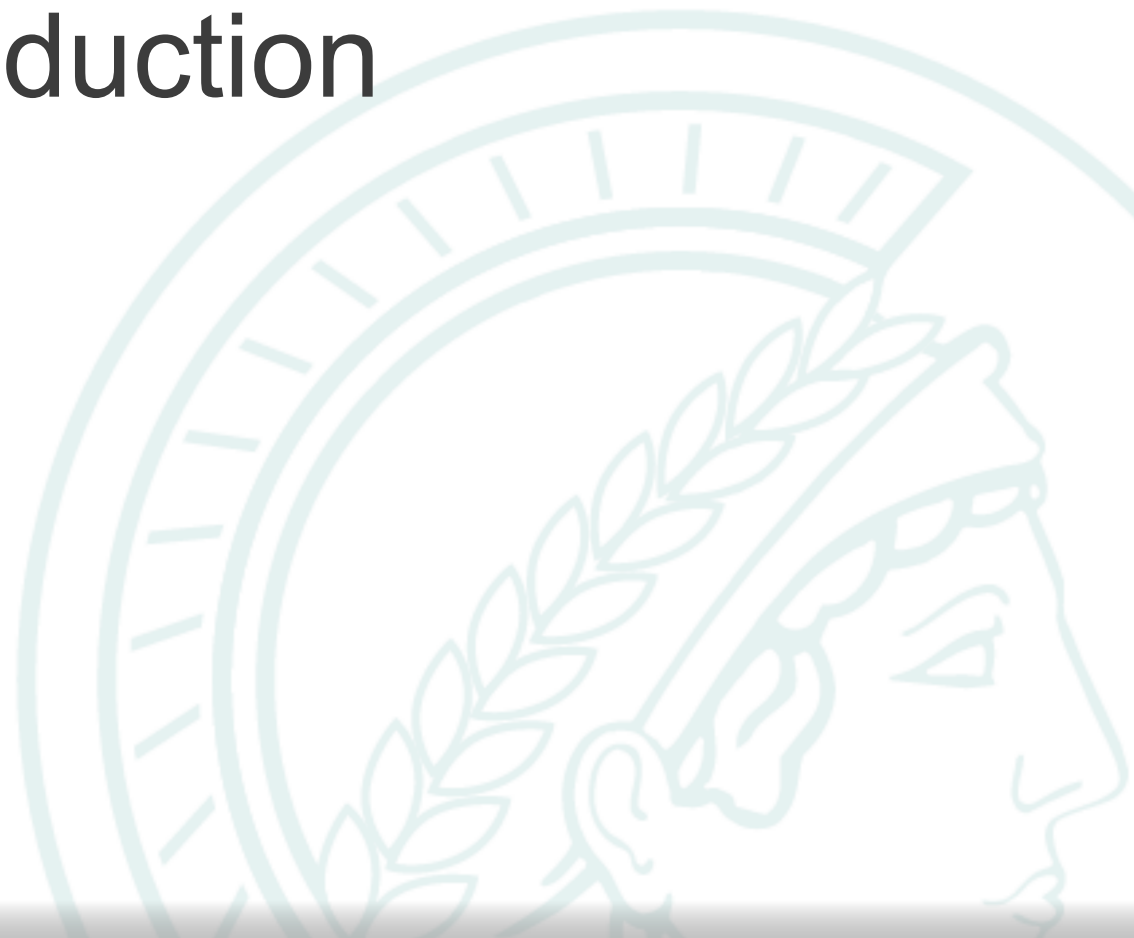# Bayesian Control Rule

## Pedro A. Ortega

Max Planck Institute for Intelligent Systems
Max Planck Institute for Biological Cybernetics

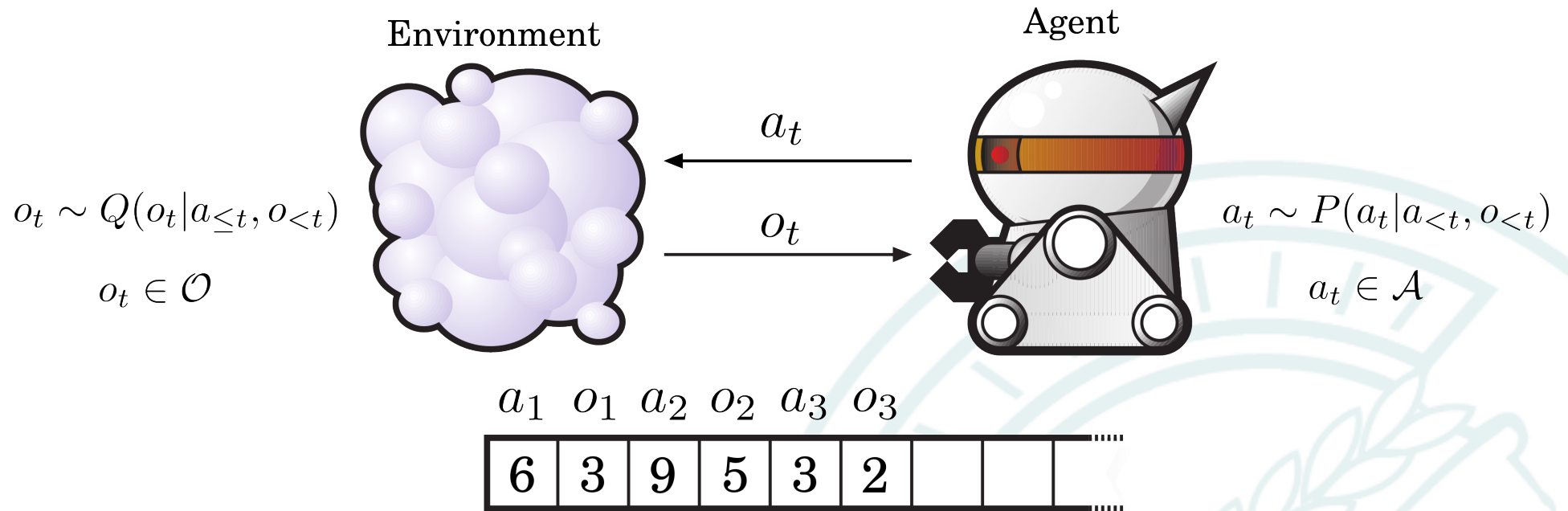# Overview

- Introduction

- Adaptation

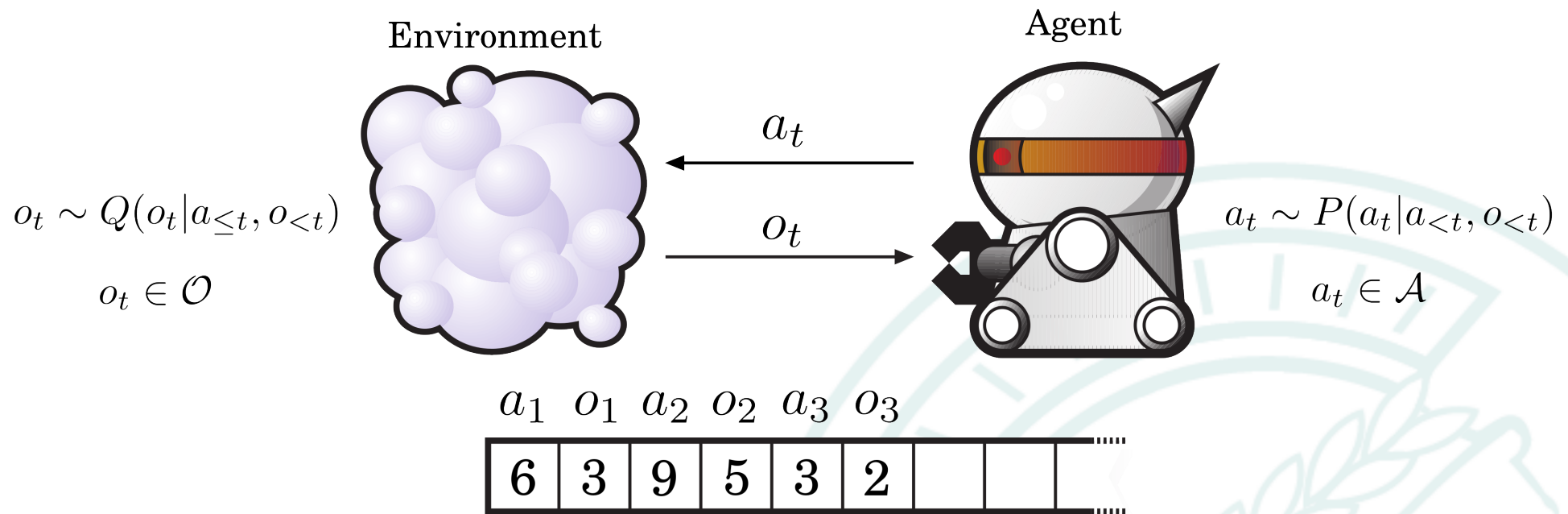- Causality

- Bayesian control rule

- Conclusions

# Introduction

# Agent-Environment Setup

Environment

Agent

$$a_t$$

$$o_t \sim Q(o_t | a_{\leq t}, o_{<t})$$

$$o_t$$

$$a_t \sim P(a_t | a_{<t}, o_{<t})$$

$$o_t \in \mathcal{O}$$

$$a_t \in \mathcal{A}$$

| $a_1$ | $o_1$ | $a_2$ | $o_2$ | $a_3$ | $o_3$ | | | |
|---|---|---|---|---|---|---|---|---|
| 6 | 3 | 9 | 5 | 3 | 2 | | | |

# Agent-Environment Setup

Environment

Agent

$$o_t \sim Q(o_t|a_{\leq t}, o_{<t})$$

$$o_t \in \mathcal{O}$$

$a_t$

$o_t$

$$a_t \sim P(a_t|a_{<t}, o_{<t})$$

$$a_t \in \mathcal{A}$$

| $a_1$ | $o_1$ | $a_2$ | $o_2$ | $a_3$ | $o_3$ | | | |
|---|---|---|---|---|---|---|---|---|
| 6 | 3 | 9 | 5 | 3 | 2 | | | |

Environment can be a bandit, MDP, POMDP or any other **controllable stochastic process**.

# Adaptive Control



$$P(a_1)$$

$$P(o_1|a_1)$$

$$P(a_2|a_1,o_1)$$

$$P(o_2|a_1,o_1,a_2)$$

$$U(a_1,o_1,a_2,o_2)$$

**In theory:**

- Choose policy maximizing **subjective expected utility**.

**In practice: intractable!**

- Policy space **grows exponentially** with planning horizon.
- Policy choice **causally precedes** interactions.

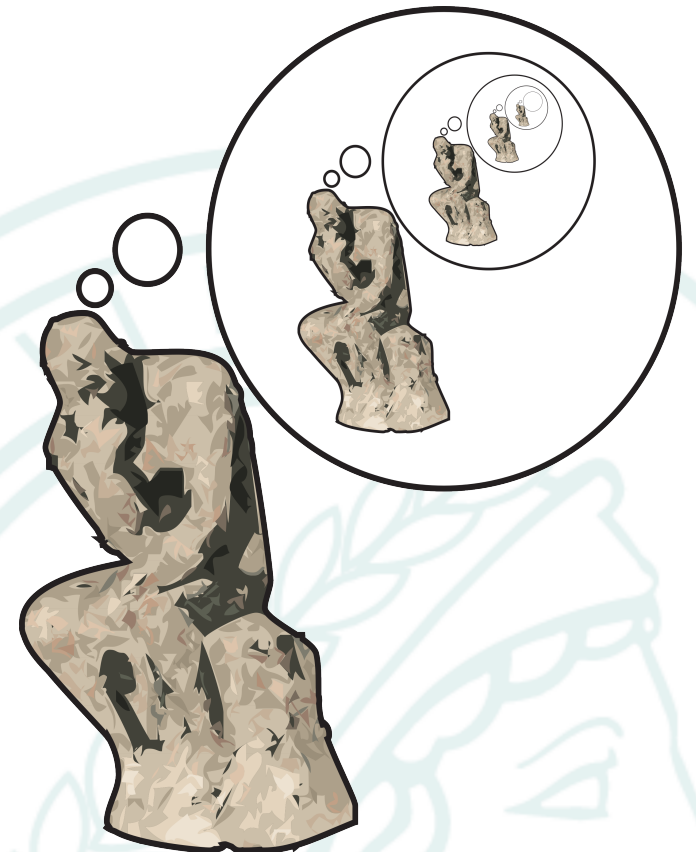# Choose policy **before** interacting?

**What if choosing the optimal policy was tractable?**

This implies:

- precomputing **all the possible lives**,

- and then picking the **optimal policy**.

However:

- Prediction has no accuracy, because it is **not supported** by any data.

- The optimal policy is **statistically meaningless in the beginning**!

# Can we choose policies **dynamically**?

- **Delay** choice of optimal policy – when **justified** by data.

- Agent is **uncertain** about the optimal policy.

- **Practical** adaptive control and RL **do this** explicitly/implicitly.
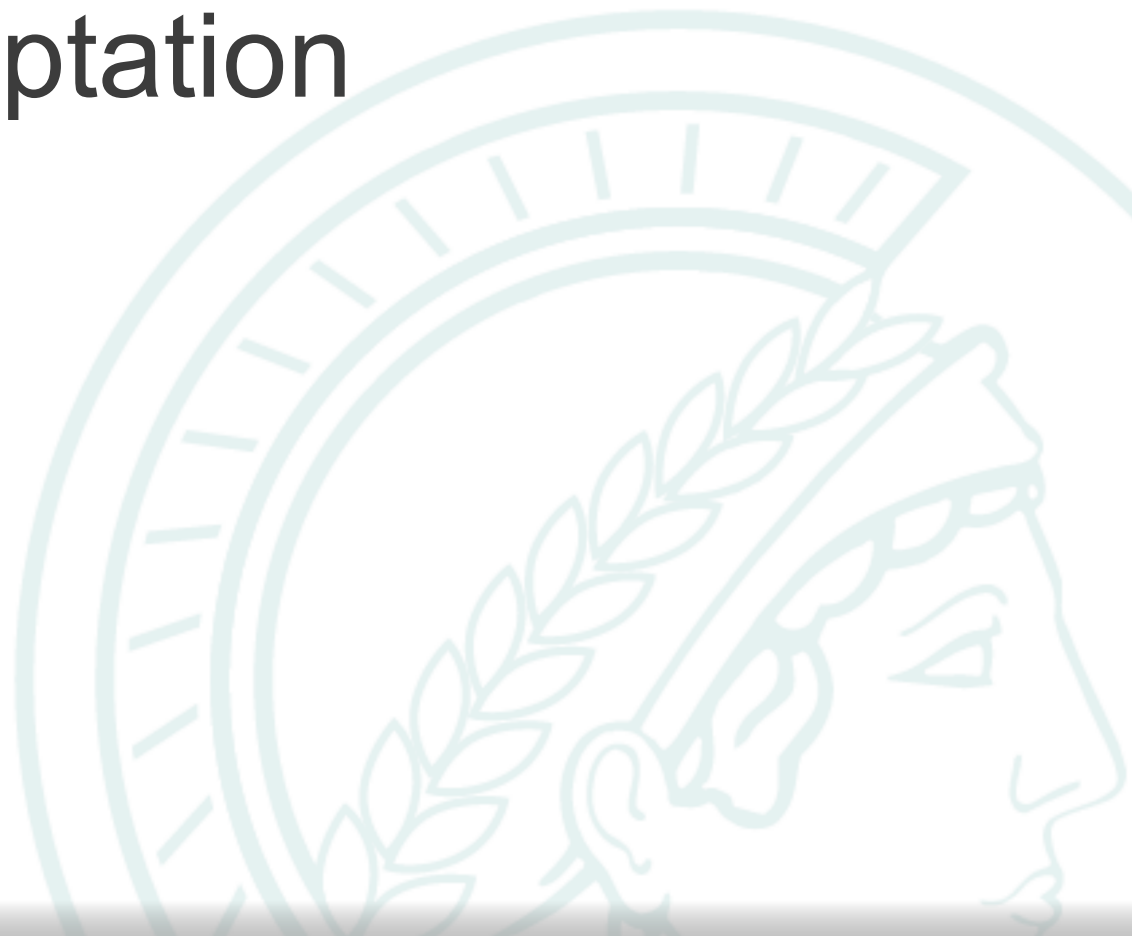
- Implementation of "**intuition**"

# Questions

How do we choose the optimal policy **dynamically**?

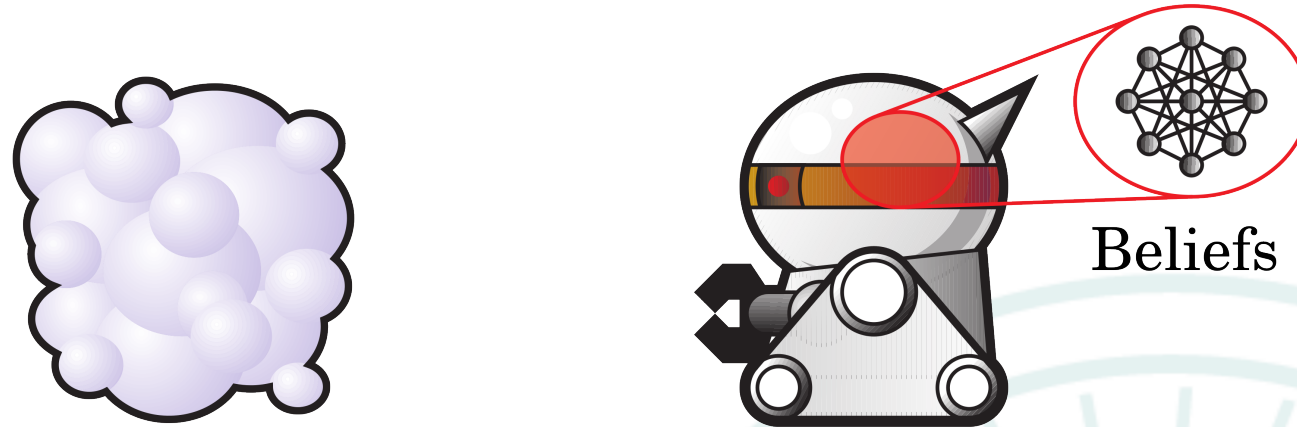How is uncertainty over the policy **represented**?

How are **actions issued** when the policy is uncertain?
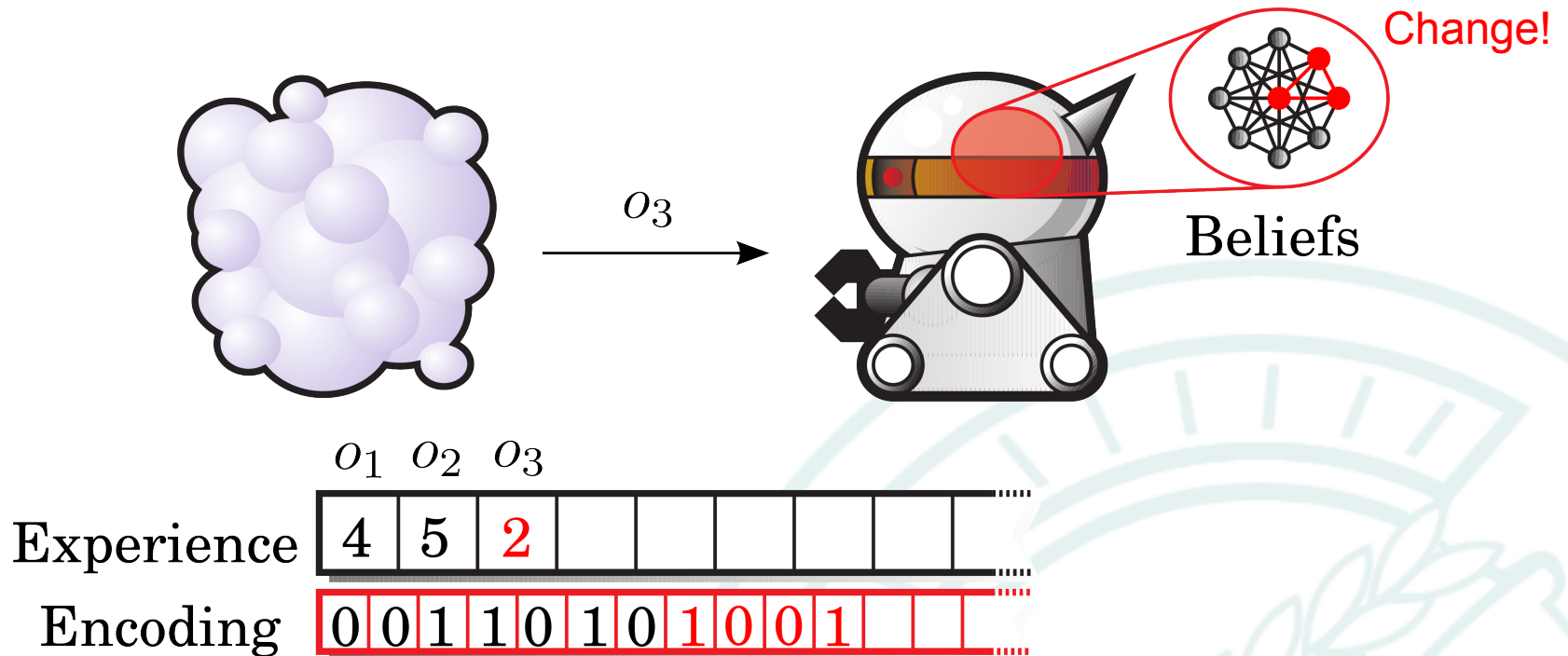
How is this uncertainty **reduced**?

# Adaptation

# The Cost of Experience



Beliefs

|  | $o_1$ | $o_2$ | $o_3$ | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Experience | 4 | 5 | | | | | | | |

| Encoding | 0 | 0 | 1 | 1 | 0 | 1 | 0 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

- Agent records observations.

# The Cost of Experience



Change!

Beliefs

$o_1$ $o_2$ $o_3$

| | Experience | 4 | 5 | 2 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Encoding | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |

- Agent records observations.
- Acquiring experience implies **changes** in the belief structure.

# The Cost of Experience



Change!

$o_3$

Beliefs

$$o_1 \quad o_2 \quad o_3$$

| Experience | 4 | 5 | 2 | | | | | | |
|------------|---|---|---|---|---|---|---|---|---|

| Encoding | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | | |
|----------|---|---|---|---|---|---|---|---|---|---|---|---|---|

- Agent records observations.

- Acquiring experience implies **changes** in the belief structure.

- **Can we minimize these changes?**

# Adaptive Compression

- When the environment is **known**, maximal compression is achieved when codeword lengths are chosen as

$$l(o_{\leq t}) := -\log Q(o_{\leq t})$$

- Conversely, every code **implies predictions**

$$P(o_{\leq t}) = 2^{-l(o_{\leq t})}$$

- The belief structure of the agent **embodies the assumptions** about the environment.

# Adaptive Compression (cont.)

- How to compress when the environment is **unknown**?

- Consider set of possible environments $\Theta$, with probabilities $P(\theta)$ and models $P(o_{\leq t}|\theta)$.

- Choose a predictor $\tilde{P}$ minimizing expected codeword length:

$$L_t[\tilde{P}] = \underbrace{\sum_\theta P(\theta)}_{\text{Choice of } \theta} \left\{ \sum_{o_{\leq t}} P(o_{\leq t}|\theta) \log \frac{\overbrace{P(o_{\leq t}|\theta)}^{\text{Environment } \theta}}{\underbrace{\tilde{P}(o_{\leq t})}_{\text{Predictor}}} \right\}$$

- Solution: **Bayesian mixture**

$$\tilde{P}(o_{\leq t}) := \sum_{\theta} P(o_{\leq t}|\theta)P(\theta) = P(o_{\leq t})$$

- Predictive distribution

$$P(o_t|o_{<t}) = \sum_{\theta} P(o_t|o_{<t})P(\theta|o_{<t})$$

- Bottom line: adaptive compression is solved by **pretending** that the Bayesian mixture is the true environment
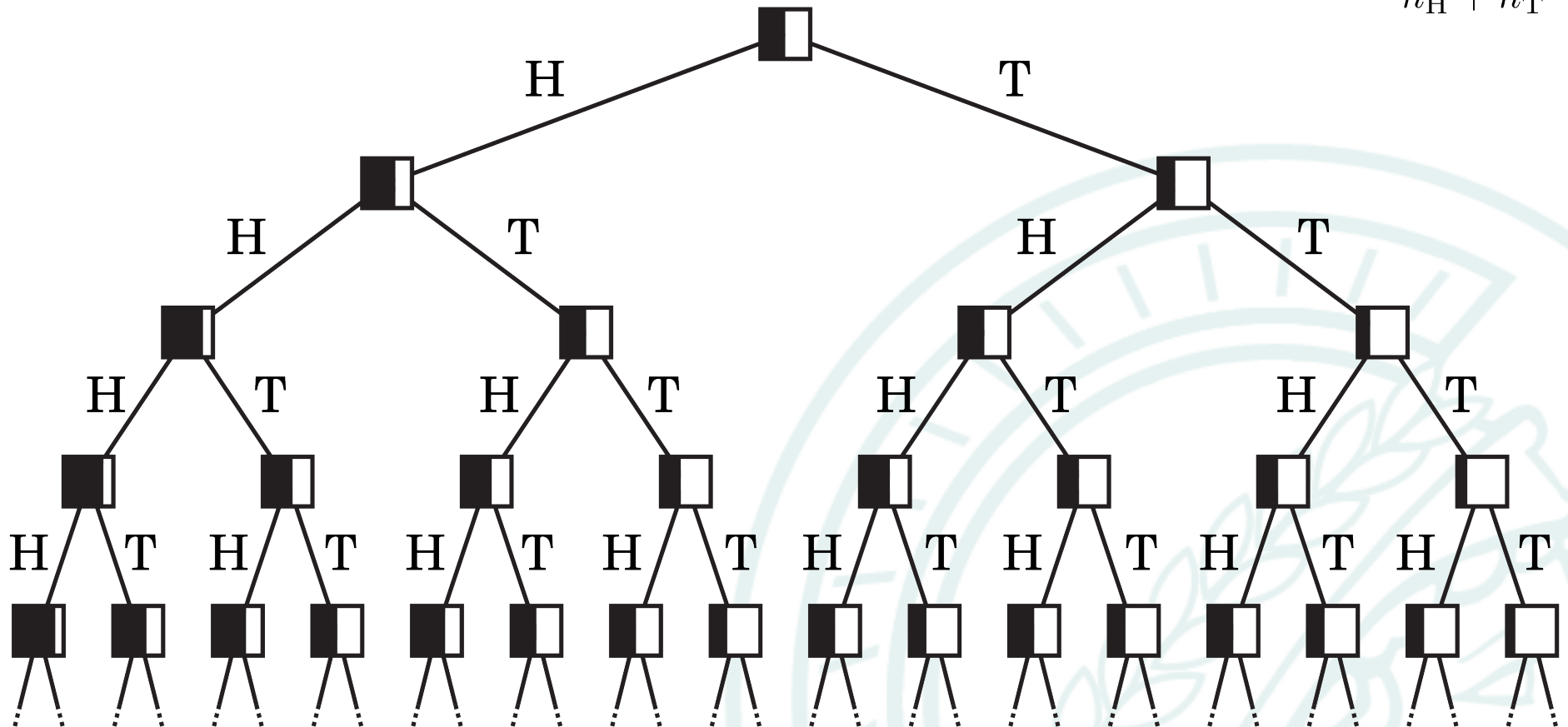
bias $\theta = \frac{3}{4}$

bias $\theta = \frac{1}{4}$

$P(\theta) = \frac{1}{2}$

mixture

$$P(\text{H}|\text{experience}) = \frac{n_\text{H} + 1}{n_\text{H} + n_\text{T} + 2}$$

mixture over all biases in [0,1]

# Summary

The Bayesian mixture is the optimal compressor of experience for an unknown environment.
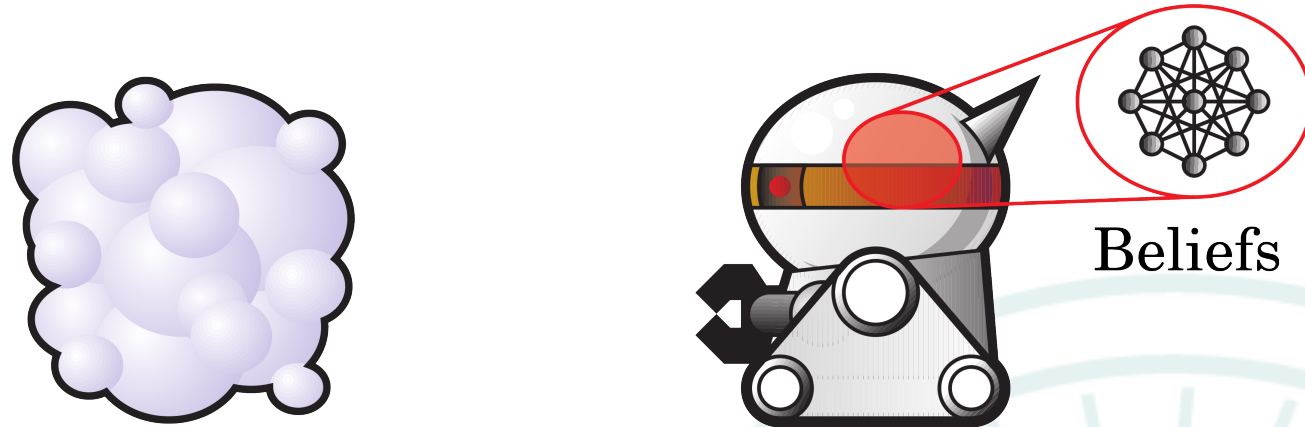
# Causality

# Extension to Actions

- Can we use this for **adaptive behavior**?

- Instead of **competing hypotheses**, we would have **competing behaviors** $(\theta, \pi) \in \Theta \times \Pi$ :

$$P(a_{\leq t}, o_{\leq t} | \theta, \pi) \qquad P(\theta, \pi)$$

- Would lead to

$$P(\text{next action} | \text{experience}) = P(a_t | a_{<t}, o_{<t})$$
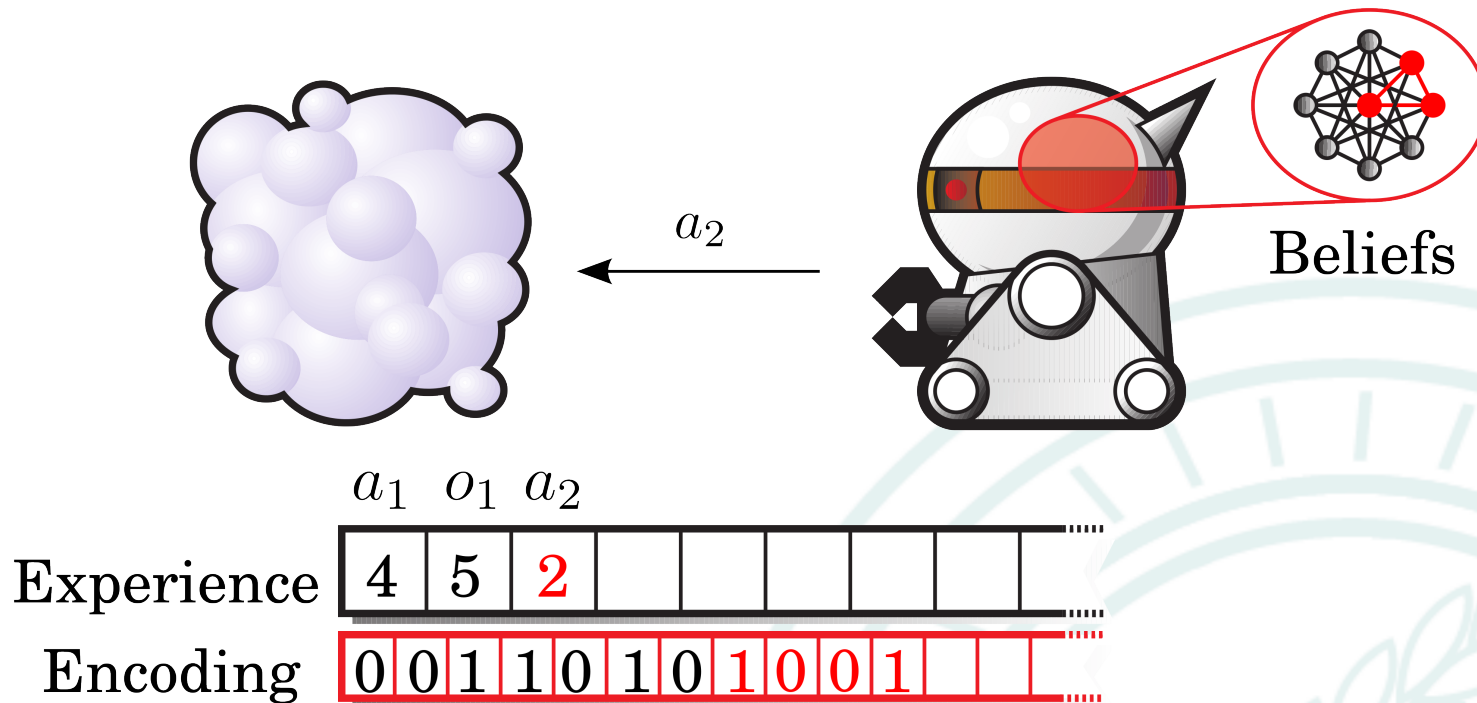
# The Cost of Experience II



$a_1$ $o_1$ $a_2$

| Experience | 4 | 5 | | | | | | |
|---|---|---|---|---|---|---|---|---|

| Encoding | 0 | 0 | 1 | 1 | 0 | 1 | 0 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

- Agent records actions & observations.

# The Cost of Experience II



Beliefs

$$a_1 \quad o_1 \quad a_2$$

| Experience | 4 | 5 | 2 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|

| Encoding | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | | |

- Agent records actions & observations.
- Again, actions **change** the belief structure.

# The Cost of Experience II



Beliefs

$a_2$

$a_1$ $o_1$ $a_2$

Experience | 4 | 5 | 2 | | | | | |

Encoding | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | | |

- Agent records actions & observations.
- Again, actions **change** the belief structure.
- However, they **do not change the beliefs**.

# Problems of the Bayesian Update

- Posterior beliefs

$$P(\theta, \pi | a_t, o_t, \ldots)$$

$$\propto \text{likelihood} \times \text{prior}$$

$$= P(o_t | \theta, a_t, \ldots) P(a_t | \pi, \ldots) \times P(\theta, \pi | \ldots)$$

# Problems of the Bayesian Update

- Posterior beliefs

$$P(\theta, \pi | a_t, o_t, \ldots)$$

$$\propto \text{likelihood} \times \text{prior}$$

$$= P(o_t | \theta, a_t, \ldots) P(a_t | \pi, \ldots) \times P(\theta, \pi | \ldots)$$

...but **our actions produce evidence**, we conclude the optimal policy from our own actions.

# Problems of the Bayesian Update

- Posterior beliefs

$$P(\theta, \pi | a_t, o_t, \ldots)$$

$$\propto \text{likelihood} \times \text{prior}$$

$$= P(o_t | \theta, a_t, \ldots) P(a_t | \pi, \ldots) \times P(\theta, \pi | \ldots)$$

...but **our actions produce evidence**, we conclude the optimal policy from our own actions.

# Problems of the Bayesian Update

- Posterior beliefs

$$P(\theta, \pi | a_t, o_t, \ldots)$$

$$\propto \text{likelihood} \times \text{prior}$$

$$= P(o_t | \theta, a_t, \ldots) P(a_t | \pi, \ldots) \times P(\theta, \pi | \ldots)$$

...but **our actions produce evidence**, we conclude the optimal policy from our own actions.

- **We cannot change events that causally precede the present.**

# Causality

- Solution: treat actions as **causal interventions**

$$P(\theta, \pi | \hat{a}_t, o_t, \ldots)$$

$$\propto \text{likelihood} \times \text{prior}$$

$$= P(o_t | \theta, \hat{a}_t, \ldots) P(\hat{a}_t | \pi, \ldots) \times P(\theta, \pi | \ldots)$$

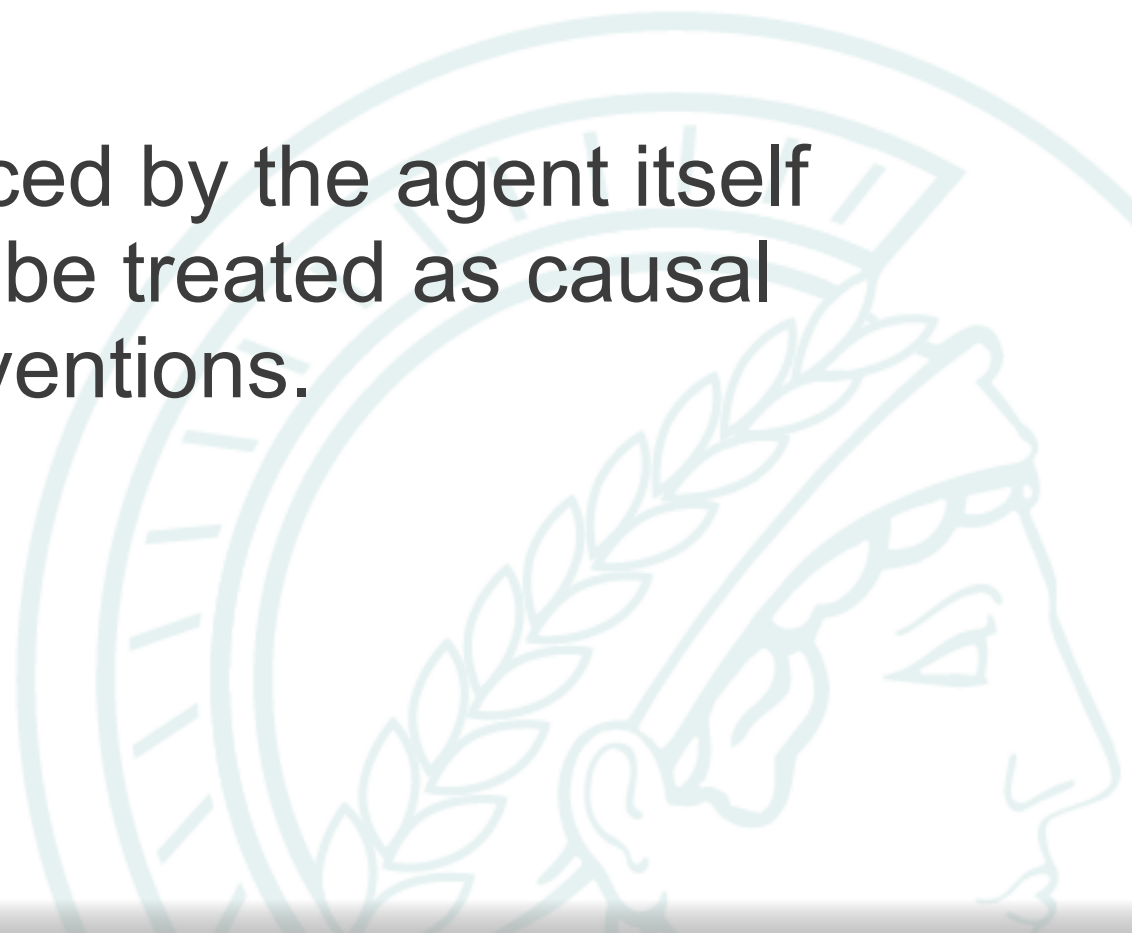$$= P(o_t | \theta, a_t, \ldots) \times P(\theta, \pi | \ldots)$$

# Causality

- Solution: treat actions as **causal interventions**

$$P(\theta, \pi | \hat{a}_t, o_t, \ldots)$$

$$\propto \text{likelihood} \times \text{prior}$$

$$= P(o_t | \theta, \hat{a}_t, \ldots) P(\hat{a}_t | \pi, \ldots) \times P(\theta, \pi | \ldots)$$

$$= P(o_t | \theta, a_t, \ldots) \times P(\theta, \pi | \ldots)$$

- Causal intervention informs us that we have to **ignore the evidence** produced by the action.

# Causality

- Solution: treat actions as **causal interventions**

$$P(\theta, \pi | \hat{a}_t, o_t, \ldots)$$

$$\propto \text{likelihood} \times \text{prior}$$
$$= P(o_t | \theta, \hat{a}_t, \ldots) P(\hat{a}_t | \pi, \ldots) \times P(\theta, \pi | \ldots)$$
$$= P(o_t | \theta, a_t, \ldots) \times P(\theta, \pi | \ldots)$$

- Causal intervention informs us that we have to **ignore the evidence** produced by the action.
- Caveat: $\qquad \pi = \pi(\theta)$

# Summary

Actions are produced by the agent itself and thus need to be treated as causal interventions.

# Bayesian Control Rule

# Bayesian Control Rule

Given a set $\Theta$ of

- behaviors $P(a_{\leq t}, o_{\leq t} | \theta)$

- prior probabilities $P(\theta)$

**sample** actions from $P(a_t | \hat{a}_{<t}, o_{<t})$

# Bayesian Control Rule (cont.)

**Time $t$**

$\theta$

$\theta^*$

$\theta$

$\theta$

Prior:

$$P(\theta|\hat{a}_{<t}, o_{<t})$$

Acting:

$$\theta^* \sim P(\theta|\hat{a}_{<t}, o_{<t})$$
$$a_t^* \sim P(a_t|\theta^*, \hat{a}_{<t}, o_{<t})$$

Observing:

$$P(\theta|\hat{a}_{\leq t}, o_{\leq t})$$
$$\propto P(\theta|\hat{a}_{\leq t}, o_{\leq t})P(o_t|\theta, a_{\leq t}, o_{<t})$$

**Time $t+1$**

$\theta$

Posterior:

$$P(\theta|\hat{a}_{\leq t}, o_{\leq t})$$

# Example: 2-Armed Bandit

- Bernoulli-distributed rewards, unknown biases.

- Hypotheses: $\Theta = [0, 1] \times [0, 1]$

- Prior: $P(\theta) = U(0, 1) \times U(0, 1)$

- Observations: $P(o|\theta, a) = B(o; \theta_a)$
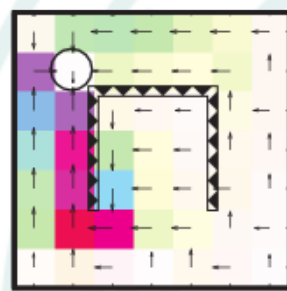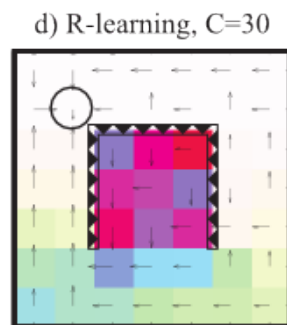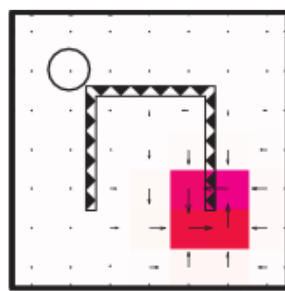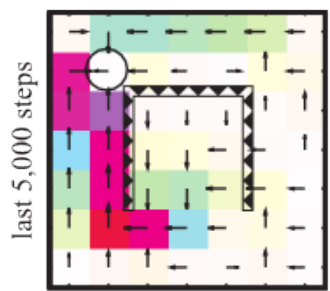
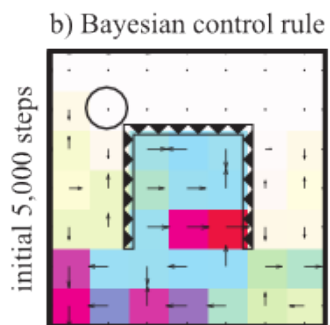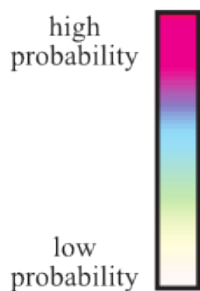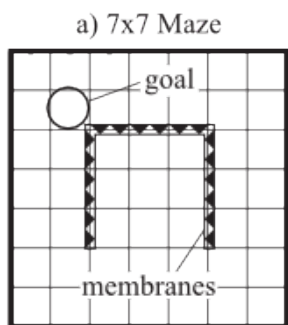- Actions: $P(a|\theta) = \delta_a^{\arg\max_i \theta_i}$

# Example: 2-Armed Bandit

- Bernoulli-distributed rewards, unknown biases.
- Hypotheses: $\Theta = [0, 1] \times [0, 1]$
- Prior: $P(\theta) = U(0, 1) \times U(0, 1)$
- Observations: $P(o|\theta, a) = B(o; \theta_a)$
- Actions: $P(a|\theta) = \delta_a^{\arg\max_i \theta_i}$

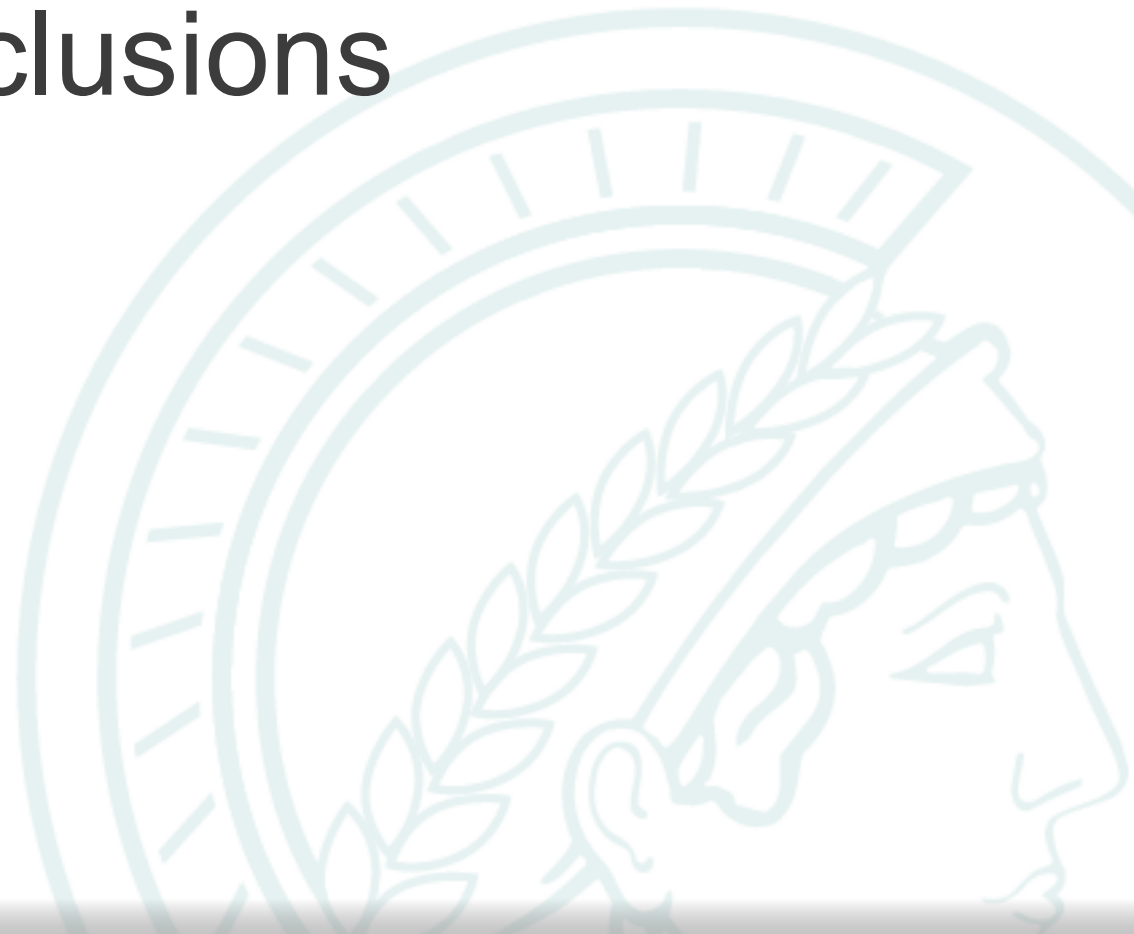- Recently proven to be **asymptotically optimal** [Kaufmann, Korda, Munos 2012].

# Results for 10-Armed Bandit

# Markov Decision Processes



a) 7x7 Maze — goal, membranes

b) Bayesian control rule

c) R-learning, C=5

d) R-learning, C=30

initial 5,000 steps / last 5,000 steps

high probability / low probability

e) Average Reward — R-learning, C=30; R-learning, C=5; Bayesian control rule

x1000 time steps
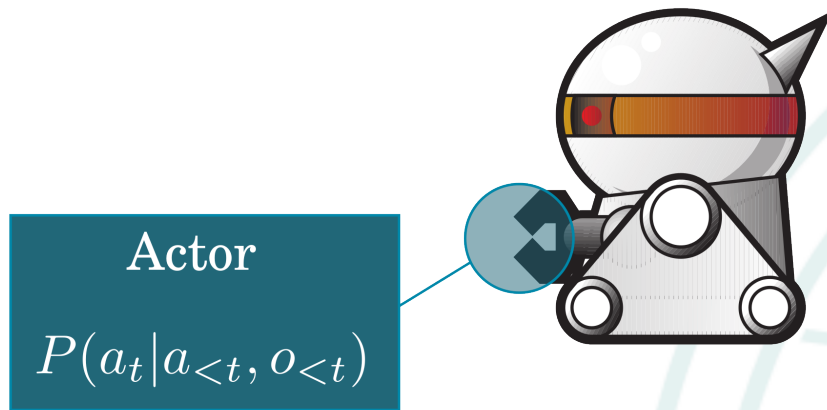
# Conclusions

# Blueprint of the Agent

# Blueprint of the Agent



Actor

$$P(a_t | a_{<t}, o_{<t})$$

# Blueprint of the Agent

**Predictor**

$$P(o_t | a_{\leq t}, o_{<t})$$

**Actor**

$$P(a_t | a_{<t}, o_{<t})$$

# Blueprint of the Agent



compiles to

**Predictor**

$$P(o_t | a_{\leq t}, o_{<t})$$

**Actor**

$$P(a_t | a_{<t}, o_{<t})$$

**Hypotheses**

$$P(o_t | \theta, a_{\leq t}, o_{<t})$$

# Blueprint of the Agent

# Blueprint of the Agent



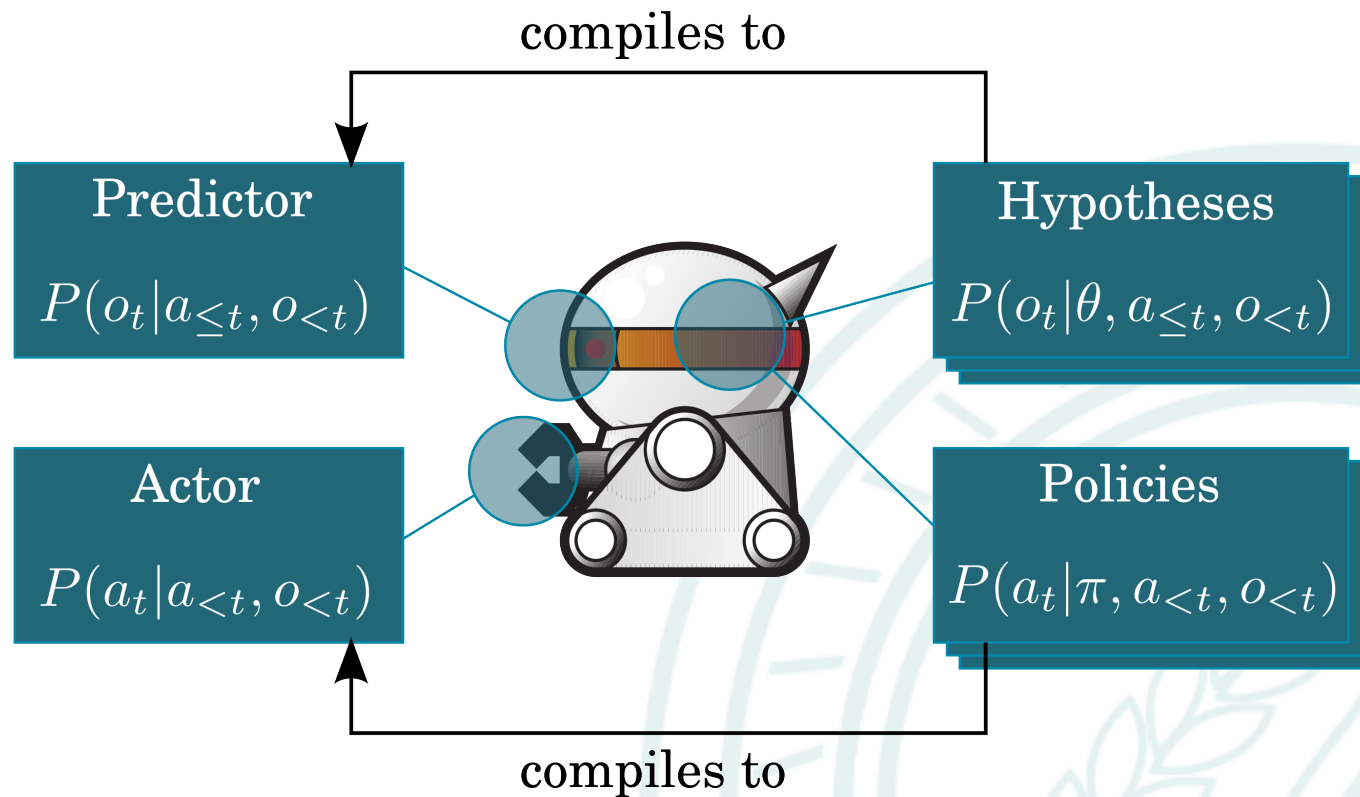compiles to

| Predictor | Hypotheses |
|---|---|
| $P(o_t \mid a_{\leq t}, o_{<t})$ | $P(o_t \mid \theta, a_{\leq t}, o_{<t})$ |

desired coupling

| Actor | Policies |
|---|---|
| $P(a_t \mid a_{<t}, o_{<t})$ | $P(a_t \mid \pi, a_{<t}, o_{<t})$ |

compiles to

# Properties

- Stochastic controller that **refines its policy with experience.**

- Ingredients: **Bayes + Causality**.

- Transforms control into inference.

- Related to **Random Beliefs** & **Thompson sampling.**

- Allows tackling **game-theoretic** problems.

- Exploits **built-in reward mechanism** of Bayes' rule.

- Works also with **complex causal models**.

# Pros and Cons

## Pros

- Simple and general.
- Converges to desired behavior in "ergodic" tasks.
- Suitable for on-line.
- Trades-off exploration versus exploitation.
- Automatic temporal credit assignment.

## Cons

- Sub-optimal in the transient.
- Does not converge in non-ergodic environments.
- Convergence speed highly depends on environment.
- Design of behaviors can be difficult.