

# **Learning in Neural Networks (week10)**

## **The role of exploration, novelty, and surprise in RL**

### **Objectives for today:**

- how to encourage exploration**
- understand surprise**
- understand difference of novelty and surprise**
- use of surprise to modulate learning rate**
- use of novelty to guide exploration**

Previous slide.

## Background reading:

[An analysis of model-based Interval Estimation for Markov Decision Processes](#)

Strehl and Littman, 2008

<https://www.sciencedirect.com/science/article/pii/S0022000008000767>

[Novelty is not Surprise: Human exploratory and adaptive behavior in sequential decision-making](#)

HA Xu\*, A Modirshanechi\*, MP Lehmann, W Gerstner, MH Herzog, PLOS Comput. Biol. E1009070, (2021)

[Learning in Volatile Environments with the Bayes Factor Surprise](#)

V Liakoni\*, A Modirshanechi\*, W Gerstner, J Brea  
Neural Computation 33 (2), 269-340 (2021)

[A taxonomy of surprise definitions](#)

A Modirshanechi, J Brea, W Gerstner  
Journal of Mathematical Psychology 110, 102712 (2022)

# Novelty and Surprise

Q1: Why does an agent need to explore?

Q2: What is novelty?

Q3: What is surprise?

Q4: What is the difference between the two?

Q5: Why are they useful?

Q6: Why should we talk about novelty in this class?

Previous slide.

Today we will ask 6 questions:

Why is exploration important?

What is novelty, What is surprise, What is the difference, Why are they useful.

And why should we talk about it in this class?

# Learning in Neural Networks

Wulfram Gerstner

EPFL, Lausanne, Switzerland

## The role of exploration, novelty, and surprise in RL

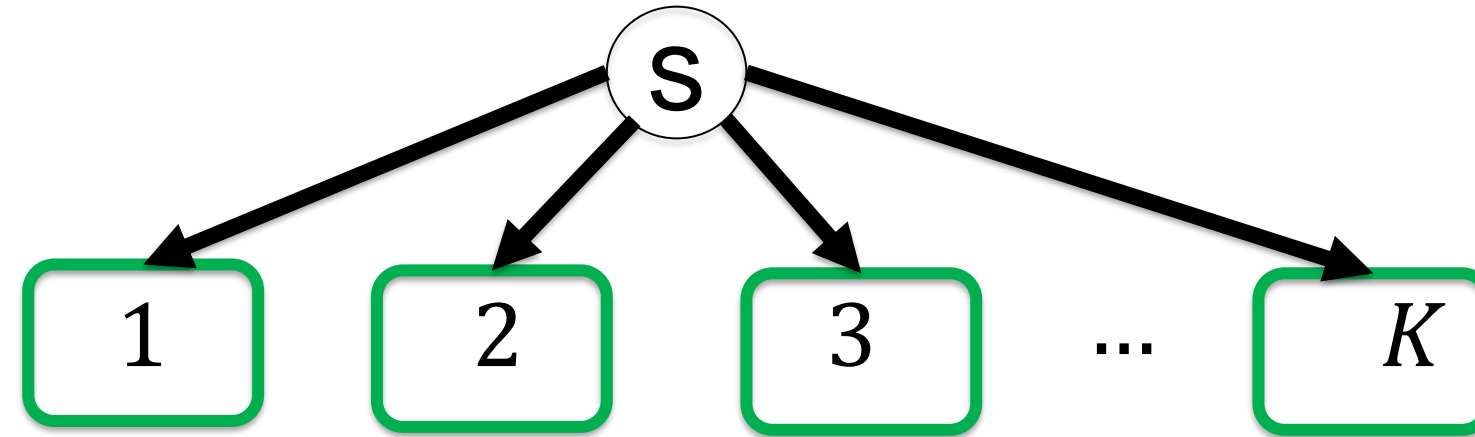
### 1. Exploration and intrinsic 'rewards'

Previous slide.

We start with recalling the problem of exploration-exploitation

# Review: Multi-armed Bandits: MAB (1-step horizon)

- Single state. We have  $K$  possible actions:



Which action to choose at time  $t$ ?

- With true average reward:

$$\mu_i = E[r|a = i]$$

 $\mu_1$  $\mu_2$  $\mu_3$  $\dots$  $\mu_K$ 

Optimal policy:  $a_t = \mathop{\text{arg max}}_i \mu_i$

- Naïve estimates of averages:

$$\hat{\mu}_i^{(t)} = \frac{\sum_{\tau \in T_i^{(t)}} r_\tau}{|T_i^{(t)}|}$$

 $\hat{\mu}_1^{(t)}$  $\hat{\mu}_2^{(t)}$  $\hat{\mu}_3^{(t)}$  $\dots$  $\hat{\mu}_K^{(t)}$ 

$$T_i^{(t)} = \{\tau \leq t : a_\tau = i\}$$

Not optimal:  $a_t = \mathop{\text{arg max}}_i \hat{\mu}_i^{(t)}$

Solutions based on random exploration:

- Epsilon-greedy
- Softmax

## Comments for the previous slide:

- If we knew the exact average reward  $\mu_i = E[r|a = i]$  of each arm, then the optimal solution would trivially be to choose the arm with highest average reward:  $a_t = \arg \max_i \mu_i$
- A naïve approach is to estimate the average reward by the empirical averages and greedily choose the action with maximum estimated average reward:  $a_t = \arg \max_i \hat{\mu}_i^{(t)}$
- The naïve greedy policy is prone to fail in finding the best action.
- You have seen epsilon-greedy and the softmax policy as two approaches for dealing with this problem by adding randomness to the action-selection.
- Our focus will be on “directed exploration” by using exploration bonuses.



# Example: average rewards in MAB (1-step horizon)

- MAB with 4 possible actions (Example):

$$\mu_i = E[r|a = i]$$

rewards are stochastic (binomial)

$$P(r_t = 2\mu_i | a = i) = 0.5 = P(r_t = 0 | a = i)$$

$$\mu_1 = 1$$

$$\mu_2 = 0.9$$

$$\mu_3 = 9.9$$

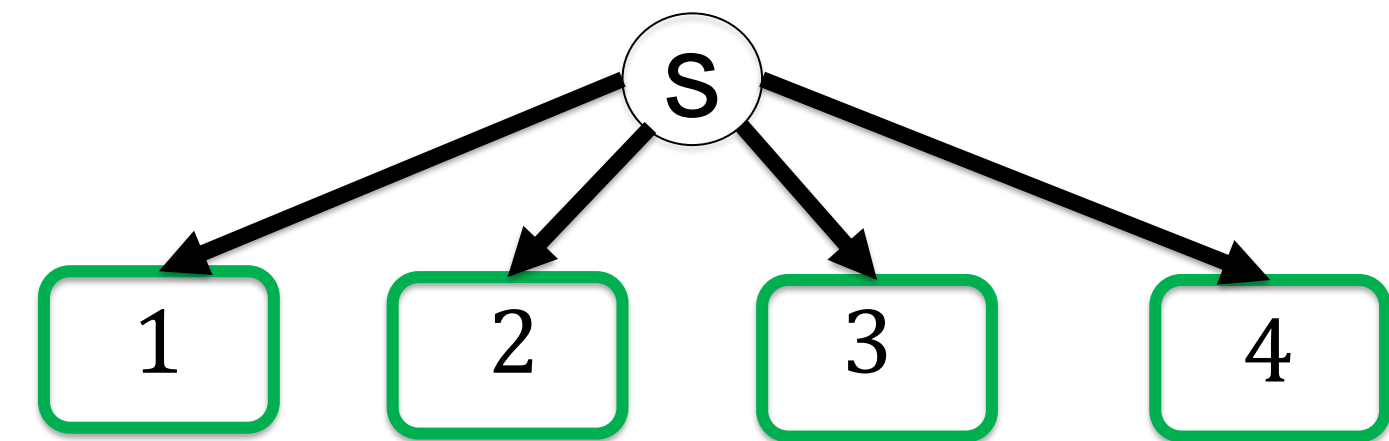
$$\mu_4 = 10.0$$

$$\hat{\mu}_i^{(t)} = \frac{\sum_{\tau \in T_i^{(t)}} r_\tau}{|T_i^{(t)}|}$$

**Idea:** explore while tails of distributions overlap

1

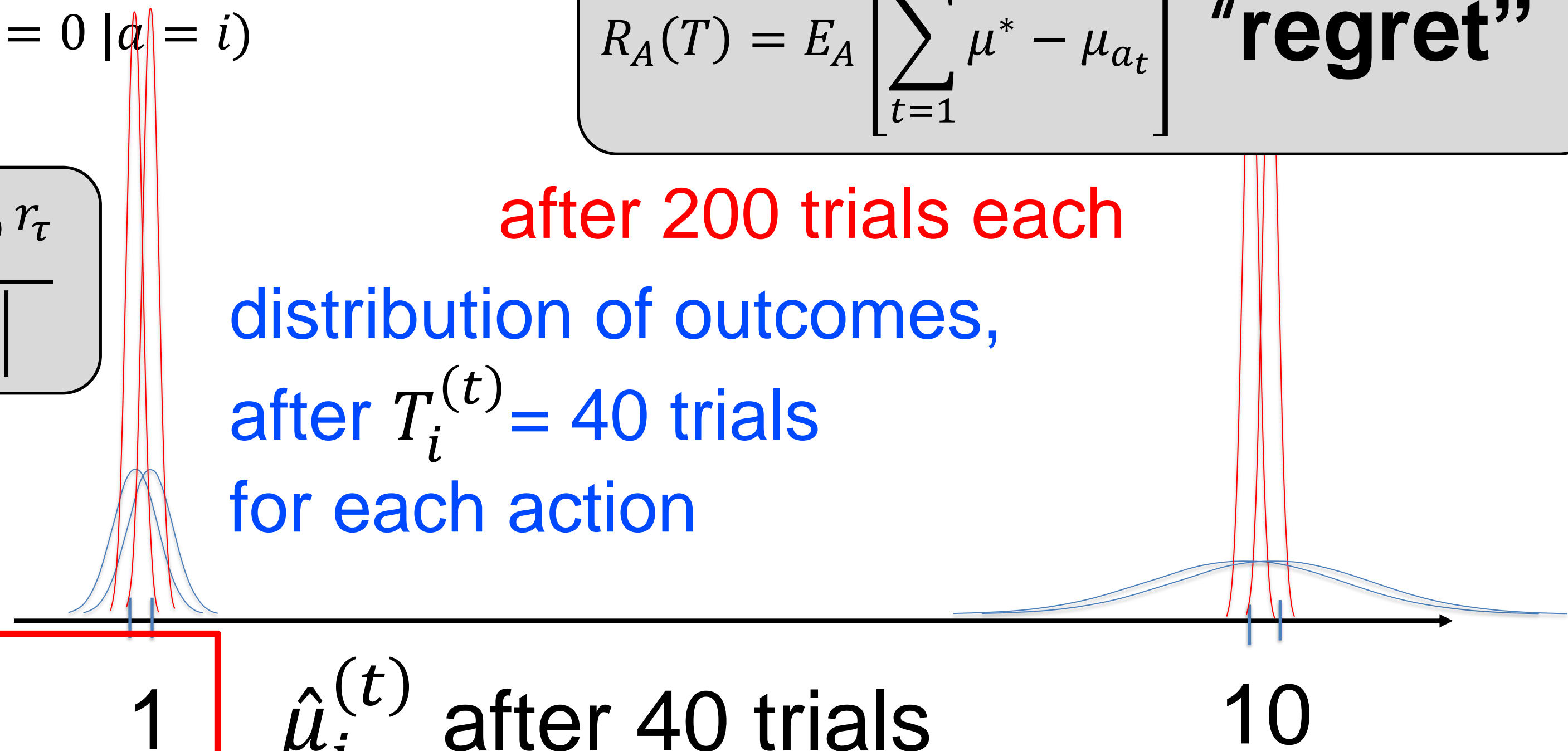
$\hat{\mu}_i^{(t)}$  after 40 trials for each action



$$R_A(T) = E_A \left[ \sum_{t=1}^T \mu^* - \mu_{a_t} \right] \text{ "regret"}$$

after 200 trials each

distribution of outcomes,  
after  $T_i^{(t)} = 40$  trials  
for each action



10

- Comments for the previous slide:
- Example of MAB with 4 actions. Each action yields a reward with 50 percent probability.
- Two actions have low rewards (about 1); the two other have high rewards about 20.
- Imagine that at the beginning you played each action 40 times and evaluate the mean return.
- If you repeated the game many times, each time starting with playing each action 40 times, you would get a distribution (hand-drawn here).
- As long as the distributions overlap, we continue to play all actions. Hence, after  $t=160$ , we should continue to play actions 3 and 4, while actions 1 and 2 can be safely dropped as a possibility.
- Note that the width of the distributions can be pre-calculated with a model of stochasticity

# Exploration Bonus (multi-step horizon)

Assuming we know the true  $P(s'|s, a)$  and  $R(s, a)$ , the Bellman equation is

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q^*(s', a'); \quad a_t = \arg \max_a Q^*(s_t, a)$$

If we do not know the true  $P(s'|s, a)$  and  $R(s, a)$ , the Bellman equation can be replaced by

$$\hat{Q}_{\text{MB}}^{(t)}(s, a) = \hat{R}^{(t)}(s, a) + \text{Bonus}(s, a) + \gamma \sum_{s'} \hat{P}^{(t)}(s'|s, a) \max_{a'} \hat{Q}_{\text{MB}}^{(t)}(s', a'); \quad a_t = \arg \max_a \hat{Q}_{\text{MB}}^{(t)}(s_t, a)$$

One of the choices is 
$$\text{Bonus}(s, a) = \frac{\beta}{\sqrt{T_{s,a}^{(t)}}}$$

- **Adding exploration bonus provably improves the performance of RL algorithms.**
- **Hence, to optimally seek a reward, best seek a ‘modified reward’ .**

- Comments for the previous slide:

$\hat{P}^{(t)}(s'|s, a)$  is the estimated transition probability

The theory itself is not treated in class. The main point of the slide to show that in CS and RL, the problem of exploration is known to be important and has been addressed.

There are several ‘good’ choices to add an exploration bonus.

The exploration bonus is integrated into the standard Bellman equation.

# Exploration:

## Standard theory of exploration assumes:

- stationary environment (but the world is not stationary!)
- many trials (but does a biological agent has the time for many repetitions?)
- model-based RL with update of Bellman equation in background  
(but is that plausible for biological agents?)
- adds an exploration bonus into the **general Bellman equation** for reward  
(but why treat search for novelty as the 'same kind of' reward?)

- Comments for the previous slide:

The assumptions of the theory do not fit with what we know about biological agents!



# Review: Neuromodulators

- 4 or 5 neuromodulators
- near-global action
- internally created signals

Dopamine/reward/TD:  
*Schultz et al., 1997,*  
*Schultz, 2002*

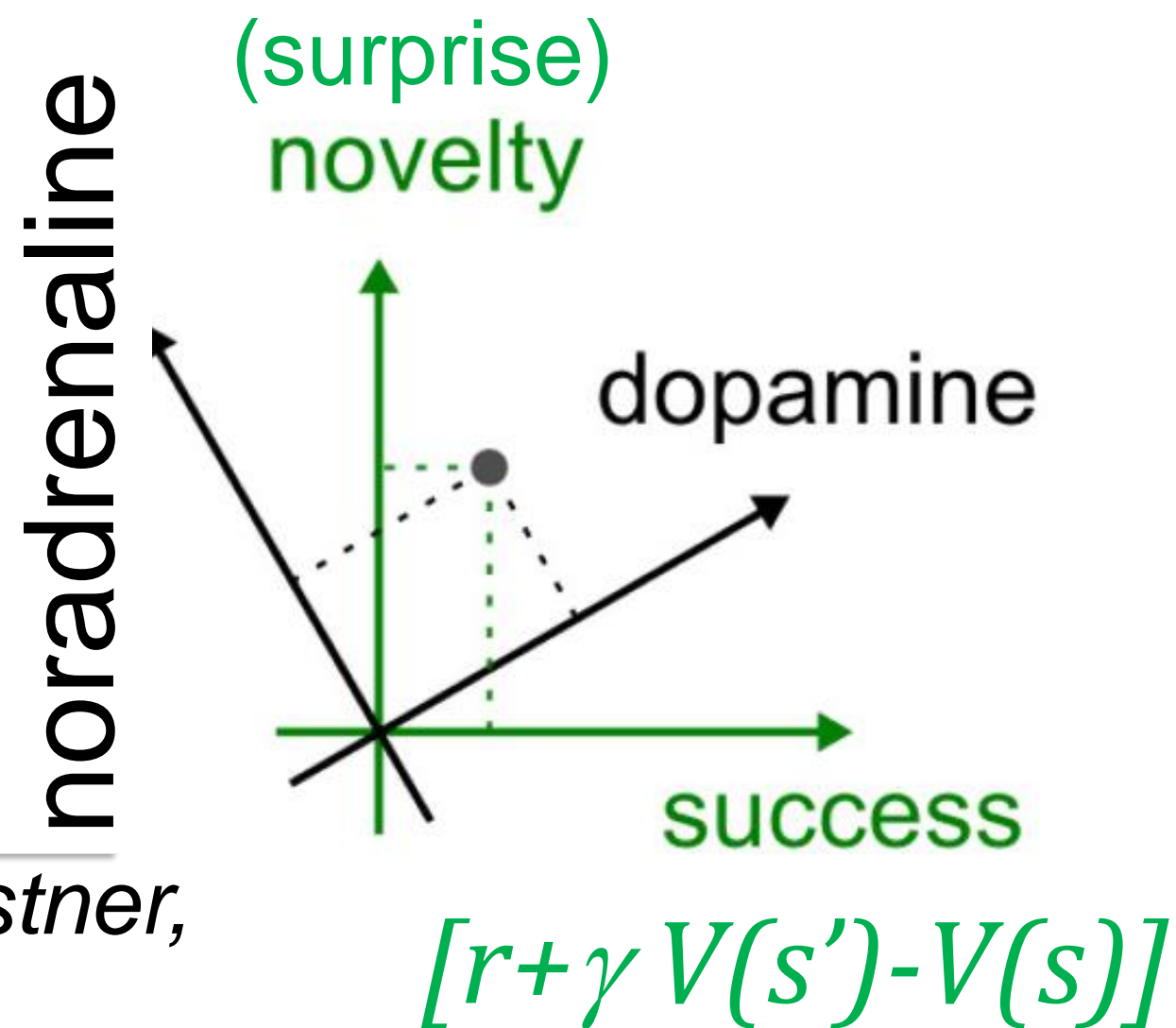
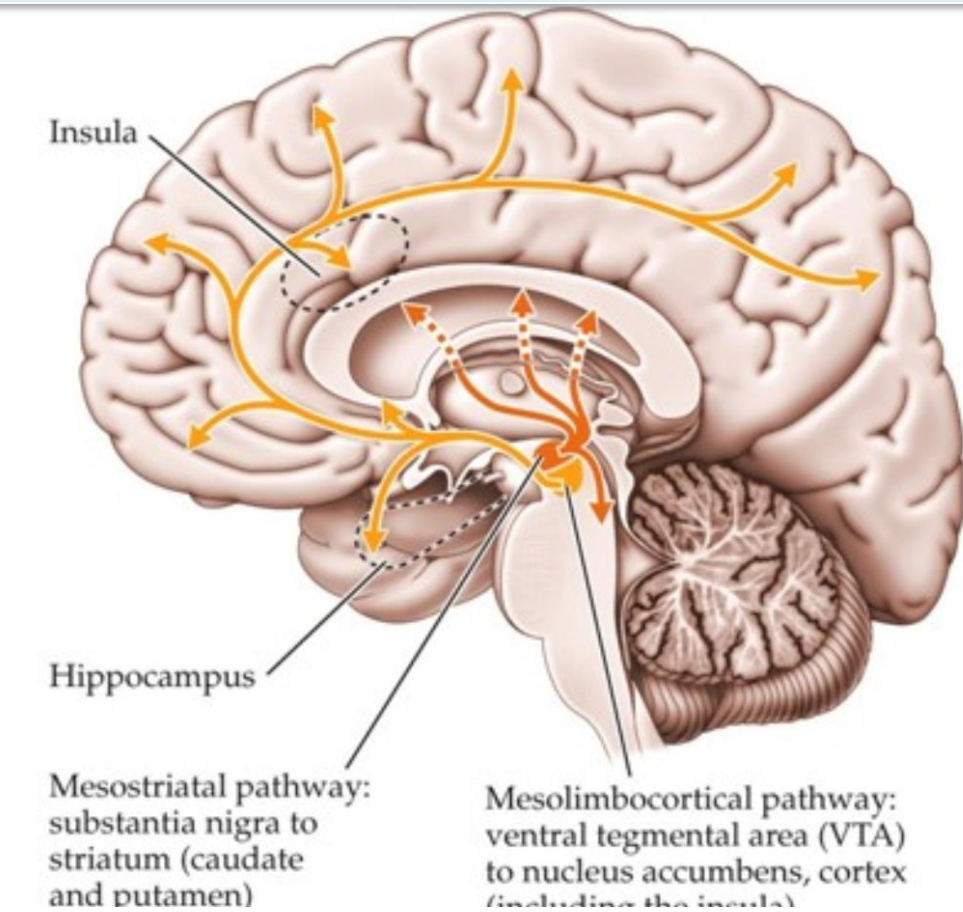


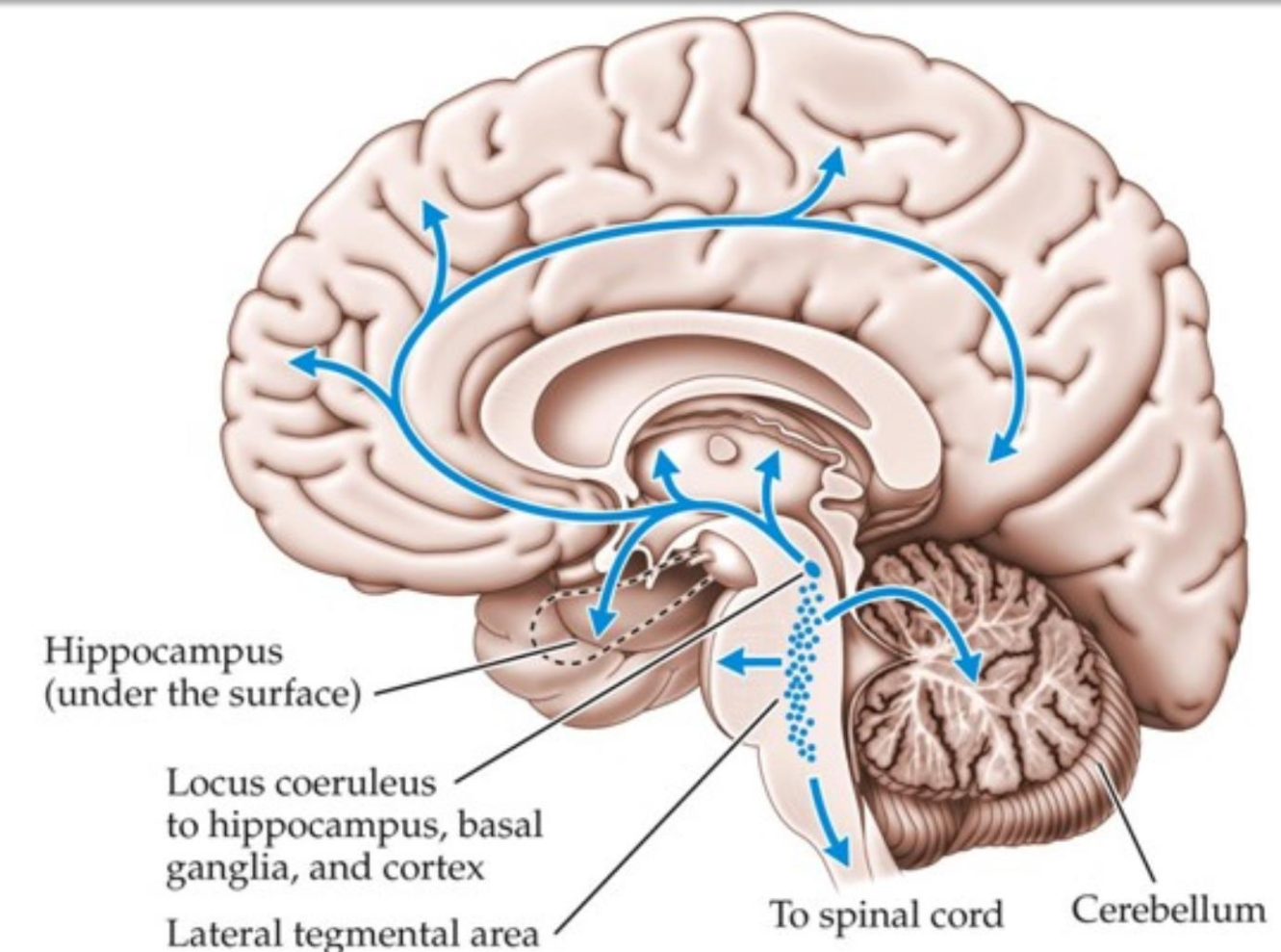
Image:  
*Fremaux and Gerstner,*  
*Frontiers (2016)*

Image: *Biological Psychology, Sinauer*

## Dopamine (DA)



## Noradrenaline (NE)



Previous slide. Review

The most famous neuromodulator is dopamine (DA) which is related to reward, as we will see.

But there are other neuromodulators such as noradrenaline (also called norepinephrine, NE) which is related to surprise.

Left: the mapping between neuromodulators and functions is not one-to-one.

Indeed, dopamine also has a 'surprise' component.

Inversely, noradrenaline also has a reward component.

Right: most neuromodulators send axons to large areas of the brain, in particular to several cortical areas. The axons branch out in thousands of branches.

Thus the information transmitted by a neuromodulator arrives nearly everywhere. In this sense, it is a 'global' signal, available in nearly all brain areas.

Note that the TD error is an internally created signal. The TD can be positive at time  $t$  even if no explicit reward is given at time  $t$ .

Similarly, surprise is an internally generated signal indicating model mismatch.



## Review: Neuromodulators (NM)

- several neuromodulators: not just dopamine but also others
- each has near-global action
- internally created signals, could be used as 'internal rewards'
- they are related, but **different** signals
- combinations of neuromodulators relay different functions
  - (i) changes in environment → surprise signal
  - (ii) exploration of environment → novelty signal
  - (iii) unexpected reward → success signal

'general emotional brain signals are related to neuromodulators'

No reason to combine them into a single Bellman equation!

Previous slide.

Since there are different neuromodulator signals that code for different functionalities, there is no need to combine them into a SINGLE internal reward signal! (In contrast to the CS/RL theory of exploration!)

Wulfram Gerstner

EPFL, Lausanne, Switzerland

# **Artificial Neural Networks and RL**

## **The role of exploration, novelty, and surprise in RL**

**1. Exploration Bonus**

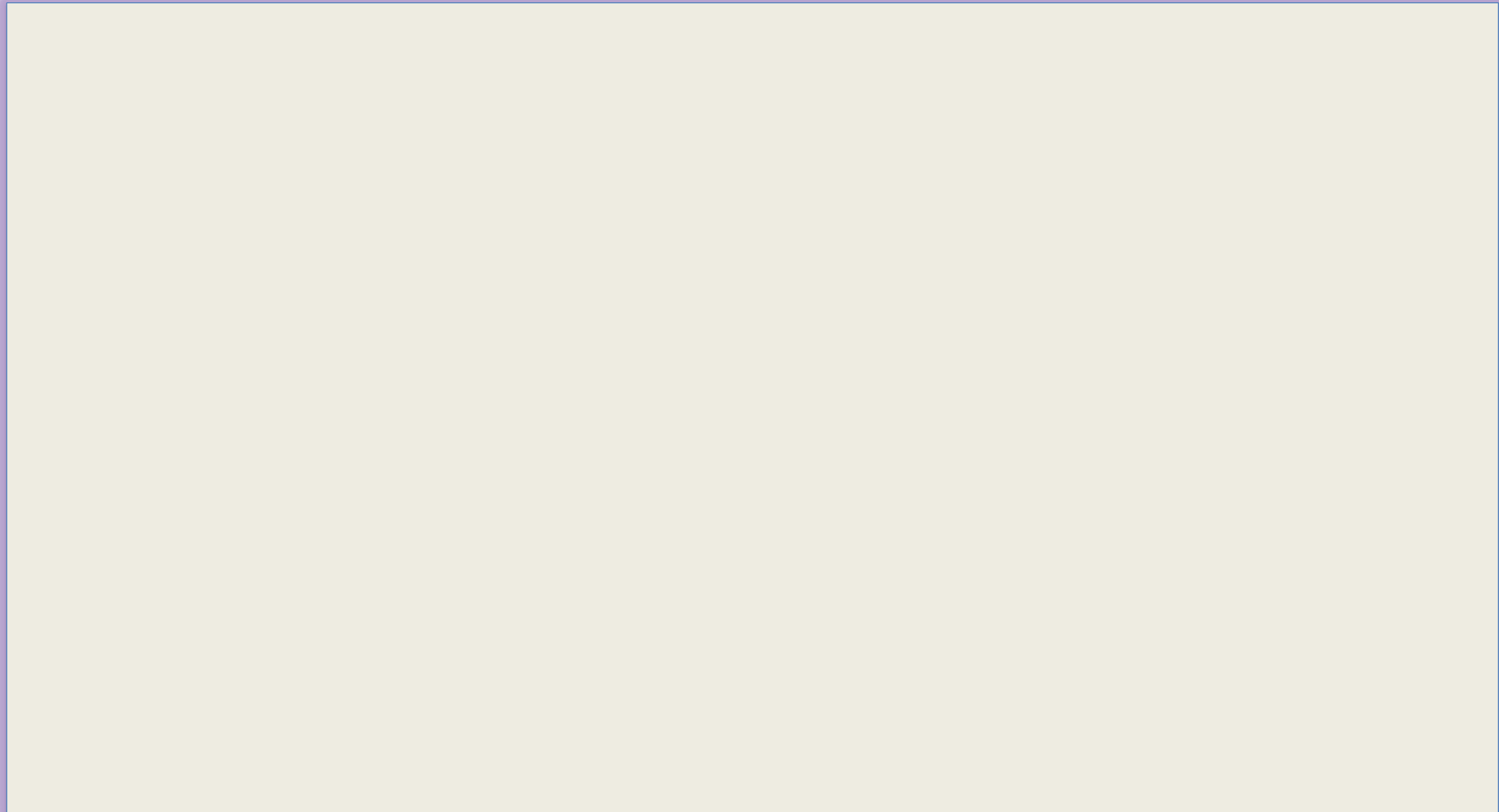
**2. Definitions of Novelty and Surprise (tabular environment)**

Previous slide.

Searching for something 'novel' could be a good heuristic exploration bonus.

We now turn to our intuitions of novelty and surprise.

*Enjoy the images!*



Novelty is not Surprise

Surprise is against models (beliefs)

Previous slide.

The video contains a sequence of about 15 flashed images.

Which ones are ‘novel’?

Which ones are ‘surprising’?

# Novelty and Surprise

Q4: *What is the difference between the two?*

First answer – **novelty and surprise are not the same.**

Second answer (more precise):

Surprise is ‘against beliefs’ or ‘against expectations’  
whereas novelty is not.

Previous slide.



# Novelty and Surprise

Surprise is 'against expectations': an example

...



Previous slide.

# Novelty in a tabular environment: discrete states

events = states  $s$  (e.g., one image). Total number is  $|s|$

**Novelty  $n$ :**

1) count events of type  $s$  up to time  $t$ :  $C^t(s)$

2) a higher count gives lower novelty.

3) the agent has spent a time  $t$  in the environment

4) the empirical observation frequency is  $p_N(s) = \frac{C^t(s) + 1}{t + |s|}$

Definition: The '**Novelty**' of a state  $s$  at time  $t$  is

$$n_t(s) = -\log p_N(s)$$

Previous slide.

Novelty can be defined empirically as the negative logarithm of the empirical frequency.

This definition gives

- At the beginning ( $t=0$ ), all states have the same high novelty (related to the total number of known states).
- The novelty of state  $s$  goes down if it has been observed several times, since its count increases.
- If a state has not been observed for a long time, it will slowly become novel again as time increases – and during that time other states have been observed.

# Surprise in a tabular environment: discrete states and actions

events = transitions  $(s, a \rightarrow s')$  given action  $a$  in state  $s$ .

**Surprise  $S$ :**

- 1) count events of type  $(s, a \rightarrow s')$  up to time  $t$ :  $C^t(s, a \rightarrow s')$
- 2) a higher count gives lower surprise.
- 3) the agent has spent a time  $t$  in the environment
- 4) the empirical observation frequency is

$$p^t(s_{t+1} = s' | s_t, a_t) = \frac{C^t(s, a \rightarrow s') + 1}{\tilde{C}^t(s, a) + |s|}$$

Definition: The '**Surprise**' of a transition is

$$S_{BF}^{t+1}(s') = \frac{\text{prior}}{p_s^t(s_{t+1} = s' | s_t, a_t)}$$

*Bayes  
Factor  
Surprise*

Previous slide.

Surprise is related to expectation – if you do not expect something, then you cannot be surprised. Hence surprise needs contexts and experience that enable an agent to build a belief. Expectations arise from the belief.

While novelty is derived from observation counts of states, surprise is derived from observation counts of transitions.

There are several definitions of surprise.

The specific surprise considered here is the Bayes Factor Surprise.

# Definitions of Novelty and Surprise

Q1: What is novelty?

Definition: The **‘Novelty’** of a state  $s$  is

$$n^t(s) = -\log p_N(s)$$

Q2: What is surprise?

Definition: The **‘Surprise’** of a transition is

$$S_{BF}^{t+1}(s') = \frac{\text{prior}}{p_s^t(s_{t+1}=s' | s_t, a_t)}$$

There are 17 different definitions of surprise.  
This here is the Bayes-Factor surprise.

*Modirshanechi et al.  
(2022)*

Previous slide. Summary.

Note that there are also other definitions of surprise.



# Teaching monitoring – monitoring of understanding

[ ] up to here, at least 60% of material was new to me.

[ ] I have the feeling that I have been able to follow  
(at least) 80% of the lecture up to here.

Previous slide. Summary.

We now turn to Question 4. Why is surprise (or novelty) useful?

We start with surprise.

# **Artificial Neural Networks and RL**

## **The role of exploration, novelty, and surprise in RL**

- 1. Formal Exploration Bonus**
- 2. Definitions of Novelty and Surprise (tabular environment)**
- 3. Why is Surprise useful?**

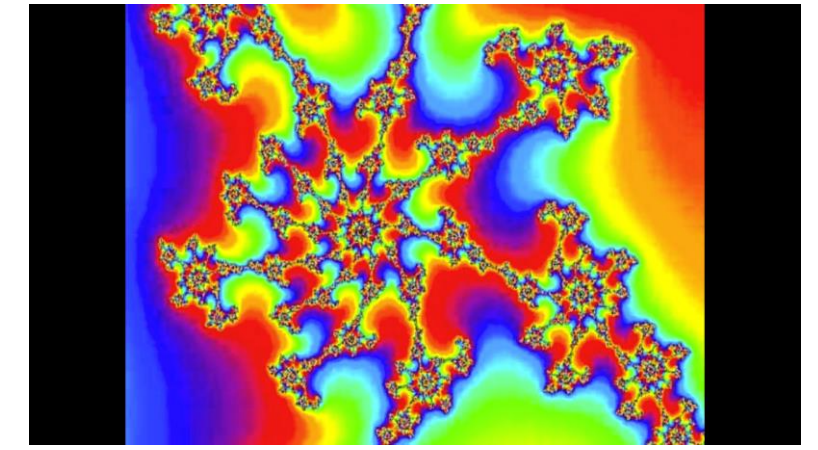
Previous slide. Summary.

We now turn to Question 4. Why is surprise (or novelty) useful?

We start with surprise.

# When are we surprised?

3 9 7 3 9 7 3 9 7 3 9 7 3 9 4 3 9 7



## Surprise against expectations from your current belief

- Expectations arise from models of the world
- We always make models
- We know that the models are not perfect
- **Surprise enables us to adapt the models**

→ Hypothesis:

**Surprise boosts plasticity (3<sup>rd</sup> factor)/ increases the learning rate**

**Note: no reward!!!!**

Previous slide. Review

Similar to the video with the fractals, the series of numbers has a surprising element.

The world around us is incredibly complex. We try to understand it by making models. However, our brain is prewired (inference prior set by evolution) so that we know that our models are simplified and wrong.

At the moment when our expectations arising from our world model is wrong we get a surprise signal. The use of the surprise signal is to increase the learning rate so that we can rapidly re-adapt our model.



# Review: Neuromodulators

- 4 or 5 neuromodulators
- near-global action
- internally created signals

Dopamine/reward/TD:  
*Schultz et al., 1997,*  
*Schultz, 2002*

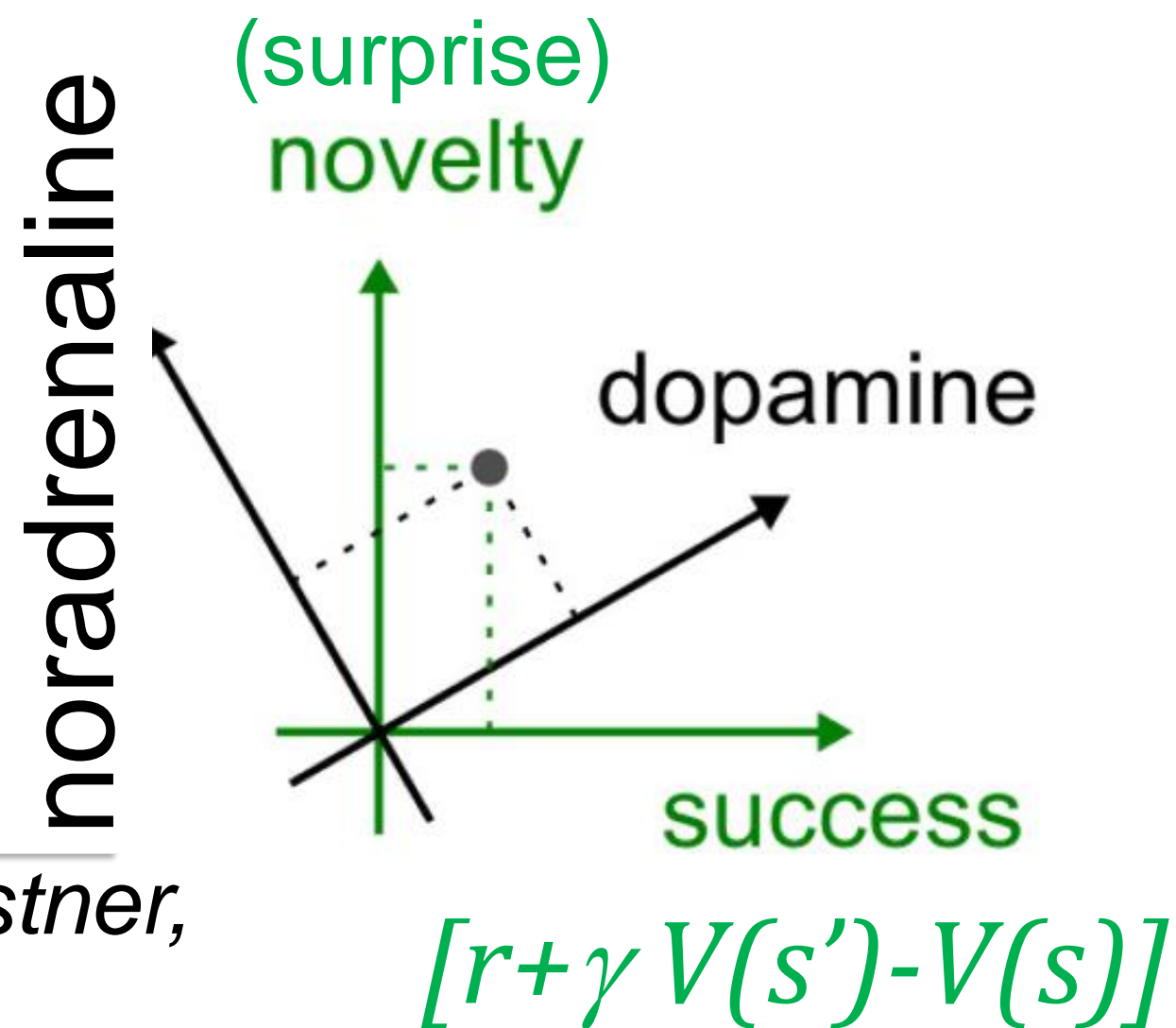
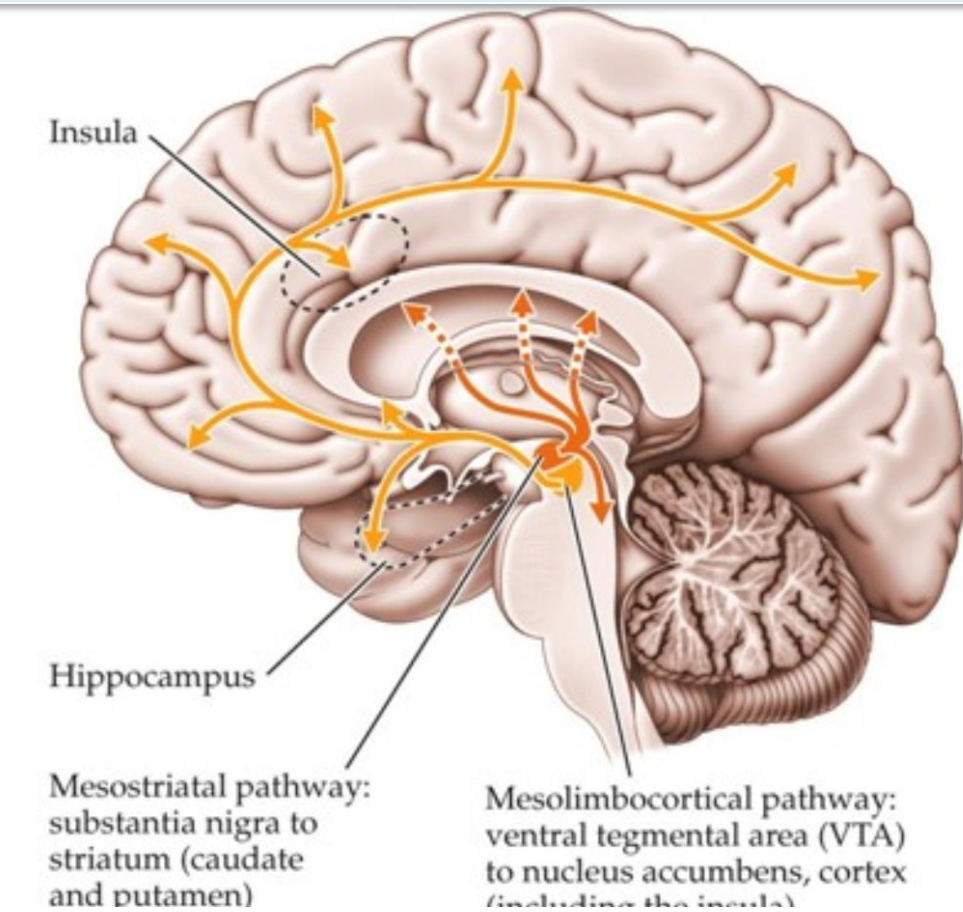


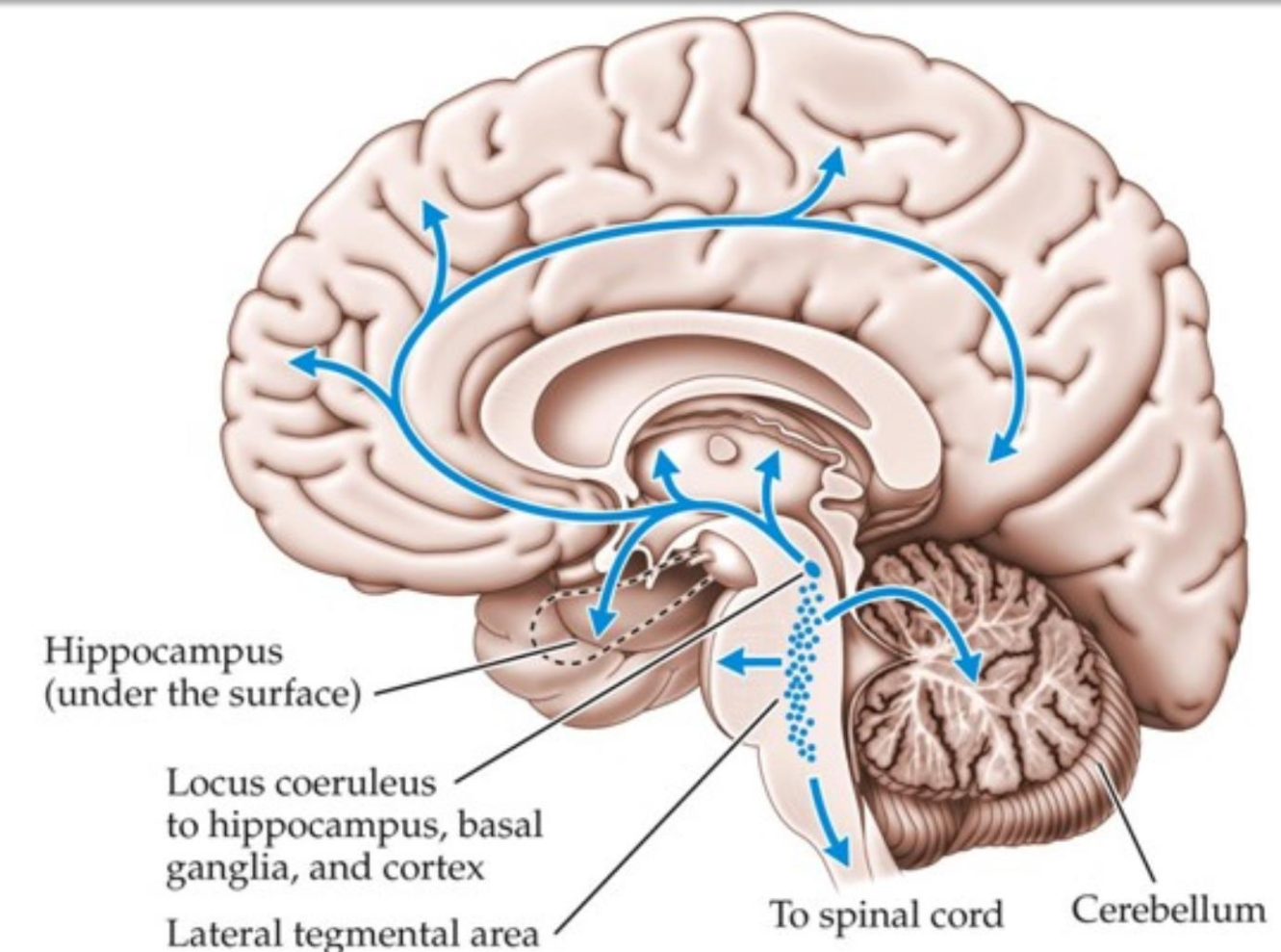
Image:  
*Fremaux and Gerstner,*  
*Frontiers (2016)*

Image: *Biological Psychology, Sinauer*

## Dopamine (DA)



## Noradrenaline (NE)



Previous slide. Review

The most famous neuromodulator is dopamine (DA) which is related to reward, as we will see.

But there are other neuromodulators such as noradrenaline (also called norepinephrine, NE) which is related to surprise.

Left: the mapping between neuromodulators and functions is not one-to-one.

Indeed, dopamine also has a 'surprise' component.

Inversely, noradrenaline also has a reward component.

Right: most neuromodulators send axons to large areas of the brain, in particular to several cortical areas. The axons branch out in thousands of branches.

Thus the information transmitted by a neuromodulator arrives nearly everywhere. In this sense, it is a 'global' signal, available in nearly all brain areas.

Note that the TD error is an internally created signal. The TD can be positive at time  $t$  even if no explicit reward is given at time  $t$ .

Similarly, surprise is an internally generated signal indicating model mismatch.



# Review: Formalism of Three-factor rules with eligibility trace

$x_j$  = activity of presynaptic neuron

$\varphi_i$  = activity/state of postsynaptic neuron

Step 1: co-activation sets eligibility trace

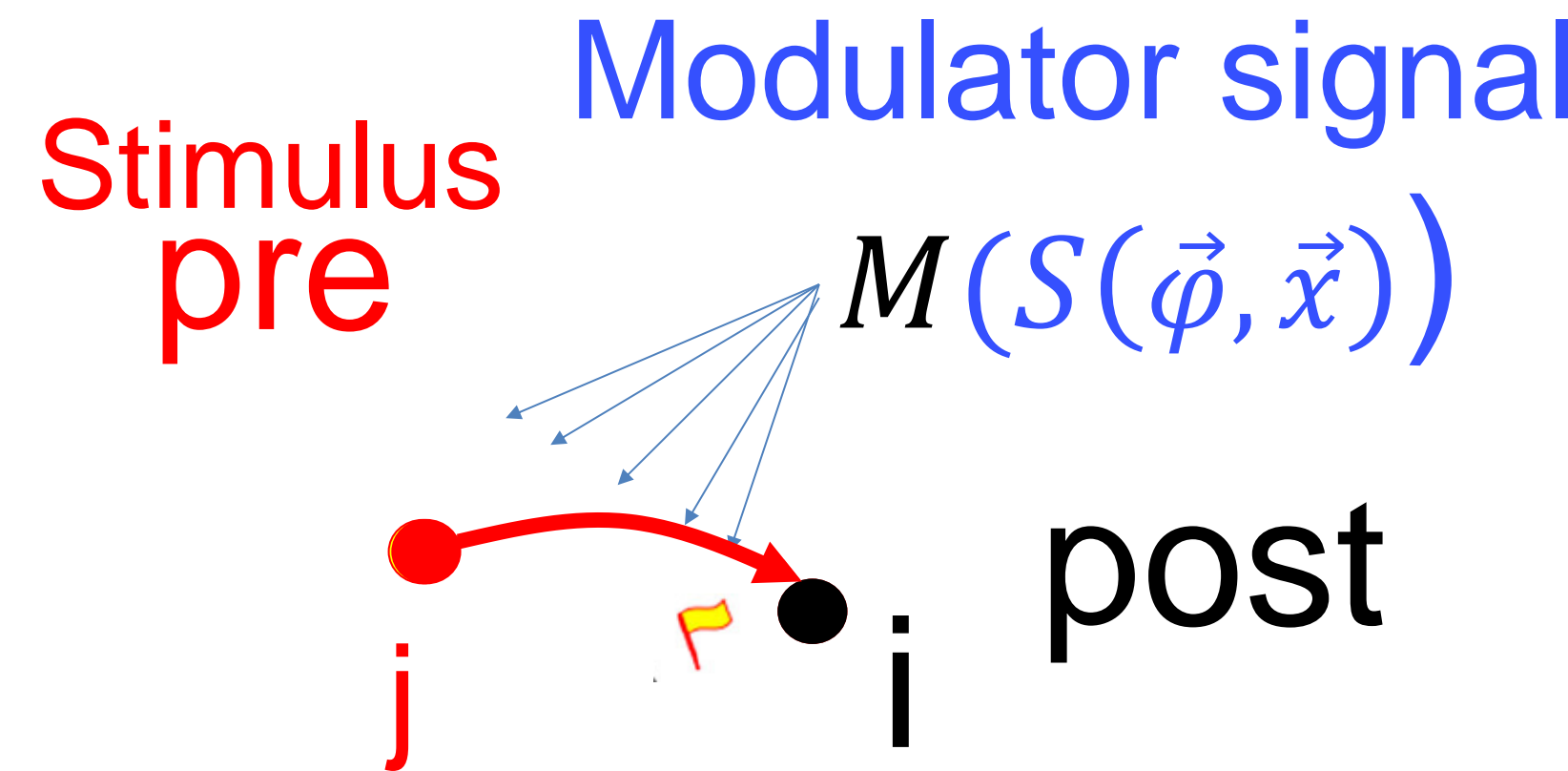
$$\Delta z_{ij} = \eta f(\varphi_i) g(x_j)$$

Step 2: eligibility trace decays over time

$$z_{ij} \leftarrow \lambda z_{ij}$$

Step 3: eligibility trace translated into weight change

$$\Delta w_{ij} = \eta M(S(\vec{\varphi}, \vec{x})) z_{ij}$$



$M(S)$ :

- TD-error
- surprise

Previous slide.

Three-factor rules are implementable with eligibility traces.

1. The joint activation of pre- and postsynaptic neuron sets a 'flag'. This step is similar to the Hebb-rule, but the change of the synapse is not yet implemented.

2. The eligibility trace decays over time

3. However, if a neuromodulatory signal  $M$  arrives before the eligibility trace has decayed to zero, an actual change of the weight is implemented.

The change is proportional to

- the momentary value of the eligibility trace
- the value of the neuromodulator signal

The neuromodulator could signal the

- TD-error
- or Surprise

Usefulness of Surprise? It modulates(similar to the TD error) the learning rate of RL! Surprising events increase the learning rate.

Wulfram Gerstner

EPFL, Lausanne, Switzerland

# Artificial Neural Networks and RL

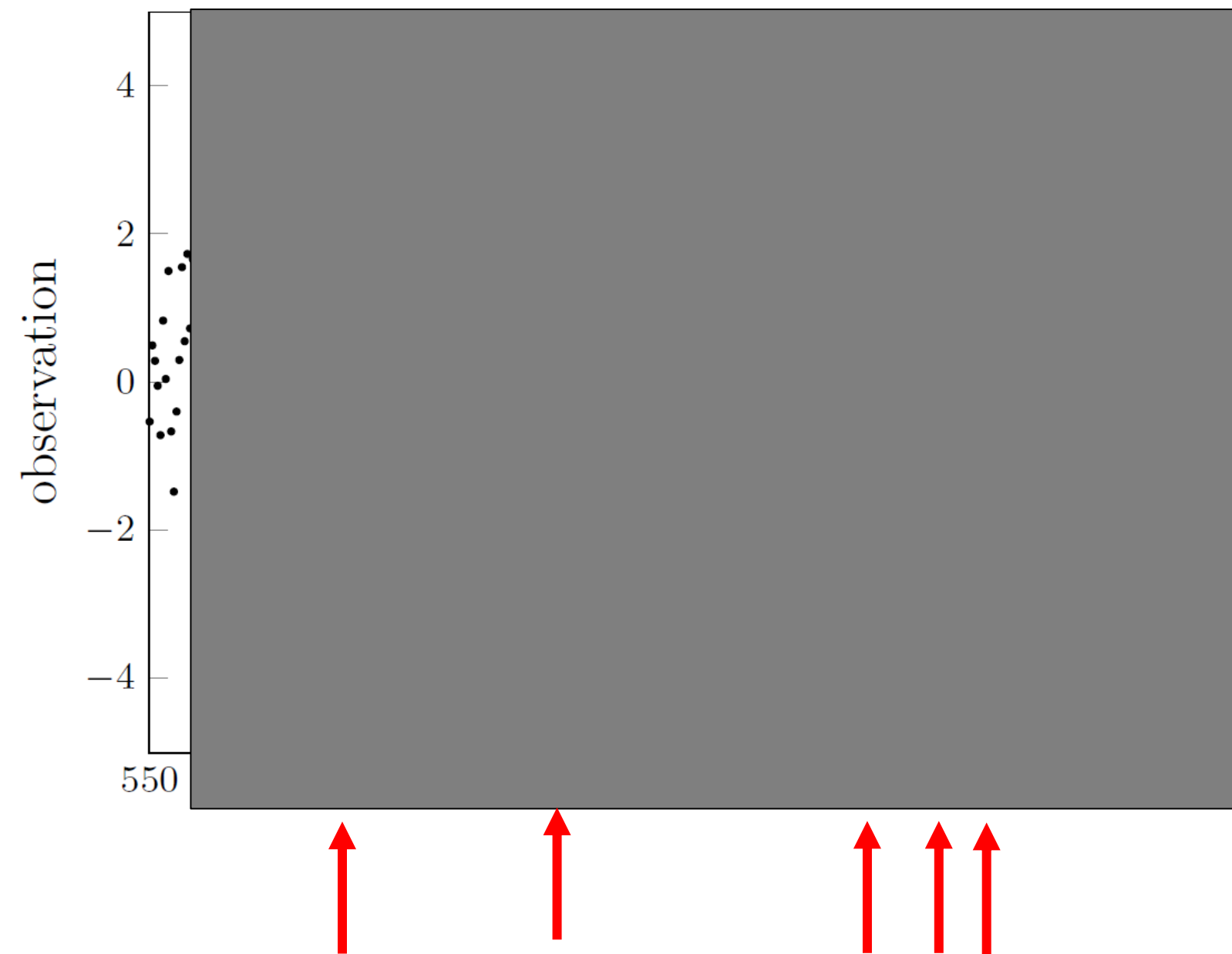
## The role of exploration, novelty, and surprise in RL

1. Definitions of Novelty and Surprise (tabular environment)
2. Why is Surprise useful?
3. **Change-point detection by Bayes-Factor Surprise**

Previous slide.

Our claim is that the Bayes-Factor surprise is ideal for detecting change points.

# Surprise boosts plasticity in volatile environments



Volatile environment:  
abrupt changes with small probability  
→ 'change points'

→ you have to **reset** model after a **change point**

**generative model = nonstationary stochastic process**

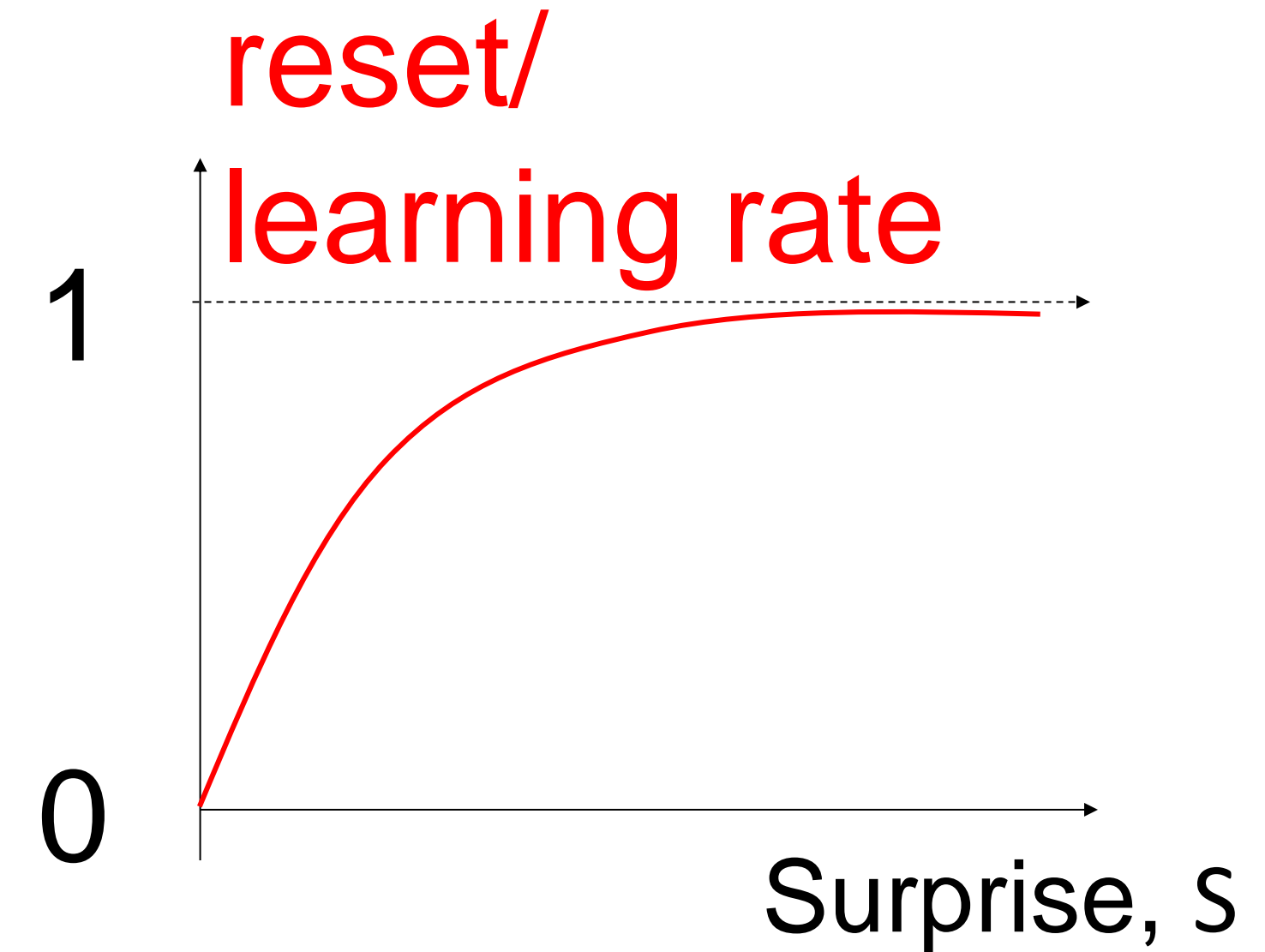
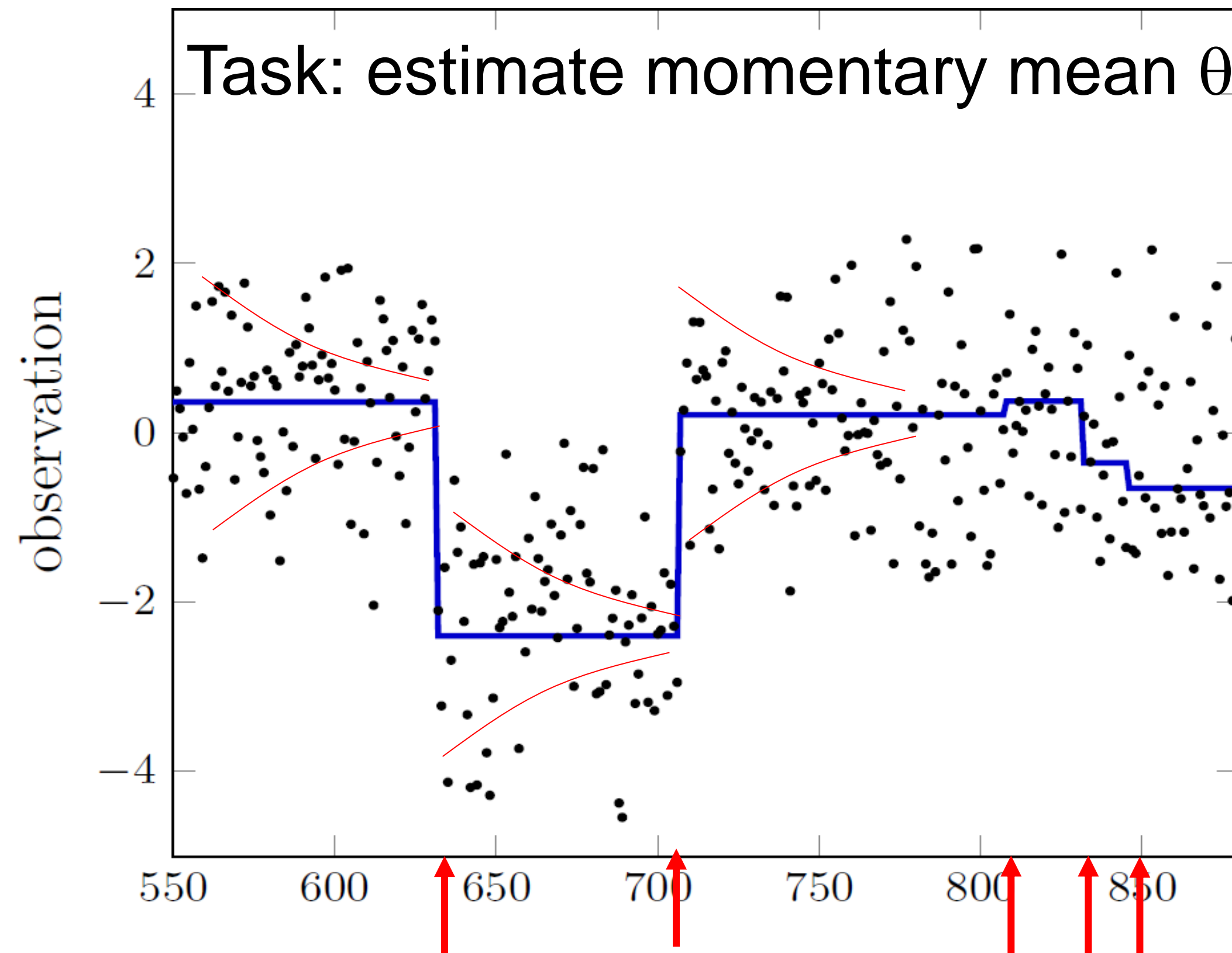
- here:
- mean of Gaussian is fixed for many steps
  - mean jumps at 'change points' : probability  $\ll 1$
  - variance is fixed
  - task is to estimate **momentary mean** of Gaussian

Previous slide.

The volatile environment has stationary segments, interrupted by unpredictable 'change points' that occur at low probability.

If you want to make predictions about the next stimulus (or here: its mean), then the best strategy is to reset your model completely if you have detected a change point.

# Surprise boosts plasticity in volatile environments



in volatile environment, best approach (Bayesian):

- reset your belief to prior, if observation does not make sense
- plasticity of system must increase if 'surprising observation'

Previous slide.

The volatile environment has stationary segments, interrupted by unpredictable change points that occur at low probability.

During the stationary segment your belief gets more precise, and your predictions (regarding the mean of the distribution) get therefore better.

But the best strategy is to reset your model completely if you have detected a change point. So the challenge is to detect the change points.

The optimal way of doing this is the Bayes-Factor surprise.

Plasticity of the model must then increase when you detect a change point, so that you reset to the prior and integrate new data points starting from the prior.

Plasticity (learning rate) of the model must then increase when you detect a change point, so that you reset to the prior and integrate new data points starting from the prior.



# Surprise boosts plasticity in volatile environments

$$S_{\text{BF}}(y_{t+1}; \pi^{(t)}) = \frac{P(y_{t+1}; \pi^{(0)})}{P(y_{t+1}; \pi^{(t)})}.$$

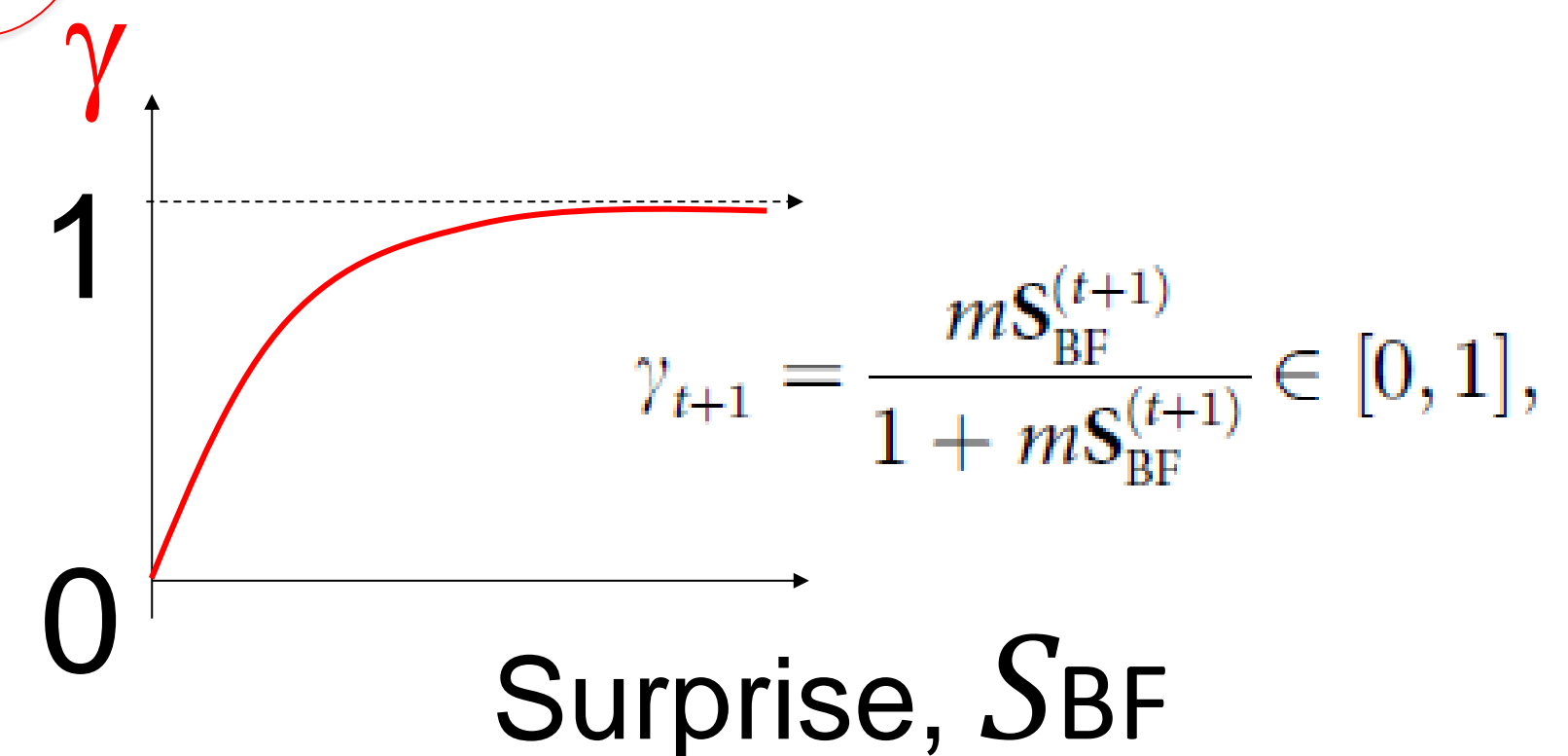
Probability of observation  $y$   
under prior belief  $\pi^{(0)}$

Probability of observation  $y$   
under current belief  $\pi^{(t)}$

→ reset your belief to prior, if observation  $y$  does not make sense

$$\pi^{\text{new}}(\theta) = (1 - \gamma) \pi^{\text{integration}}(\theta | y^{\text{new}}, \pi^{\text{old}}) + \gamma \pi^{\text{reset}}(\theta | y^{\text{new}}, \pi^{(0)}).$$

→ 'exact Bayesian inference'  
in volatile environment modulates  
update with factor  $\gamma$



Previous slide.

We claimed that plasticity (learning rate) of the model must increase when you detect a change point, so that you reset to the prior and integrate new data points starting from the prior.

This is formalized in the long equation in the middle.

Using a careful analysis of the statistical estimation in the presence of change points you find that:

If it unlikely (small  $\gamma$ ) that there was a change point between the previous data and the current data point (observation  $y^{\text{new}}$ ), then you should use standard statistical updates of your estimates to INTEGRATE the new data into your current belief.

If it is likely ( $\gamma$  close to 1) that there was a change point, then you should reset to your prior and integrate the new data point using statistical updates starting with the prior as your current belief.

Moreover, this factor  $\gamma$  depends monotonically on the Bayes-Factor Surprise  $S_{\text{BF}}$

# Surprise boosts plasticity in volatile environments

$$S_{\text{BF}}(y_{t+1}; \pi^{(t)}) = \frac{P(y_{t+1}; \pi^{(0)})}{P(y_{t+1}; \pi^{(t)})}.$$

Probability of observation  $y$   
under prior belief  $\pi^{(0)}$

Probability of observation  $y$   
under current belief  $\pi^{(t)}$

→ reset your belief to prior, if observation  $y$  does not make sense

**Exact update rule not implementable, but**

Bayes-Factor Surprise plays crucial role in approximate methods:

- Particle Filter with  $N$  particles,
- Message-Passing with  $N$  messages,
- Published approximations

Previous slide.

The general theoretical framework cannot be integrated out over several time steps. Therefore approximations are necessary.

However, what is important is the gist of the argument:  
A high surprise indicates that the learning rate should be increased.

# Summary: Definitions of Novelty and Surprise

What is novelty?

$$p_N(s) = \frac{C^t(s) + 1}{t + |s|}$$

Definition: The '**Novelty**' of a state **s** is

$$n^t(s) = -\log p_N(s)$$

What is surprise?

$$p^t(s_{t+1} = s' | s_t, a_t) = \frac{C^t(s, a \rightarrow s') + 1}{\tilde{C}^t(s, a) + |s|}$$

Definition: The '**Surprise**' of a transition is

$$S_{BF}^{t+1}(s') = \frac{\text{prior}}{p_s^t(s_{t+1} = s' | s_t, a_t)}$$

There are 17 different definitions of surprise.  
This here is the Bayes-Factor surprise.

*Modirshanechi et al.*  
(2022)

# Summary: Why is surprise useful?

- **Detect change points** in environment statistics
- **Adapt learning rate** after change point.
- **Bayes-Factor Surprise** is a good surprise measure for this

Wulfram Gerstner

EPFL, Lausanne, Switzerland

# Artificial Neural Networks and RL

## The role of exploration, novelty, and surprise in RL

1. Definitions of Novelty and Surprise (tabular environment)
2. Why is Surprise useful?
3. Change-point detection by Bayes-Factor Surprise
4. **Why is Novelty useful?**

Previous slide.

We are done with surprise and turn now to the second part of Question 4.  
Why is novelty useful?

We start with a detour in order to review well-known results from RL, in particular TD learning and eligibility traces.



# Why is Novelty useful?

## → helps to explore

### Exercise 1. How fast can we find the goal state with a stationary policy?

Consider an environment with the state space  $\mathcal{S}$ , a goal (terminal) state  $G \in \mathcal{S}$ , and an action space  $\mathcal{A}$  in non-goal states (i.e.,  $\mathcal{S} - \{G\}$ ). After taking action  $a \in \mathcal{A}$  in state  $s \in \mathcal{S}$ , the agent moves to state  $s' \in \mathcal{S}$  with the transition probability  $p(s'|s, a)$ . These transition probabilities are unknown to the agent. We use  $T$  to denote the first time an agent find the goal state  $G$ , i.e.,  $s_T = G$ . If we assume that the agent uses a stationary policy  $\pi$ , then we can define the average of  $T$  given each initial state  $s \in \mathcal{S}$  as

$$\mu_\pi(s) := \mathbb{E}_\pi[T|s_0 = s],$$

where  $s_0$  is the state at time  $t = 0$ . In this exercise, we study  $\mu_\pi(s)$  in its most general case.

- a. What is the value of  $\mu_\pi(G)$ ?

*Hint:* Note that  $T$  is equal to the smallest  $t \geq 0$  when we have  $s_t = G$ .

- b. What is the relationship between  $\mathbb{E}_\pi[T|s_1 = s]$  and  $\mu_\pi(s)$ ?

*Hint:* Note that  $\mu_\pi(s)$  is the average of  $T$  if the agent starts in state  $s$  at time  $t = 0$ , whereas  $\mathbb{E}_\pi[T|s_1 = s]$  is the average of  $T$  if the agent starts in state  $s$  at time  $t = 1$ .

- c. Find a system of linear equations for finding  $\mu_\pi(s)$  for  $s \in \mathcal{S} - \{G\}$ .

*Hint:* Use the fact that  $p_\pi(s'|s) = \sum_{a \in \mathcal{A}} \pi(a|s)p(s'|s, a)$ .

### Exercise 2. The magic of seeking novelty.

Previous slide.

We are done with surprise and turn now to the second part of Question 4.  
Why is novelty useful?

We start with a detour in order to review well-known results from RL, in particular TD learning and eligibility traces.

# Review: TD-learning in the general sense

$$Q(s, a) = \sum_{s'} P_{s \rightarrow s'}^a \left[ R_{s \rightarrow s'}^a + \gamma \sum_{a'} \pi(s', a') Q(s', a') \right]$$

SARSA

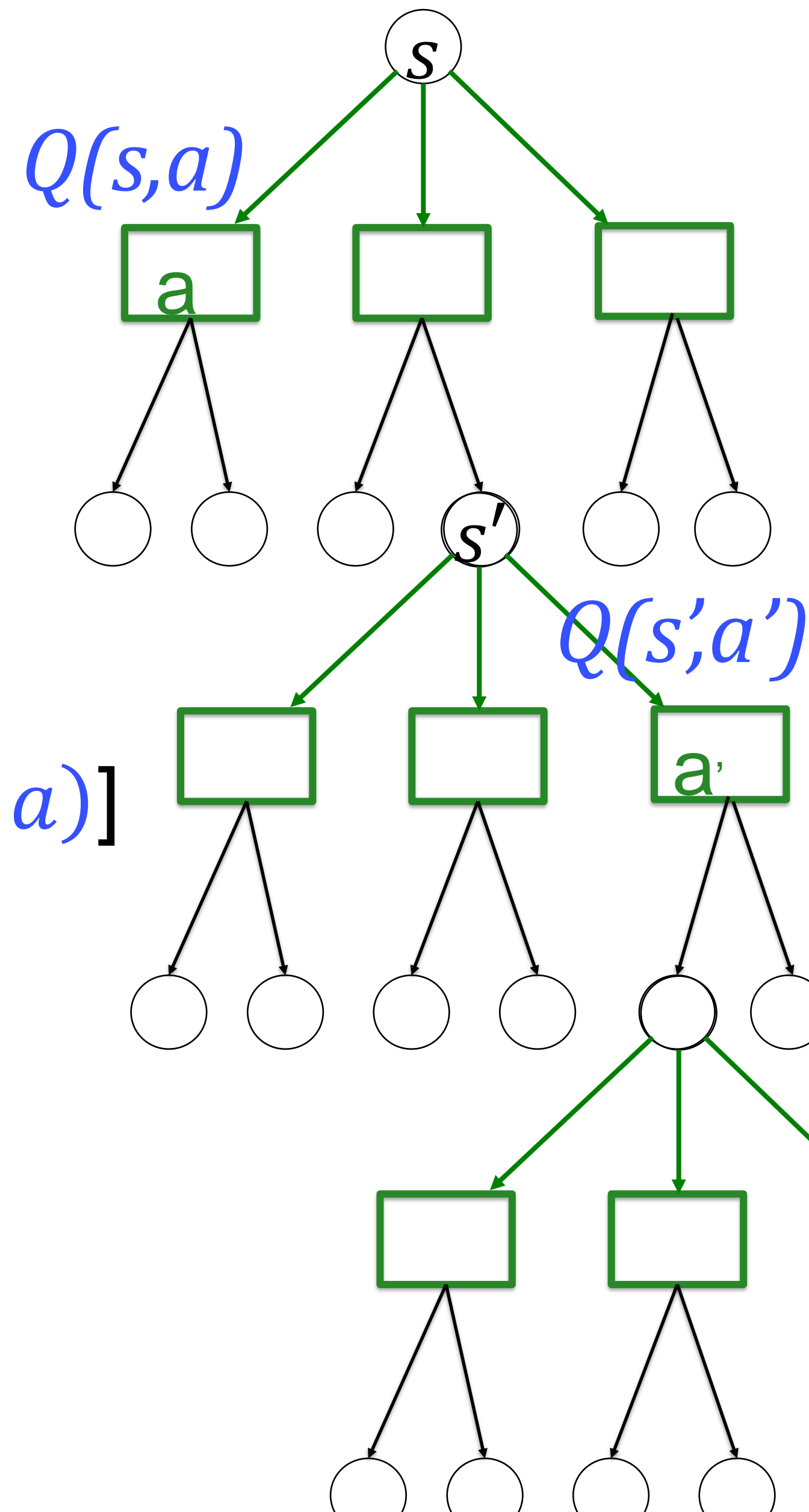
$$\Delta Q(s, a) = \eta [r_t + \gamma Q(s', a') - Q(s, a)]$$

Expected SARSA

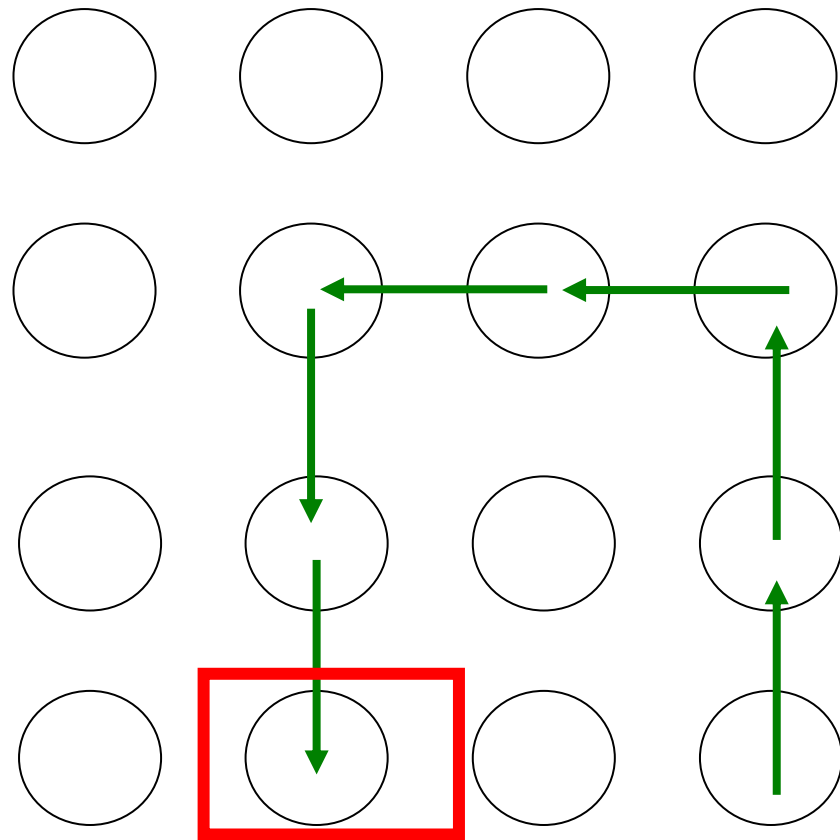
$$\Delta Q(s, a) = \eta [r_t + \gamma \{ \sum_{a'} \pi(s', a') Q(s', a') \} - Q(s, a)]$$

Q-learning

$$\Delta Q(s, a) = \eta [r_t + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$



# Review: Eligibility Traces, SARSA( $\lambda$ )

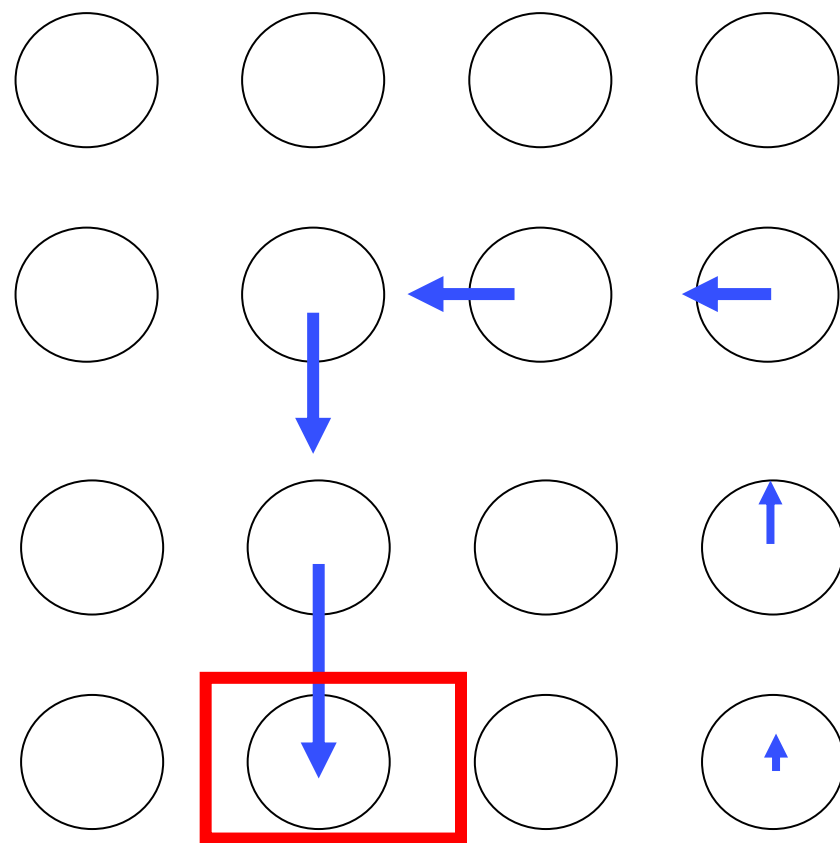


Idea:

- keep memory of previous state-action pairs
- memory decays over time
- update eligibility trace for **all** state-action pairs

$$e(s, a) \leftarrow \lambda e(s, a) \quad \text{decay of **all** traces}$$

$$e(s, a) \leftarrow e(s, a) + 1 \quad \text{if action } a \text{ chosen in state } s$$



- update **all** Q-values at **all time steps**  $t$ :

$$\Delta Q(s, a) = \eta \underbrace{[r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]}_{\text{RPE = TD error } \delta_t} e(s, a)$$

RPE = TD error  $\delta_t$

Note:  $\lambda=0$  gives standard SARSA

# Review: Model-based

versus

# Model-free

- learns model of environment  
‘transition matrix’
- knows ‘rules’ of game
- planning ahead is possible
- can update Bellman equation  
in ‘background’ without action
- can simulate action sequences  
(without taking actions)
- is not

- does not
- does not
- cannot plan ahead
- cannot
- cannot
- Eligibility traces and V-values  
keep memory of past
- completely online, causal,  
forward in time.

# Reward-based learning versus Novelty-based learning

rewards

$r_t$

Q-values

$Q_R^{(t)}(s, a)$

Bellman eq.  
estimation/update

Model-based

prioritized  
sweeping

$Q_{MB,R}^{(t)}(s, a)$

Model-free

eligibility  
traces

$Q_{MF,R}^{(t)}(s, a)$

novelty

$n_t$

Q-values

$Q_N^{(t)}(s, a)$

Bellman eq.  
estimation/update

Model-based

prioritized  
sweeping

$Q_{MB,N}^{(t)}(s, a)$

Model-free

eligibility  
traces

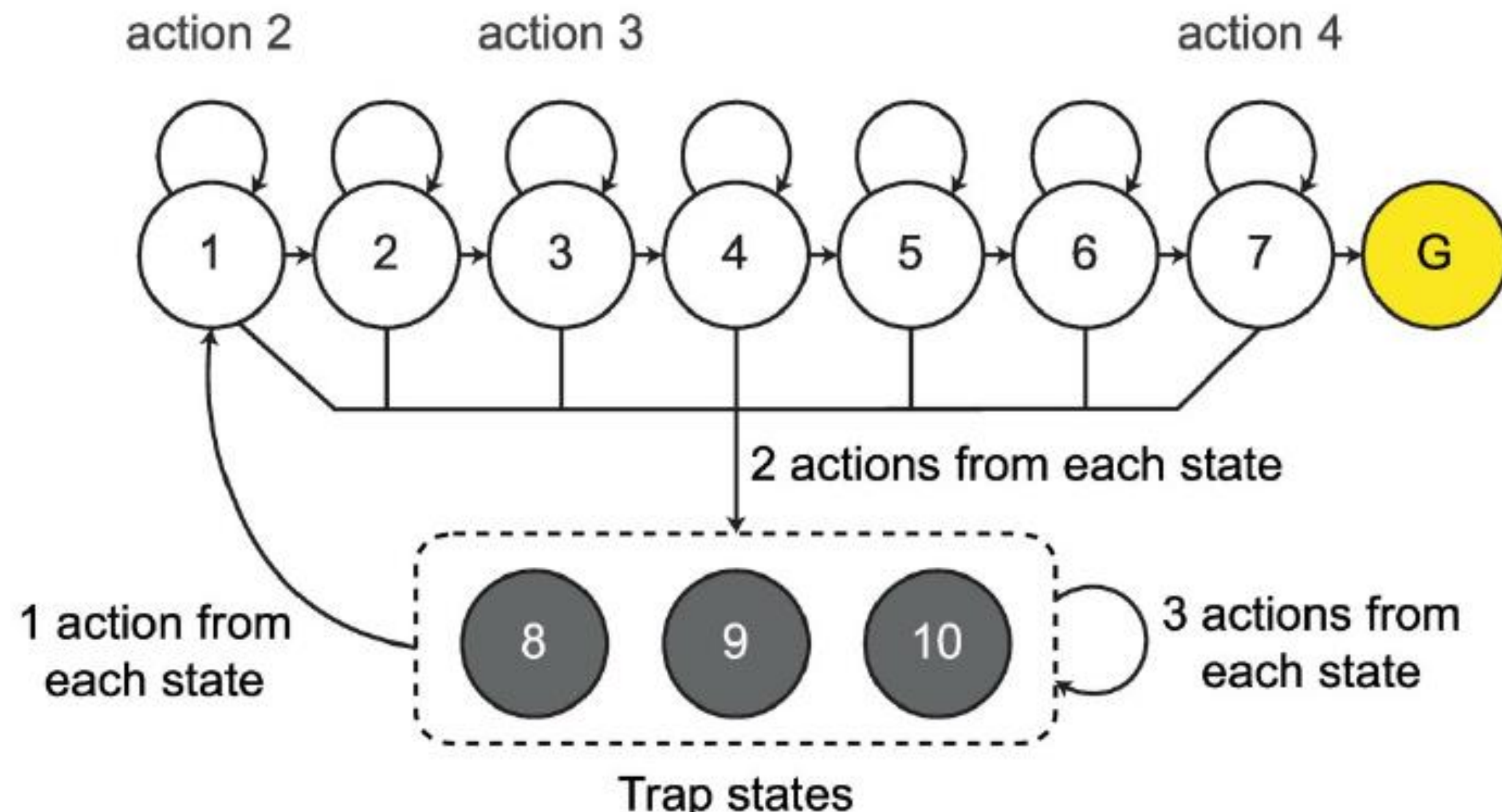
$Q_{MF,N}^{(t)}(s, a)$



# Initial exploration of an environment

Environment with 10 states (+ goal)

4 actions per state



Start in state 1:  
With random policy,  
how many actions  
on average before  
finding goal?

[ ] 100-500

[ ] 1000 – 5000

[ ] more than 10000

Actions are deterministic.  
Fixed random assignment.

Previous slide.

With random exploration, how long would it take on average to find the goal?  
There are only 10 states with four actions each, plus the goal.



# Improve exploration of an environment

Focus on 1<sup>st</sup> episode, before any reward.

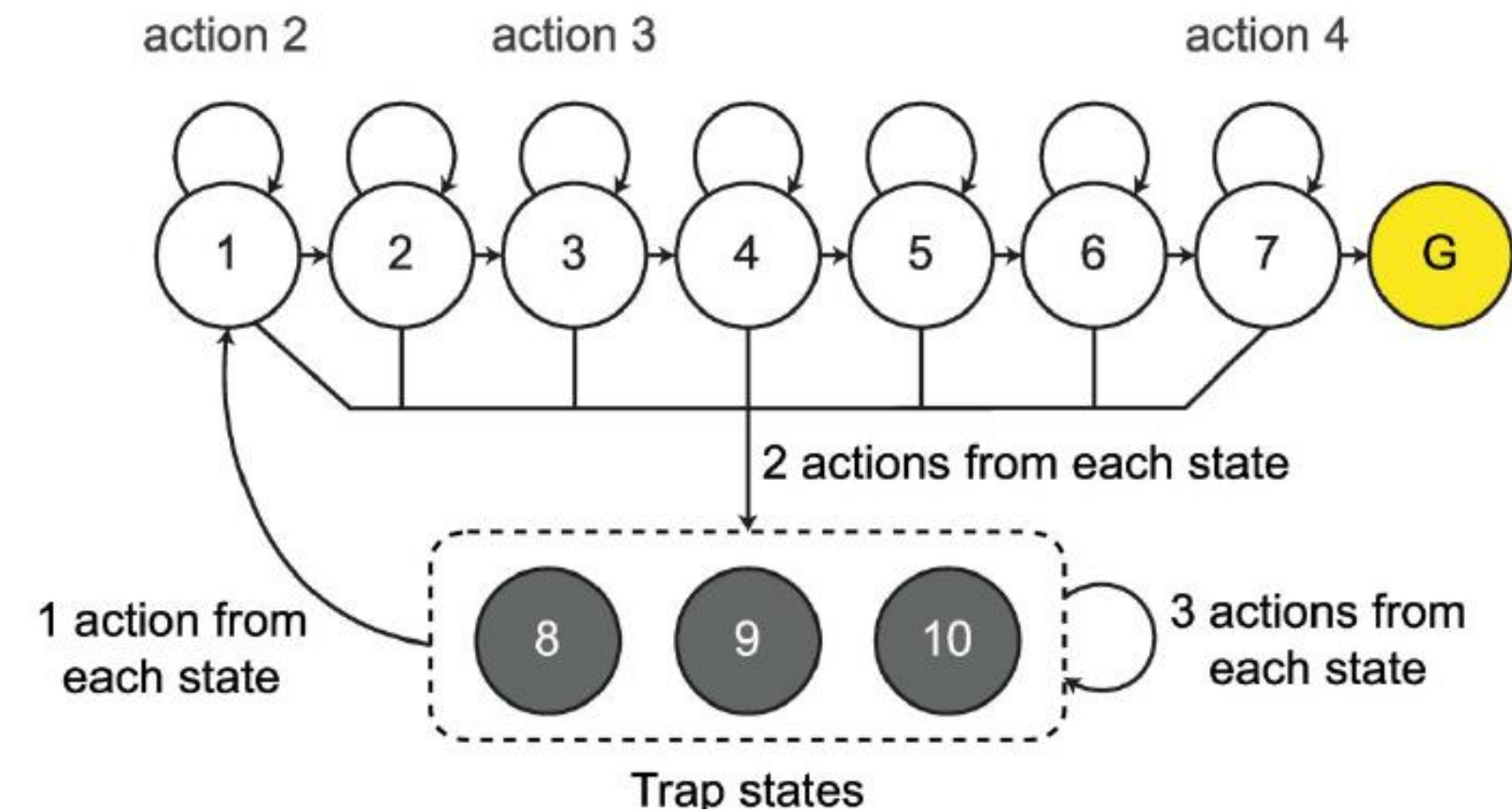
Improve exploration! Solutions?

1. Optimistic initialization?

Initialize  $Q_R(s, a) = 10$  for all  $s, a$

$$\Delta Q_R(s, a) = \eta[r_t + \gamma \max_{a'} Q_R(s', a') - Q_R(s, a)]$$

- Possible but comparatively slow.
- Does not generalize well for episode 2.



Previous slide.

Optimistic initialization is not sufficient to drive exploration.

# Novelty encourages exploration of an environment

Focus on 1<sup>st</sup> episode, before any reward.

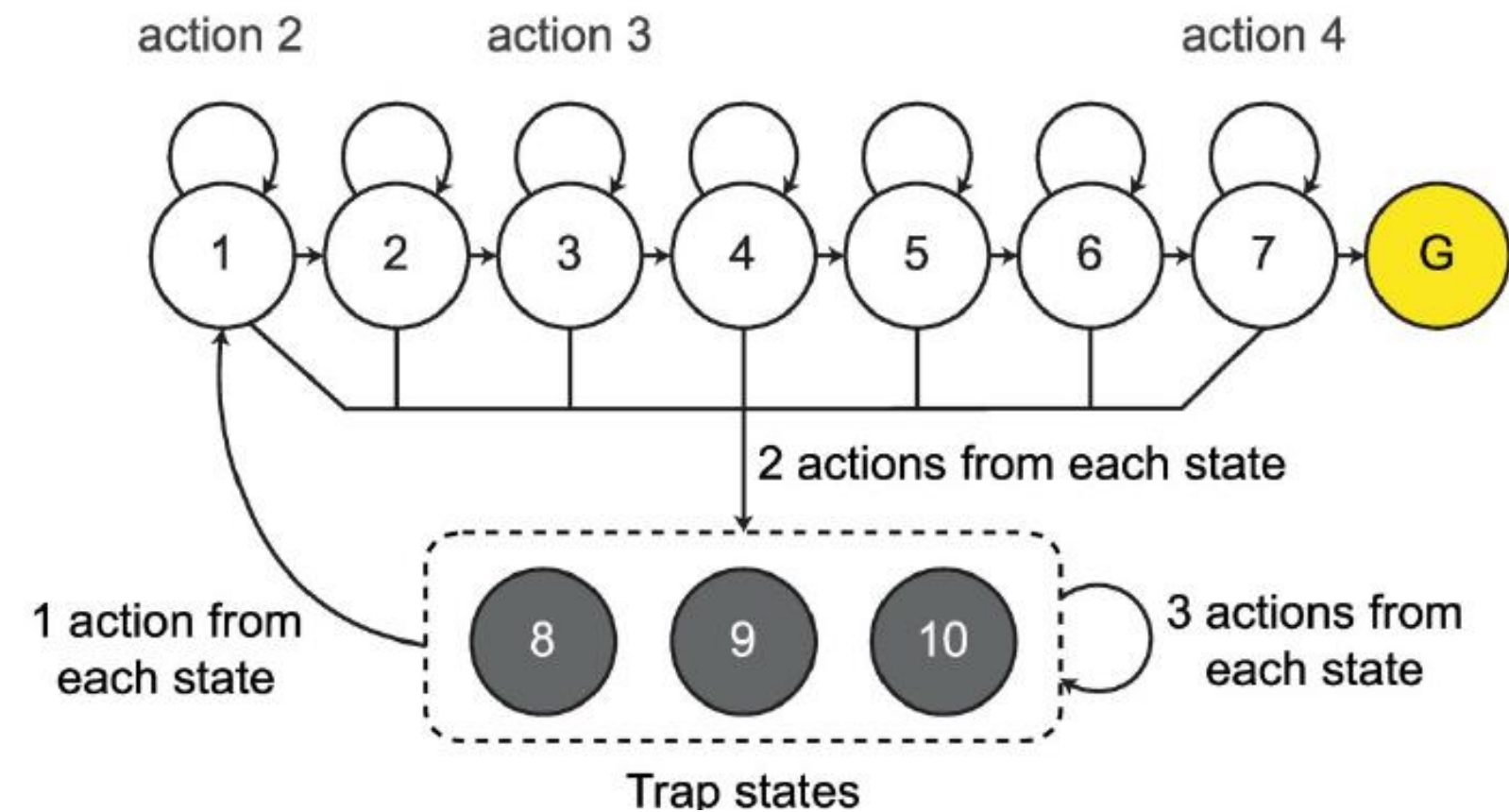
Improve exploration! Solutions?

2. Novelty at time  $t$  is  $n_t$

Novelty Prediction Error (NPE)

$$\Delta Q_N(s, a) = \eta [n_t + \gamma \max_{a'} Q_N(s', a') - Q_N(s, a)]$$

→ Separate Q-value for novelty!



Previous slide.

We now use the novelty-Q-values.

Note that every state has some level of novelty. So the novelty prediction error NPE gives non-zero values for most transitions.

Does this lead to good novelty values? To answer this let us look at the next slide.

RPE: Reward Prediction Error

NPE: Novelty Prediction Error



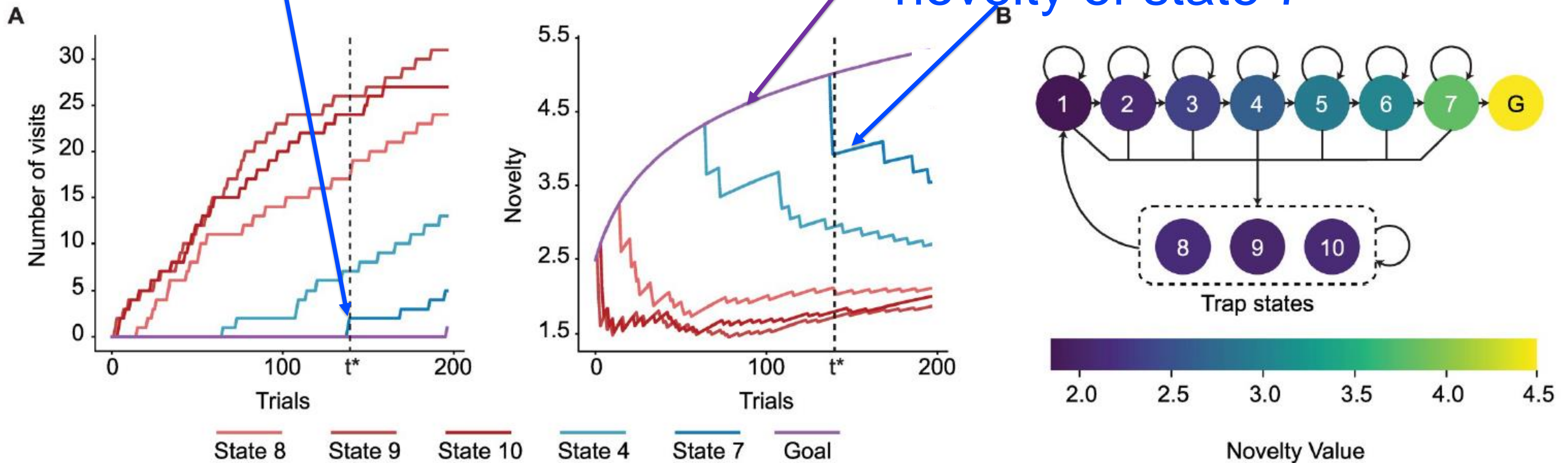
# Novelty encourages exploration of an environment

Focus on 1<sup>st</sup> episode, before any reward; with some policy

first encounter of state 7

novelty of goal

novelty of state 7



→ use novelty values  $Q_N^{(t)}(s, a)$  for action policy!

Previous slide.

The novelty of state 7 or of the goal state increases over time during episode 1.

The plot on the right shows novelty Q-values at the moment when state 7 was found for the first time. There is a nice gradient of increasing novelty towards the goal.

This suggests that novelty Q-values are useful to guide exploration

**Fig 3. Novelty in episode 1 of block 1.** A. The number of state visits (left panel) and novelty (right panel) as a function of time for one representative participant: The number of visits increases rapidly for the trap states and remains 0 for a long time for the states closer to the goal. Novelty of each state is defined as the negative log-probability of observing that state (see Eqs [1](#) and [2](#)) and, hence, increases for states which are not observed as time passes. The first time participants encounter state 7 (the state before the goal state) is denoted by  $t^*$ . B. Average (over participants) novelty (color coded) at  $t^*$ : Novelty of each state is a decreasing function of its distance from the goal state.

# Artificial Neural Networks and RL

## The role of exploration, novelty, and surprise in RL

1. Definitions of Novelty and Surprise (tabular environment)
2. Why is Surprise useful?
3. Change-point detection by Bayes-Factor Surprise
4. Why is Novelty useful?
5. Hybrid Model with Novelty, Surprise, and Reward

Previous slide.

Now we study a specific model that combines many aspects.

Reminder:

RPE: Reward Prediction Error = TD error of reward-consistency

NPE: Novelty Prediction Error = TD error of novelty consistency



# Combine Novelty and Reward: ideas

→ use separate novelty values  $Q_N^{(t)}(s, a)$  for action policy!

→ **exploration**

→ use separate reward values  $Q_R^{(t)}(s, a)$  for action policy!

→ **exploitation**

→ Combine the two and switch relative importance

→ Switch from **exploration to exploitation (and back)**

Note: do not simply add exploration bonus!

~~$$\hat{Q}_{\text{MB}}^{(t)}(s, a) = \hat{R}^{(t)}(s, a) + \frac{\beta}{\sqrt{T_{s,a}^{(t)}}} + \gamma \sum_{s'} \hat{P}^{(t)}(s'|s, a) \max_{a'} \hat{Q}_{\text{MB}}^{(t)}(s', a')$$~~

Previous slide.

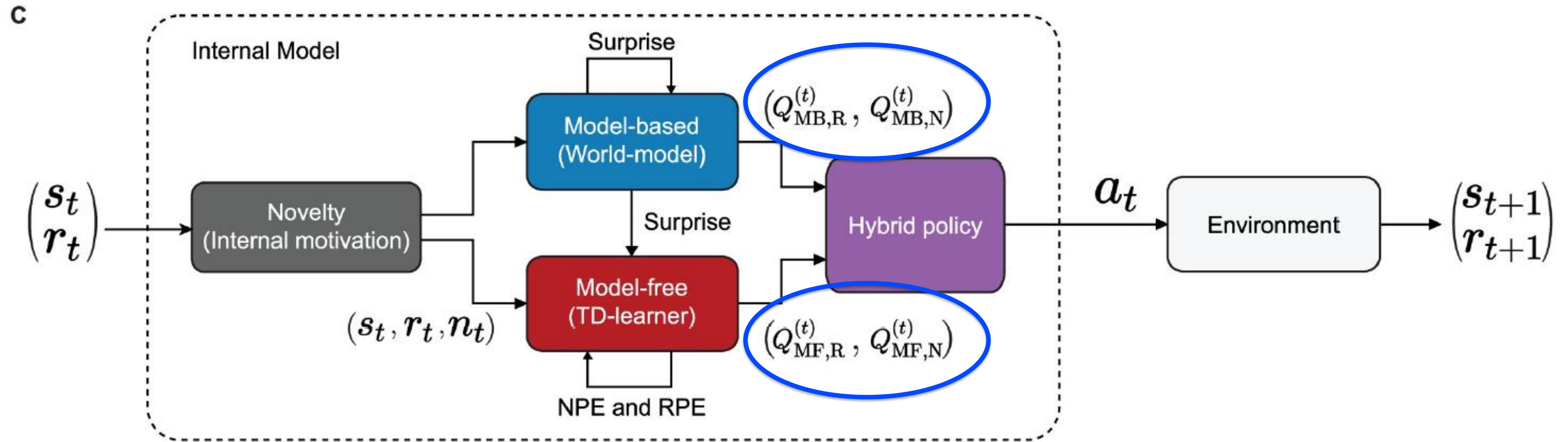
Now we study a specific model that combines many aspects.

Reminder:

RPE: Reward Prediction Error = TD error of reward-consistency

NPE: Novelty Prediction Error = TD error of novelty consistency

# Hybrid model with separate paths for Novelty and Reward (learning rate controlled by Surprise)



$$\text{RPE} = [r_t + \gamma \max_{a'} Q_R(s', a') - Q_R(s, a)]$$

$$\text{NPE} = [n_t + \gamma \max_{a'} Q_N(s', a') - Q_N(s, a)]$$

4 separate  
sets of  
Q-values!

Previous slide.

In total we have in this Hybrid model 4 sets of Q-values:

Reward-driven Q-values, in the versions model-free and model based.

Novelty-driven Q-values, in the versions model-free and model based.

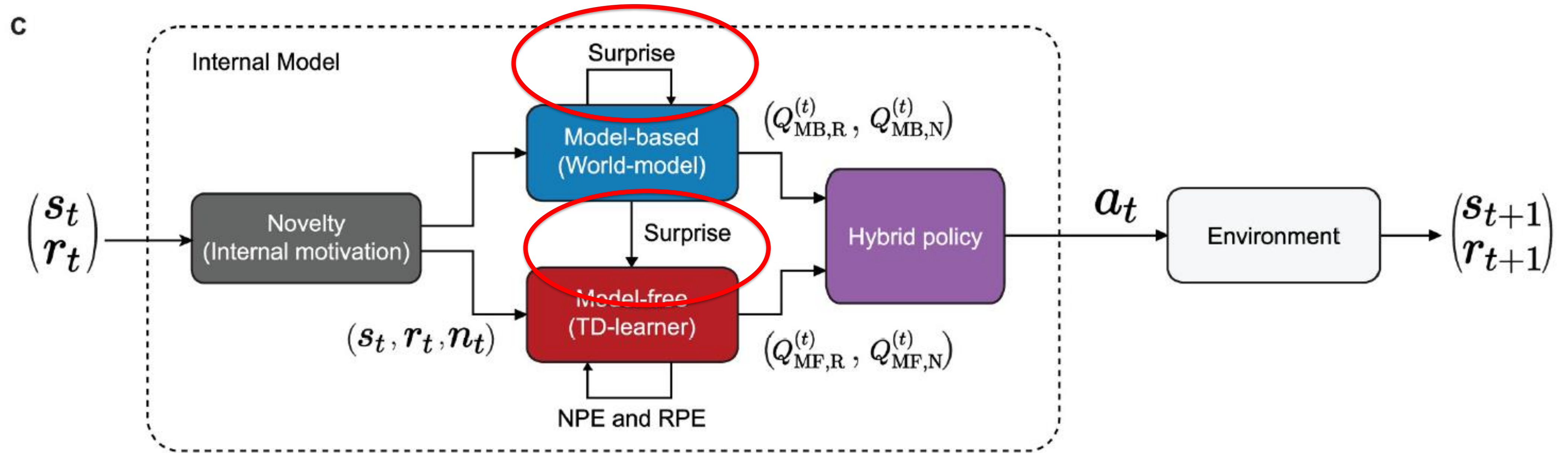
All 4 Q-values are then combined in a softmax fashion to choose the best action.

The relative weighting factors can be changed.

Before the first episode, it might be good to give more importance to novelty,  
and after the first episode more importance to rewards.

algorithm: Information of state  $s_t$  and reward  $r_t$  at time  $t$  is combined with novelty  $n_t$  (grey block) and passed on to the world-model (blue block, implementing the model-based branch of SurNoR) and TD learner (red block, implementing the model-free branch). The surprise value computed by the world-model modulates the learning rate of both the TD-learner and the world-model. The output of each block is a pair of Q-values, i.e, Q-values for estimated reward  $Q_{MF,R}$  and  $Q_{MB,R}$  as well as for estimated novelty  $Q_{MF,N}$  and  $Q_{MB,N}$ . The hybrid policy (in purple) combines these values.

# Hybrid model with separate paths for Novelty and Reward (learning rate controlled by Surprise)



World model: estimated transition matrix  $\hat{P}^{(t)}(s'|s, a)$

- used in model-based Q-learning for background updates
- used to evaluate surprise (Bayes Factor Surprise)
- surprise influences learning rate



Previous slide.

In total we have in this Hybrid model 4 sets of Q-values:

Note that the model-based version need a 'world model'.

- The world model can be used for background updates.
- The world model contains estimated transition probabilities
- The world model can then also used to evaluate 'surprise'
- We use the Bayes Factor surprise
- Surprise will influence the learning rates of ALL four RL algorithms

# **Artificial Neural Networks and RL**

## **The role of exploration, novelty, and surprise in RL**

- 1. Definitions of Novelty and Surprise (tabular environment)**
- 2. Why is Surprise useful?**
- 3. Change-point detection by Bayes-Factor Surprise**
- 4. Why is Novelty useful?**
- 5. Hybrid Model with Novelty, Surprise, and Reward**
- 6. Two Experiments (Markov Decision Problem for humans!)**

Previous slide.

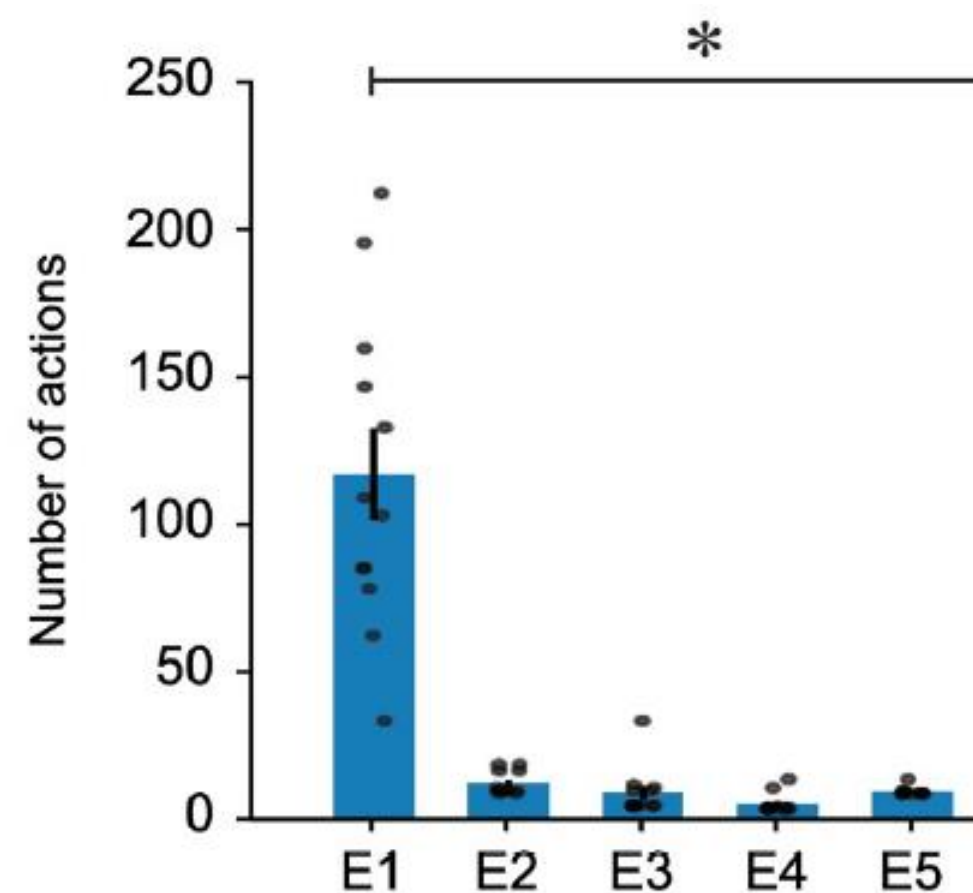
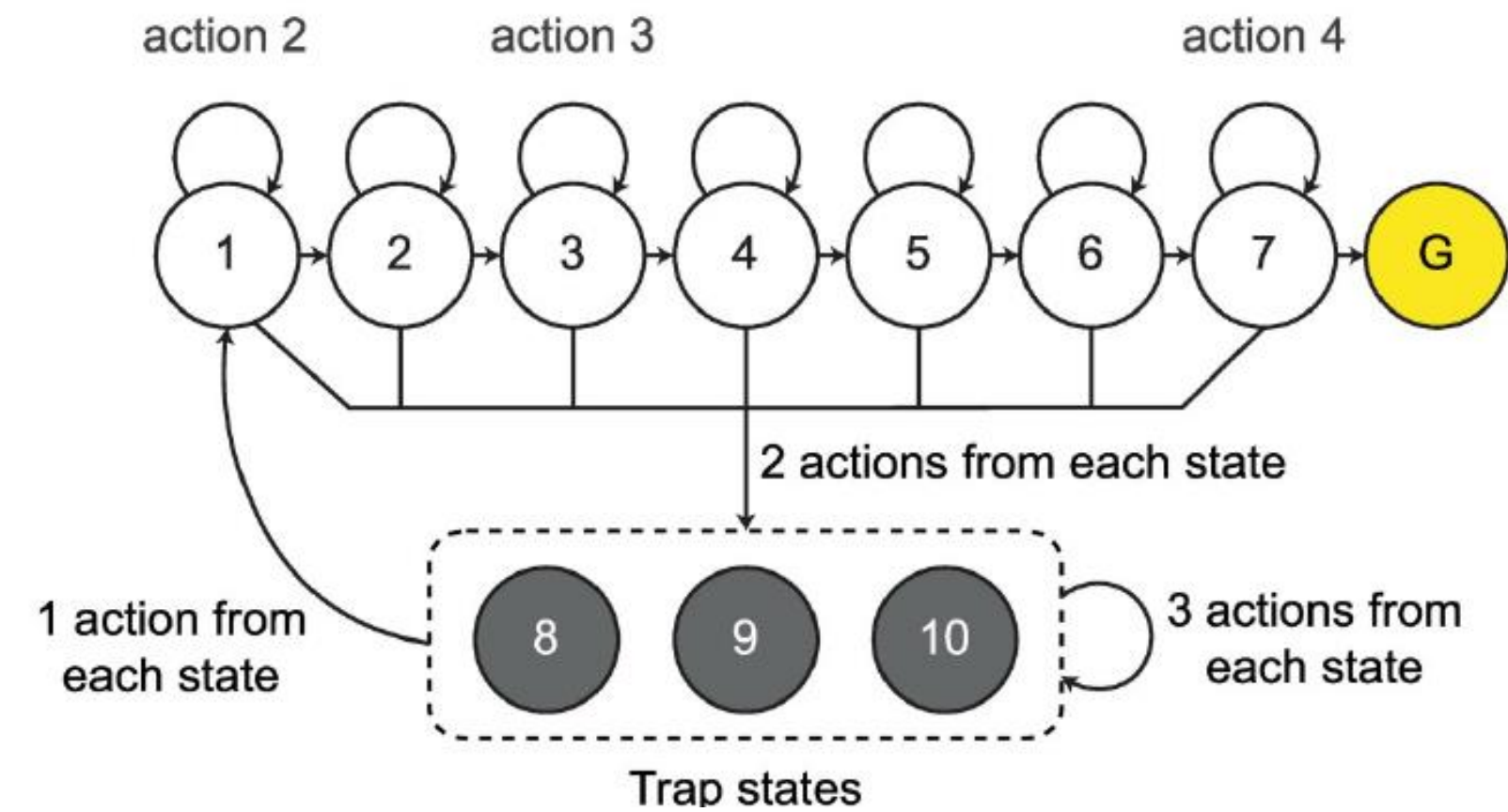
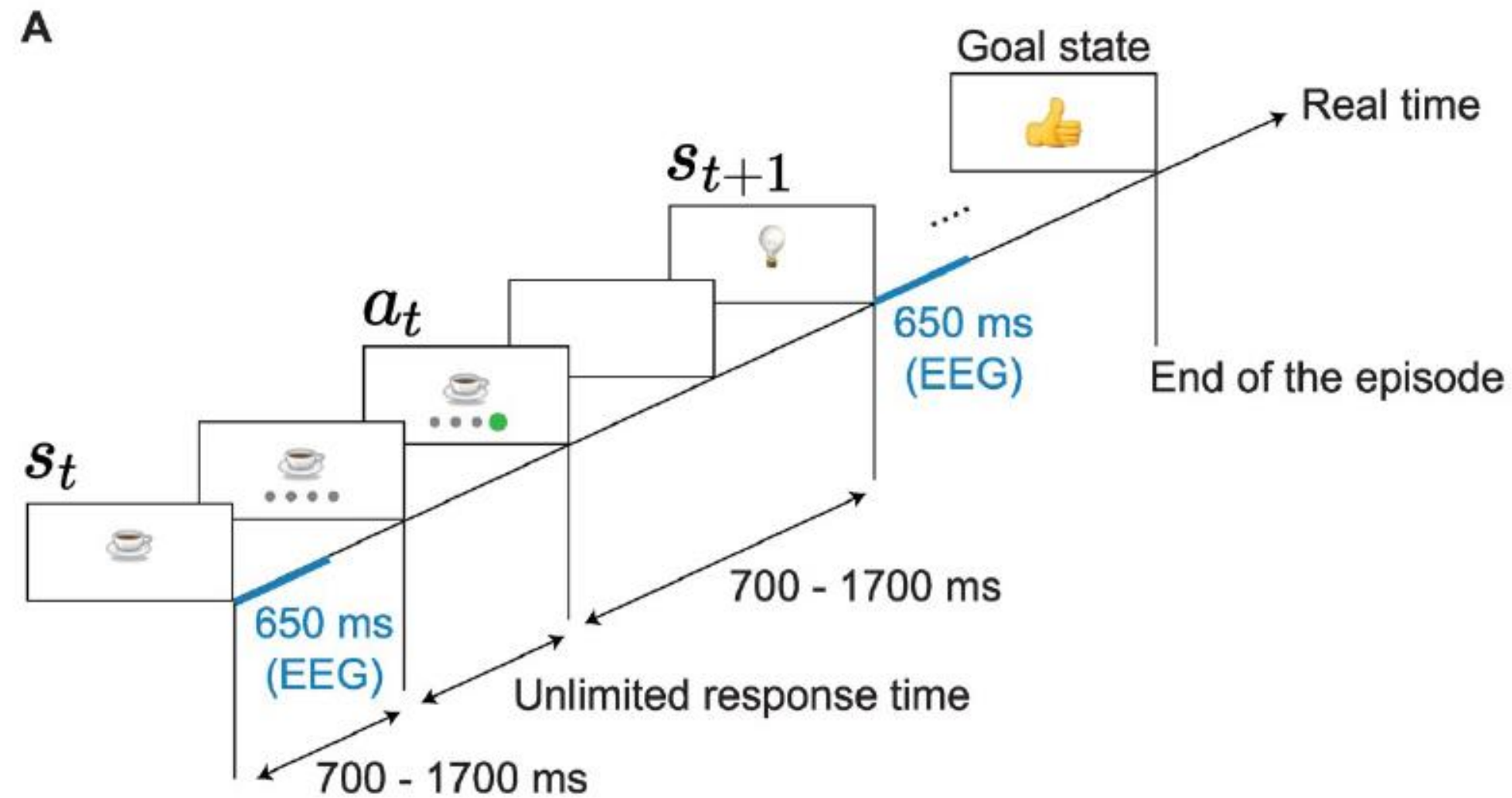
RL algorithms are inspired by human and animal behavior.

Thus, sometimes it is a good idea, how humans would perform in a given environment.

Markov Decision Processes are ideal testbeds for tabular RL algorithms.  
So, let us test humans in such an environment!



# Environment: Markov Decision Process



Finding 1)

Participants need about 150 actions in episode 1

Finding 2)

In episode 2, participants go straight to goal

Previous slide.

Human participants are put into a Markov Decision Process.

They have four action buttons to navigate from one image to the next.

They have been told before the experiment that there are 10 states and one goal state, each identified by an image. The 11 images (including goal) have been shown once.

Until image onset, participants have to wait for a time of about 1s until four grey disks were present – these are the action buttons.

The goal image in this example is the thumb-up image.

Right: Structure of the environment for the first 5 episodes (block 1).

Finding 1) humans are MUCH faster than the random exploration strategy to find the goal for the first time.

Finding 2) humans are extremely good in episodes 2-5 to return to the goal. The starting condition is not always state 1, but can also be a different state (varies across episodes, but the same starting state for all participants).

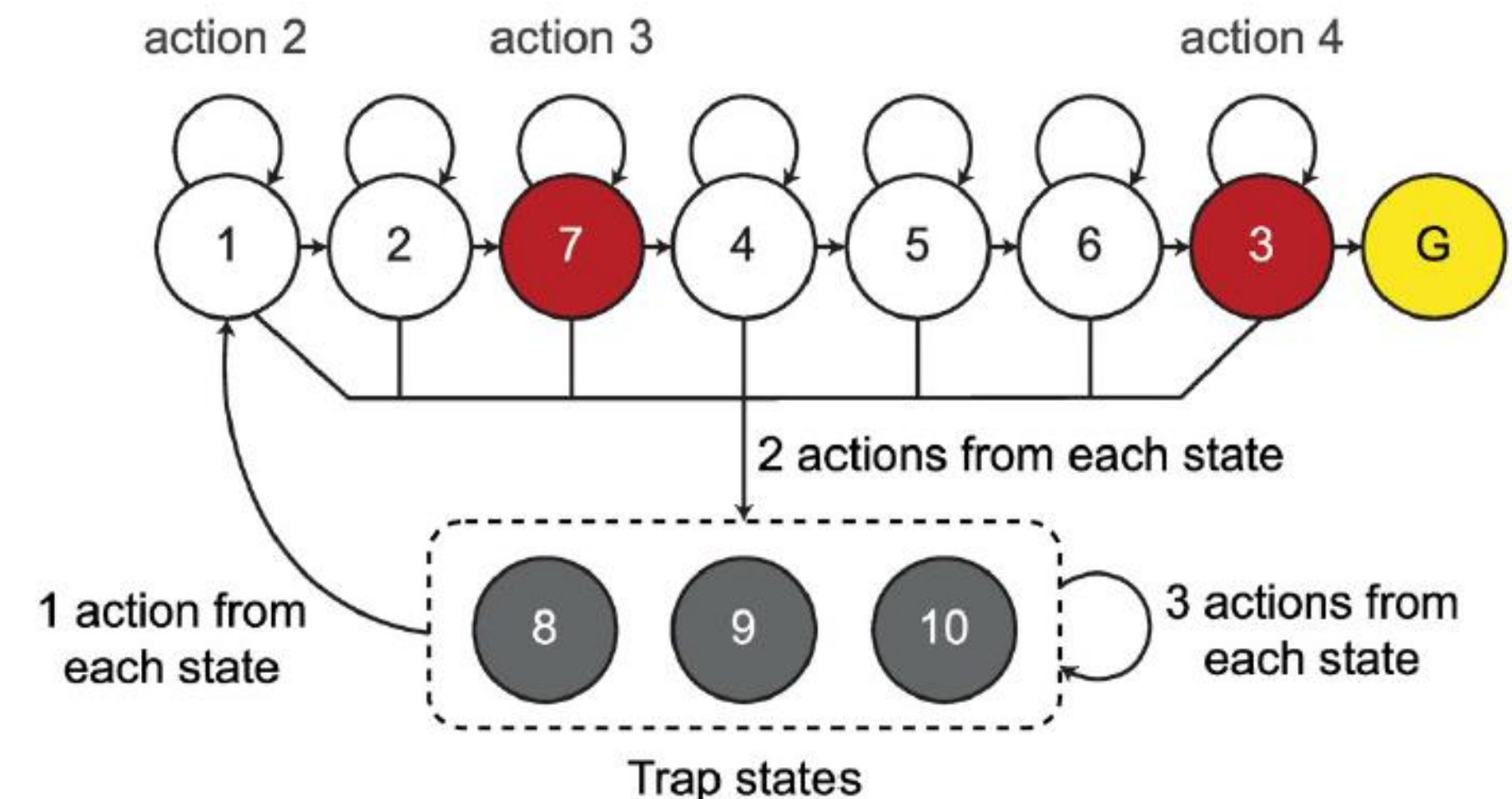
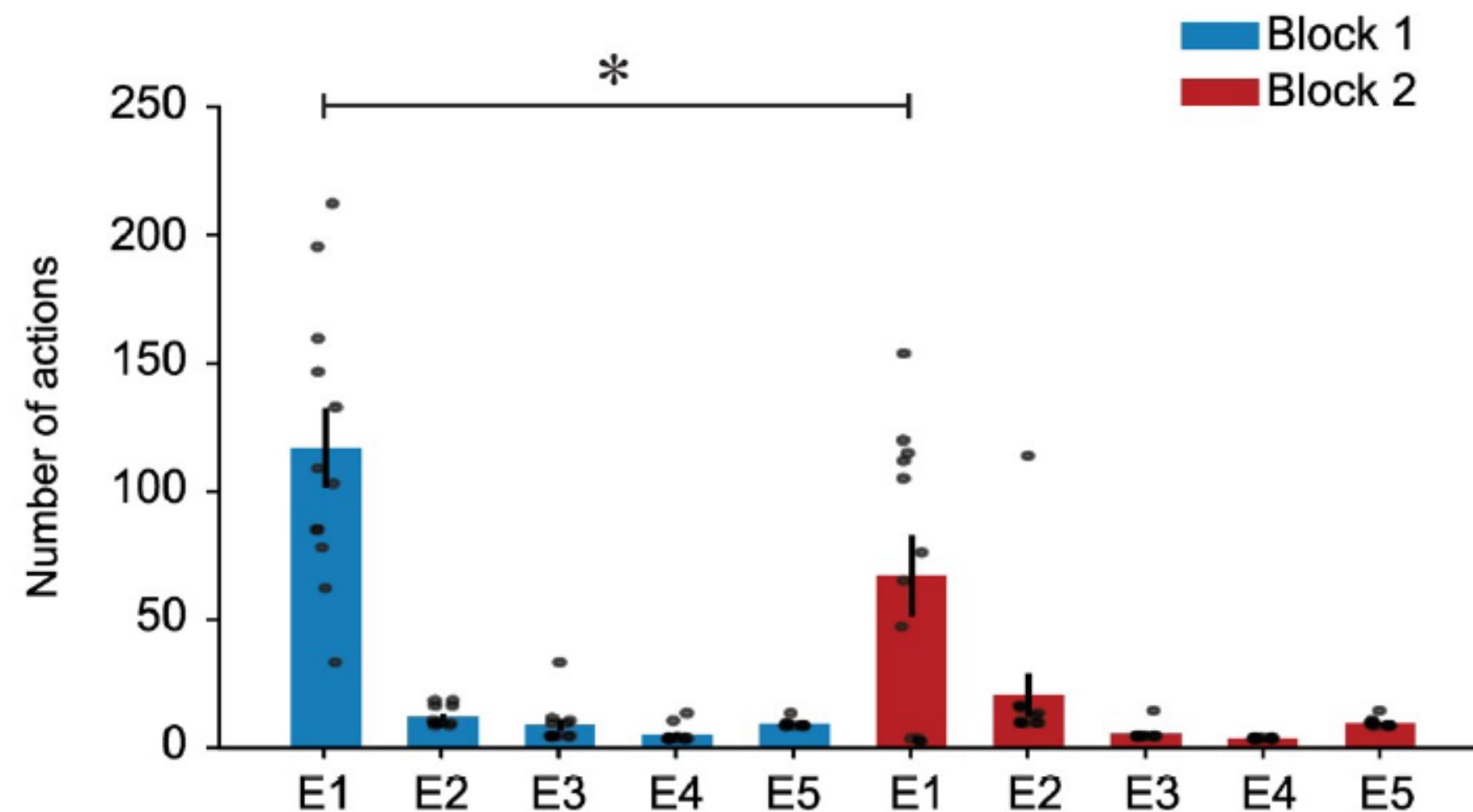
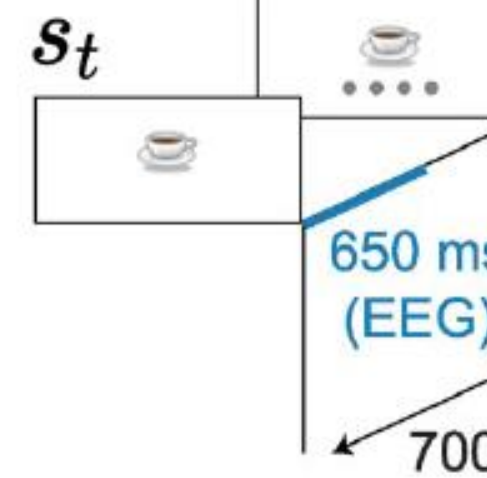
# Volatile Environment: Switch after episode 5

Finding 3)

In episodes 5 and 6, participants rapidly relearn!

Questions:

- Is Surprise necessary to explain relearning?
- Are humans model-based or model-free?
- Is novelty a good explanation of results?



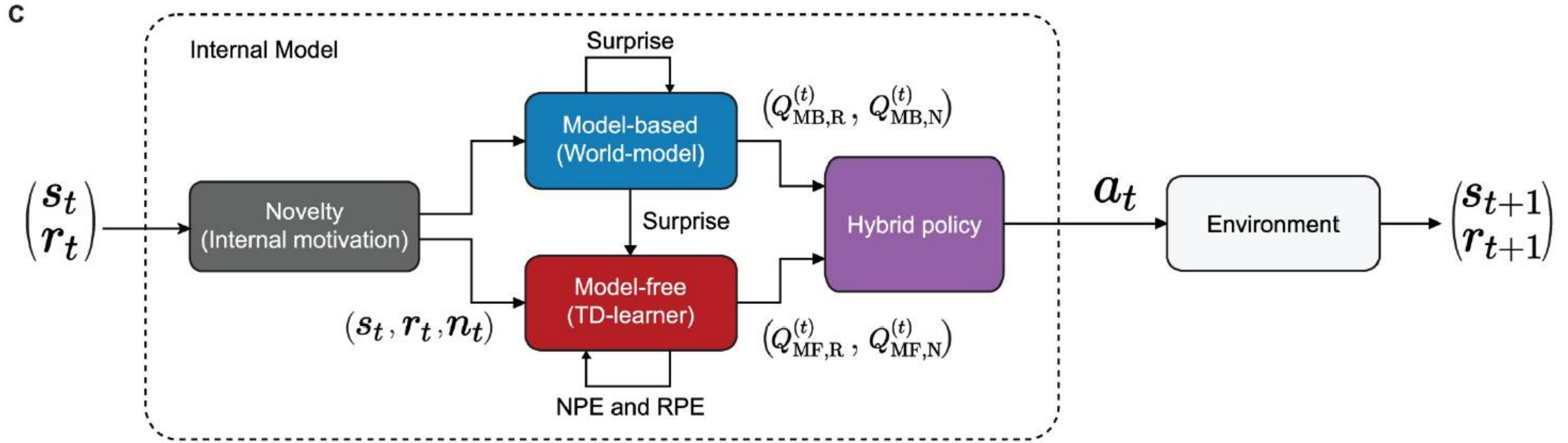
Previous slide.

After episode 5, states 3 and 7 have been swapped. Thus the environment is not stationary (volatile environment).

Humans rapidly readapt.  
Would algorithms also re-adapt?



# Review: Hybrid model with separate paths Surprise, Novelty, Reward (SurNoR)



$$\text{RPE} = [r_t + \gamma \max_{a'} Q_R(s', a') - Q_R(s, a)]$$

$$\text{NPE} = [n_t + \gamma \max_{a'} Q_N(s', a') - Q_N(s, a)]$$

4 separate  
sets of  
Q-values!

Previous slide.

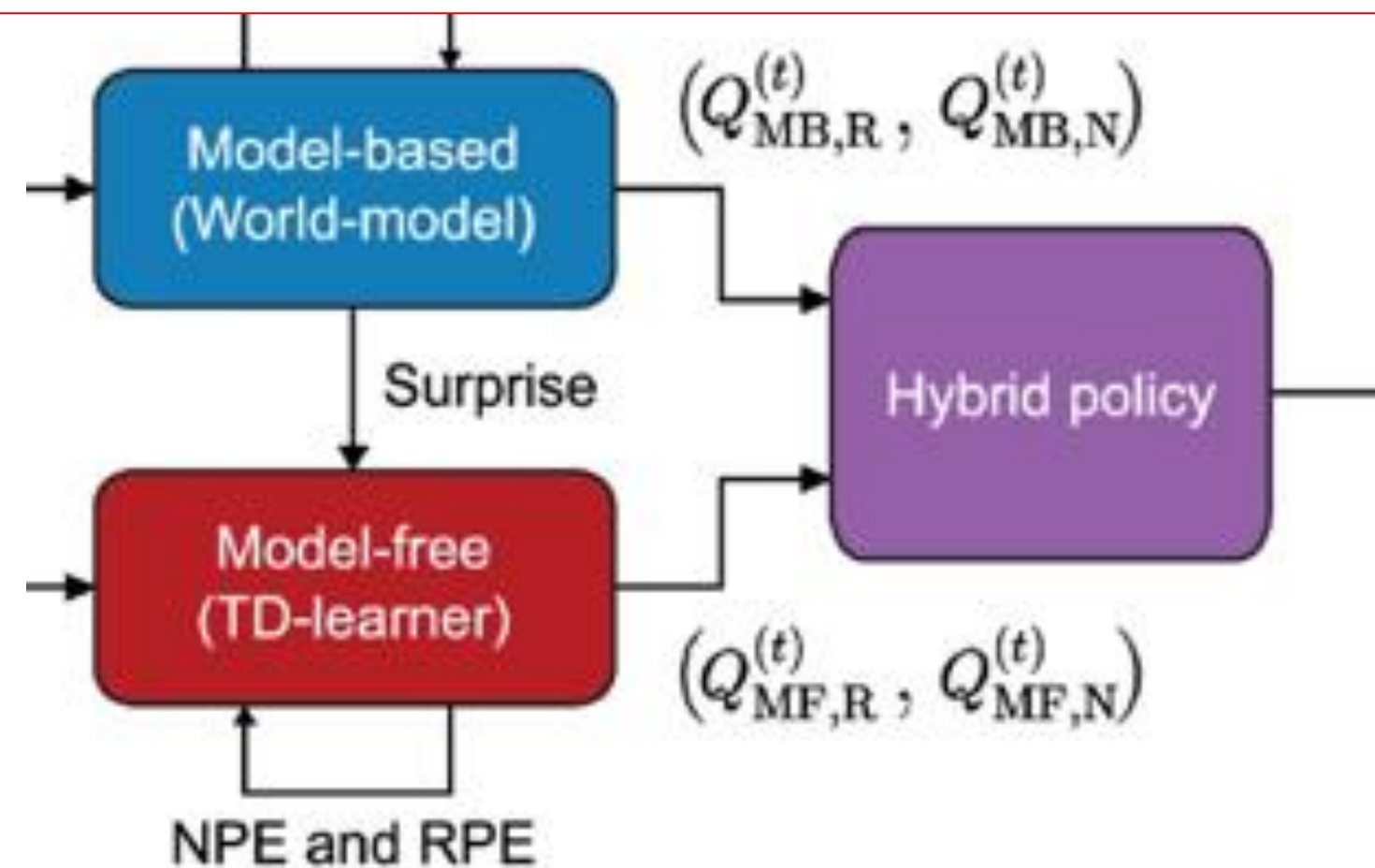
Note that in the formal theory of exploration bonus, we simply added the bonus in the Bellman equation.

However, here we claim that it is useful to develop two separate Bellman equations, one for novelty and one for reward. Each one has separate Q-values.

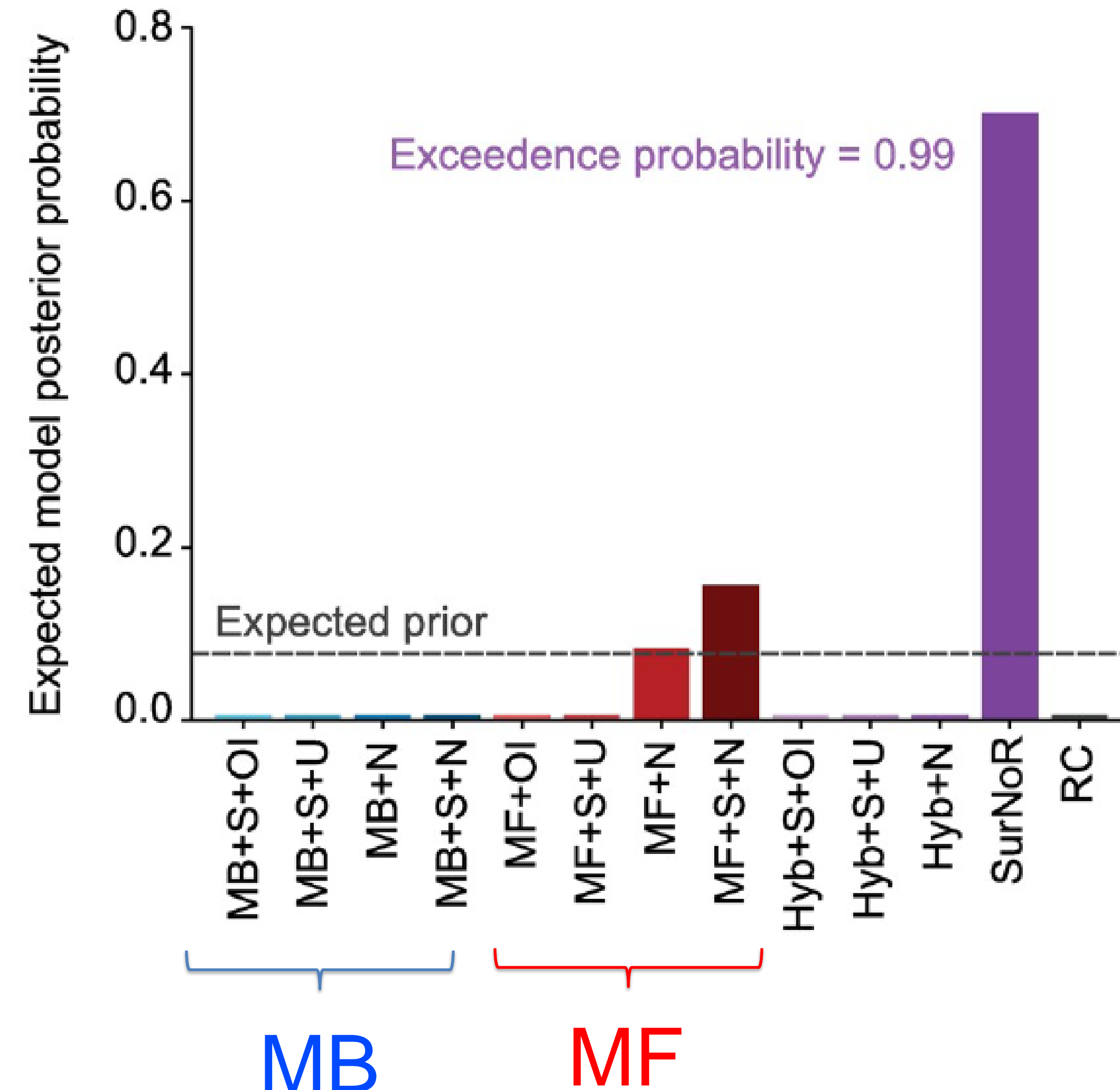
Each one of these, can be implemented as model-free or model-based.

# Comparison of Models: Surprise, Novelty, Reward

Finding 4)  
Rapid relearning needs surprise



- Turn off novelty
- Turn off surprise
- Turn off model-based  $\rightarrow$  MF
- Turn off model-free  $\rightarrow$  MB
- OI = Optimistic Initialization



Previous slide.

The best model is the combination of Surprise, Novelty and Reward (SuRNoR).

The second best model is model-free (MF) RL with surprise (S), novelty (N), and reward.

Turning off surprise lowers the performance (Hybrid model and surprise).

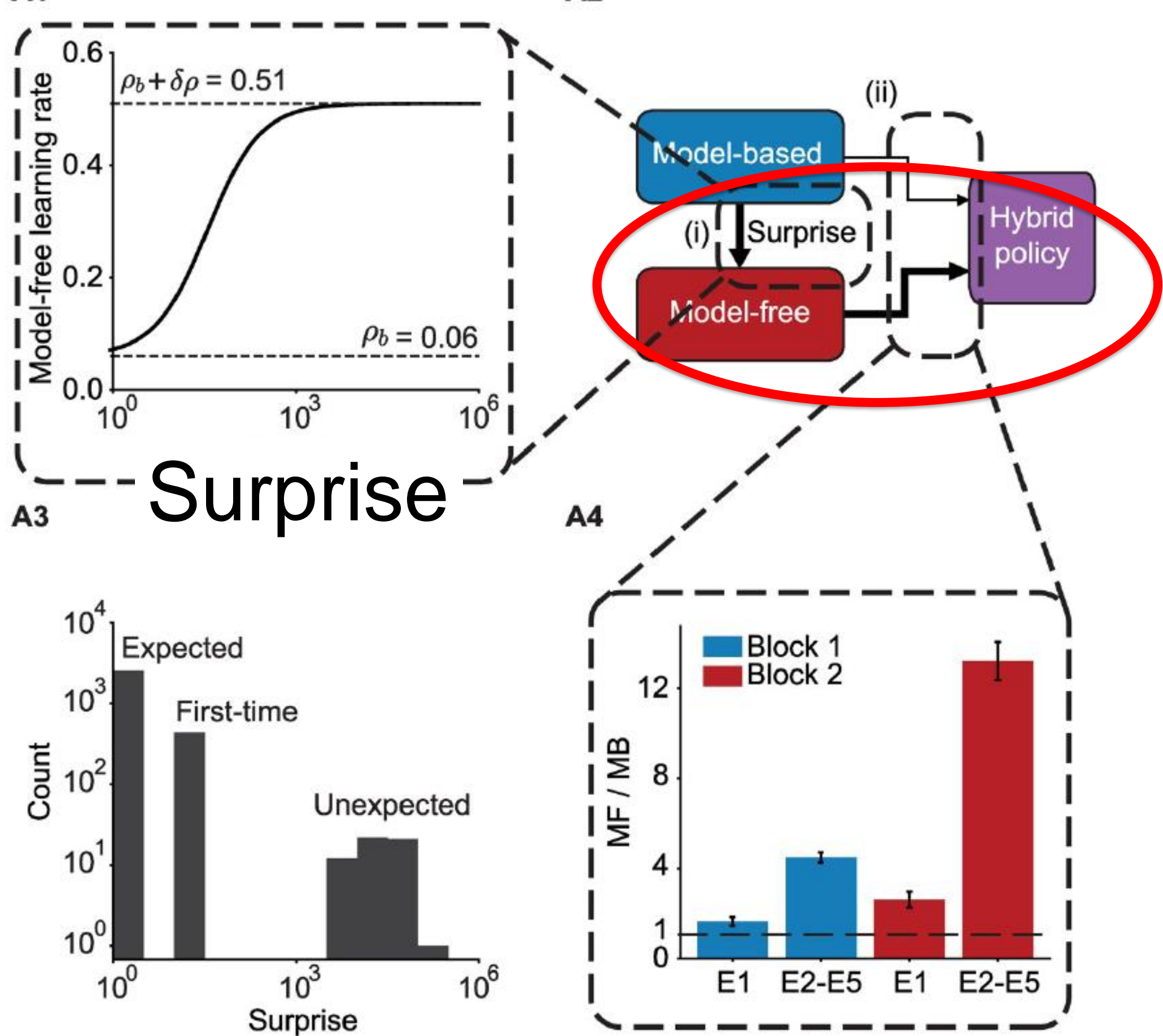
Model-based compares less well with human data than model-free. The combination of model-based and model-free (Hybrid model/SuRNoR) explains the data best.



# Relative importance of model-based versus model-free

Finding 5)  
Model-free dominates  
Human behavior!

surprise-modulated learning rate



Previous slide.

One can separately analyze the relative importance of the model-free and the model-based pathway to the hybrid policy in the SuRNoR model.

One finds that model-based never dominates, so that we conclude that human participants are best described by model-free algorithms with surprise.

# Surprise is used modulate learning in RL

Finding 6)

Surprise is against expectations.

Hence surprise needs a **world model**.

However, world model is

- Not used to do planning!
- Only used to extract surprise!

## **World-model not used for planning!**

Previous slide.

Surprise needs a world model, but we said that the model-free algorithm better explains the behavior.

The interpretation is that human participants develop a model of the world, but they only use it to detect surprise (change points) which allows them to re-adapt the model.

But they do not use it to plan ahead or do updates of the Bellman equation in the background.

# Surprise, World models, and Planning

Finding 6)

World model is available to humans

- But not used to do planning!
- Only used to extract surprise!

For humans:

- Planning is hard (not intuitive/natural)
- Exception: Planning in 2-dim or 3-dim environments
- Planning needs 'paper and pencil': "let's work this out"

Humans are not 'optimal'. Humans use heuristics.

Heuristics is mostly good for natural tasks.

Markov Decision Problems are 'not natural'

Previous slide.

Planning is simple for humans in 2-dim or 3-dim environments.

But not for Markov Decision Problems.

Abstract problems require (for most humans) a slow process of math-like solution process: whenever you feel, it would be easier to work something out with paper or pencil, you try to use a 'world model' that is non-intuitive for humans.

# Reward-based learning versus Surprise-based learning

Reward-Prediction Error → Surprise

defined as  
TD error

→ defined as  
Bayes Factor Surprise

stimulated by  
chocolate, money,  
praise, ...

→ stimulated by observations  
not consistent with momentary  
model of environment

modulates  
learning rate

→ modulates  
learning rate

Previous slide.

Summary: Comparison of Reward Prediction Error and Surprise.



# Second experiment: Detect Brain Signals during Reinforcement Learning with Surprise Trials



fMRI machine  
(standard image from WEB)

*V. Liakoni et al. (2022), Brain signals of  
a Surprise-Actor-Critic model.  
NeuroImage 246 (2022) 118780*

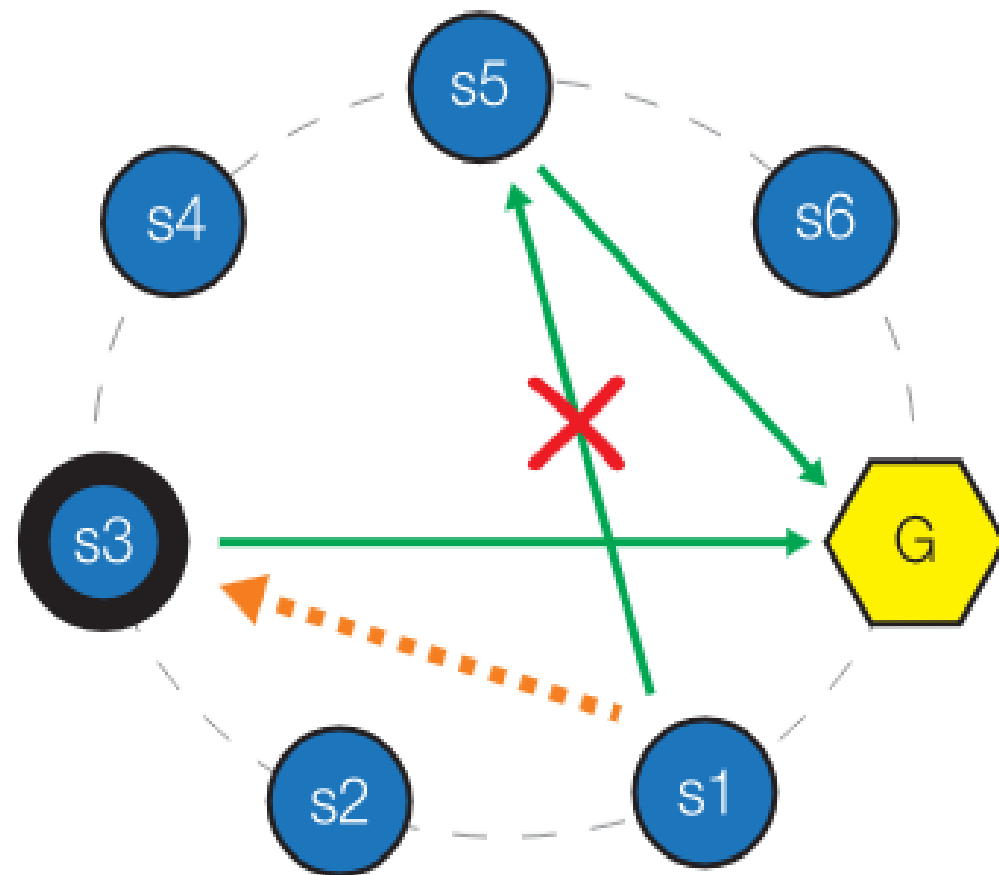
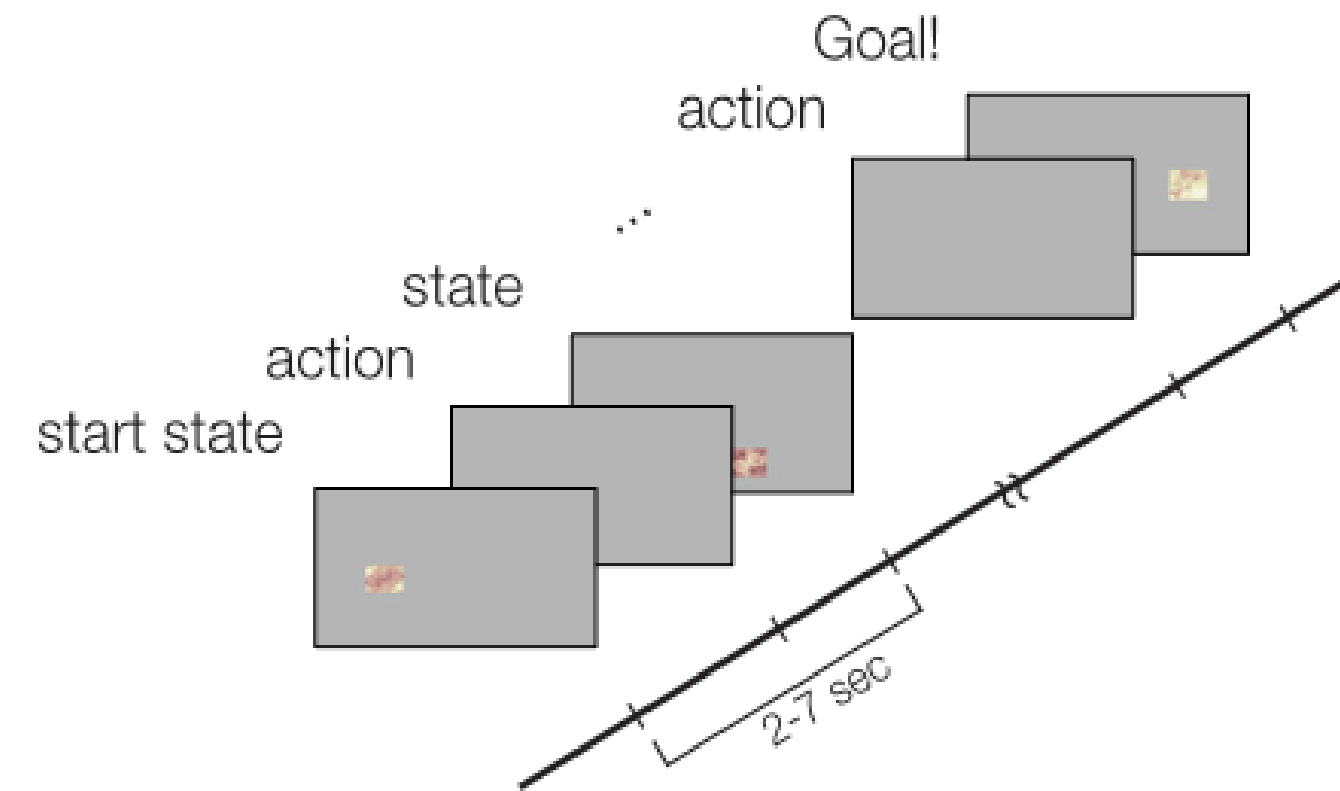
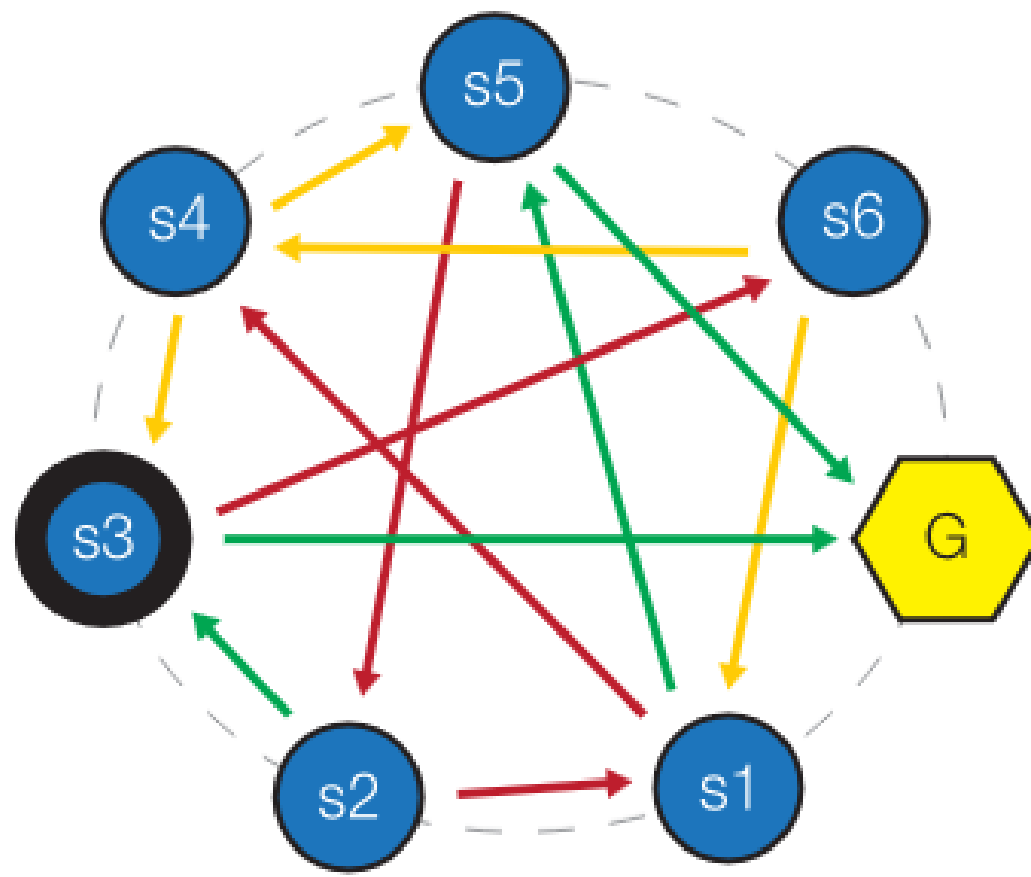
Previous slide.

We now turn to the second experiment. It is similar in spirit to the first one, but simpler and can be used in a brain-imaging device for functional Magnetic Resonance Imaging (fMRI).

In the fMRI, activated brain areas are visible by increased oxygen indicating increased blood circulation (called BOLD signal).

Normally, the BOLD signal is compared across two conditions A and B and the difference in the BOLD signal is plotted and projected onto a cut through the brain.

# Behavioral Task used to detect Brain Signals



*V. Liakoni et al. (2022), Brain signals of a Surprise-Actor-Critic model. NeuroImage 246 (2022) 118780*

Previous slide.

The task is related to the one on the previous slides, but simplified to make it usable in fMRI devices. The aim is to record brain signals correlated with an RL model.

The task consists of 7 states and one goal state. In each state two actions are possible. One of these (green) brings the participant closer to the goal.

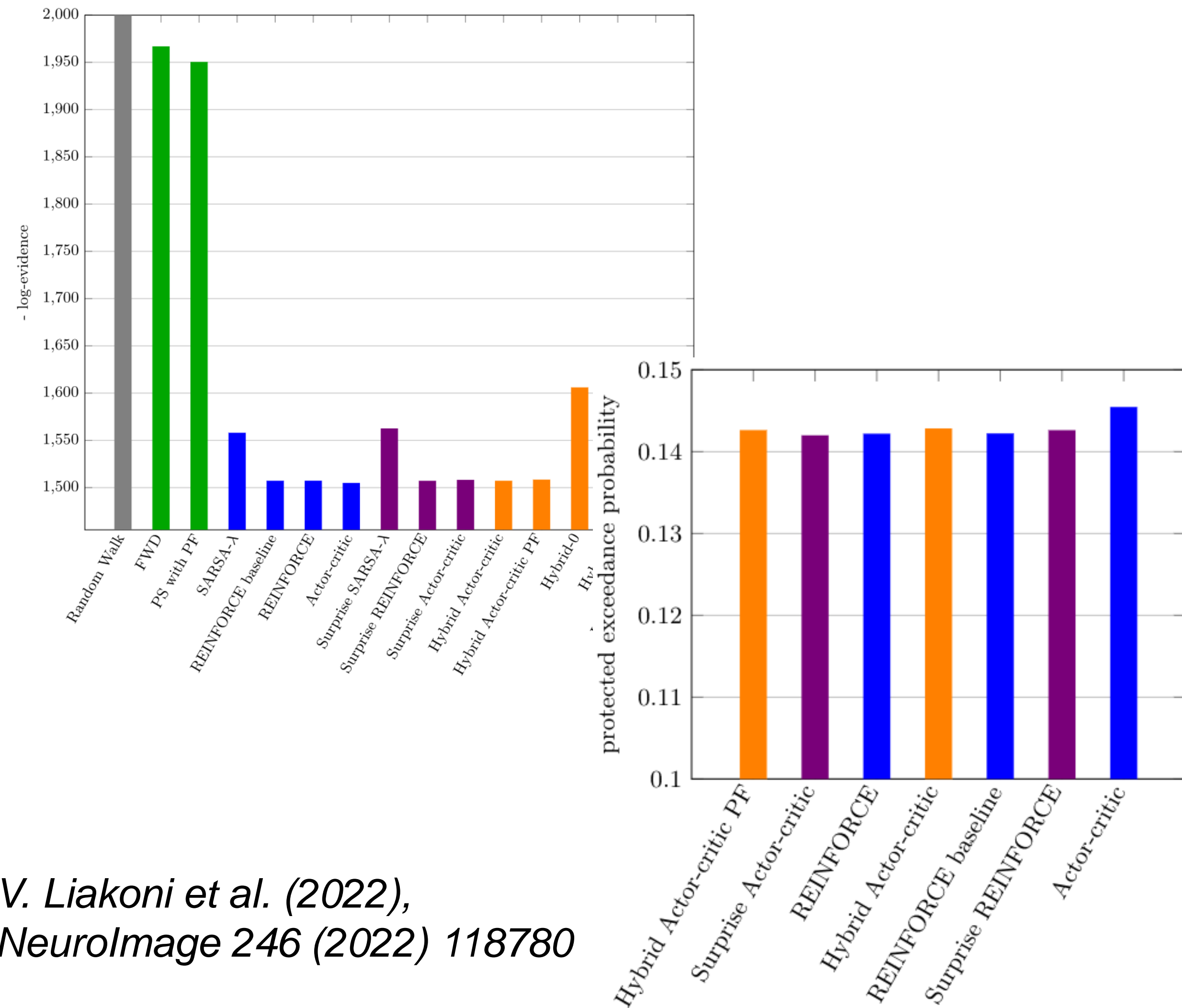
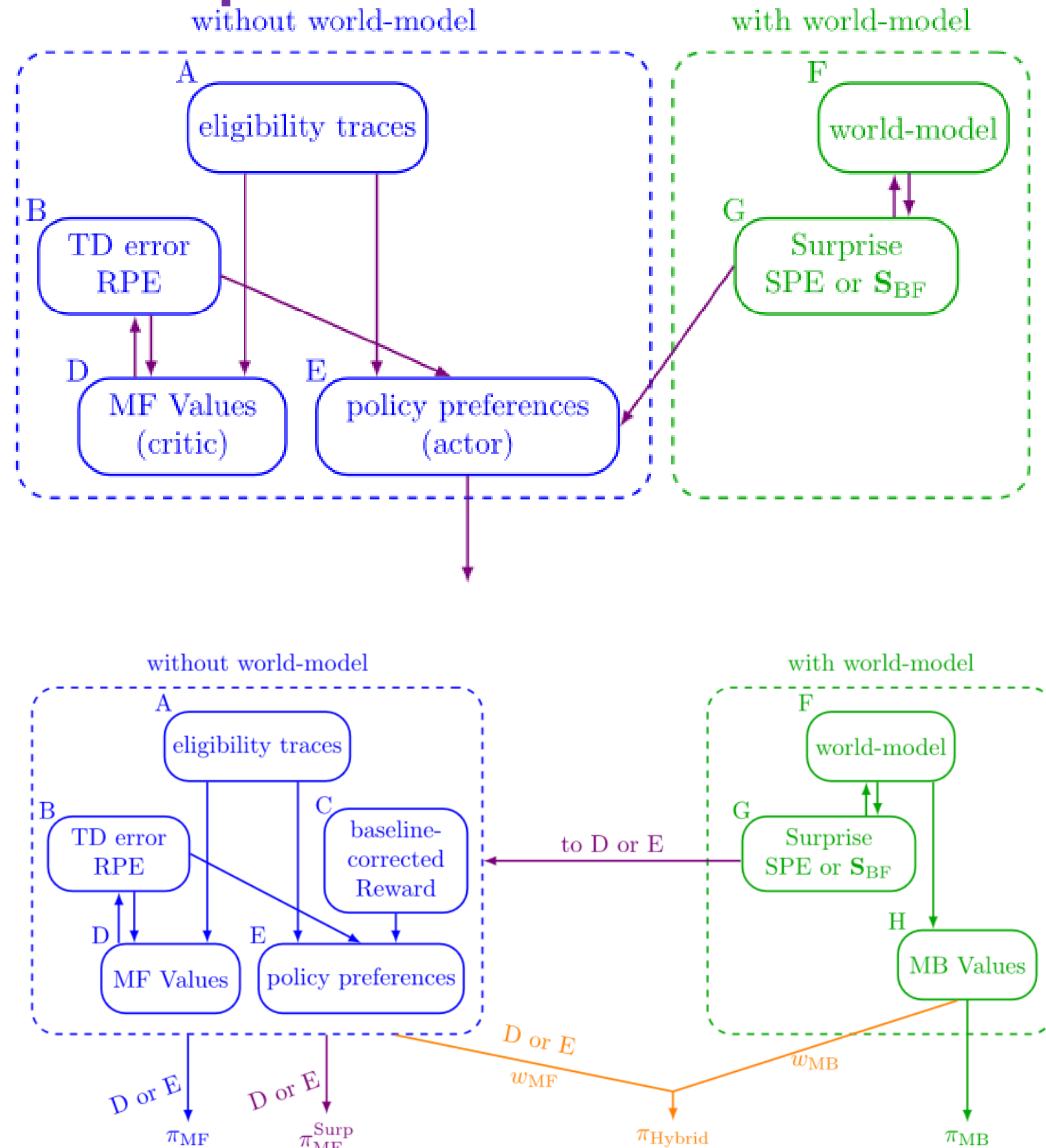
Importantly, several states can be classified as 'two steps from goal' and others as 'three steps from goal'.

This is important because if after  $N$  trials, the states in the same category (e.g., 'two steps from goal' have been visited an equal number of times), then they are expected to have the same  $V$ -values.

Bottom: Occasionally an action that would normally lead closer to the goal leads to a completely different state (surprise trials). Note that the both state  $S3$  and  $S5$  have the same distance from goal (one step from goal), and therefore the same  $V$ -values. Hence a potential surprise signal is not influenced by an additional difference in  $V$ -values.

# Surprise Actor-Critic and other RL algorithms

## Surprise Actor-Critic



V. Liakoni et al. (2022),  
NeuroImage 246 (2022) 118780

Previous slide.

LEFT:

One specific algorithm is the Surprise-modulated Actor-Critic algorithm (top). The world model (right box) is only used in order to modulate the learning rate of the actor-critic, but not for planning. The actor-critic is model-free with eligibility traces.

Bottom: RL algorithms come in different flavors:

- 1) Model based (right) versus model-free (left) or hybrid therefore (orange, bottom)
- 2) TD-type algorithm versus policy gradient algorithm (actor, right with-in left box) or actor-critic that combines TD (value estimation) with policy gradient (actor).

The Surprise Actor-Critic algorithm is one specific combination of these ingredients, but many other algorithms can be formulated in this framework.

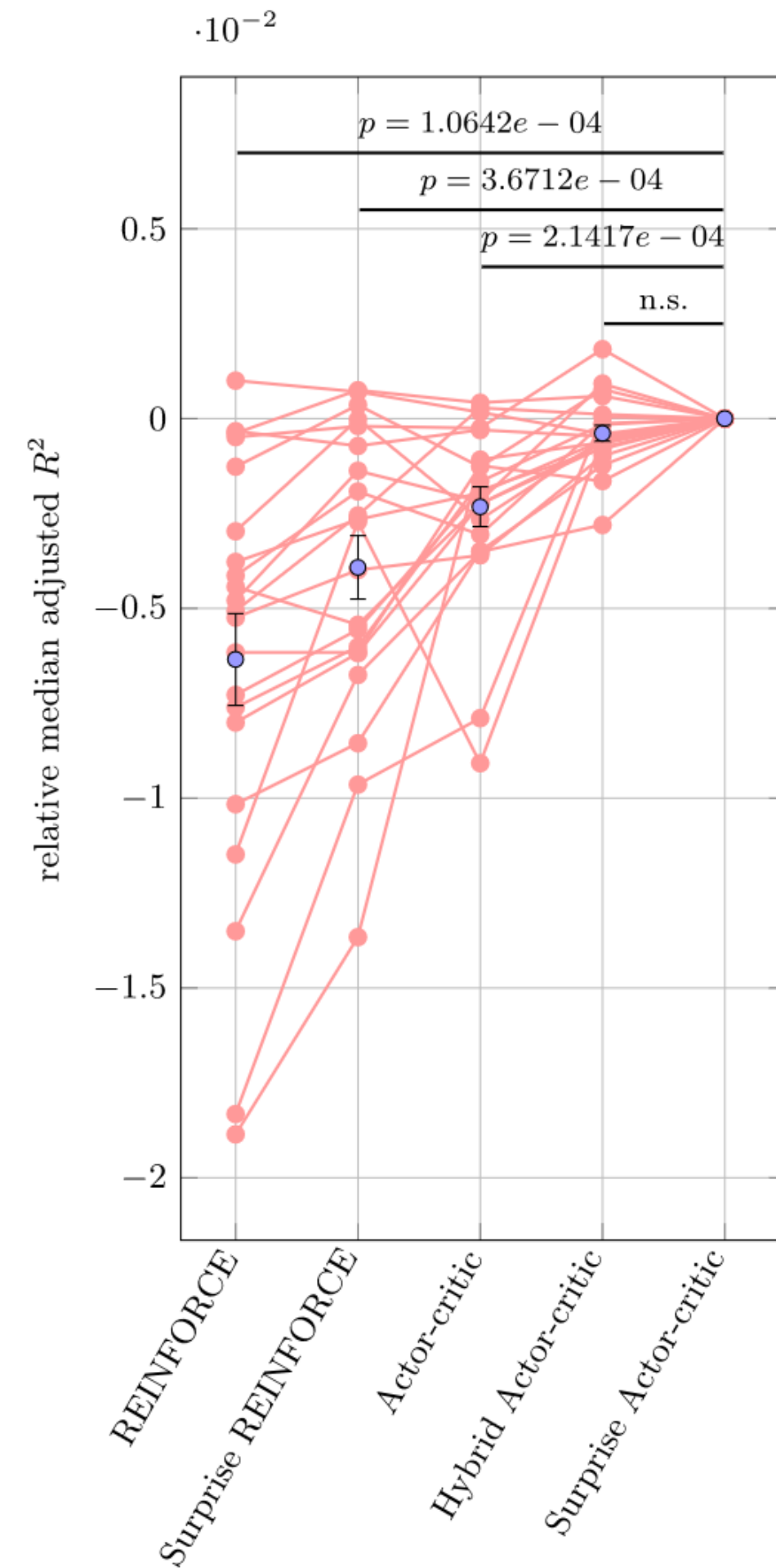
RIGHT: Performance in terms of negative log-evidence of many algorithms in explaining the behavior data of human participants (lowest is best). The 7 best ones amongst these are indistinguishable in terms of log-evidence and also indistinguishable in terms of 'protected exceedance probability'.



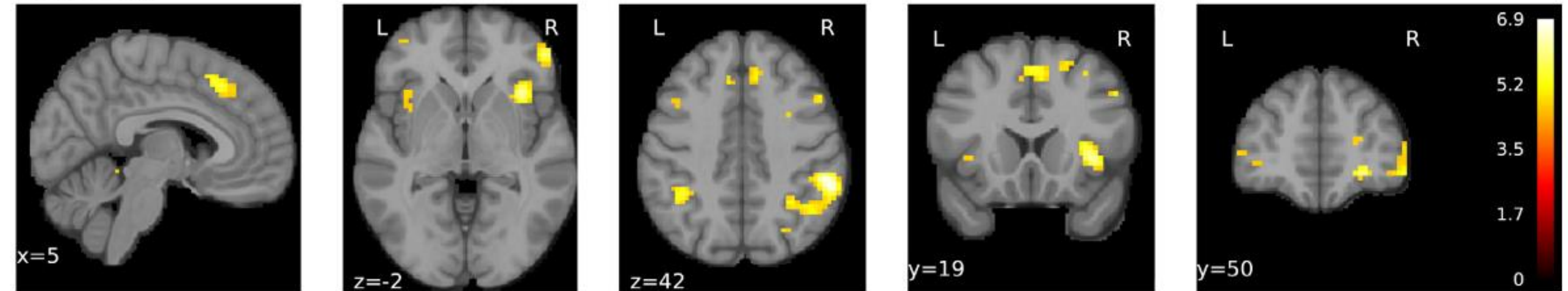
# Model Performance on explaining brain activity

## Surprise Actor-Critic

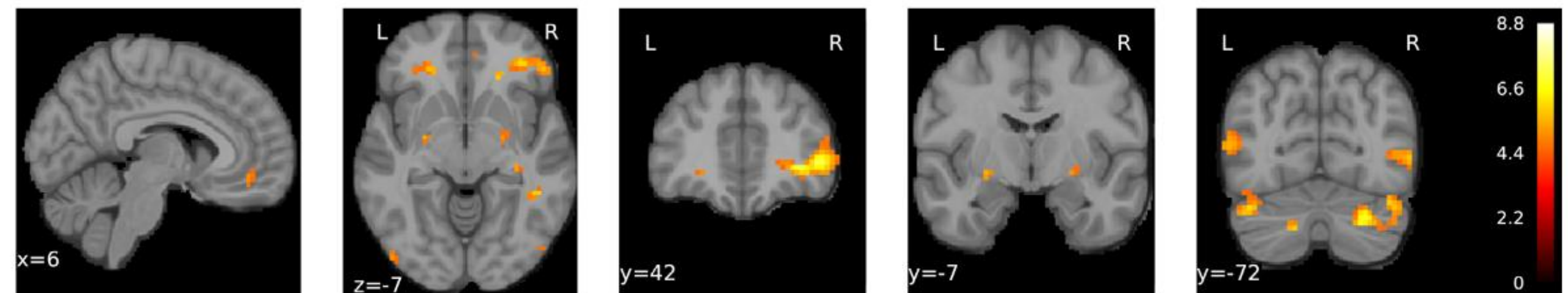
V. Liakoni et al. (2022),  
NeuroImage 246 (2022) 118780



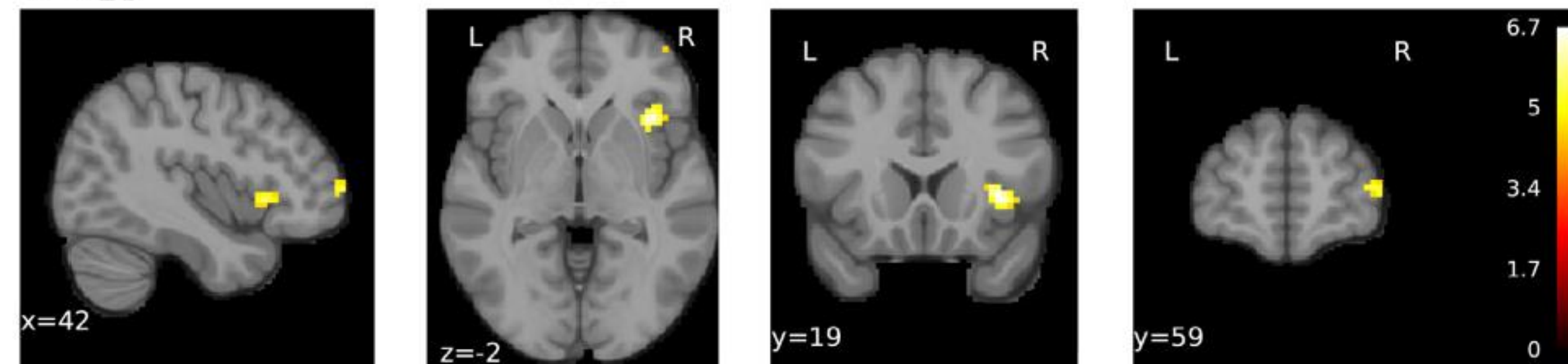
A  $SPE_{BF}$



B RPE



C  $S_{BF}$



Previous slide.

**Fig. 5. Neural model comparison.** A. Difference in median adjusted  $R^2$  across the whole brain for the winning computational models and the Surprise Actor-critic. Each red line corresponds to a participant and is centered with respect to the Surprise Actor-critic. The median adjusted  $R^2$  of the Surprise Actor-critic is significantly larger from the one of the REINFORCE, Surprise REINFORCE and the Actor-Critic (Wilcoxon signed rank test  $p < .001$ , i.e. passing a Bonferonni corrected threshold of 0.0125 for the 4 comparisons performed). The performance of the Surprise Actor-critic and of the Hybrid Actor-critic were not significantly differen

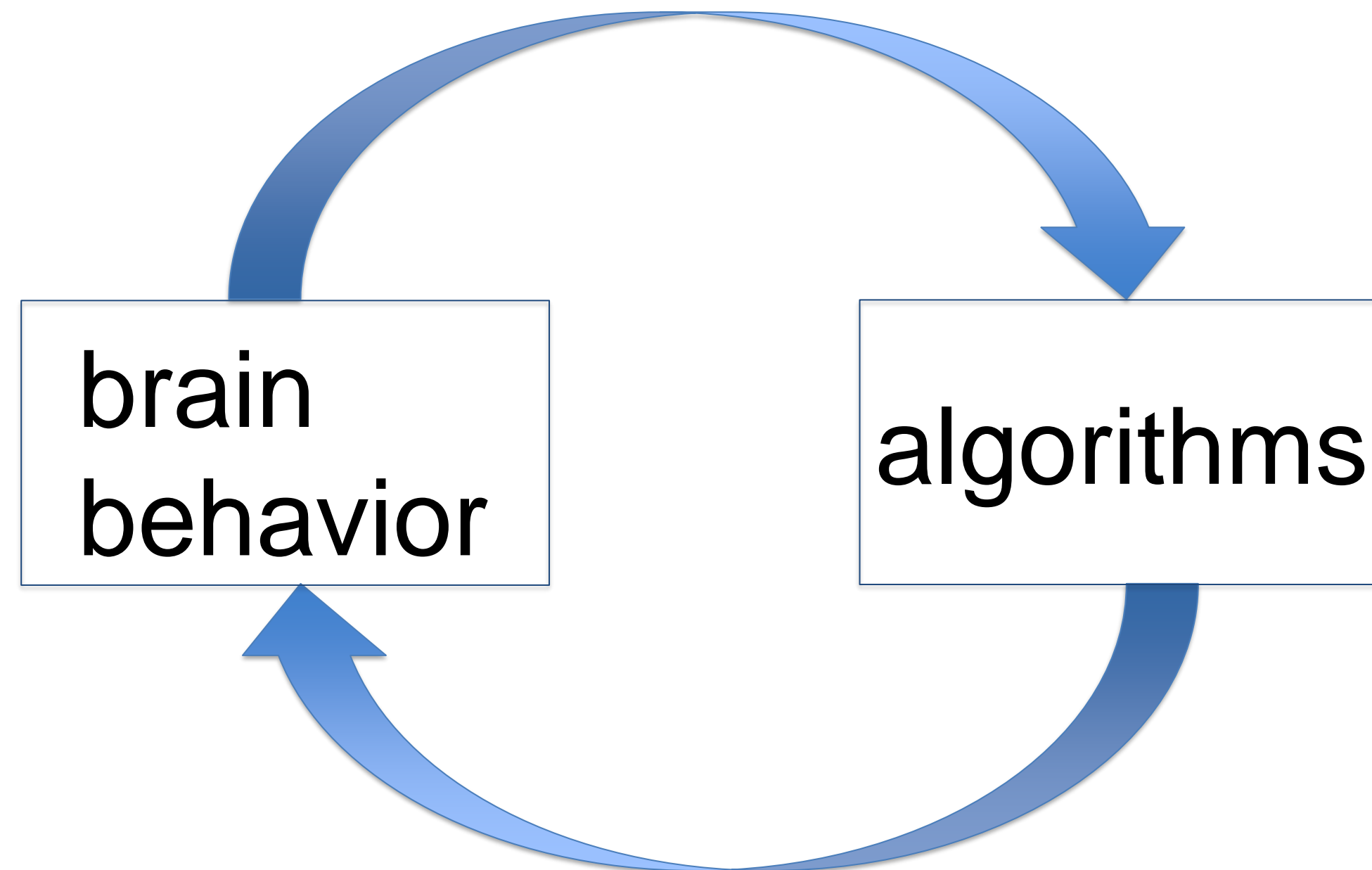
**Fig. 6. Neural correlates of learning signals of the Surprise Actor-critic – GLM<sub>4</sub>.** T-statistic maps (21 subjects, random effects whole brain analysis, cluster-wise correction with a cluster-defining threshold (CDT) of  $p = 10^{-4}$  and a FWE-corrected threshold of  $p = .05$ , nonparametric permutation test with maximum statistic approach) of A.  $SPE_{BF}$ . We find significant correlation in SMA, insula, middle frontal gyrus, angular gyrus, supramarginal gyrus and in the superior frontal gyrus. B. RPE. We find significant correlation in the inferior frontal and orbitofrontal gyrus, the striatum (putamen and pallidum), the vmPFC, and the inferior occipital gyrus. C.  $S_{BF}$ . We find significant correlation in the right insula and the right middle frontal gyrus.

Left: Message of Fig. 5: Brain data is best explained by the Surprise Actor-Critic (model-free) and equally well by the hybrid actor-critic. For Fig. 6 we use the model-free Surprise Actor-Critic since it is simpler.

Right: Message of Fig. 6: Precise brain areas can be identified that have activity significantly correlated with the State-Prediction Error (SPE), the Reward-Prediction Error (RPE), or the Bayes-Factor Surprise ( $S_{BF}$ ). The areas identified for RPE in this specific experiment are similar to those identified in many other studies. The state prediction error is the update signal for the transition matrix  $\hat{P}^{(t)}(s'|s, a)$ .



# Current Research in Reinforcement Learning and in Brain Sciences:



- Exploration → not exploration bonus, but separate modules
- Novelty → Novelty supports exploration
- Surprise → Surprise detects changes/adapts learning rate

Previous slide. Review from previous lectures.

RL has two roots: optimization for Markov Decision Problems and Brain sciences/psychology

The interaction has not stopped. Modern RL still takes up influences from Brain Sciences. Examples are the role of novelty, surprise, and their roles for exploration and in volatile environments.

At the same time RL has strongly influence the brain sciences!

The END

Previous slide.

# Learning in Neural Networks

Wulfram Gerstner

EPFL, Lausanne, Switzerland

## The role of exploration, novelty, and surprise in RL

### Appendix: More on Formal Exploration Bonus

(Thanks to Dr. Alireza Modirshanechi)

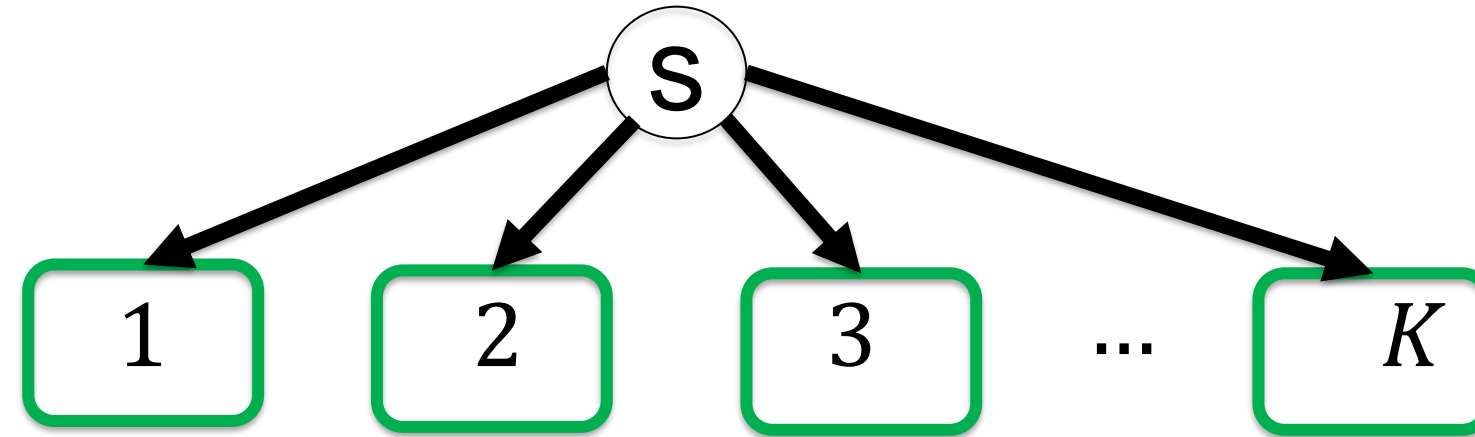
Previous slide.

We start with some results from the formal theory of exploration.

- a) For multi-armed bandits (1-step horizon)
- b) For full Markov Decision Problem (multi-step horizon)

# Review: Multi-armed Bandits: MAB (1-step horizon)

- Single state. We have  $K$  possible actions:



Which action to choose at time  $t$ ?

- With true average reward:

$$\mu_i = E[r|a = i]$$

 $\mu_1$  $\mu_2$  $\mu_3$  $\dots$  $\mu_K$ 

Optimal policy:  $a_t = \arg \max_i \mu_i$

- Naïve estimates of averages:

$$\hat{\mu}_i^{(t)} = \frac{\sum_{\tau \in T_i^{(t)}} r_\tau}{|T_i^{(t)}|}$$

 $\hat{\mu}_1^{(t)}$  $\hat{\mu}_2^{(t)}$  $\hat{\mu}_3^{(t)}$  $\dots$  $\hat{\mu}_K^{(t)}$ 

$$T_i^{(t)} = \{\tau \leq t : a_\tau = i\}$$

Not optimal:  $a_t = \arg \max_i \hat{\mu}_i^{(t)}$

Solutions based on random exploration:

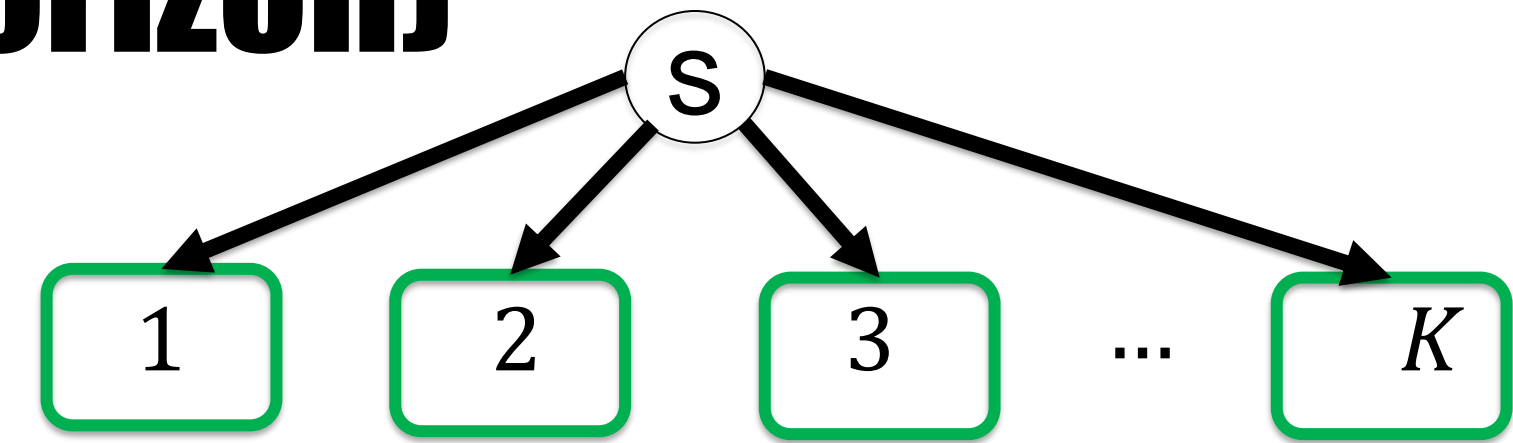
- Epsilon-greedy
- Softmax

## Comments for the previous slide:

- If we knew the exact average reward  $\mu_i = E[r|a = i]$  of each arm, then the optimal solution would trivially be to choose the arm with highest average reward:  $a_t = \arg \max_i \mu_i$
- A naïve approach is to estimate the average reward by the empirical averages and greedily choose the action with maximum estimated average reward:  $a_t = \arg \max_i \hat{\mu}_i^{(t)}$
- The naïve greedy policy is prone to fail in finding the best action.
- You have seen epsilon-greedy and the softmax policy as two approaches for dealing with this problem by adding randomness to the action-selection.
- Our focus will be on “directed exploration” by using exploration bonuses.



# Regret in Multi-armed Bandits (1-step horizon)



- MAB with  $K$  possible actions:

$$\mu_i = E[r|a = i]$$

Highest reward rate:  $\mu^* = \max_i \mu_i$

- “Regret” of algorithm  $A$   
(e. g.,  $\epsilon$ -greedy):

$$R_A(T) = E_A \left[ \sum_{t=1}^T \mu^* - \mu_{a_t} \right]$$

with best  
action you  
can choose

with your  
actual choices

- Consistent algorithms:

$$\lim_{T \rightarrow \infty} \frac{R_A(T)}{T} = 0 \quad \Rightarrow \quad \lim_{T \rightarrow \infty} \frac{E_A[\sum_{t=1}^T \mu_{a_t}]}{T} = \mu^*$$

- Theorem 1 of Lai and Robbins 1985:

Under specific conditions, if algorithm  $A$  is consistent, then, loosely speaking,  $R_A(T)$  is at least proportional to  $\log T$ .

a loose notion of optimality

**Idea:** you need to play other actions, even if that means that  $R_A(T)$  increases

## Comments for the previous slide:

- Before discussing how to deal with exploration-exploitation dilemma, we discuss a common method for evaluating different algorithms in multi-armed bandits.
- A key notion to evaluate an algorithm  $A$  is regret  $R_A(T)$  measuring the **expected** difference between the choices of the algorithm and the best possible actions, summed over the first  $T$  steps.
- An algorithm is called consistent, if its average regret  $\frac{R_A(T)}{T}$  vanishes over time.
- It is proven (under certain conditions; see Lai and Robbins 1985 in Advances in Applied Mathematics) that the regret  $R_A(T)$  of a consistent algorithm scales at least logarithmically with time  $T$ .
- At the same time, consistency requires that the regret  $R_A(T)$  increase slower than  $T$ . The statement therefore is  $\log T$  is the best you can do.
- This framework introduces a loose notion of optimality: An optimal algorithm is a consistent algorithm whose regret scales logarithmically with time  $T$ .

# Example: average rewards in MAB (1-step horizon)

- MAB with 4 possible actions (Example):

$$\mu_i = E[r|a = i]$$

rewards are stochastic (binomial)

$$P(r_t = 2\mu_i | a = i) = 0.5 = P(r_t = 0 | a = i)$$

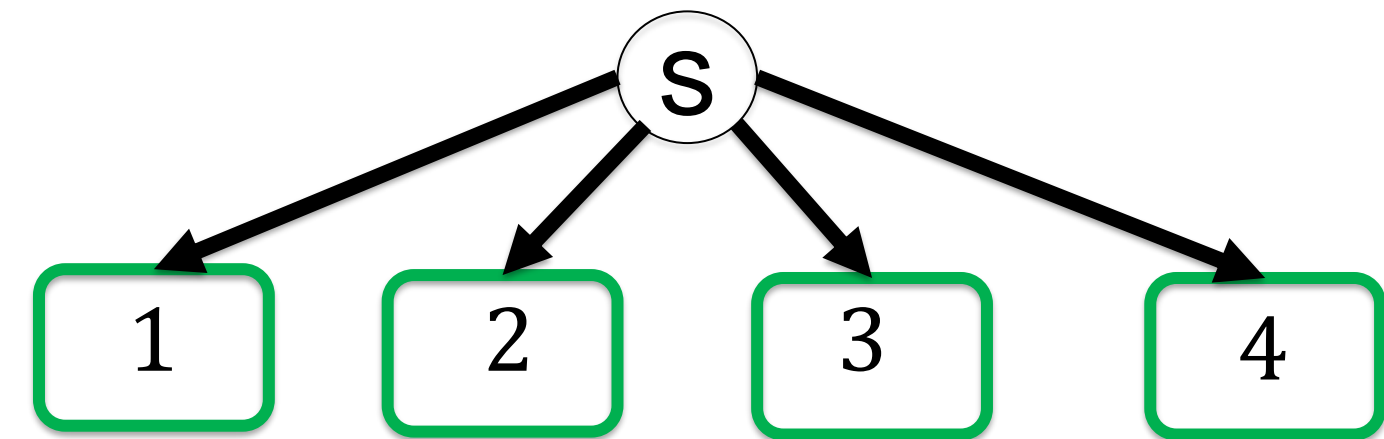
$$\mu_1 = 1$$

$$\mu_2 = 0.9$$

$$\mu_3 = 9.9$$

$$\mu_4 = 10.0$$

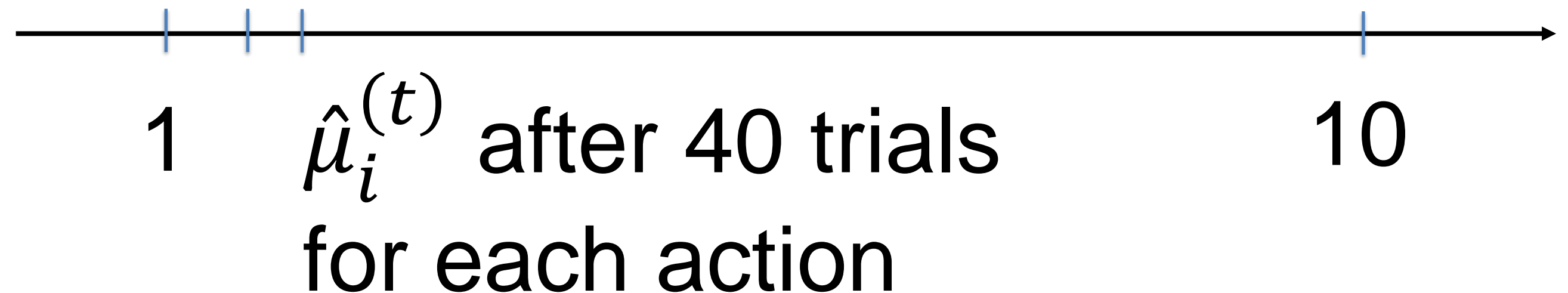
$$\hat{\mu}_i^{(t)} = \frac{\sum_{\tau \in T_i^{(t)}} r_\tau}{|T_i^{(t)}|}$$



What is the probability that  $\hat{\mu}_4^{(160)} = 2$ ?

[ ]  $\binom{N}{k} 2^{(-40)}$

[ ] between  $10^{(-7)}$  and  $10^{(-8)}$



- Comments for the previous slide:
- We assume that each action is played 40 times ( $t=160$  total time).
- How likely is it in the above example that the 'best' action with mean reward 10 would have after 40 trials a value of 2?

# Example: average rewards in MAB (1-step horizon)

- MAB with 4 possible actions (Example):

$$\mu_i = E[r|a = i]$$

rewards are stochastic (binomial)

$$P(r_t = 2\mu_i | a = i) = 0.5 = P(r_t = 0 | a = i)$$

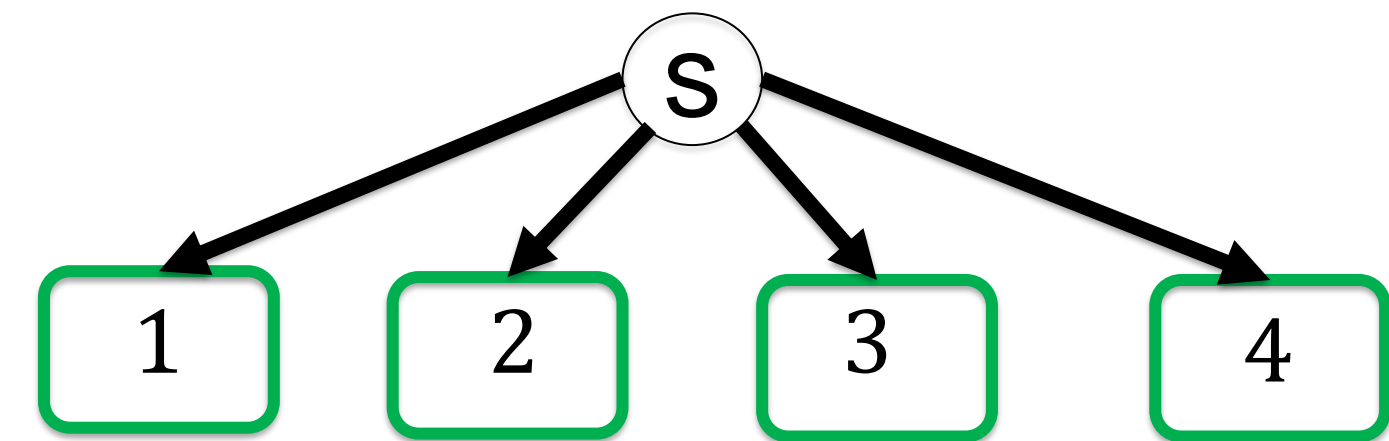
$$\mu_1 = 1$$

$$\mu_2 = 0.9$$

$$\mu_3 = 9.9$$

$$\mu_4 = 10.0$$

$$\hat{\mu}_i^{(t)} = \frac{\sum_{\tau \in T_i^{(t)}} r_\tau}{|T_i^{(t)}|}$$



$$R_A(T) = E_A \left[ \sum_{t=1}^T \mu^* - \mu_{a_t} \right]$$

after 200 trials each

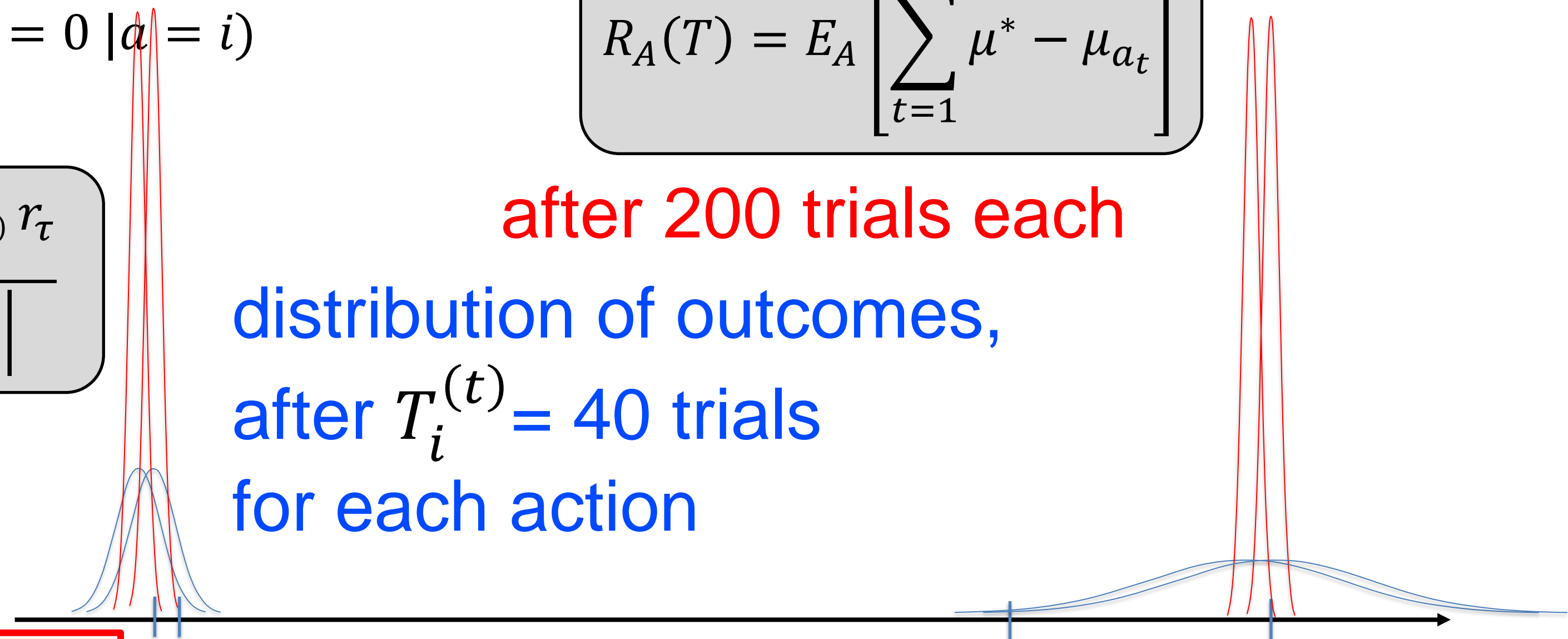
distribution of outcomes,  
after  $T_i^{(t)} = 40$  trials  
for each action

**Idea:** play other actions if tails of distribution overlap

1

$\hat{\mu}_i^{(t)}$  after 40 trials  
for each action

10

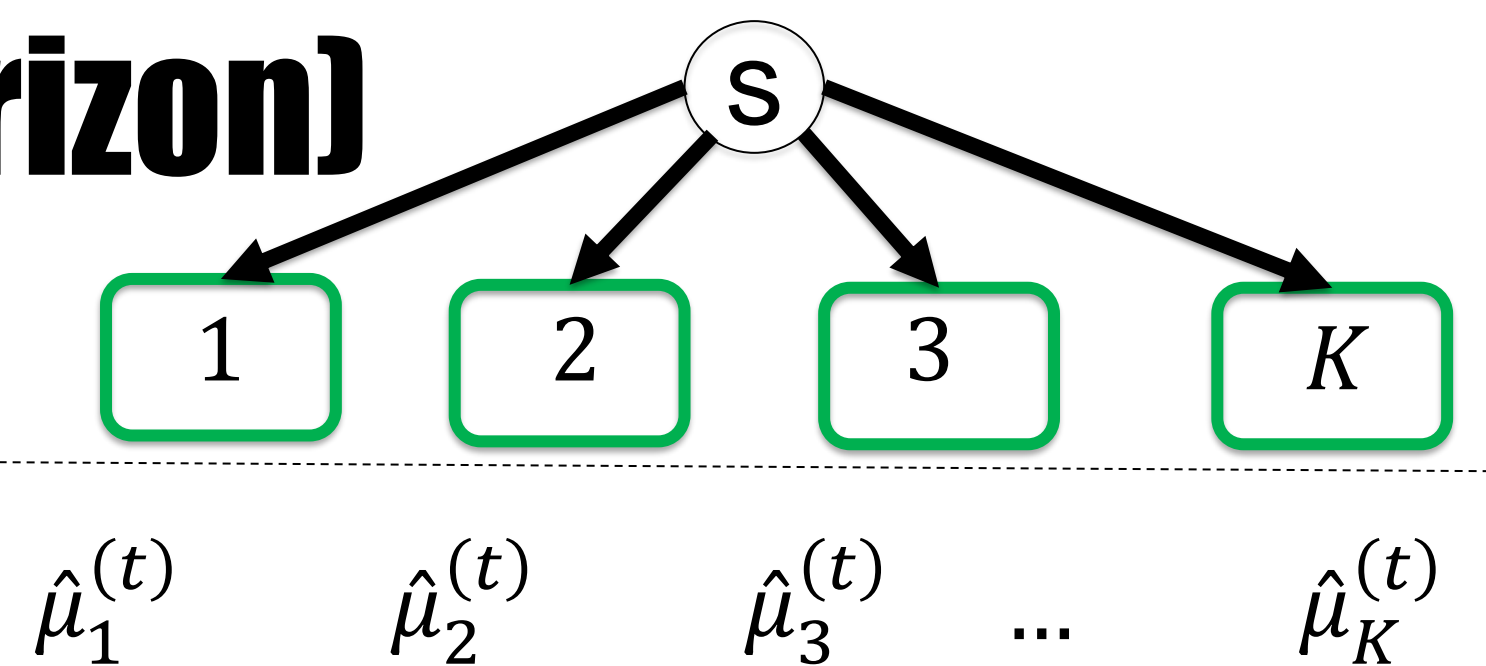


- Comments for the previous slide:
- Example of MAB with 4 actions. Each action yields a reward with 50 percent probability.
- Two actions have low rewards (about 1); the two other have high rewards about 20.
- Imagine that at the beginning you played each action 40 times and evaluate the mean return.
- If you repeated the game many times, each time starting with playing each action 40 times, you would get a distribution (hand-drawn here).
- As long as the distributions overlap, we continue to play all actions. Hence, after  $t=160$ , we should continue to play actions 3 and 4, while actions 1 and 2 can be safely dropped as a possibility.

# Exploration Bonus for MAB (1-step horizon)

- MAB with  $K$  possible actions:
- Reminder: greedy algorithm

$$\hat{\mu}_i^{(t)} = \frac{\sum_{\tau \in T_i^{(t)}} r_\tau}{|T_i^{(t)}|}$$



$$a_t = \arg \max_i \hat{\mu}_i^{(t)}$$

- Upper Confidence Bound (UCB1 in Auer et al. 2002):

$$U_i^{(t)} = \hat{\mu}_i^{(t)} + \sqrt{\frac{2 \log t}{|T_i^{(t)}|}}$$

The naïve estimate of  
average reward

Bonus for exploration  
(compare: Monte Carlo Tree Search)

$$U_1^{(t)} \quad U_2^{(t)} \quad U_3^{(t)} \quad \dots \quad U_K^{(t)}$$

$$a_t = \arg \max_i U_i^{(t)}$$

Theorem 1 of Auer et al. 2002:  
 $R_{\text{UCB1}}(T) \propto \log T + \text{const.}$

**Play greedy, but with a modified ‘value’  $U_k$**   
**→ Add exploration bonus to empirical average of reward**

- Comments for the previous slide:
- A smart optimal algorithm is Upper Confidence Bound (UCB; proposed by Auer et al. 2002 in Machine Learning) that computes a confidence bound index  $U_i^{(t)}$  for each action and chooses the one with highest index.
- The index is equal to the naïve estimate average reward  $\hat{\mu}_i^{(t)}$  plus an exploration bonus that is (i) a decreasing function of how many times an arm has been chosen  $|T_i^{(t)}|$  but (ii) an increasing function of how many actions have been taken in total (i.e.  $t$ ).
- The regret for the UCB algorithm scales logarithmically with  $T$ , hence it is an “optimal” algorithm. The constants of the regret can be fine-tuned by some variations of the algorithm (see Auer et al. 2002).



# Quiz: exploration Bonus (1-step horizon)

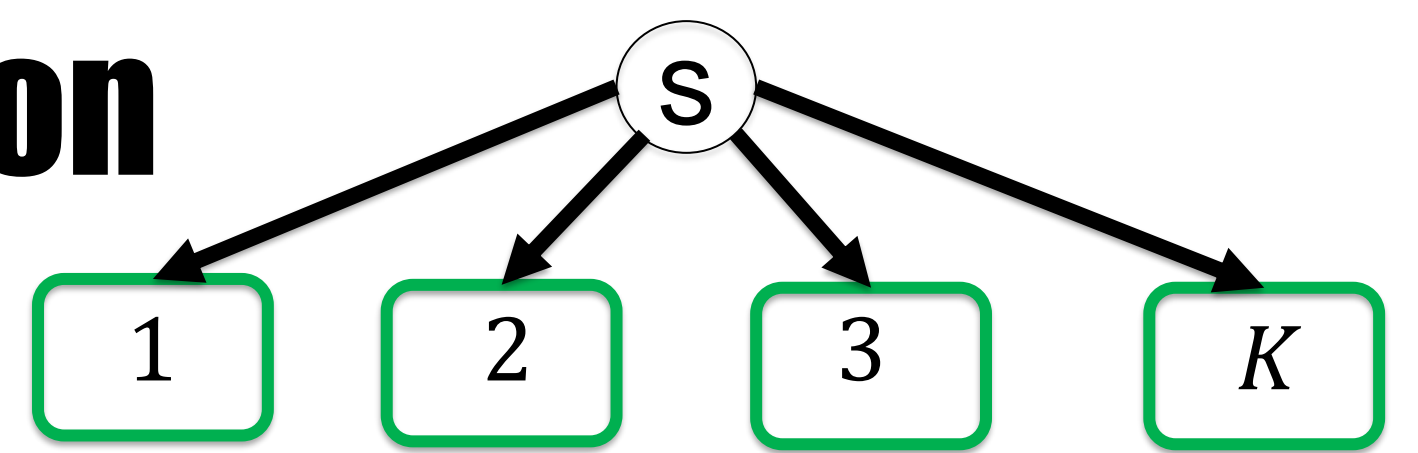
- A consistent learning algorithm eventually achieves a *zero average* regret in Multi-Armed Bandits (MAB).
- An optimal algorithm in MABs achieves a *constant total* regret.
- A good exploration bonus is  $\frac{\beta}{T_i^{(t)}}$ .
- A good exploration bonus is  $\frac{\log(t)}{\sqrt{T_i^{(t)}}}$ .

# Teaching monitoring – monitoring of understanding

[ ] up to here, at least 60% of material was new to me.

[ ] I have the feeling that I have been able to follow  
(at least) 80% of the lecture up to here.

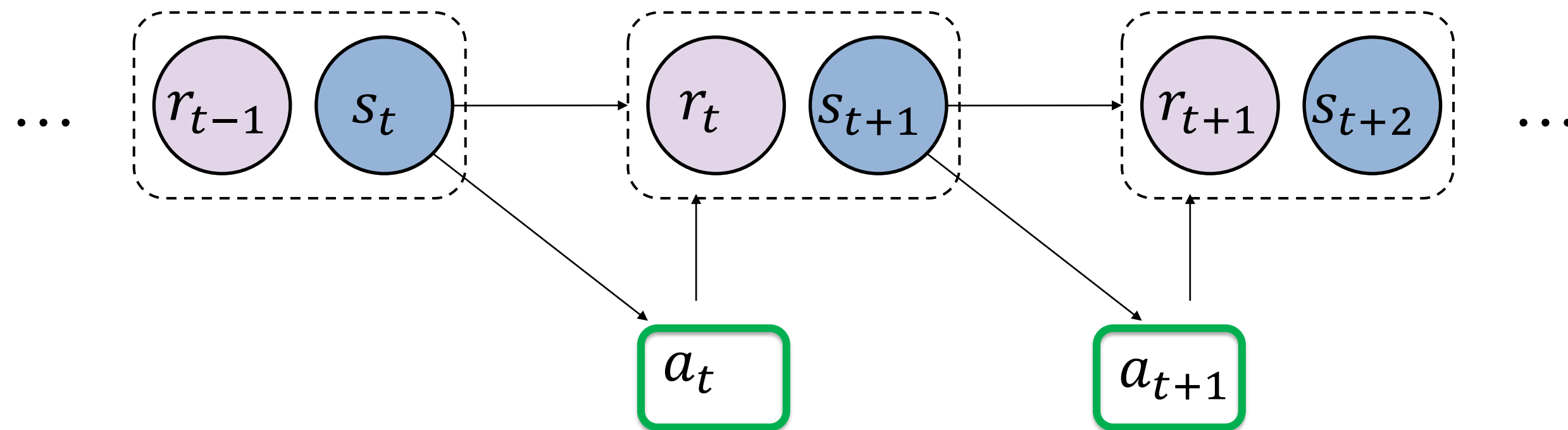
# Exploration Bonus for multi-step horizon



- MAB with  $K$  possible actions:

- 
- Markov Decision Processes (MDP):

- $P$ : transition probabilities, e.g.  $P(s'|s, a)$
- $R$ : expected reward, e.g.  $R(s, a)$



- Comments for the previous slide:
- We now want to extend from 1-step horizon (MAB) to multi-step horizon. The Multistep horizon leads to the Markov Decision Problem (MDP).

# Exploration Bonus for multi-step horizon

## Bellman equation (optimal action choice)

- Dynamic programming with true  $P(s'|s, a)$  and  $R(s, a)$  :

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q^*(s', a')$$

$$a_t = \arg \max_a Q^*(s_t, a)$$

- 
- Naïve model-based (MB) RL:

$$\hat{Q}_{\text{MB}}^{(t)}(s, a) = \hat{R}^{(t)}(s, a) + \gamma \sum_{s'} \hat{P}^{(t)}(s'|s, a) \max_{a'} \hat{Q}_{\text{MB}}^{(t)}(s', a')$$

$$\hat{R}^{(t)}(s, a) = \frac{\sum_{\tau \in T_{s,a}^{(t)}} r_{\tau}}{|T_{s,a}^{(t)}|}$$

$$\hat{P}^{(t)}(s'|s, a) = \frac{|T_{s,a,s'}^{(t)}|}{|T_{s,a}^{(t)}|}$$

$$T_{s,a}^{(t)} = \{\tau \leq t : a_{\tau} = a, s_{\tau} = s\}$$

$$T_{s,a,s'}^{(t)} = \{\tau \leq t : a_{\tau} = a, s_{\tau} = s, s_{\tau+1} = s'\}$$

$$a_t = \arg \max_a \hat{Q}_{\text{MB}}^{(t)}(s_t, a)$$

The exploration-exploitation trade-off is even more serious in MDPs than MABs.

Any trick similar to UCB?

Comments for the previous slide:

- Similar to the bandit setting, if we have access to the true transition probabilities and reward functions, then the optimal policy would be to use Dynamic Programming, solve the optimal Bellman equations, and use a greedy policy on the resulting Q-values:  $a_t = \arg \max_a Q^*(s_t, a)$
- In the absence of the complete knowledge of the environment, a naïve model-based approach is to approximate the transition probabilities and the reward values, solve the optimal Bellman equations by using these estimates, and use a greedy policy on the resulting Q-values:  $a_t = \arg \max_a \hat{Q}_{\text{MB}}^{(t)}(s_t, a)$
- The naïve model-based approach is prone to be stuck in some parts of the environment and never find the optimal policy. You have seen epsilon-greedy and the softmax policy as to approaches to deal with this issue by adding randomness to the action-selection. Here, we ask whether we can find a directed exploration approach like UCB for MDPs. What is a good exploration bonus?

# Exploration Bonus for multi-step horizon

- Dynamic programming with true  $P(s'|s, a)$  and  $R(s, a)$  :

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q^*(s', a')$$

---

- Naïve model-based (MB) RL:

$$\hat{Q}_{\text{MB}}^{(t)}(s, a) = \hat{R}^{(t)}(s, a) + \gamma \sum_{s'} \hat{P}^{(t)}(s'|s, a) \max_{a'} \hat{Q}_{\text{MB}}^{(t)}(s', a')$$

---

- Model-based interval estimation with exploration bonus (MBIE+EB in Strehl and Littman 2008):

$$\hat{Q}_{\text{MB}}^{(t)}(s, a) = \hat{R}^{(t)}(s, a) + \frac{\beta}{\sqrt{T_{s,a}^{(t)}}} + \gamma \sum_{s'} \hat{P}^{(t)}(s'|s, a) \max_{a'} \hat{Q}_{\text{MB}}^{(t)}(s', a')$$

The naïve estimate  
of average reward

Bonus for exploration (different from UCB regarding  $\log t$ )

- Comments for the previous slide:
- Model-based interval estimation with exploration bonus (MBIE+EB; proposed by Strehl and Littman 2008 in the Journal of Computer and System Sciences) uses the exact same procedure as the naïve model-based approach except that it adds an exploration bonus to the reward function.
- The exploration bonus is a decreasing function of how many times a specific action is taken in a specific state, so it encourages to take actions that have been taken less frequently in the past.



# Exploration Bonus for multi-step horizon

- Model-based interval estimation with **exploration bonus** (MBIE+EB in Strehl and Littman 2008):

$$\hat{Q}_{\text{MB}}^{(t)}(s, a) = \hat{R}^{(t)}(s, a) + \frac{\beta}{\sqrt{T_{s,a}^{(t)}}} + \gamma \sum_{s'} \hat{P}^{(t)}(s'|s, a) \max_{a'} \hat{Q}_{\text{MB}}^{(t)}(s', a')$$

- Theorem 2 in Strehl and Littman 2008:**

MBIE+EB is Probably Approximately Correct for MDPs (= it is PAC-MDP).

= loosely speaking, its choices are good enough with high probability.

- Alternative: Bayesian Exploration Bonus (BEB) by Kolter and Ng 2009**

$$\text{Bonus} = \frac{\beta}{1 + T_{s,a}^{(t)}}$$

It is **not PAC-MDP**  
but is **near-Bayesian**.

Theorem 2. Exploration based on a bonus proportional to  $\left(T_{s,a}^{(t)}\right)^{-p}$  is not PAC-MDP if  $p > 0.5$ .

- Comments for the previous slide:
- MBIE+EB is proven to be PAC-MDP (see Strehl and Littman 2008): In short and loosely speaking, this means that, with high probability, most of the actions taken by MBIE+EB are close to the actions that would have been taken by the optimal policy.
- Alternative exploration bonuses are possible, but they have different properties. For example, an exploration bonus proportional to one over  $T_{s,a}^{(t)}$  is not PAC-MDP but is “near Bayesian” (i.e., another notion of optimality; see Kolter and Ng in ICML 2009).

# Quiz: exploration Bonus (multi-step horizon)

[ ] Assuming we know the true  $P(s'|s, a)$  and  $R(s, a)$ , the Bellman equation is

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q^*(s', a'); \quad a_t = \arg \max_a Q^*(s_t, a)$$

[ ] If we do not know the true  $P(s'|s, a)$  and  $R(s, a)$ , the Bellman equation can be replaced by

$$\hat{Q}_{\text{MB}}^{(t)}(s, a) = \hat{R}^{(t)}(s, a) + \text{Bonus}(s, a) + \gamma \sum_{s'} \hat{P}^{(t)}(s'|s, a) \max_{a'} \hat{Q}_{\text{MB}}^{(t)}(s', a'); \quad a_t = \arg \max_a \hat{Q}_{\text{MB}}^{(t)}(s_t, a)$$

[ ] One of the choices is  $\text{Bonus}(s, a) = \frac{\beta}{\sqrt{T_{s,a}^{(t)}}}$

[ ] A function  $\frac{\beta}{\sqrt{T_{s,a}^{(t)}}}$  decreases more slowly than  $\frac{\beta}{1+T_{s,a}^{(t)}}$

# Summary: Exploration Bonus for multi-step horizon

- Adding exploration bonus provably improves the performance of RL algorithms.
- Hence, to optimally seek a reward, best seek a ‘modified reward’ .

- There is, however, not a single (unique) approach to
  - define an exploration bonus
  - evaluate its performance.
- For MDP a possible exploration bonus:

$$\text{Bonus} = \frac{\beta}{1 + T_{s,a}^{(t)}}$$

- These CS approaches assume: (i) stationary problem (ii) model-based RL (update of Bellman equation in the background)

# Teaching monitoring – monitoring of understanding

[ ] up to here, at least 60% of material was new to me.

[ ] I have the feeling that I have been able to follow  
(at least) 80% of the lecture up to here.

- Comments for the previous slide: