

Wulfram Gerstner

EPFL, Lausanne, Switzerland

Artificial Neural Networks and RL

The role of exploration, novelty, and surprise in RL

Objectives for today:

- 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:

[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: What is novelty?

Q2: What is surprise?

Q3: What is the difference between the two?

Q4: Why are they useful?

Q5: Why should we talk about it in an RL class?

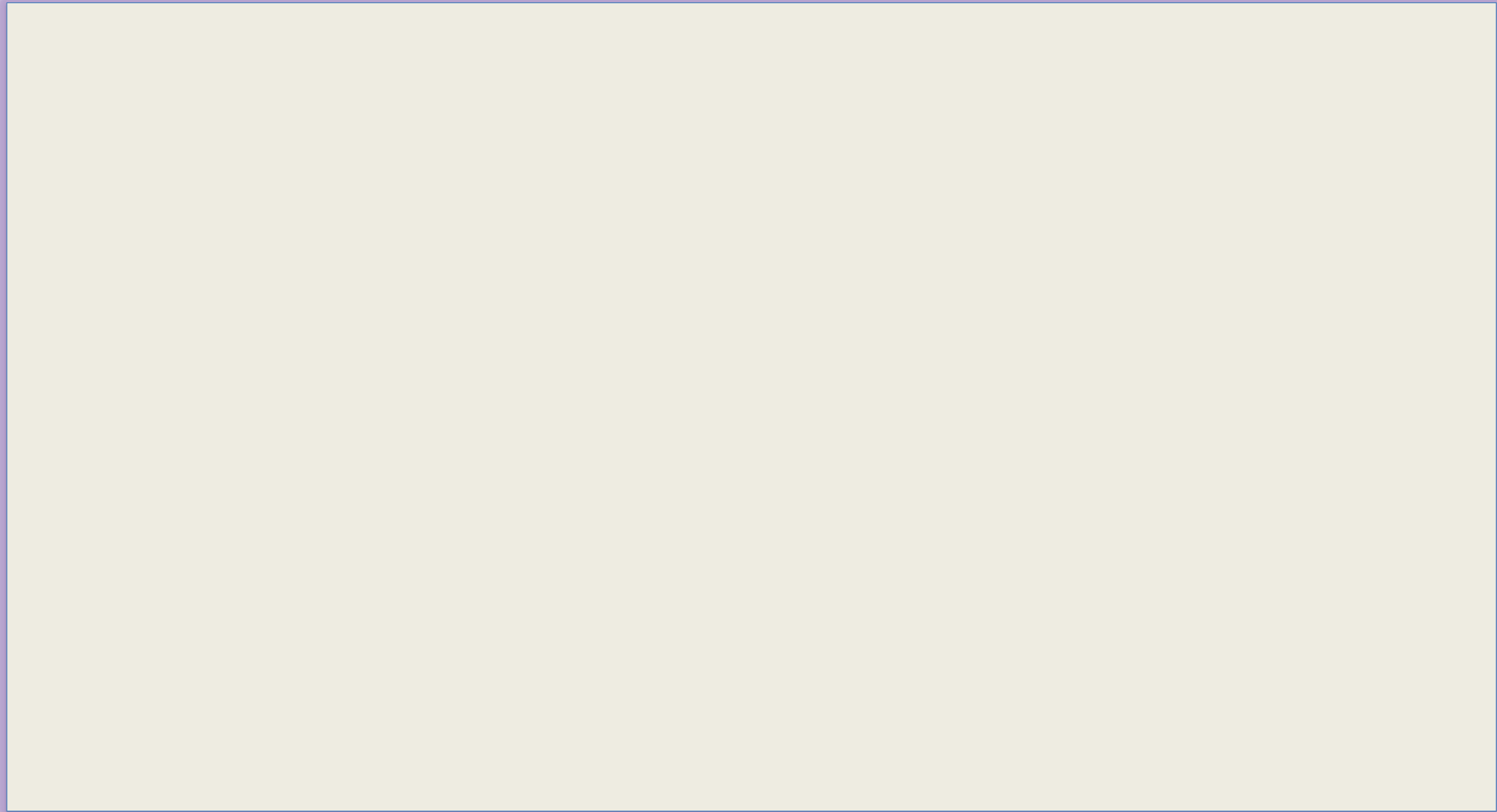
Previous slide.

Today we will ask 5 questions:

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

And why should we talk about it in a class on RL?

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

Q3: *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.

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)

Previous slide.

An example for the above statement.

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.

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?**

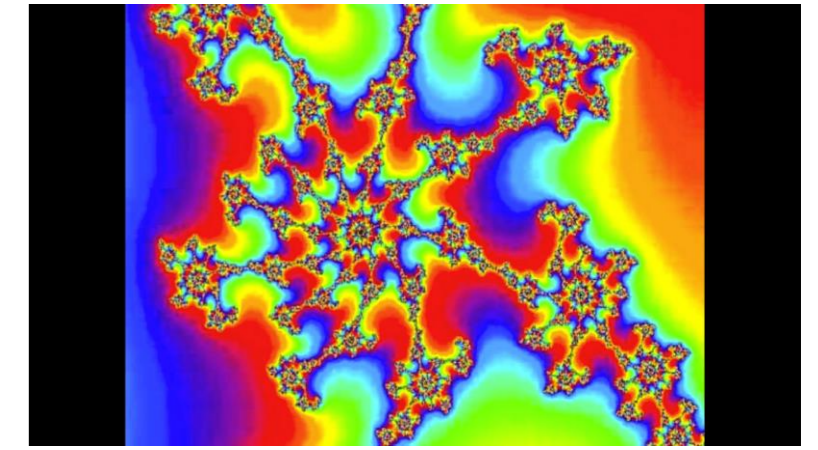
Previous slide. Summary.

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 (3rd 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

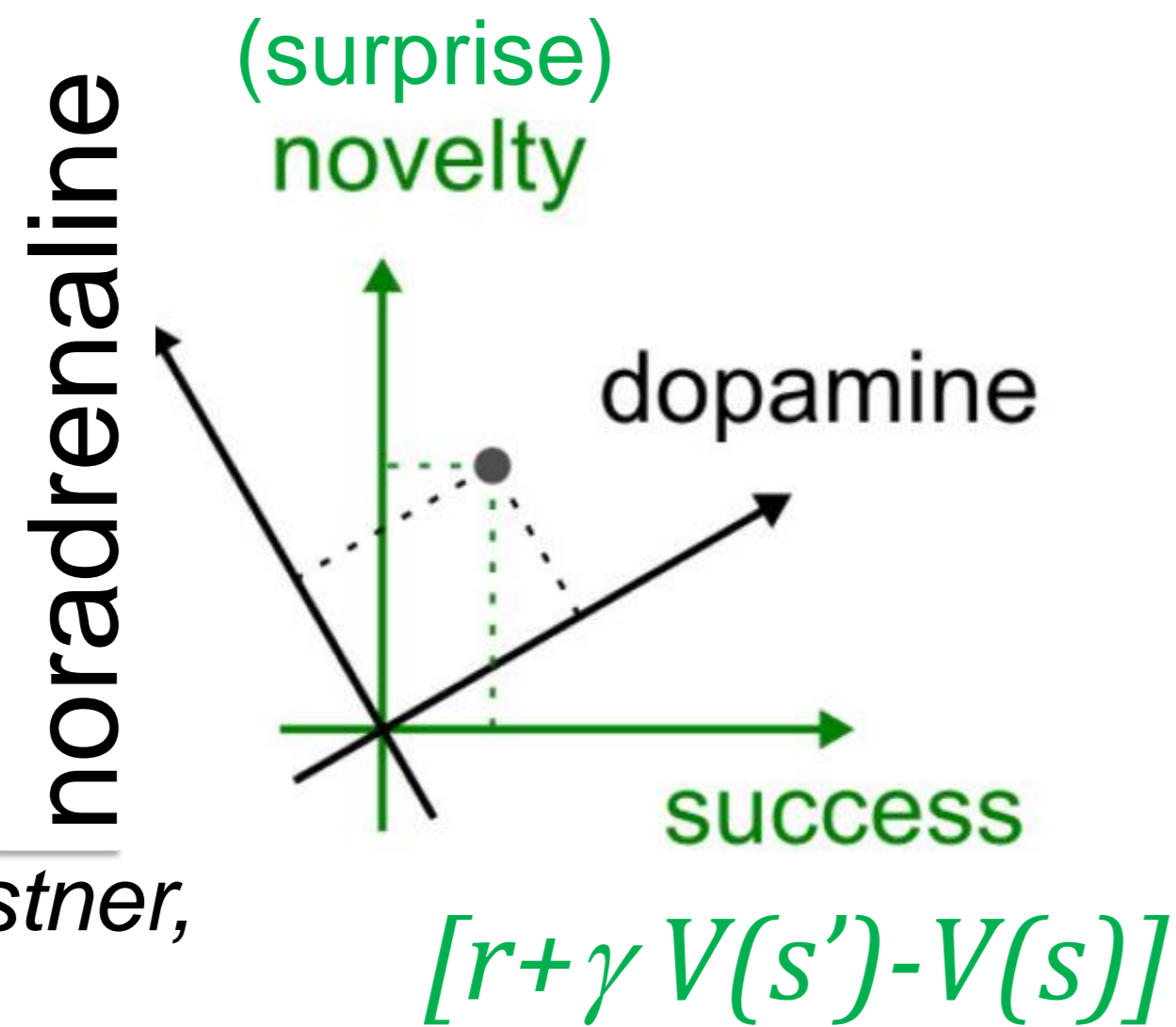
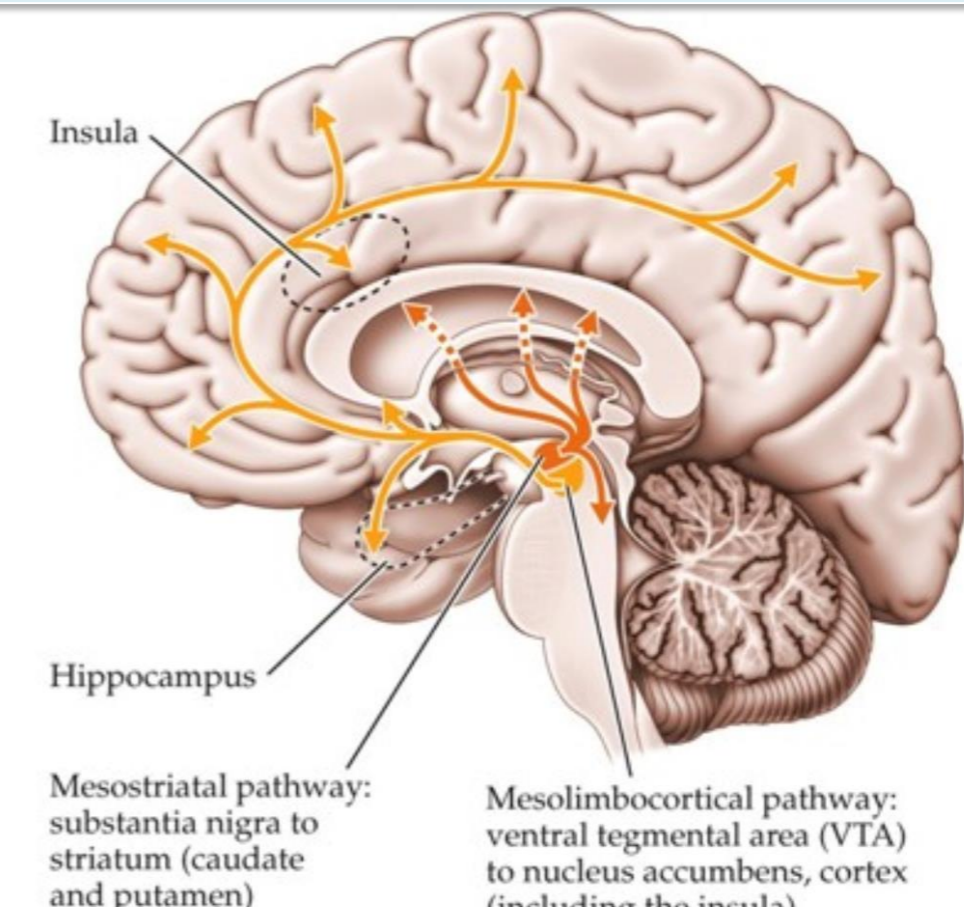


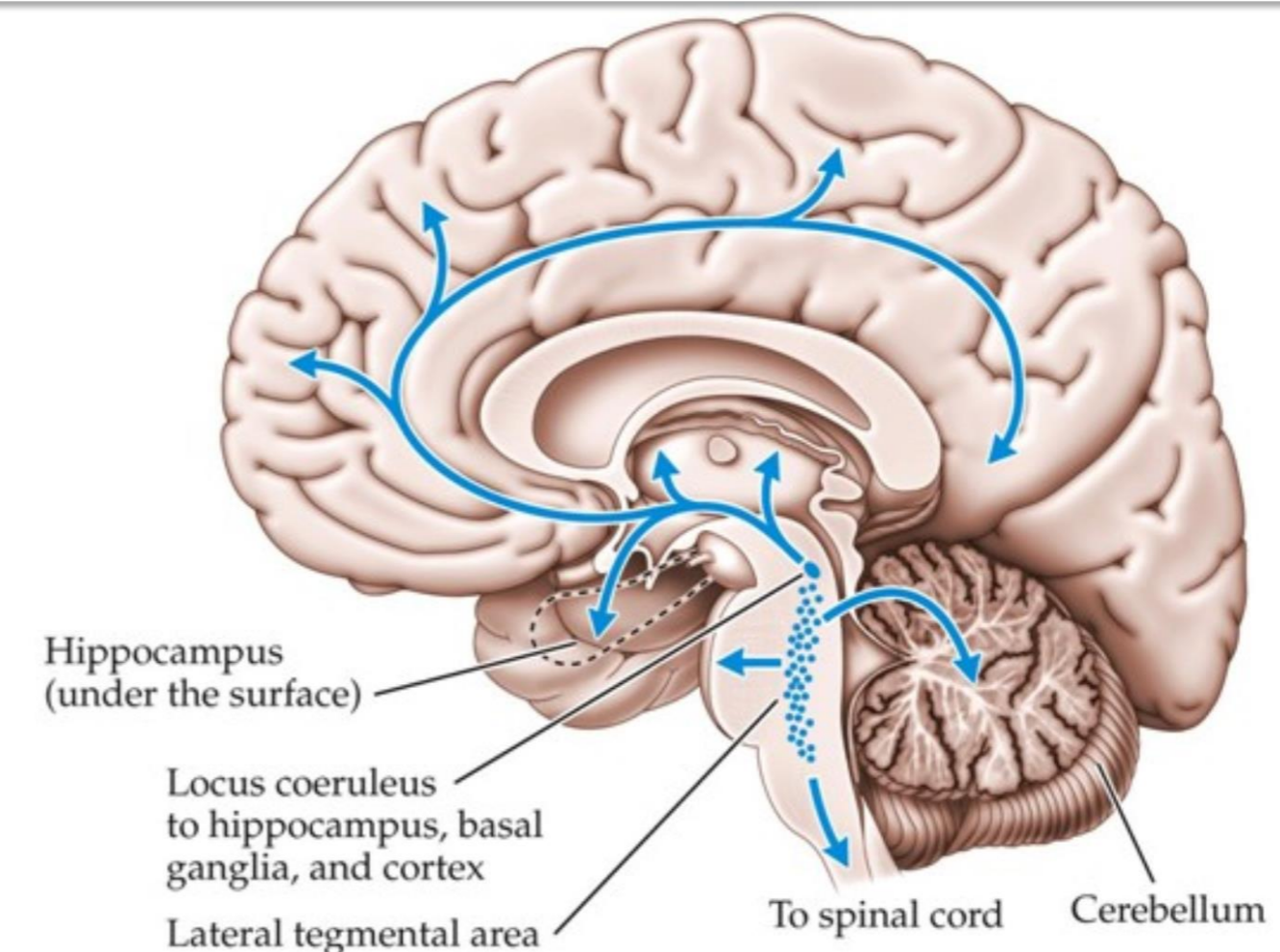
Image:
Fremaux and Gerstner,
Frontiers (2016)

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

φ_i = activity of postsynaptic neuron

Step 1: co-activation sets eligibility trace

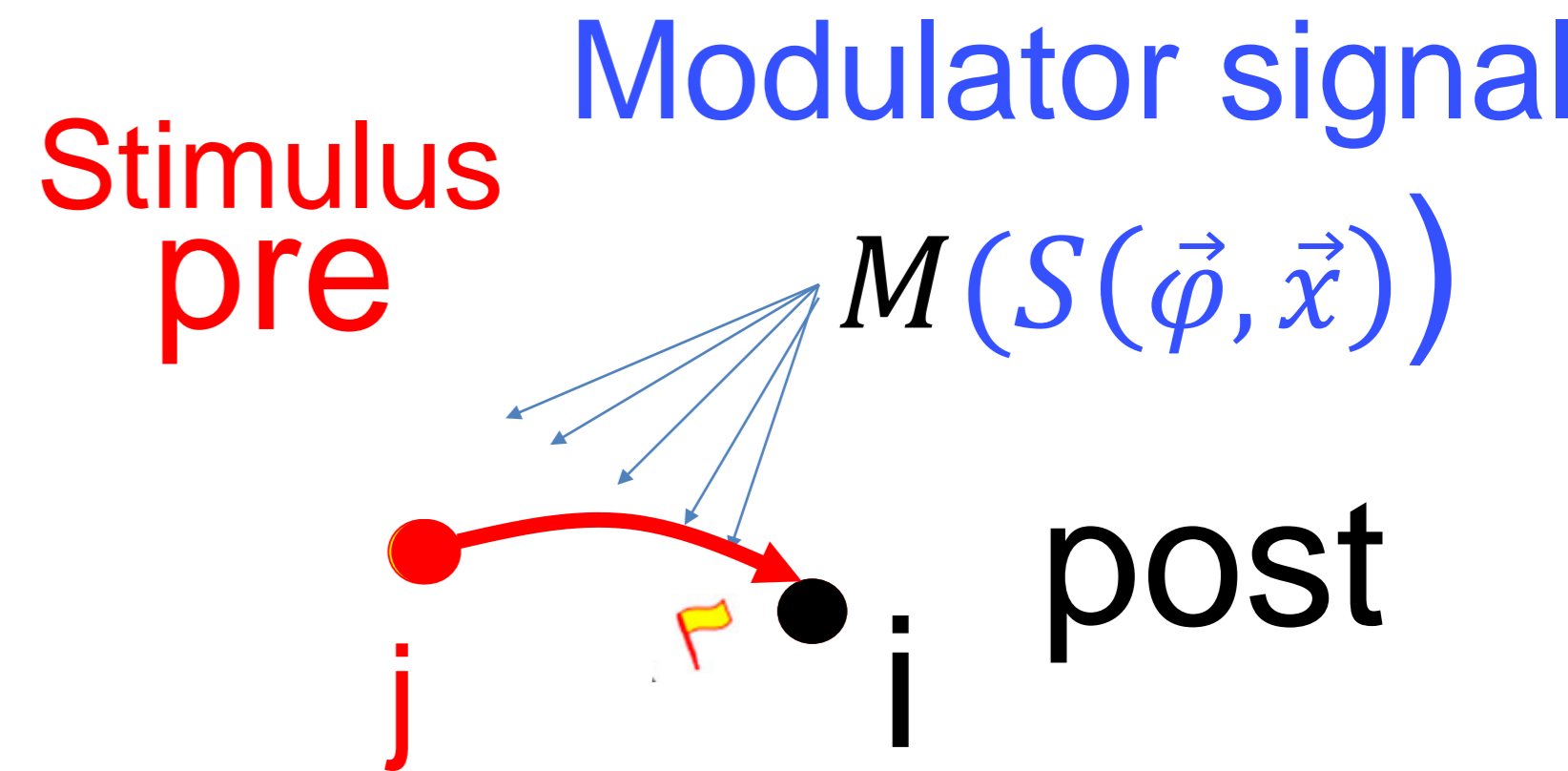
$$\Delta z_{ij} = \eta f(\varphi_i) g(x_j)$$

Step 2: eligibility trace decays over time

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

Step 3: eligibility trace translated into weight change

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



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

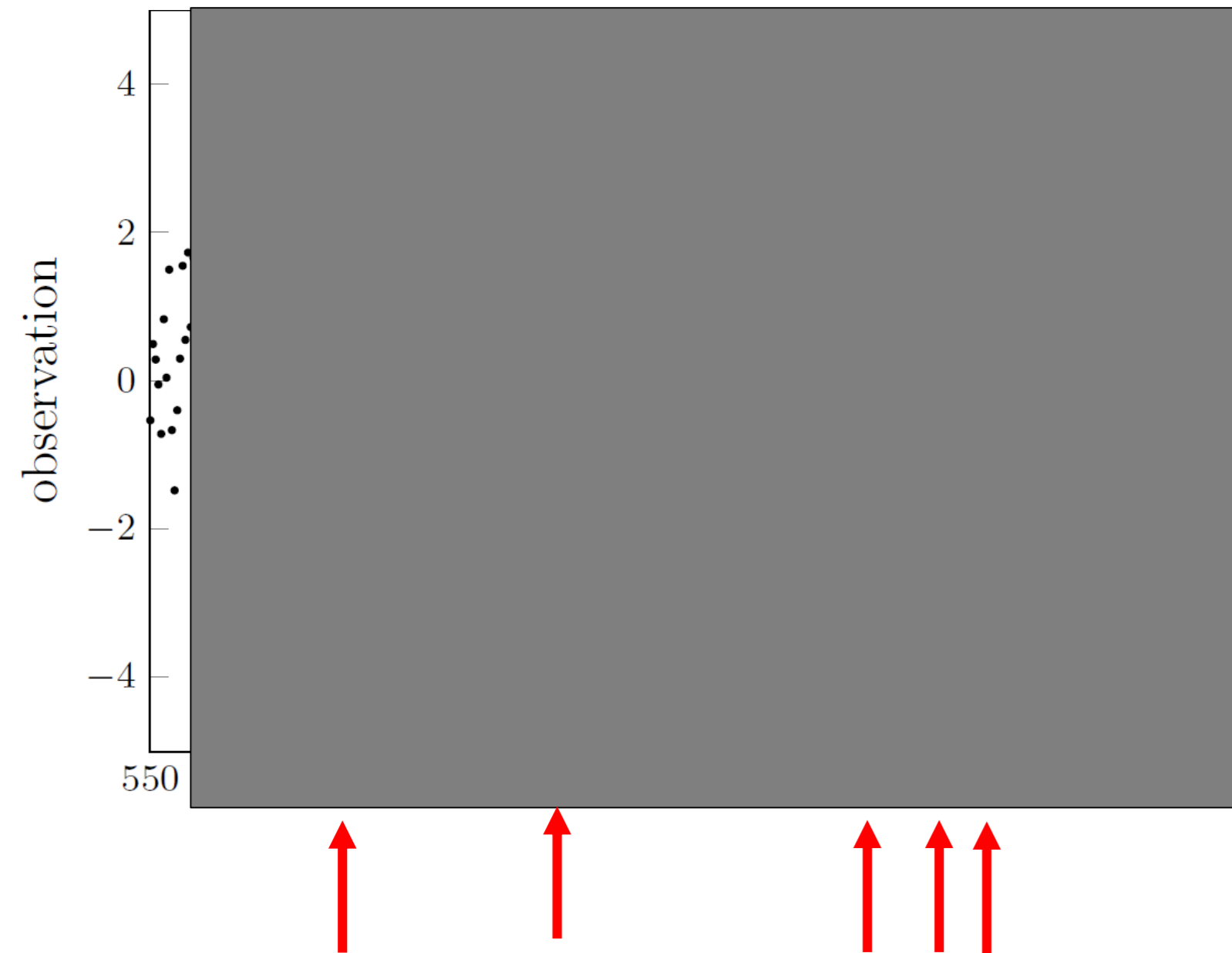
The role of exploration, novelty, and surprise in RL

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

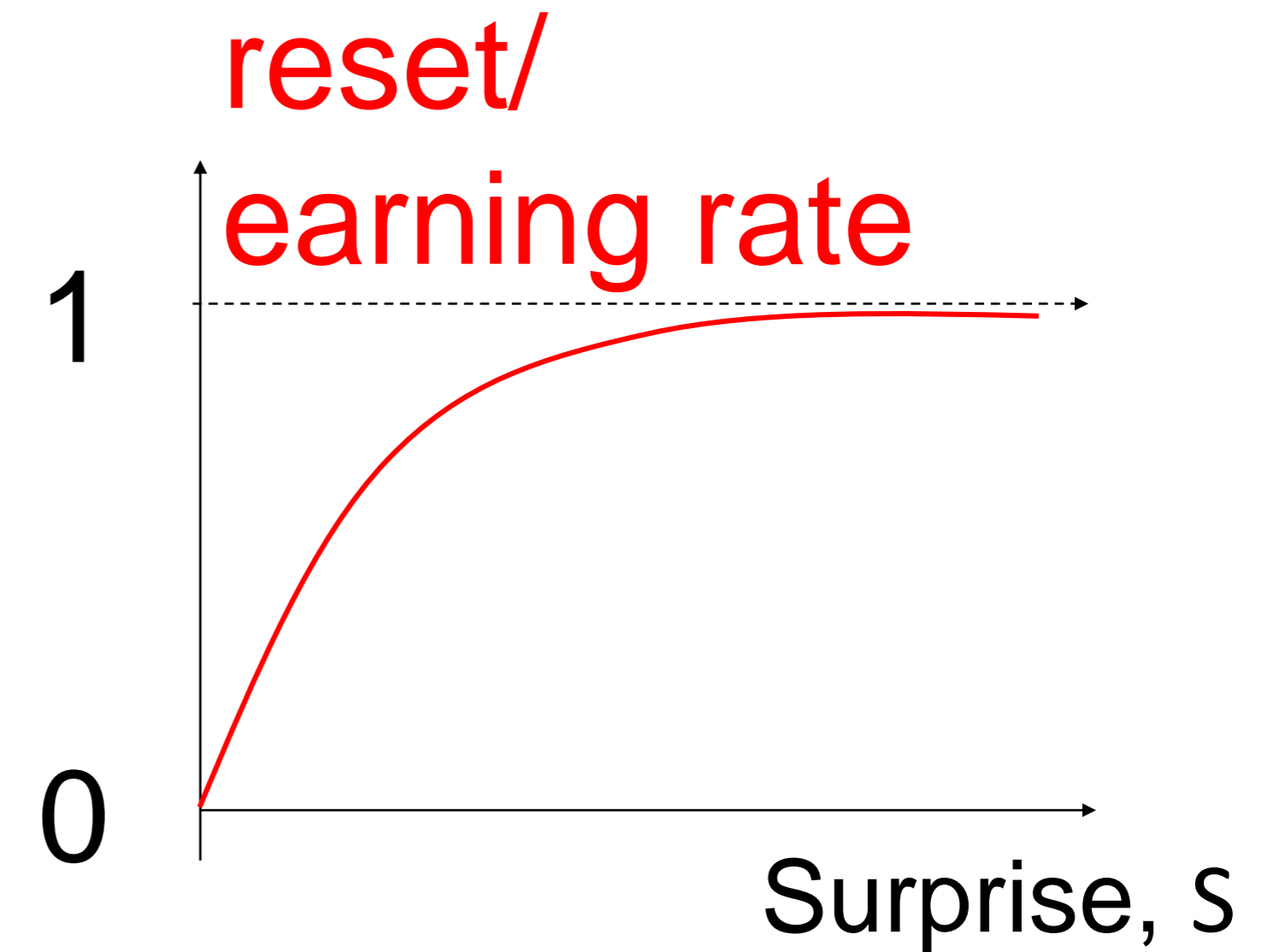
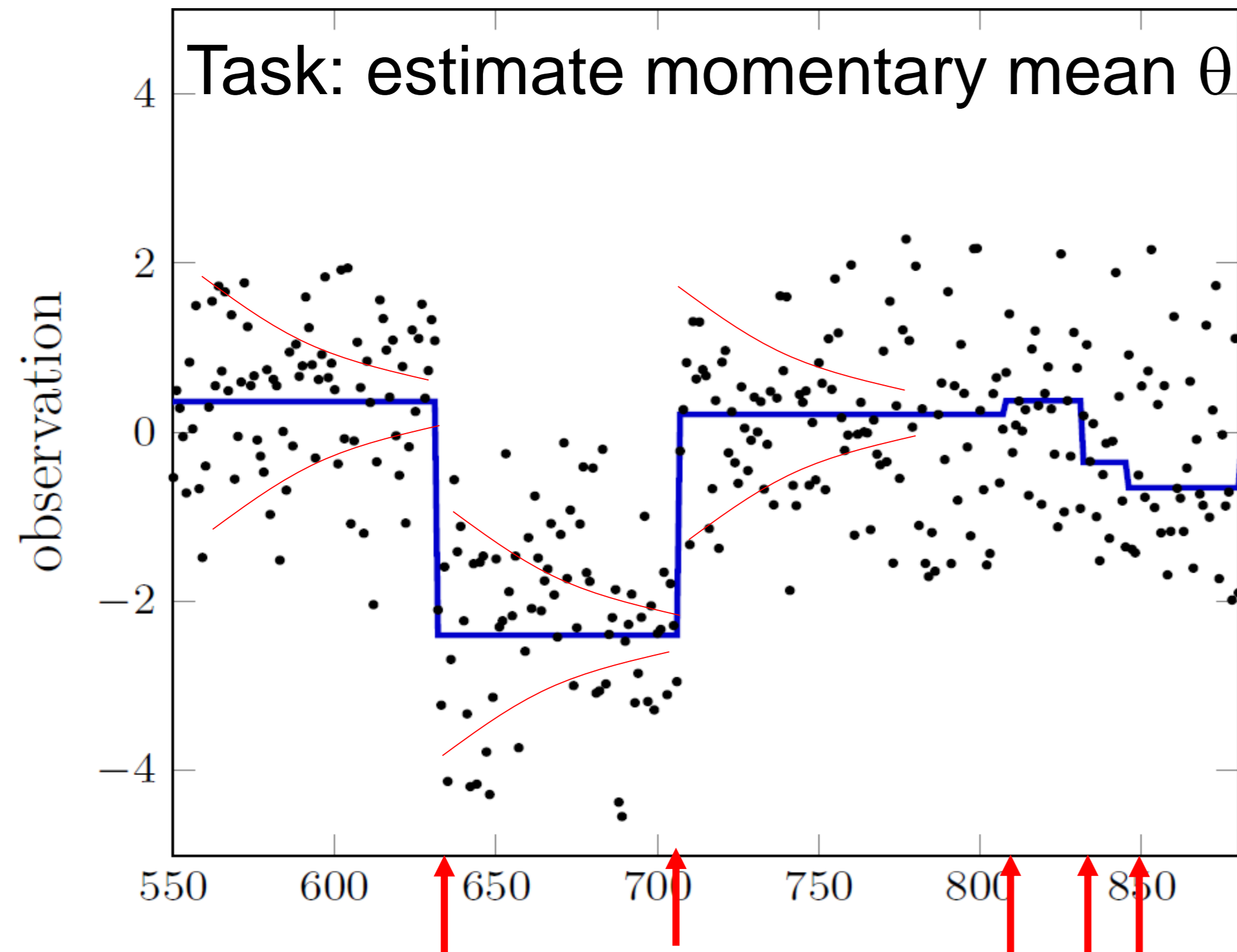
- here:
- mean of Gaussian is fixed for many steps
 - mean jumps at 'change points': probability $\ll 1$
 - variance is fixed
 - task is to estimate **momentary mean** of Gaussian

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

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

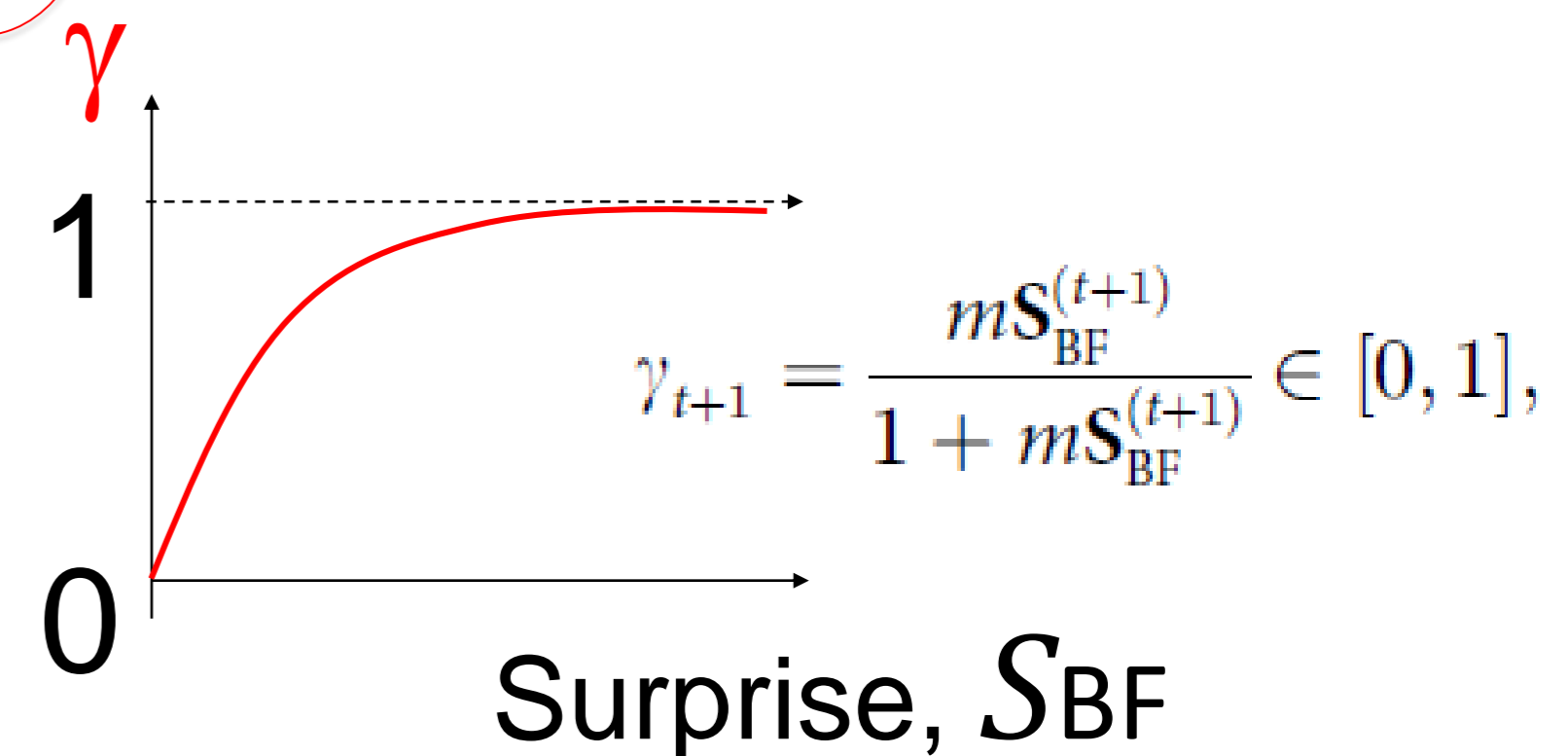
Probability of observation y
under prior belief $\pi^{(0)}$

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 γ



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

If it is likely (γ 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 γ depends monotonically on the Bayes-Factor Surprise S_{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.

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.

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

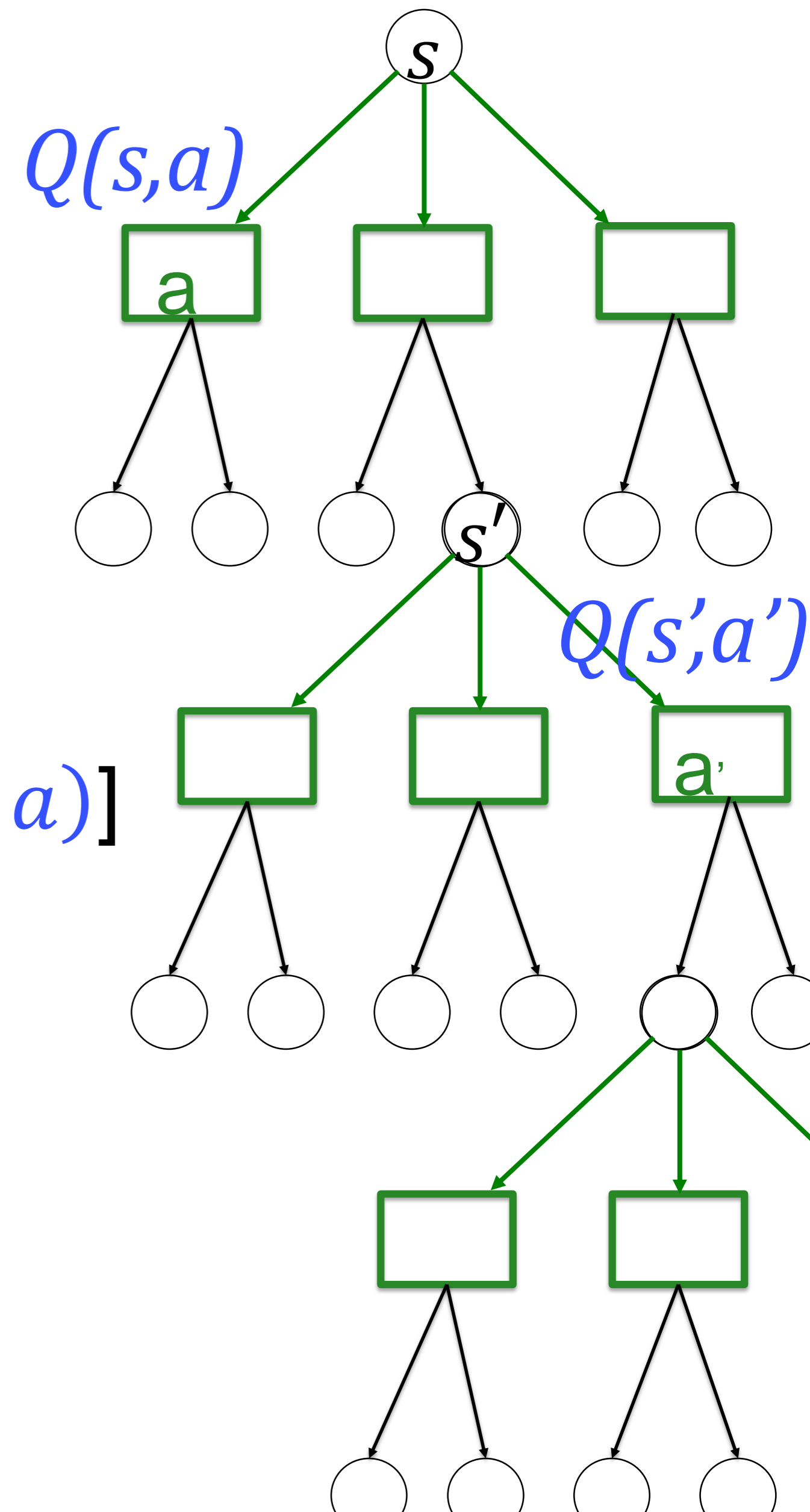
$$\Delta Q(s, a) = \eta [r_t + \gamma Q(s', a') - Q(s, a)]$$

Expected SARSA

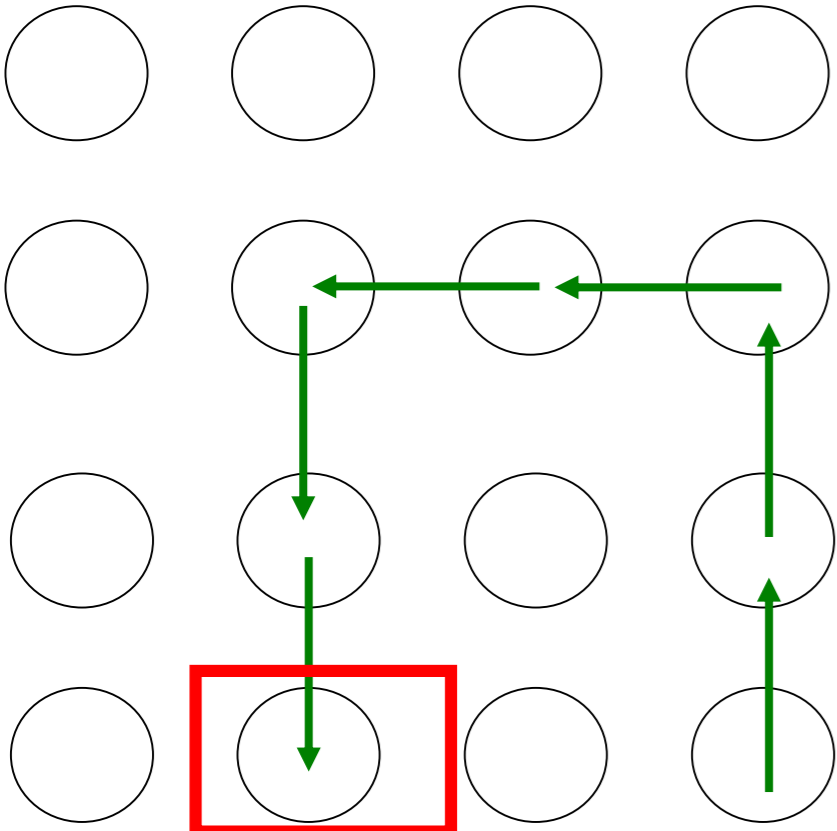
$$\Delta Q(s, a) = \eta [r_t + \gamma \{ \sum_{a'} \pi(s', a') Q(s', a') \} - Q(s, a)]$$

Q-learning

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



Review: Eligibility Traces, SARSA(λ)

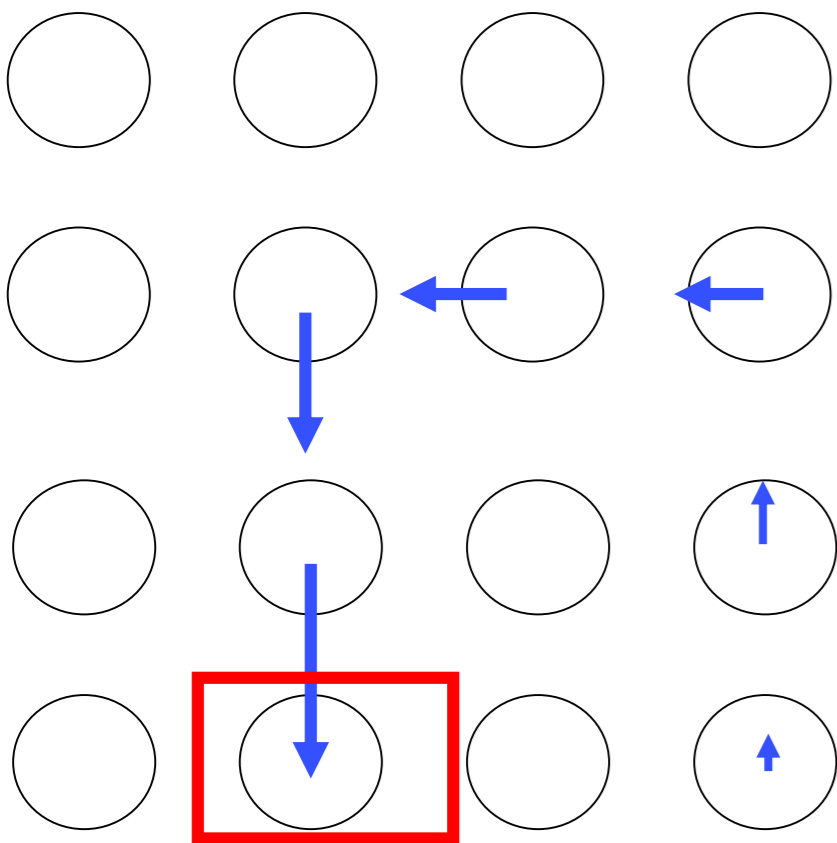


Idea:

- 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 δ_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

Model-free

prioritized
sweeping

eligibility
traces

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

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

novelty

n_t

Q-values

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

Bellman eq.
estimation/update

Model-based

Model-free

prioritized
sweeping

eligibility
traces

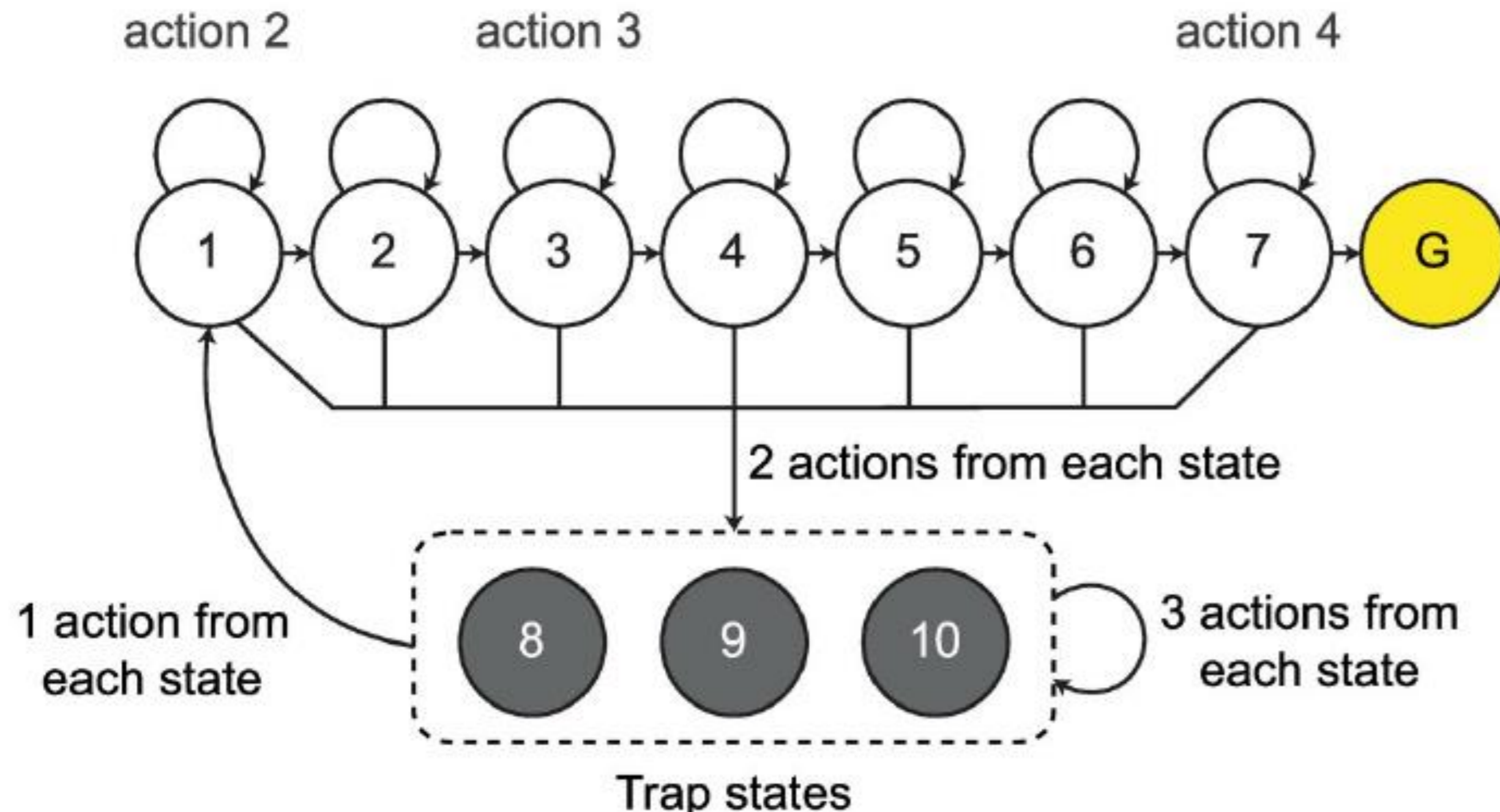
$Q_{MB,N}^{(t)}(s, a)$

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

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 1st episode, before any reward.

Improve exploration! Solutions?

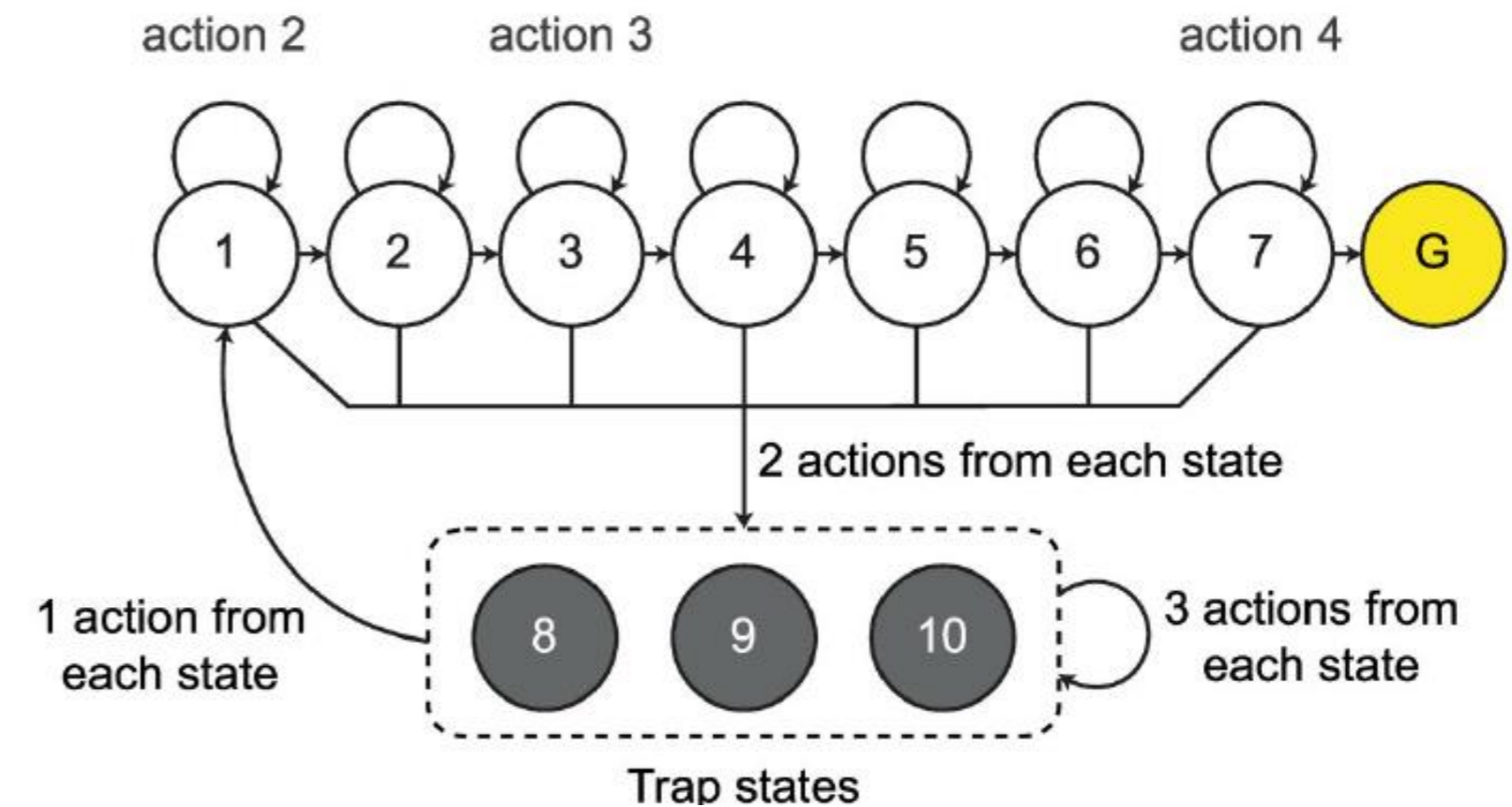
1. Optimistic initialization?

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 1st episode, before any reward.

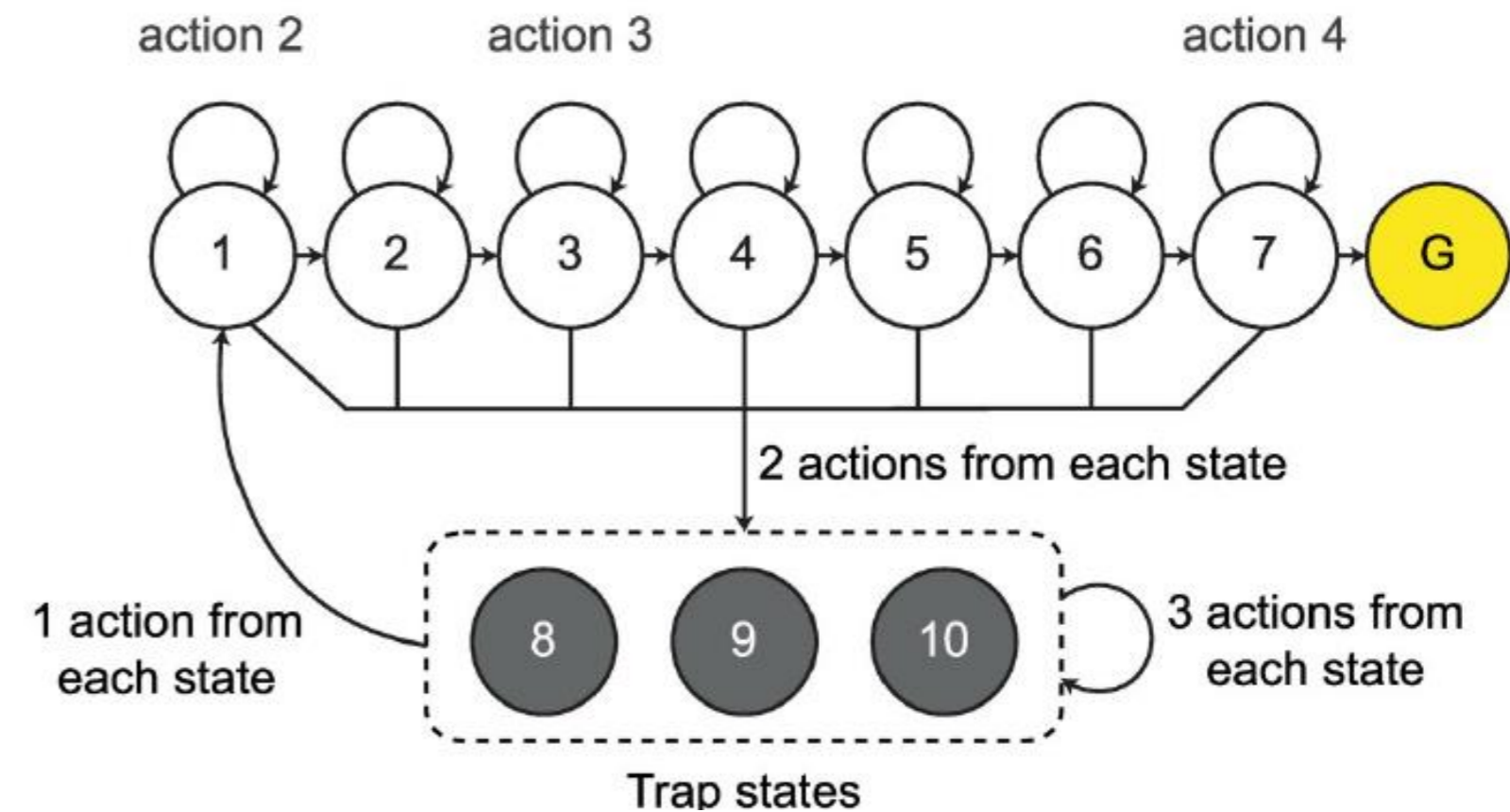
Improve exploration! Solutions?

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.

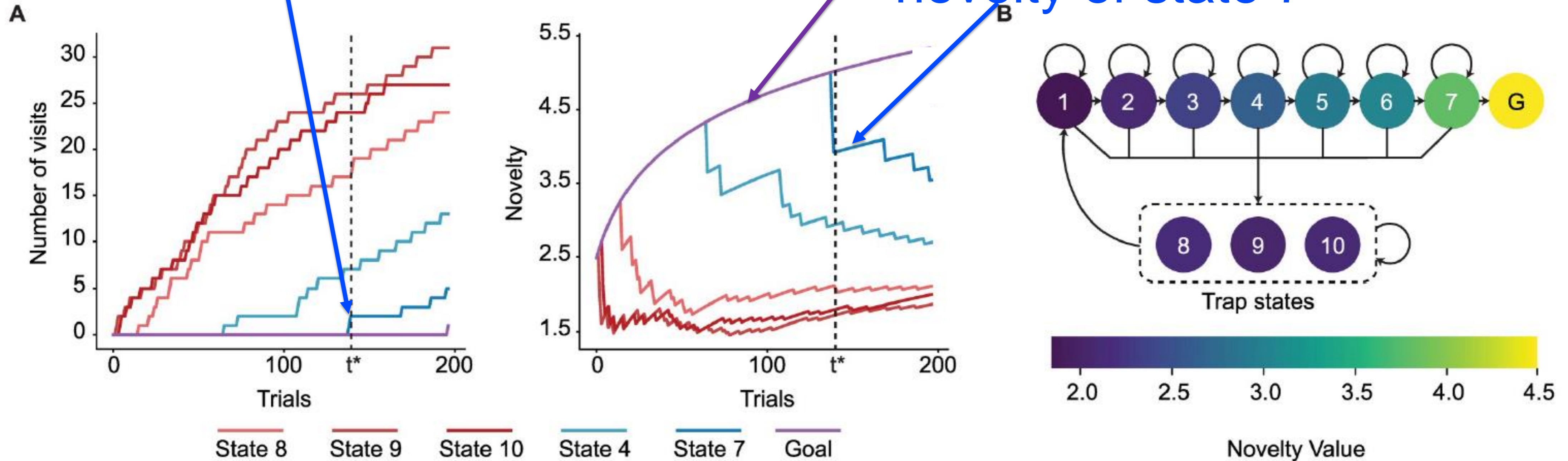
RPE: Reward Prediction Error

NPE: Novelty Prediction Error

Novelty encourages exploration of an environment

Focus on 1st 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.

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?
5. **Hybrid Model with Novelty, Surprise, and Reward**

Previous slide.

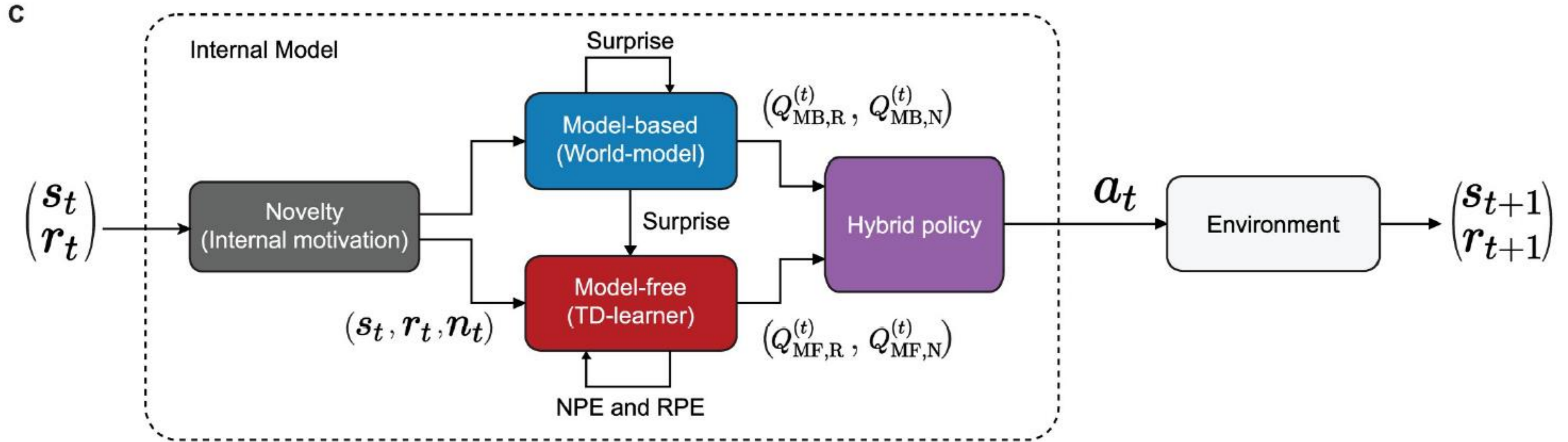
Now we study a specific model that combines many aspects.

Reminder:

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

NPE: Novelty Prediction Error = TD error of novelty consistency

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



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

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

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.

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?**
- 5. Hybrid Model with Novelty, Surprise, and Reward**
- 6. An Experiment**

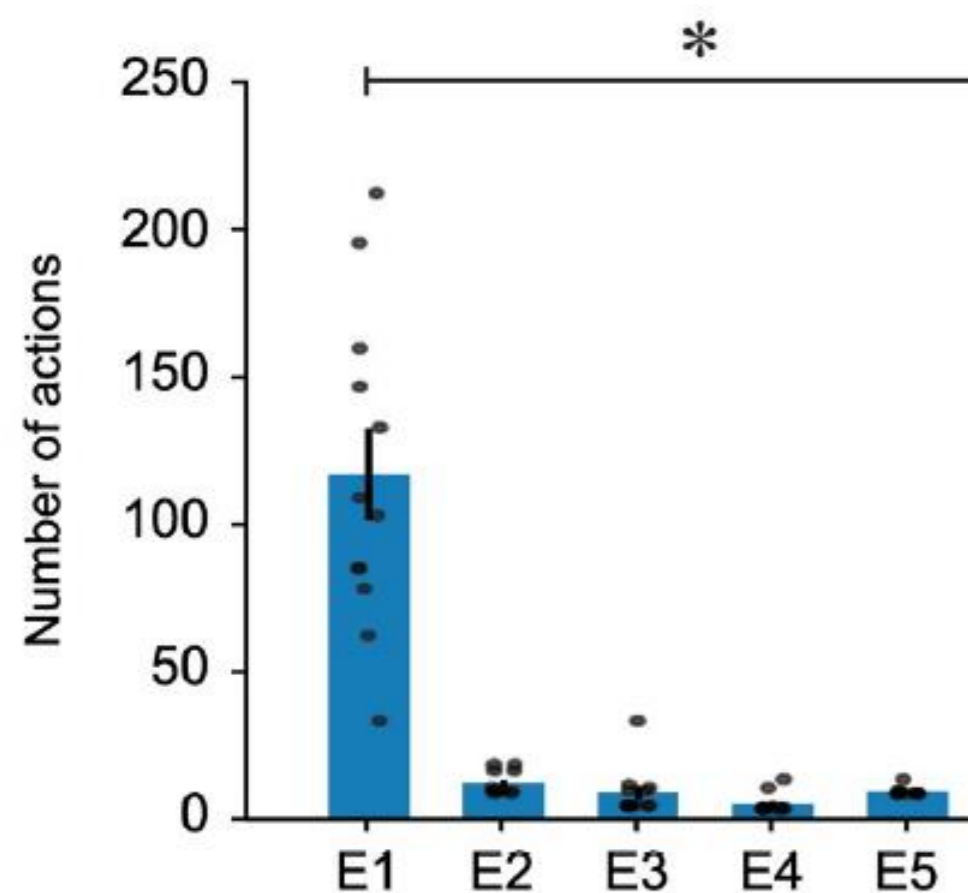
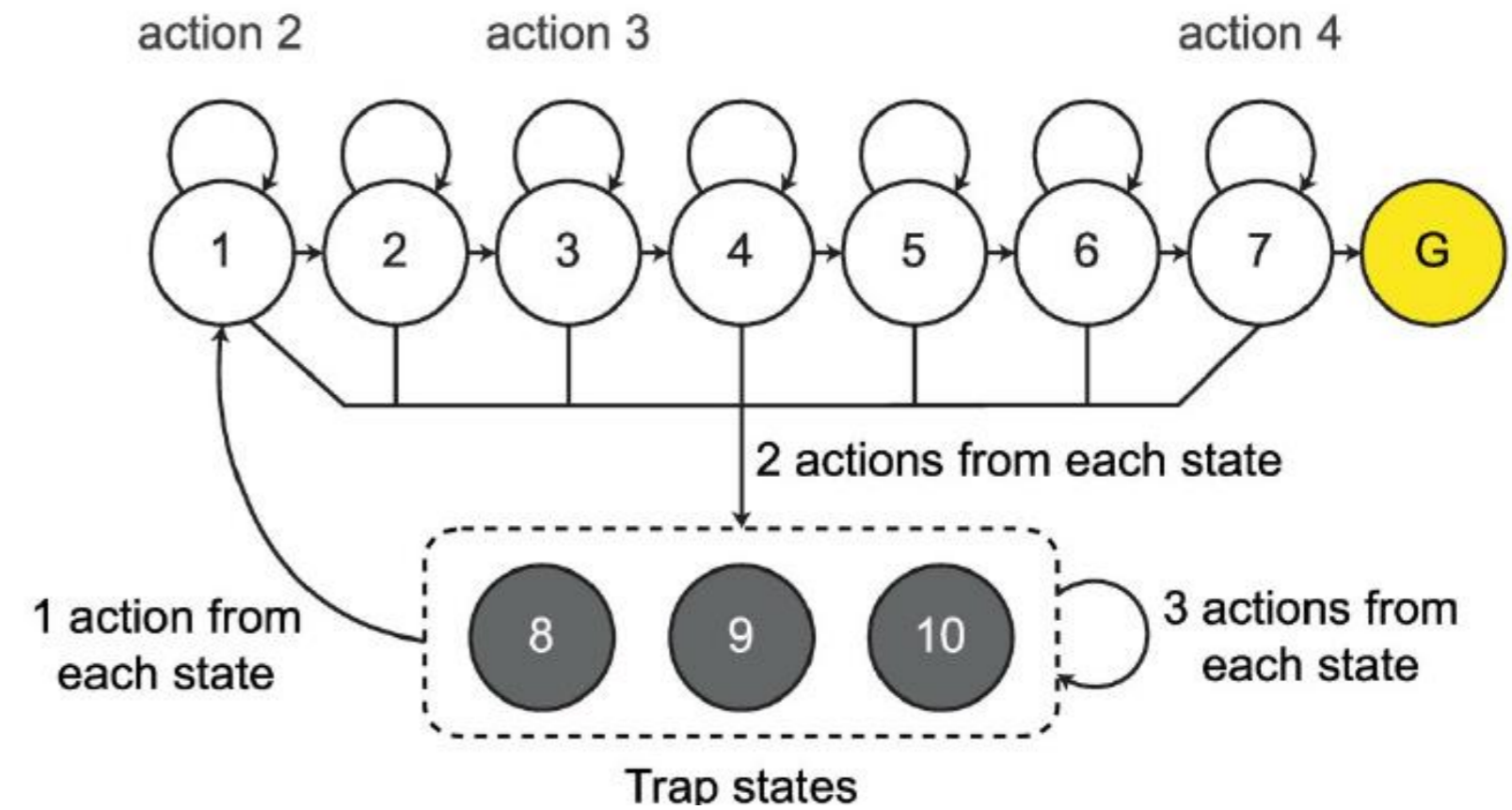
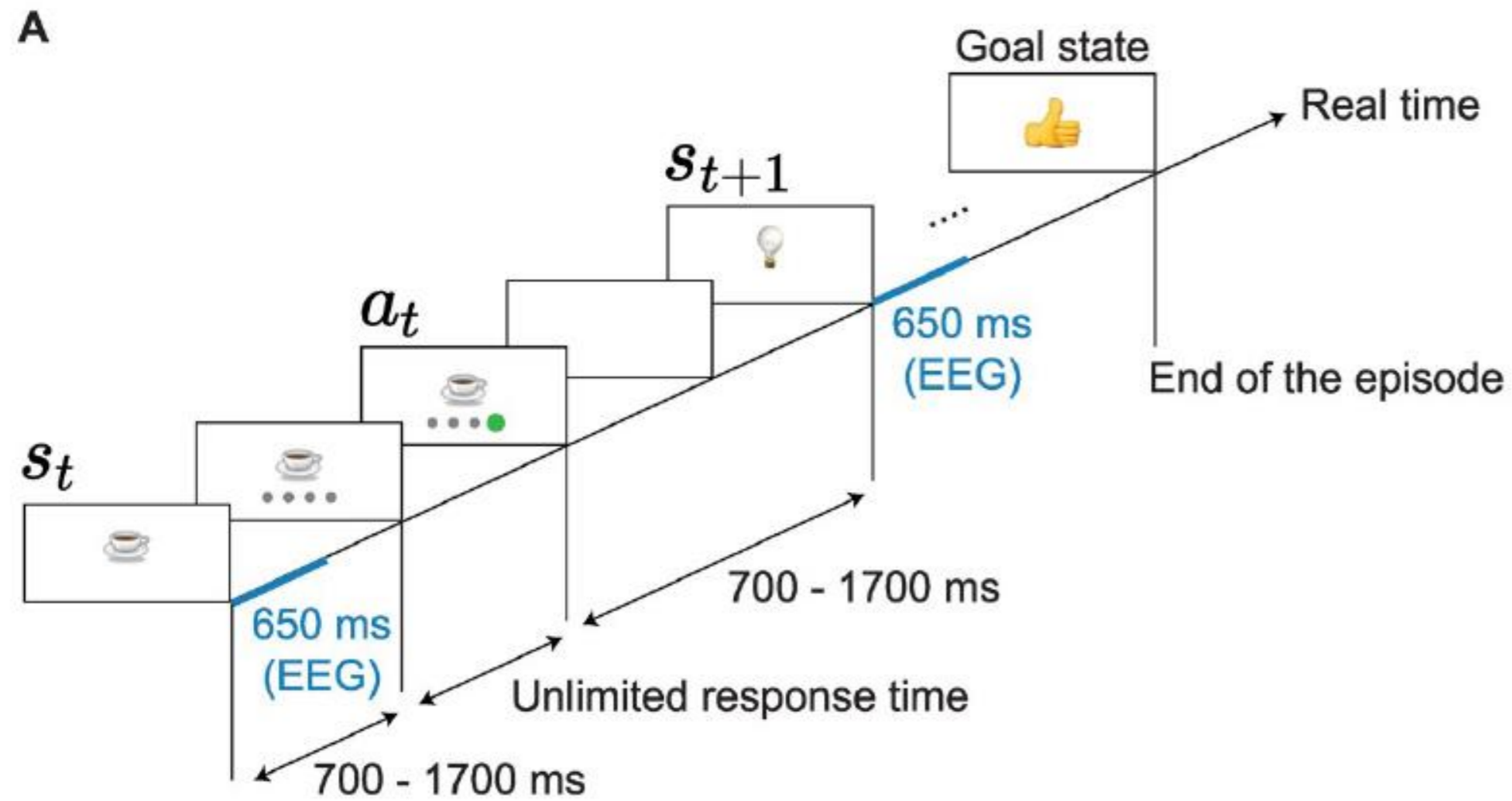
Previous slide.

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).

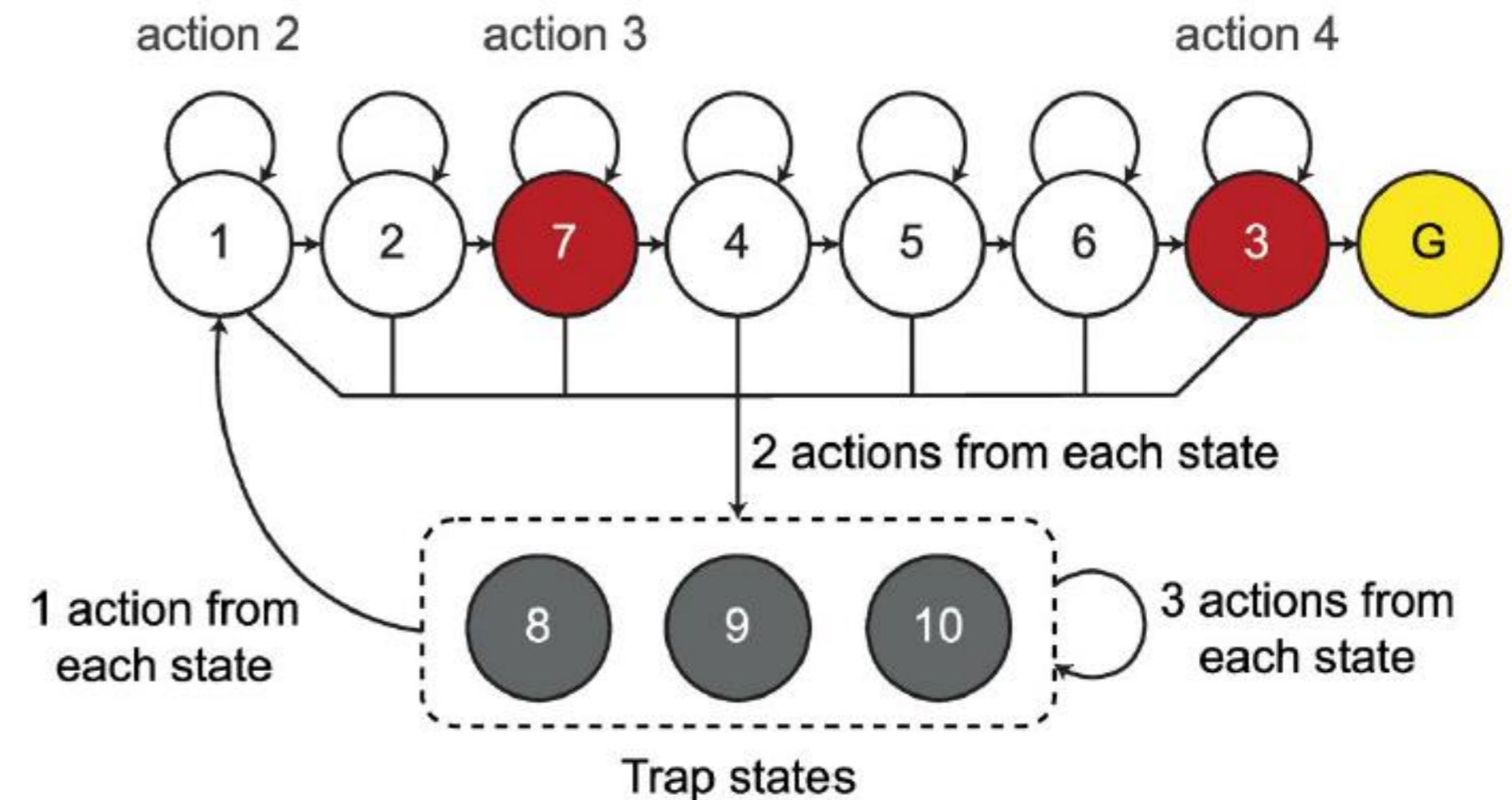
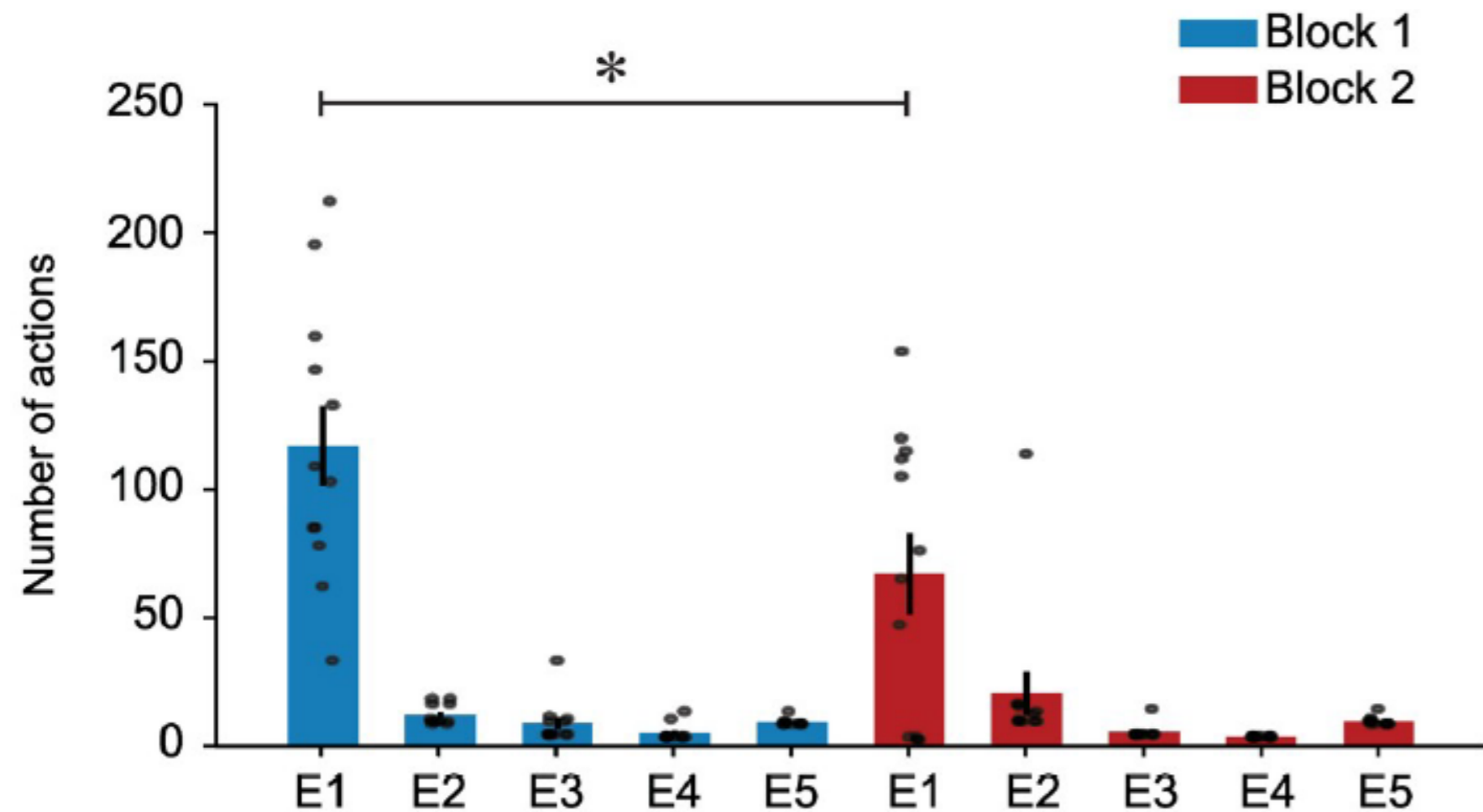
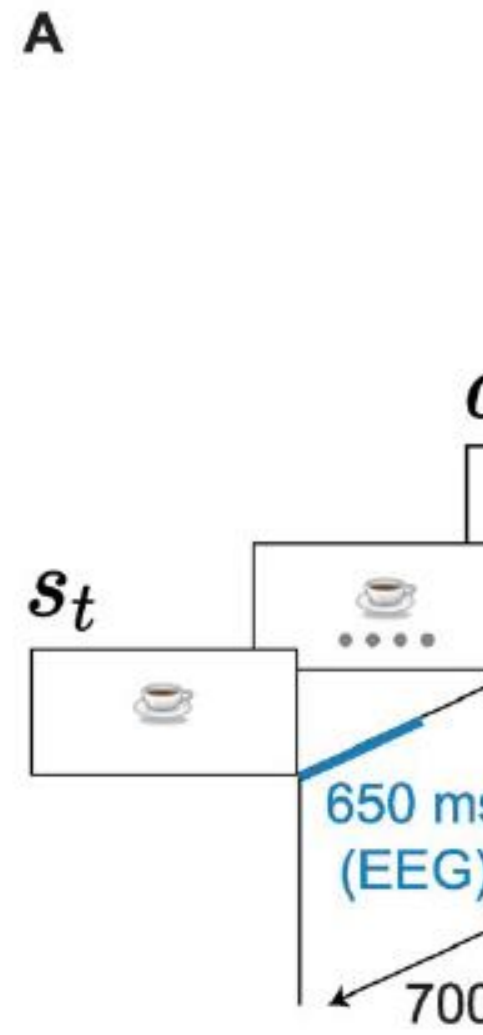
Volatile Environment: Switch after episode 5

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?



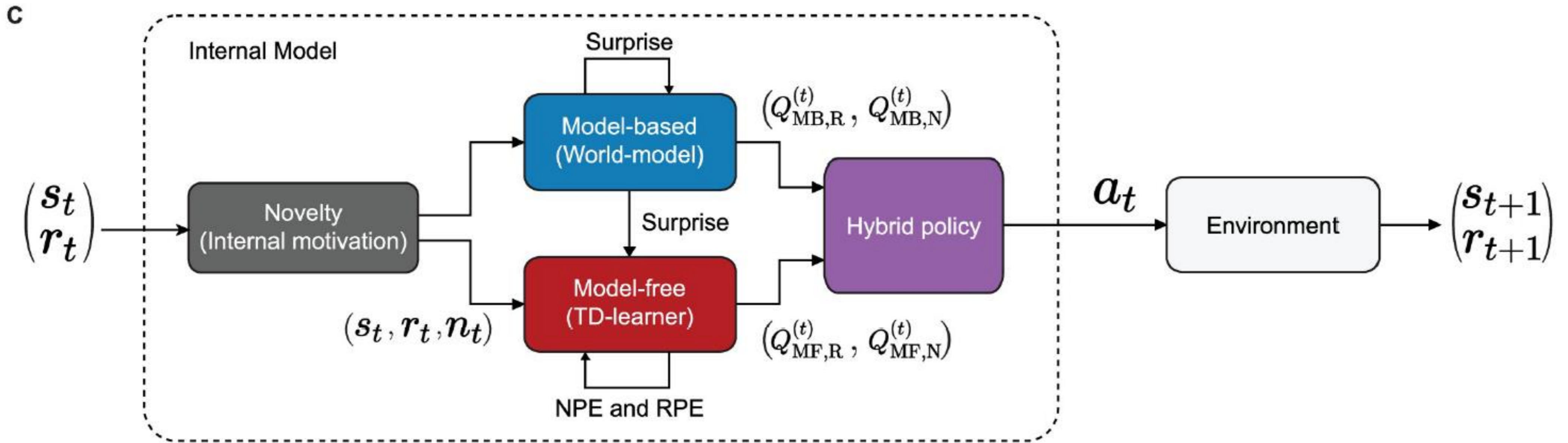
Previous slide.

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)



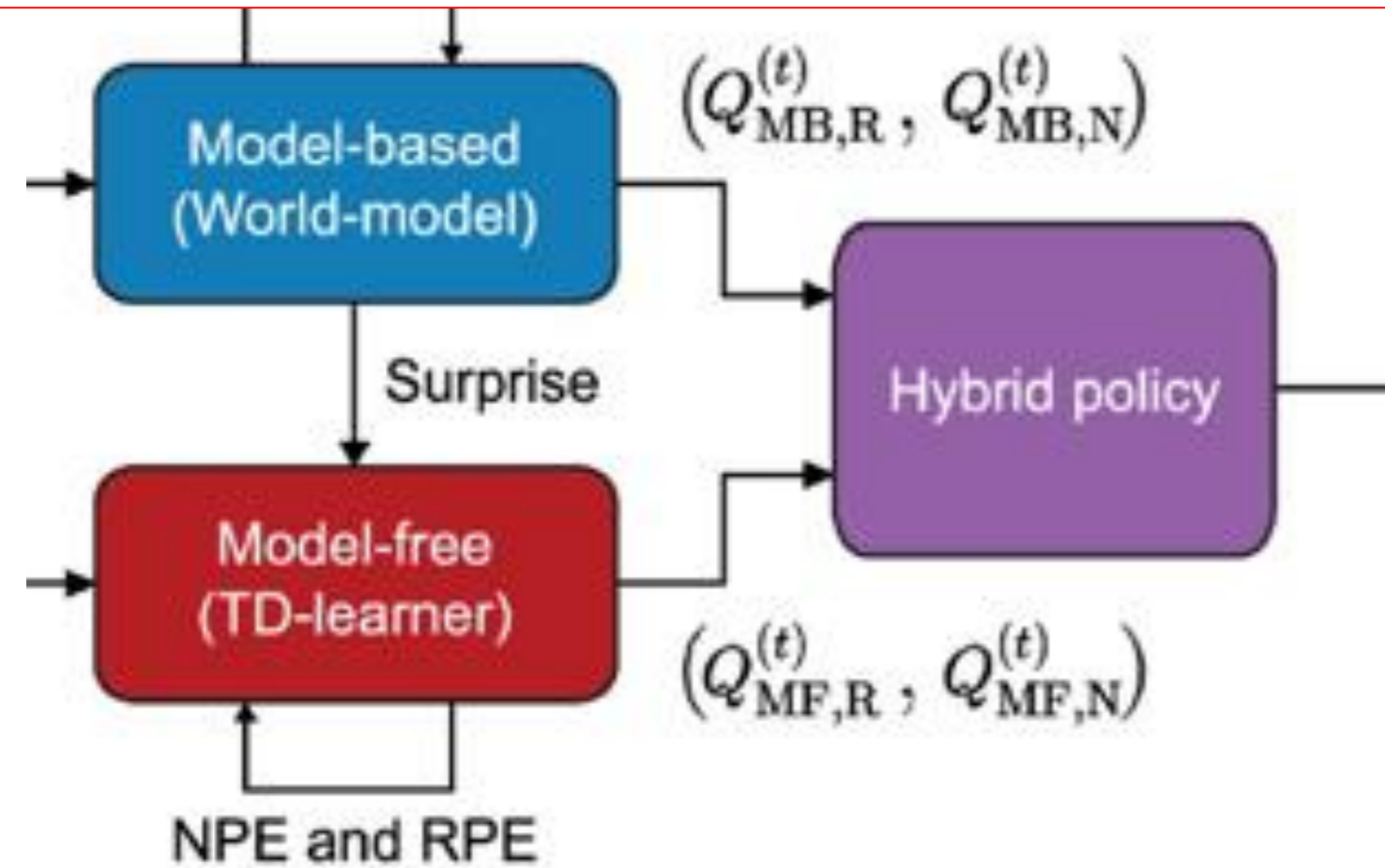
$$\text{RPE} = [r_t + \gamma \max_{a'} Q_R(s', a') - Q_R(s, a)]$$

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

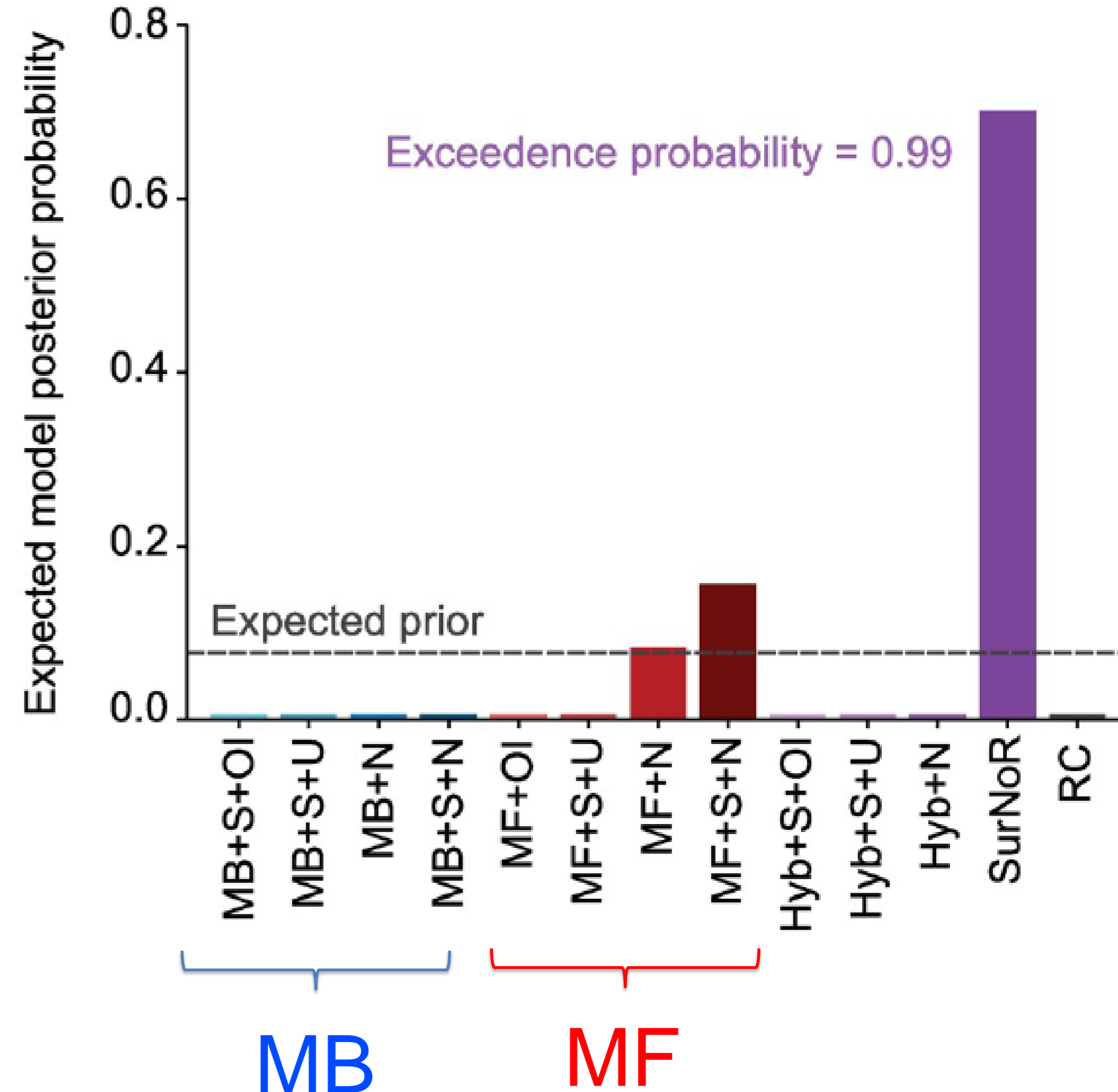
Comparison of Models: Surprise, Novelty, Reward

Finding 4)

Rapid relearning needs surprise



- Turn off novelty
- Turn off surprise
- Turn off model-based → MF
- Turn off model-free → MB
- OI = Optimistic Initialization



Previous slide.

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

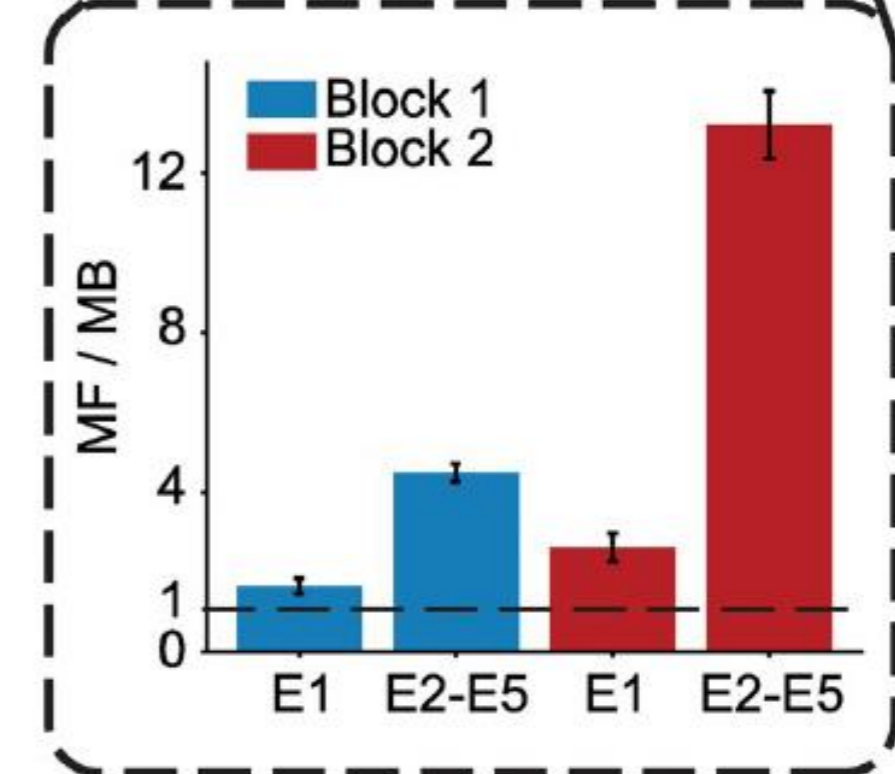
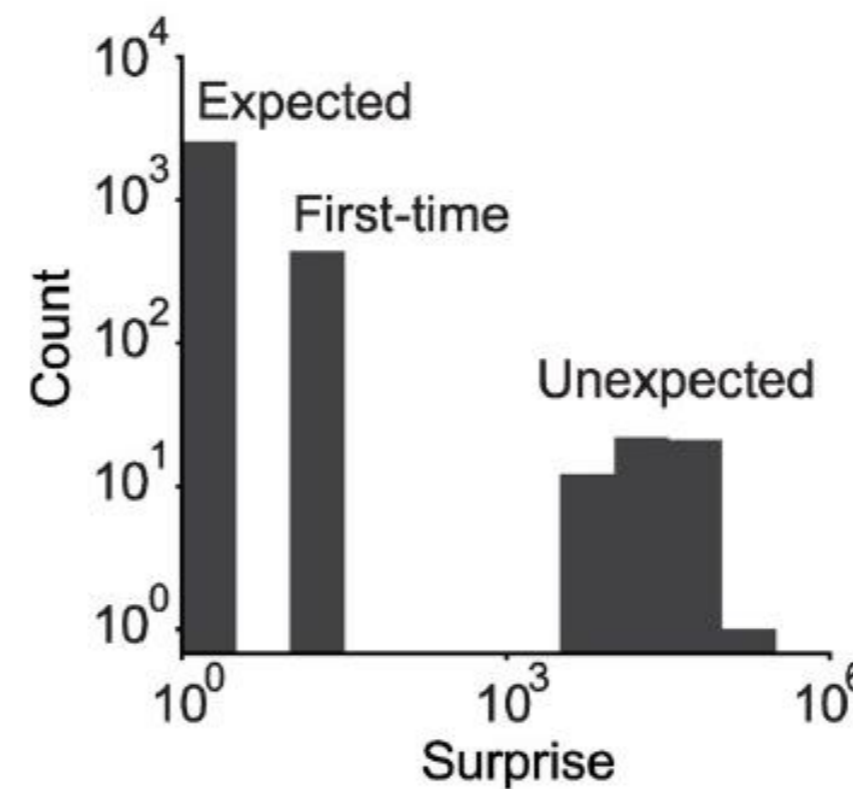
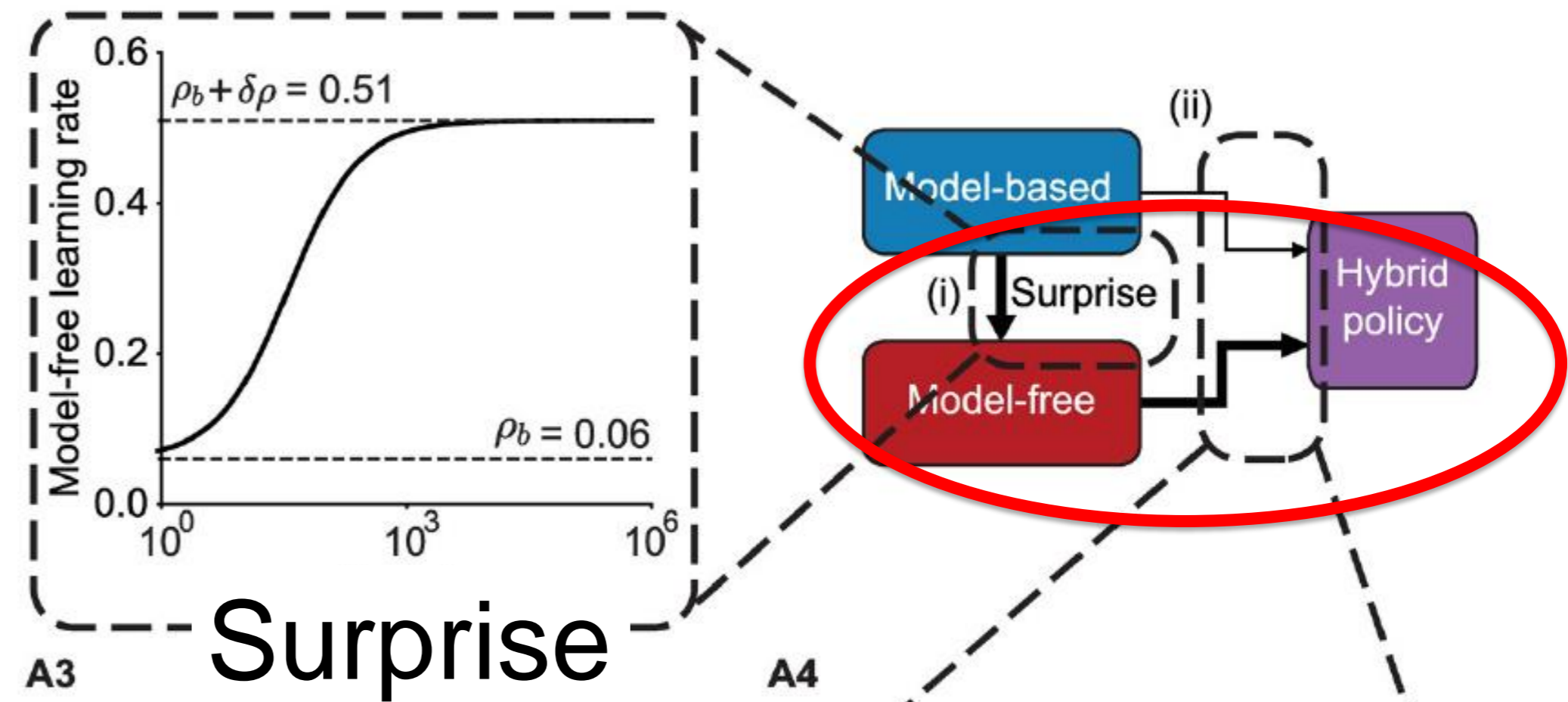
The second best model is model-free (MF) with surprise, novelty, and reward.
Turning off surprise lowers the performance (Hybrid model and surprise).

Model-based compares less well with human data than model-free.

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.

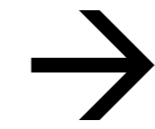
Surprise needs a world model, but we said that the model-free algorithm better explains the behavior.

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.

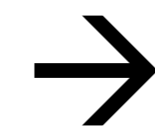
Reward-based learning versus Surprise-based learning

Reward-Prediction Error



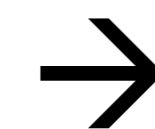
Surprise

defined as
TD error



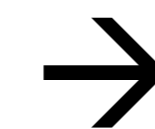
defined as
Bayes Factor Surprise

stimulated by
chocolate, money,
praise, ...



stimulated by observations
not consistent with momentary
model of environment

modulates
learning rate

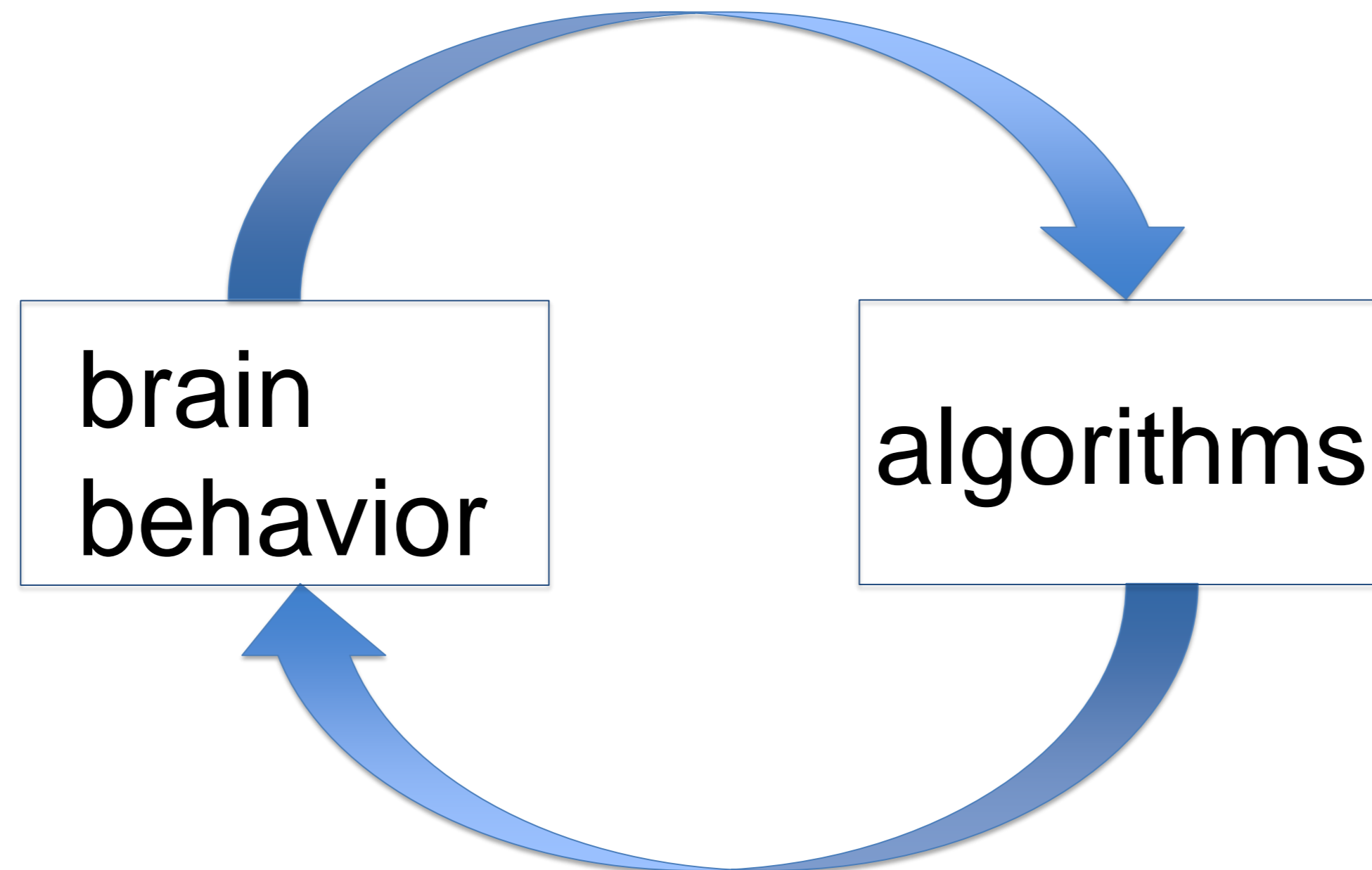


modulates
learning rate

Previous slide.

Summary: Comparison of Reward Prediction Error and Surprise.

Current Research in Reinforcement Learning:



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

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.