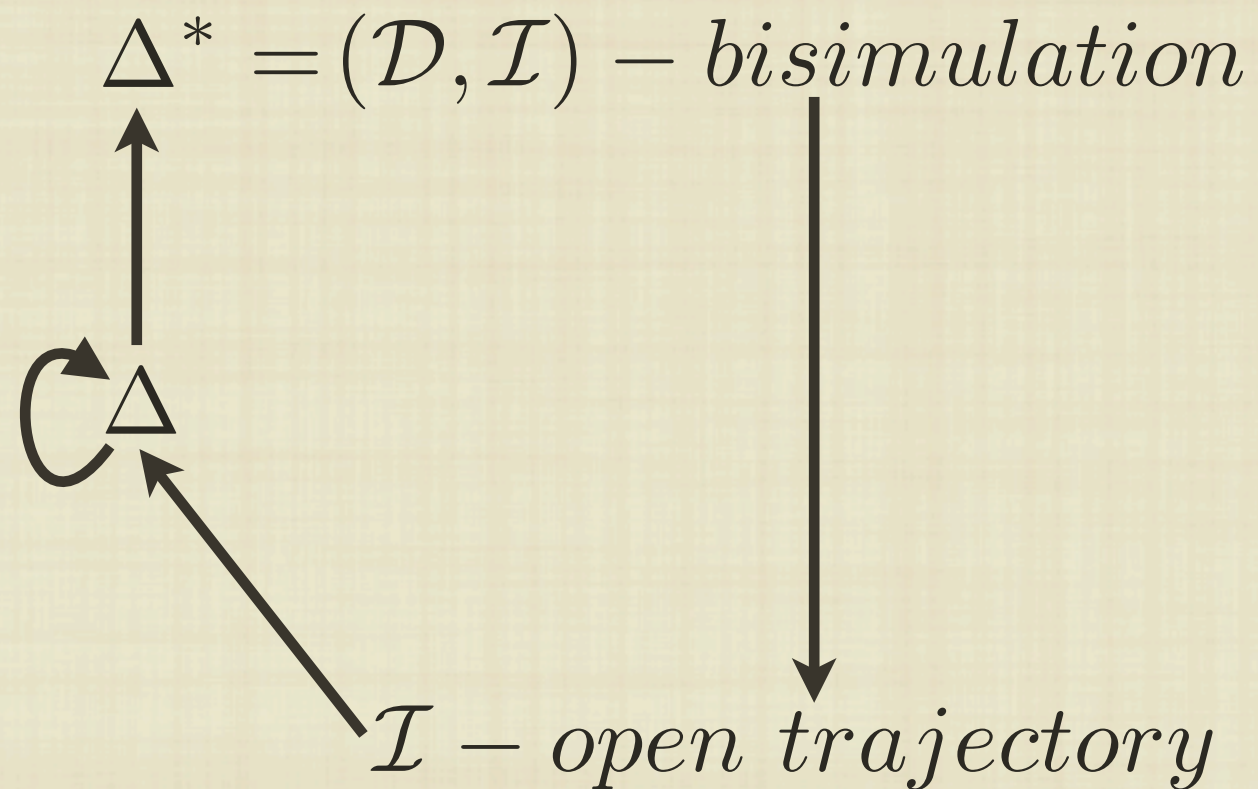
$(\mathcal{D}, \mathcal{I}) - bisimulation$

$(\mathcal{D}, \mathcal{I}) - closed\ value$

$\mathcal{I} - open\ trajectory$

$(\mathcal{D}, \mathcal{I}) - optimal\ value$

# Strengthening open trajectory

$$\mathcal{D} \not\subseteq \mathcal{I}$$

$$\Delta^* = (\mathcal{D}, \mathcal{I}) - bisimulation$$



$$\mathcal{I} - open\ trajectory$$

## From (Castro et al., 2009)

# Conclusions

■ Subsets must be chosen with care to avoid suboptimal performance

■ Open trajectory equivalence is closely related to PSRs; we showed this is not appropriate with respect to bad choices of $\mathcal{D}$ and $\mathcal{I}$.

■ In most situations we would require $\mathcal{D} \subseteq \mathcal{I}$.

■ $(\mathcal{D}, \mathcal{I})$-bisimulation is robust even when $\mathcal{D} \not\subseteq \mathcal{I}$.

$\mathcal{D} \subseteq \mathcal{I}$

$Deterministic$

$(\mathcal{D}, \mathcal{I}) - bisimulation$
$=$
$\mathcal{I} - open\ trajectory$

$\mathcal{I} - closed\ trajectory$

$(\mathcal{D}, \mathcal{I}) - bisimulation$

$(\mathcal{D}, \mathcal{I}) - closed\ value$
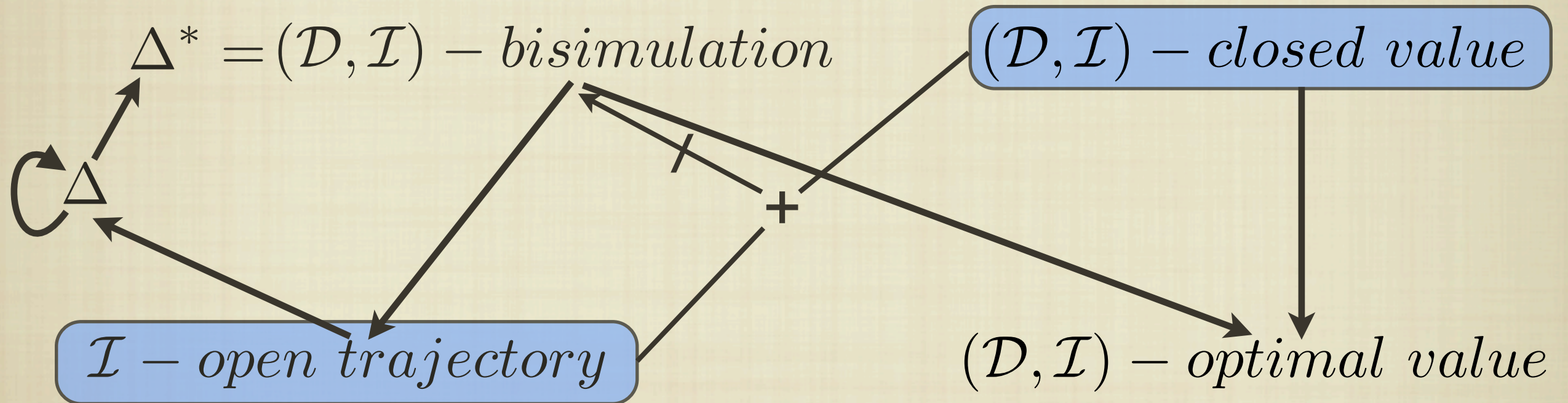
$(\mathcal{D}, \mathcal{I}) - optimal\ value$

$\mathcal{D} \not\subseteq \mathcal{I}$

*Deterministic*

$(\mathcal{D}, \mathcal{I}) - bisimulation$

$\Delta^* = (\mathcal{D}, \mathcal{I}) - bisimulation$

$(\mathcal{D}, \mathcal{I}) - closed\ value$

$\circlearrowleft \Delta$

$+$

$\mathcal{I} - open\ trajectory$

$(\mathcal{D}, \mathcal{I}) - optimal\ value$

# CURRENT WORK

- We are currently working on learning algorithms for determining $\mathcal{D}$, assuming $\mathcal{I}$ is known.

  - Start with a small $\mathcal{D}$, incrementally add more observations.

  - Start planning/learning with a small $\mathcal{D}$, use an expert/oracle to determine whether more observations are necessary

# Future work

- We project $\Omega$ onto $\Omega_{\mathcal{D}}$ and $\Omega_{\mathcal{I}}$ using binary projection matrices.

  - If we allow general projection matrices, does open trajectory equivalence yield something similar to TPSRs (Rosencratz & Gordon, 2004; Boots et al., 2010).

- Life-long learning: Many tasks to solve, different choices of $\mathcal{D}$ and $\mathcal{I}$, depending on task.

- ranking of observations to dynamically set $\mathcal{D}$ based on time requirements.

# REFERENCES

- KOZEN, D. (2007). Coinductive proof principles for stochastic Processes. Logical Methods in Computer Science 3(4:8). DOI: 10.2168/LMCS-3 (4:8) 2007.

- LITTMAN, M., R. SUTTON, AND S. SINGH (2002). Predictive representations of state. In Proceedings of the 14th Conference on Advances in Neural Information Processing Systems (NIPS-02), pp. 1555–1561.

- CASTRO, P. S., P. PANANGADEN, AND D. PRECUP (2009). Notions of state equivalence under partial observability. In Proceedings of the 21st International Joint Conference on Artificial Intelligence (IJCAI-09), pp. 1653–1658.

- BOOTS, B., S. M. SIDDIQI, AND G. J. GORDON (2010). Closing the learning-planning loop with predictive state representations. In Proc. Robotics: Science and Systems VI.

- ROSENCRANTZ, M. AND G. GORDON (2004). Learning low dimensional predictive represen- tations. In Proceedings of the International Conference on Machine Learning (ICML- 04).

$V^\pi$ equivalence

$\downarrow$

$V^*$ equivalence

**Lemma:** If $s_0$ and $t_0$ are $V^\pi$ equivalent then $V^*(s_0) \leq V^*(t_0)$.

**Proof:** Assume $V^*(s_0) > V^*(t_0)$.

$\exists V. \quad V(s,\pi) \geq V^*(s)$? **Yes! Just take** $V \equiv 1$

$V(s,\pi) \geq V^*(s) \Rightarrow \tau(V)(s,\pi) \geq V^*(s)$? **Yes!** Take any $s \neq t_0$ and $\pi \in \Pi_{CV}$

**We've shown that for any** $s \neq t_0$ **and** $\pi \in \Pi_{CV}$, $V^\pi(s) \geq V^*(s)$

**with strict inequality if** $P(s, \pi^*(s))(t_0) > 0$

$\pi(s') = \pi^*(s_0)$ if $s' = t_0$
$\pi(s') = \pi^*(s')$ otherwise

**By last corollary we know** $\exists s' \neq t_0.\ P(s, \pi^*(s))(t_0) > 0$

**Thus,** $V^\pi(s') > V^*(s')$ **contradicting optimality of** $V^*(s')$

**By contradiction,**

$= V^*(s)$

**Q.E.D.**

# Specifying data and interest

- Let $\Phi_{\mathcal{D}}$ be a projection matrix used to compute $O_{\mathcal{D}} : n \times |\Omega_{\mathcal{D}}|$:

$$O_{\mathcal{D}} = O\Phi_{\mathcal{D}}$$

- If we have $\mathcal{D}_2 \subseteq \mathcal{D}_1$, the projection $\Phi_{12}$ yields the following:

$$\Phi_{\mathcal{D}_2} = \Phi_{\mathcal{D}_1} \Phi_{12}$$

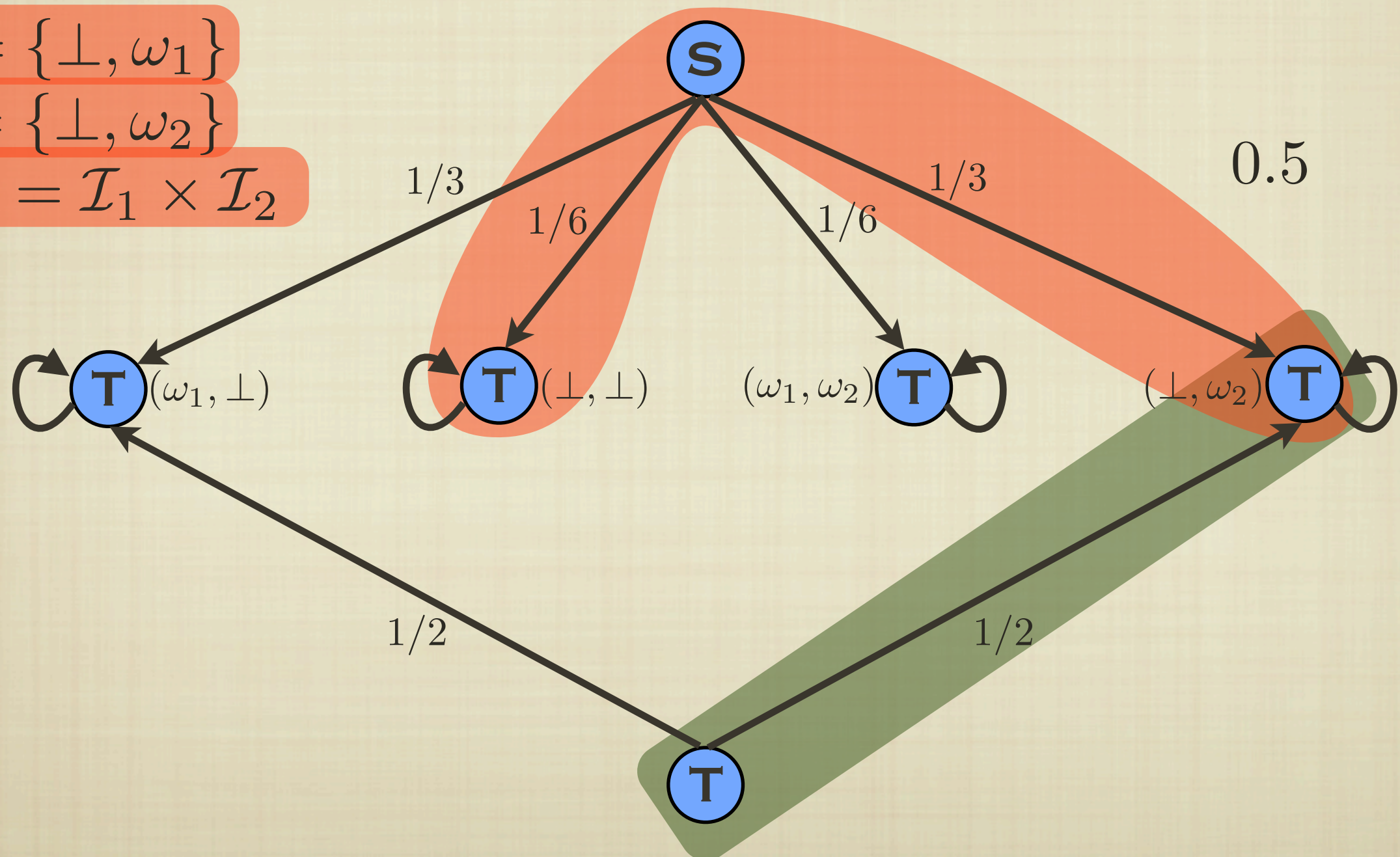$$O_{\mathcal{D}_2} = O_{\mathcal{D}_1} \Phi_{12}$$

# Approximating bisimulation

- **Proposition:** Given $\mathcal{D}, \mathcal{I}, \mu, \nu$ may be $(\mathcal{D}, \mathcal{I}_i)$-bisimilar for all $\mathcal{I}_i \subset \mathcal{I}$, but fail to be $(\mathcal{D}, \mathcal{I})$-bisimilar.

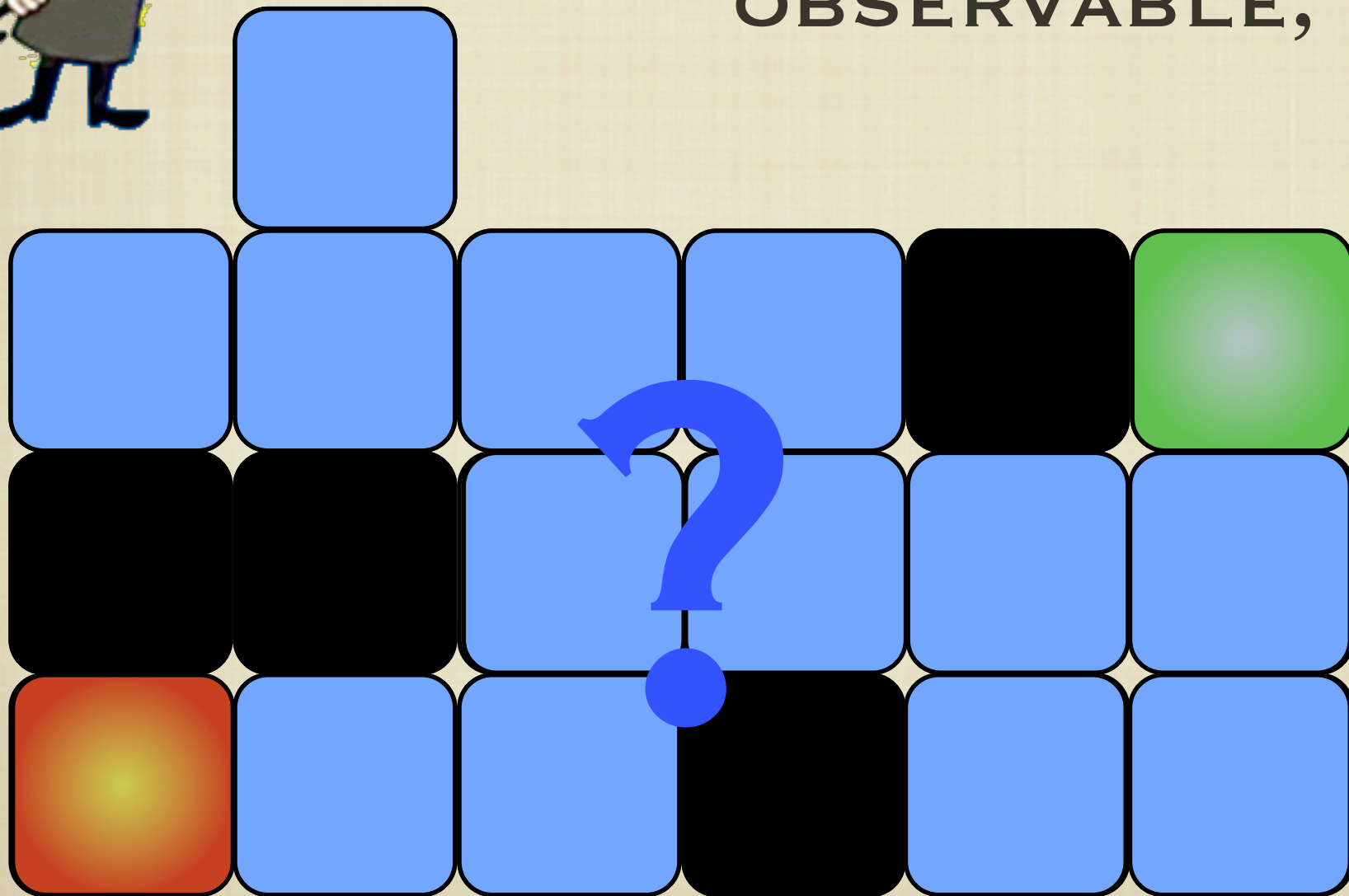$$\mathcal{I}_1 = \{\bot, \omega_1\}$$
$$\mathcal{I}_2 = \{\bot, \omega_2\}$$
$$\mathcal{D} = \mathcal{I} = \mathcal{I}_1 \times \mathcal{I}_2$$
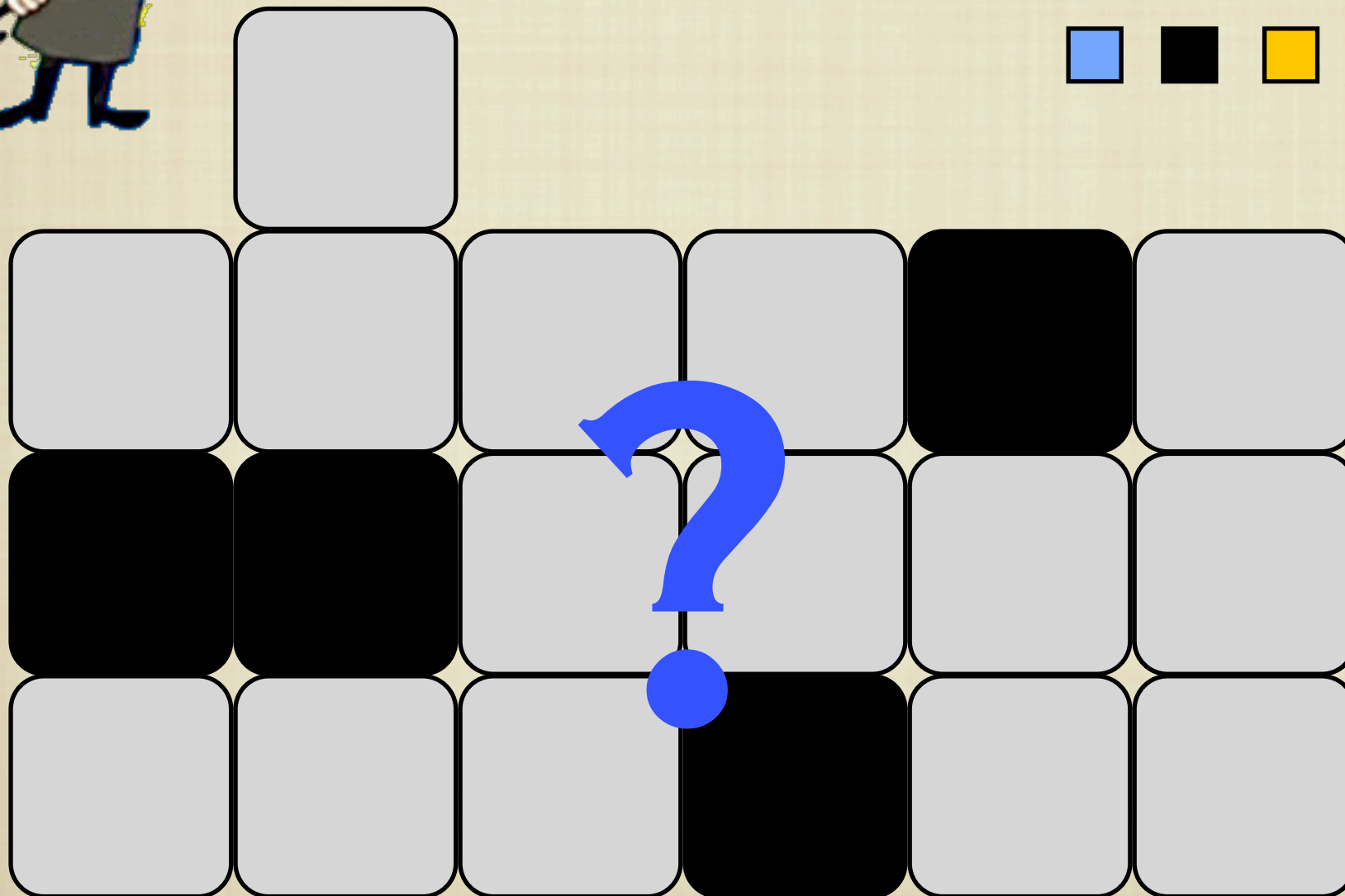
# Partially observable MDPs (POMDPs)
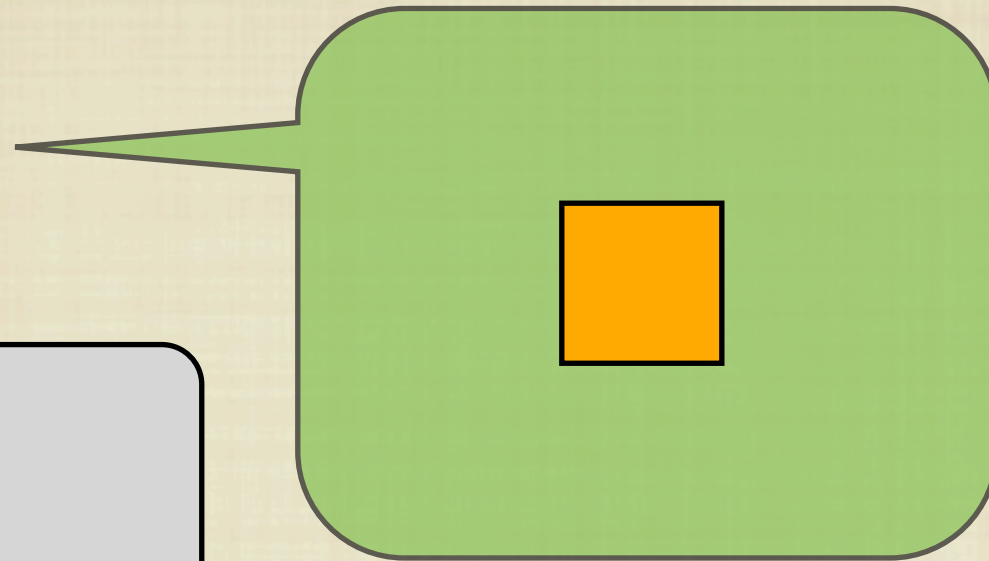
**If state is fully observable, it is a MDP**

# Partially observable MDPs (POMDPs)

**In POMDPs we only receive clues of the state**

# Partially observable MDPs (POMDPs)

Maintain a distribution over states based on clues