# COSMOS

Computational summer school on modeling social and collective behavior - Konstanz (DE) July 4th - 7th
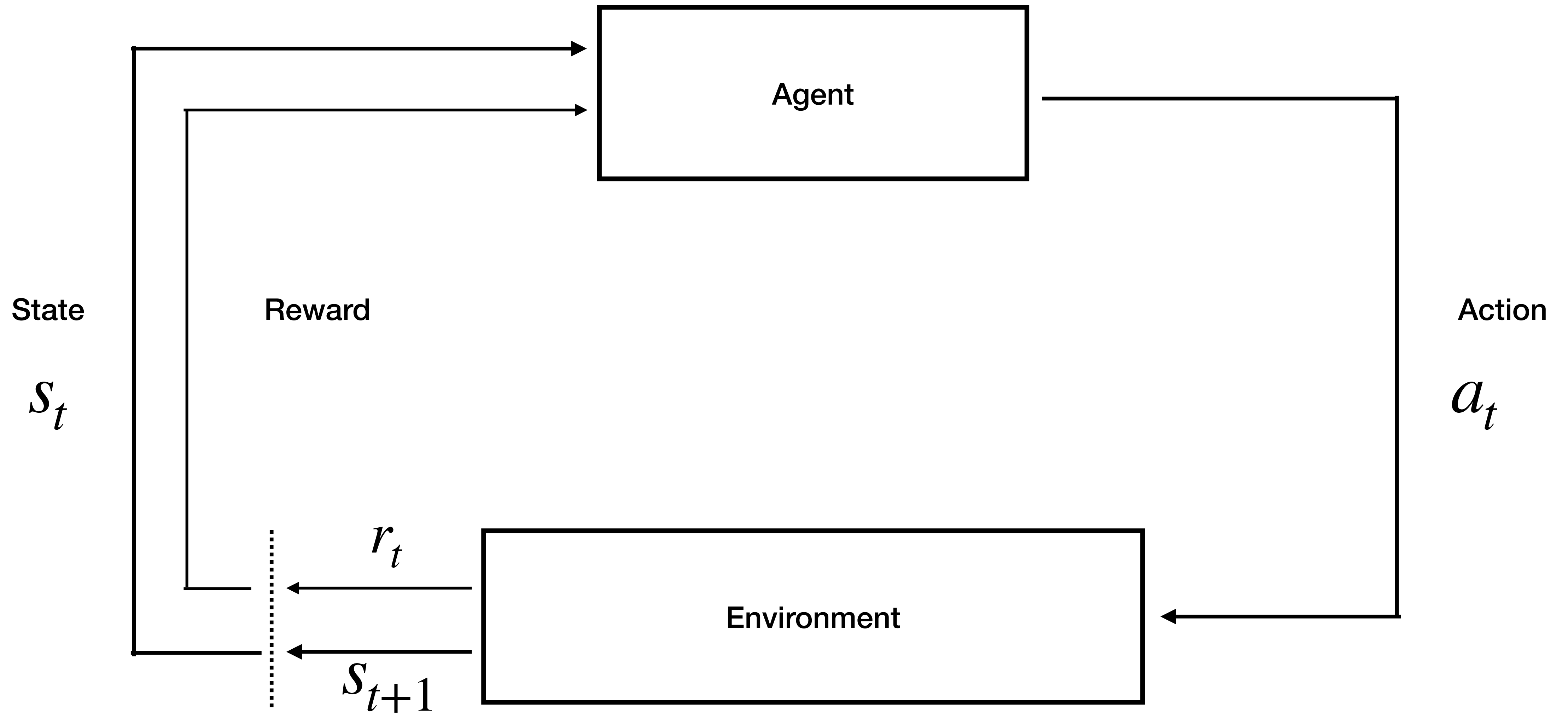
**Charley Wu & Wataru Toyokawa**
**July 5th**

# Goals of Tutorial 2:

- **Brisk introduction to asocial RL**

  - Simulating data

  - Maximum likelihood estimation (MLE) of model parameters

  - Predicting choices

- **Social learning models**

  - Imitating actions

  - Combining asocial and social learning

  - Social learning hierarchy (from imitation to Theory of Mind)

- **Scaling up to more complex problems**

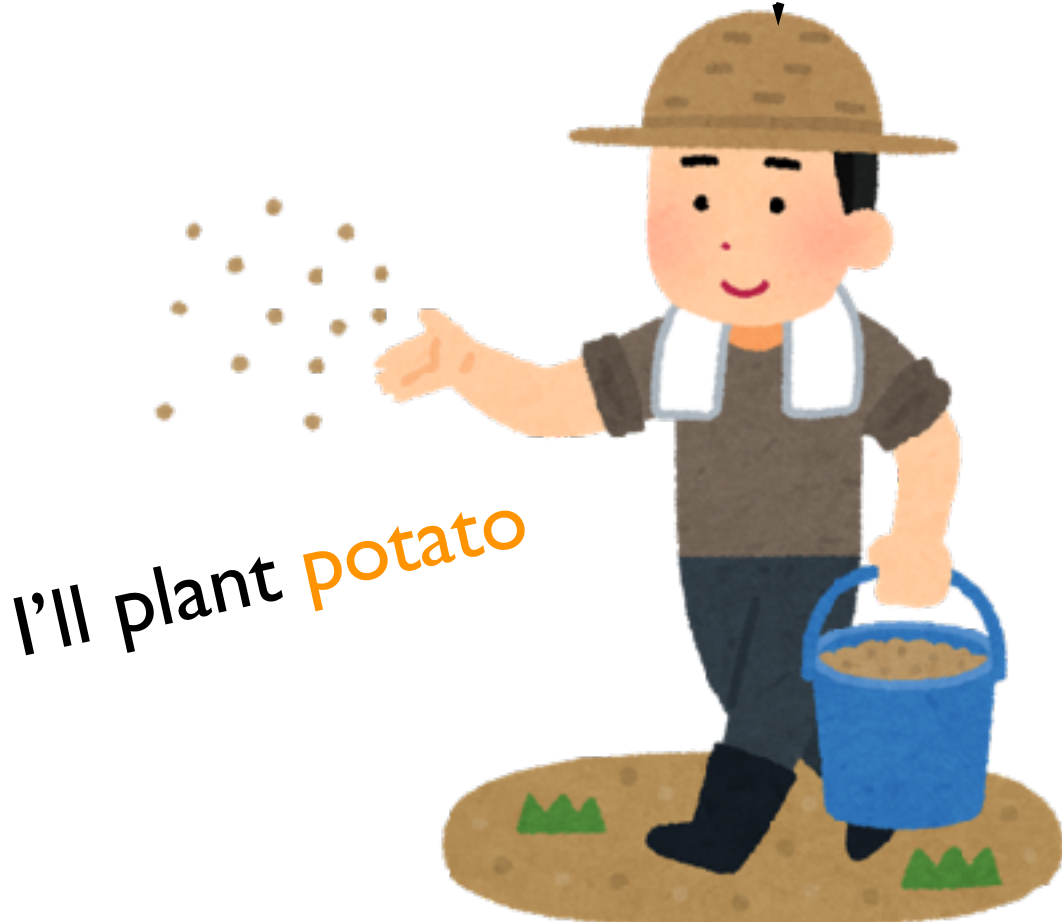- **Evolutionary simulations**

# Reinforcement Learning (RL)



Agent

Environment

State

Reward
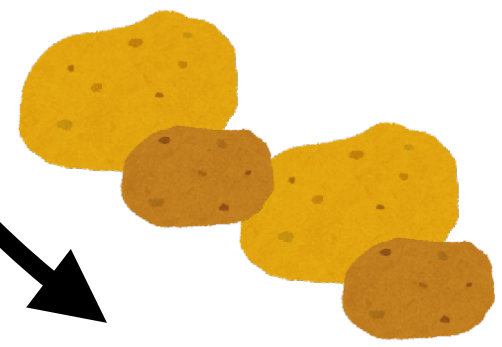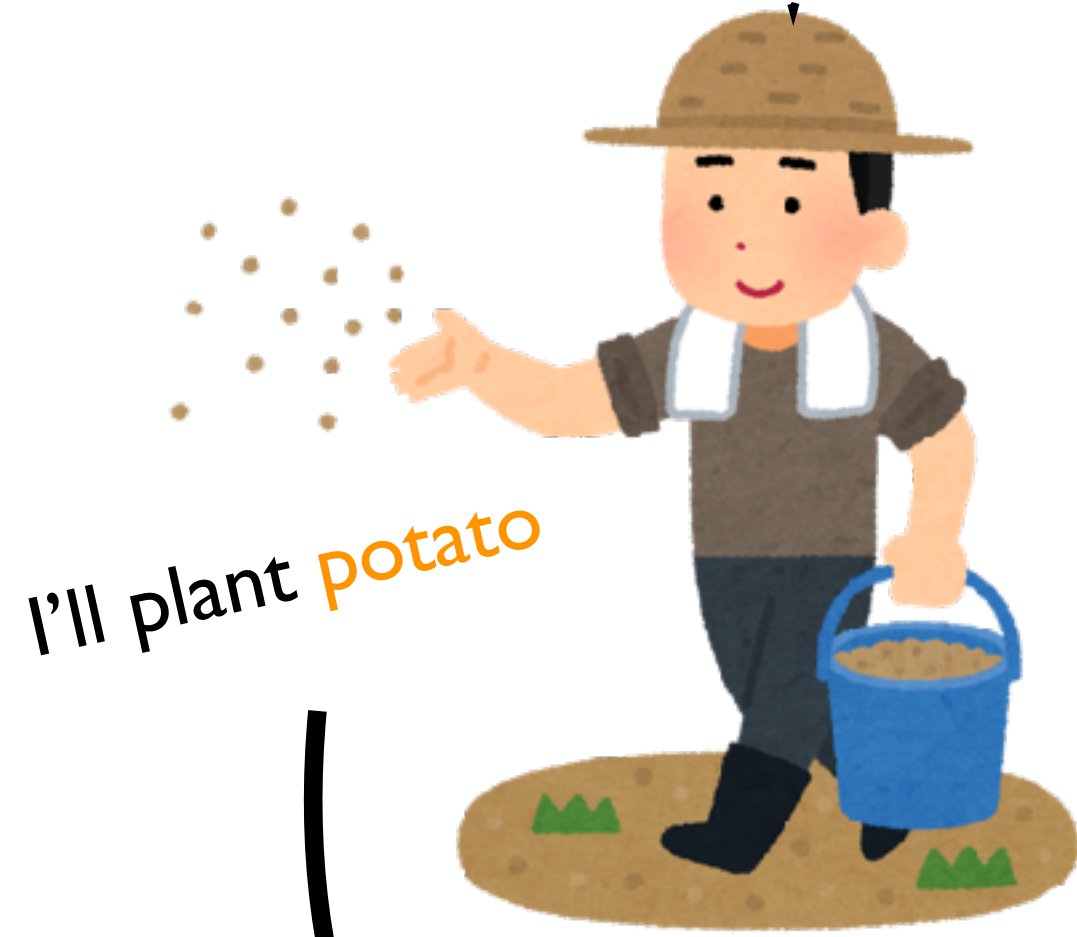
Action

$s_t$

$a_t$

$r_t$

$s_{t+1}$

Sutton & Barto (1998)

# A multi-armed bandit task

# Season 1

The farmer's preference

🥔 ≈ 🌾

I'll plant potato

**Poor yield than expected**

Season 1

The farmer's preference

🥔 ≈ 🌾

I'll plant potato

Poor yield than expected

Learning from experience

Season 2

Updated preference

🥔 < 🌾

Season 1

The farmer's preference

🥔 ≈ 🌾

I'll plant potato

Poor yield than expected

Learning from experience

Season 2

Updated preference

🥔 < 🌾

I'll plant wheat

Good yield

Season 1

The farmer's preference

🥔 ≈ 🌾

I'll plant potato
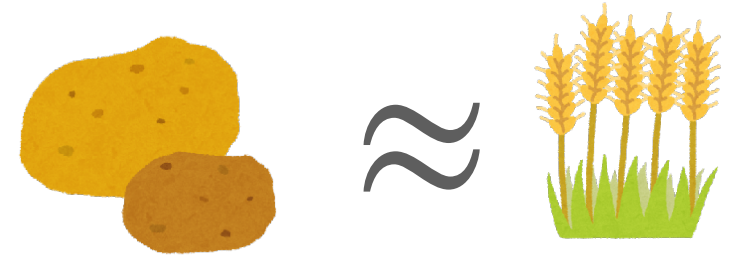
Learning from experience

Poor yield than expected

Season 2

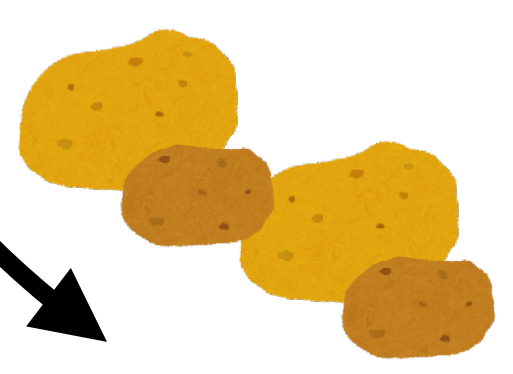Updated preference

🥔 < 🌾

I'll plant wheat

Learning from experience

Good yield

Season 3

Updated preference

🥔 << 🌾

I'll plant wheat

# RL process is not directly observable. It should be inferred statistically.

**Observable world**

t = 1

t = 2

t = 3

| Choice | | Choice | | Choice |
| --- | --- | --- | --- | --- |
| Payoff | | Payoff | | Payoff |

# RL process is not directly observable. It should be inferred statistically.

# RL process is not directly observable. It should be inferred statistically.

# Decision value updating:
# Q-Learning

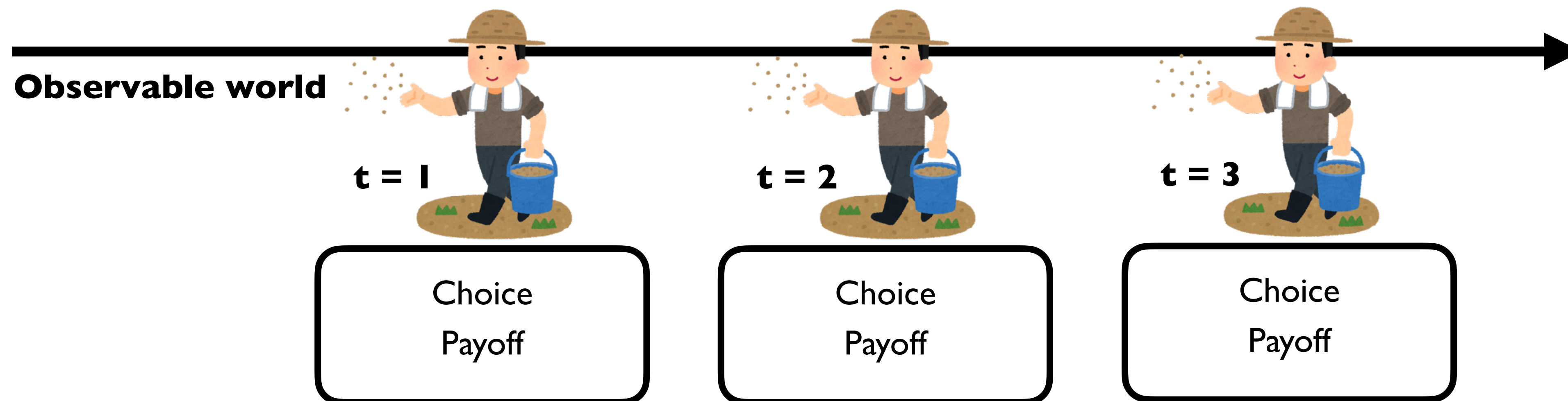t = 1



Q-values

# Decision value updating: Q-Learning

**t = 1**

**Rewards**



Q-values

I'll plant wheat

# Decision value updating: Q-Learning

Learning rate
(memory parameter; step size)

**Rewards**

t = 1

$(1-α) ×$  $+ α ×$ 

Q-values

I'll plant wheat

# Decision value updating:
# Q-Learning

Learning rate
(memory parameter; step size)

**Rewards**

**t = 2**

**t = 1**

$(1-\alpha) \times$

$+$ $\alpha \times$

0    5    10

Q-values

0    5    10

Q-values

I'll plant wheat

# Reward prediction error

$$Q_t(a)$$

**Reward prediction error**



$$(1 - \alpha)Q_t(a) + \alpha r_t(a)$$

# Reward prediction error

$$Q_{t+1}(a) \leftarrow (1 - \alpha)Q_t(a) + \alpha r_t(a)$$

**Reward prediction error**

$$Q_{t+1}(a) \leftarrow (1 - \alpha)Q_t(a) + \alpha r_t(a)$$

$$Q_{t+1}(a) \leftarrow Q_t(a) + \alpha \,[r_t(a) - Q_t(a)]$$

**reward prediction error (RPE)**

# Converting value to actions:
# Softmax policy (i.e. multinomial logistic function)

Q-values

Policy $\pi$

$Q_t(w)$

$Q_t(p)$

$\pi_t(w)$

$\pi_t(p)$

We want policy to satisfy:

$$\sum_a \pi(a) = 1$$

$$Q_1 > Q_2 \Leftrightarrow \pi_1 > \pi_2$$

# Converting value to actions:
# Softmax policy (i.e. multinomial logistic function)

Q-values

Policy $\pi$

$Q_t(w)$

$Q_t(p)$

$$\frac{\exp[\beta Q_t(a)]}{\exp[\beta Q_t(w)] + \exp[\beta Q_t(p)]}$$

$$\frac{e^{\beta Q_t(a)}}{e^{\beta Q_t(w)} + e^{\beta Q_t(p)}}$$

$\pi_t(w)$

$\pi_t(p)$

We want policy to satisfy:

$$\sum_a \pi(a) = 1$$

$$Q_1 > Q_2 \Leftrightarrow \pi_1 > \pi_2$$

# Converting value to actions:
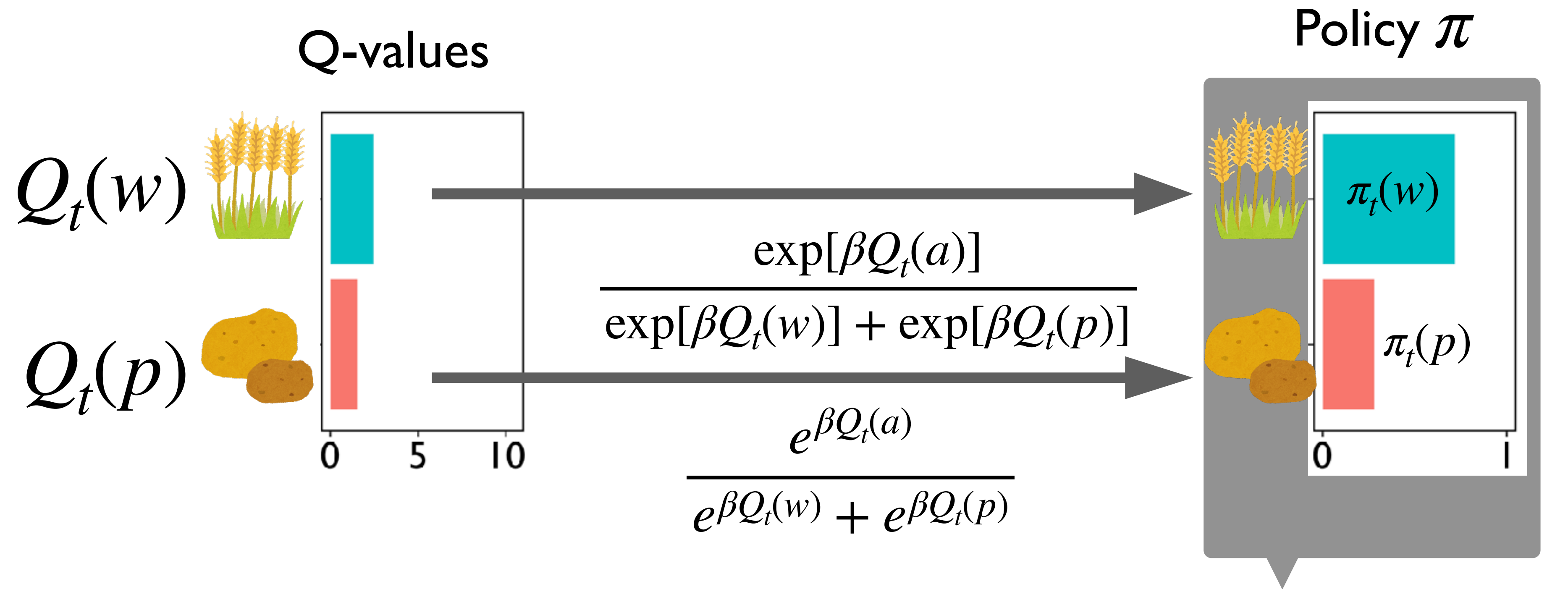# Softmax policy (i.e. multinomial logistic function)

Q-values

$Q_t(w)$

$Q_t(p)$

$$\frac{\exp[\beta Q_t(a)]}{\exp[\beta Q_t(w)] + \exp[\beta Q_t(p)]}$$

$$\frac{e^{\beta Q_t(a)}}{e^{\beta Q_t(w)} + e^{\beta Q_t(p)}}$$
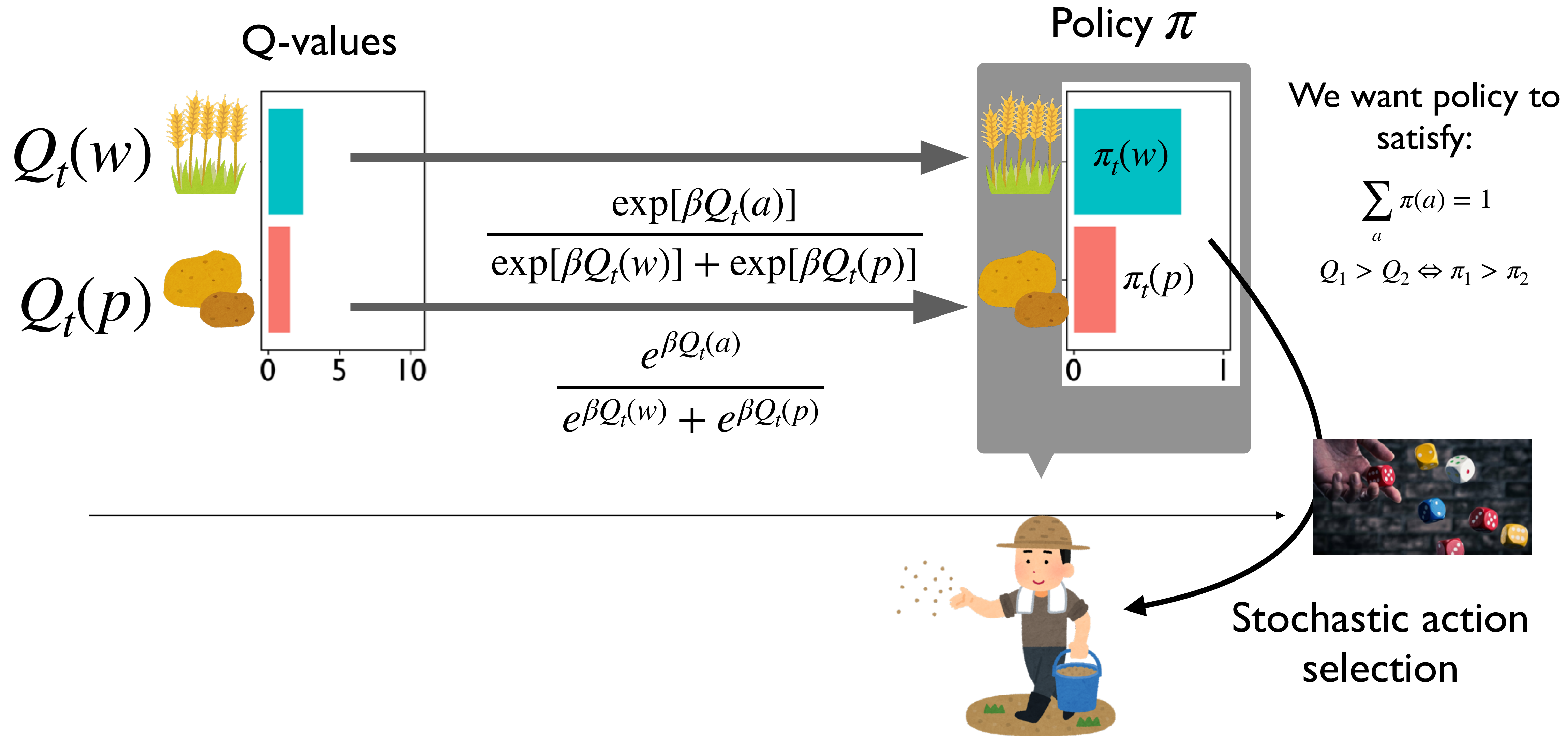
Policy $\pi$

$\pi_t(w)$

$\pi_t(p)$

We want policy to satisfy:

$$\sum_a \pi(a) = 1$$

$$Q_1 > Q_2 \Leftrightarrow \pi_1 > \pi_2$$

Stochastic action selection

Softmax policy

$$\pi_t(w) = \frac{\exp[\beta Q_t(w)]}{\exp[\beta Q_t(w)] + \exp[\beta Q_t(p)]}$$

Probability of choosing wheat

$Q_t(w) - Q_t(p)$

β
— 0.5
— 1
— 2
— 3
— 4

https://cosmos-konstanz.github.io/notebooks/tutorial-2-models-of-learning.html#simulating-data
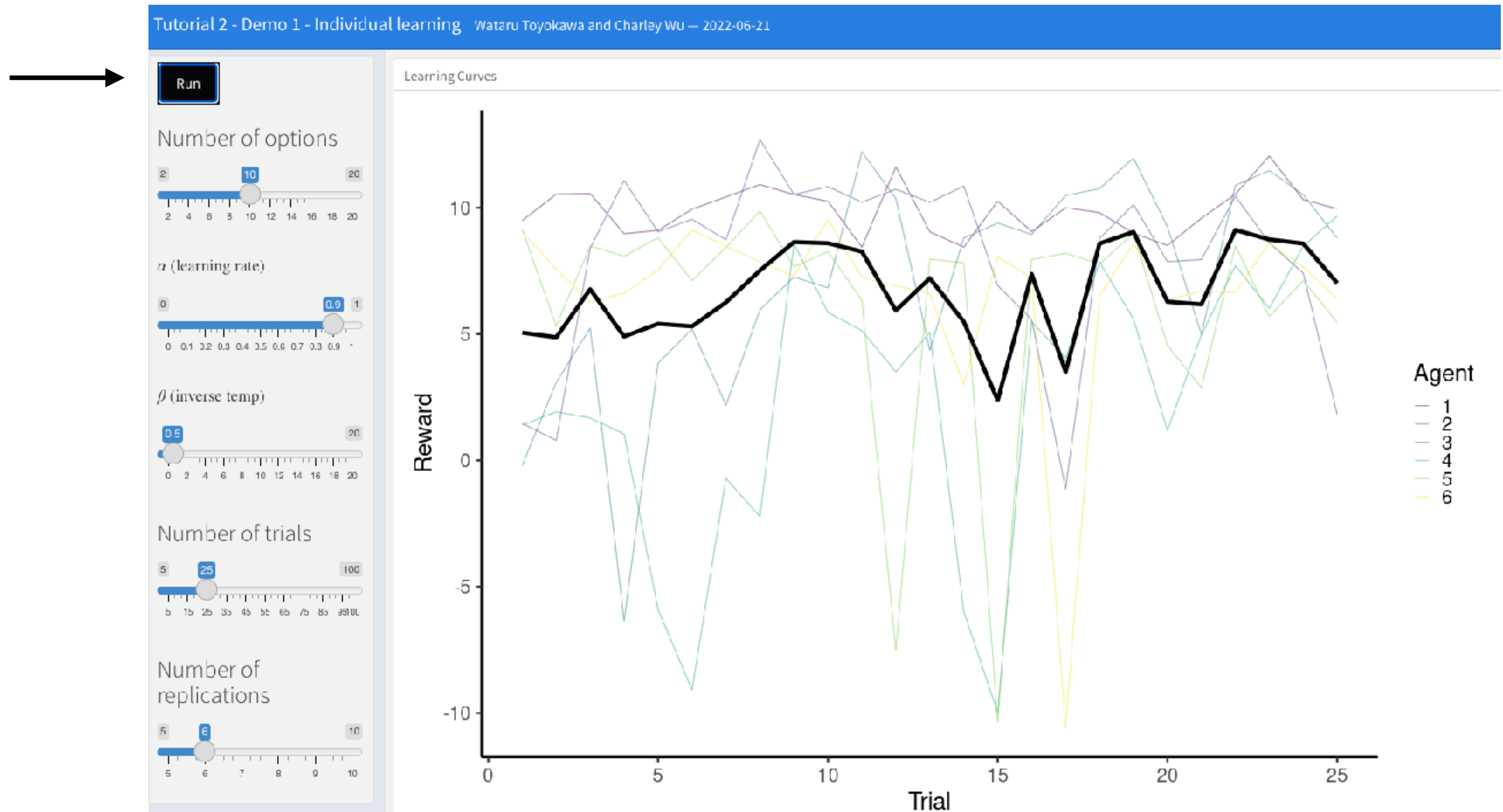
# Demo 1: Tweaking individual learning parameters



Which learning parameters $(\alpha, \beta)$ typically produce the best results?

# Likelihood function

# Likelihood function

Coin Flip Model

# Likelihood function

Coin Flip Model



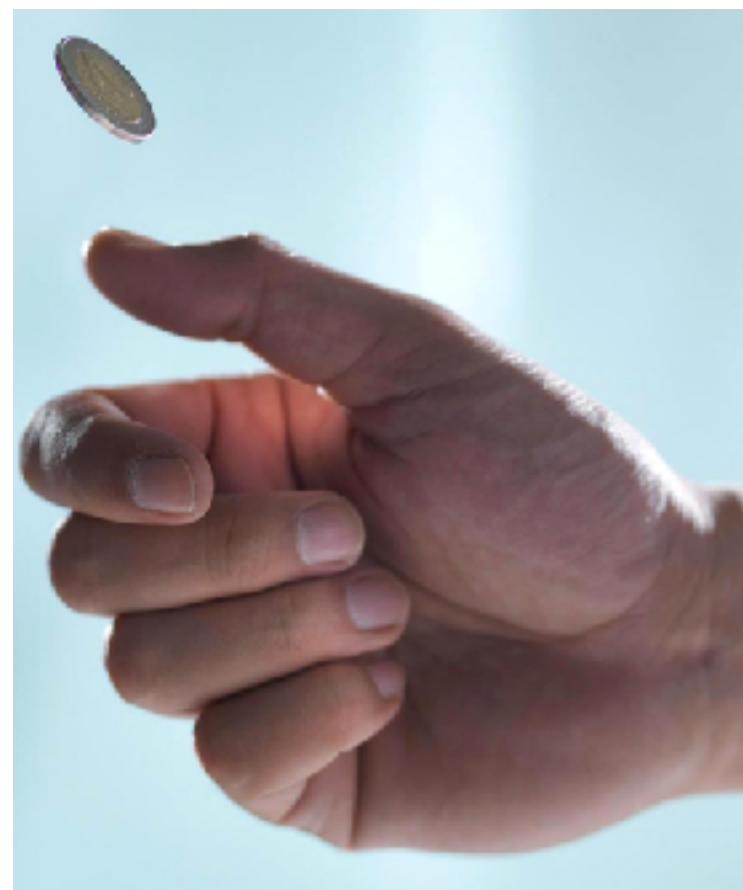Observed Data:

$$D = \{d_1 = h, d_2 = t, \ldots d_n = t\}$$

Model:

$$D \sim \textbf{Binomial}(n, \theta)$$

$$\theta = P(h)$$

# Likelihood function

Coin Flip Model



Observed Data:

$$D = \{d_1 = h, d_2 = t, \ldots d_n = t\}$$

Model:

$$D \sim \textbf{Binomial}(n, \theta)$$

$$\theta = P(h)$$

# Likelihood function

Beyond only simulating data, we also want to use models to describe experimental data.

To fit a model to data, we first need to define a **Likelihood Function:**

$$P(D\,|\,\theta)$$

describing the probability that the observed data $D$ was generated based on model parameters $\theta$
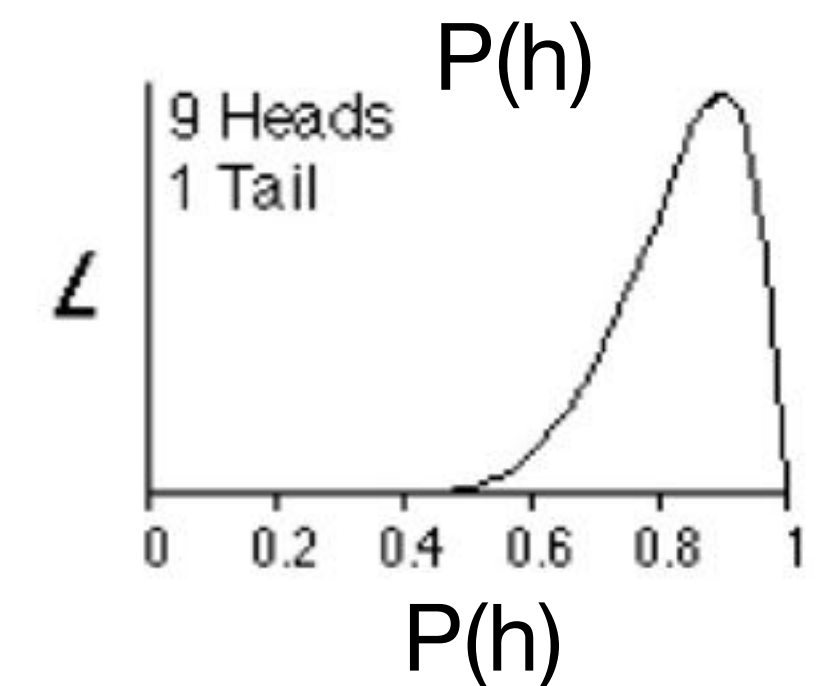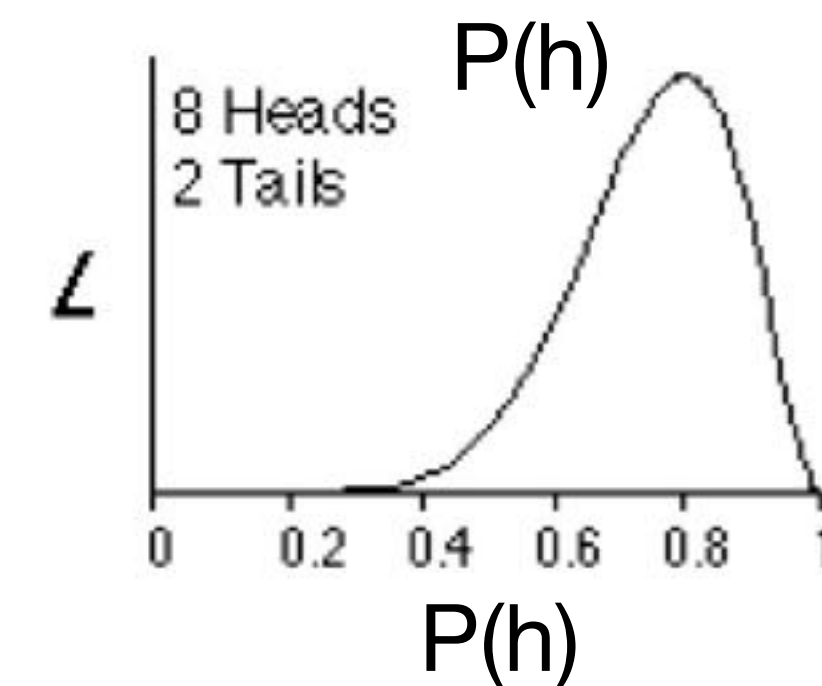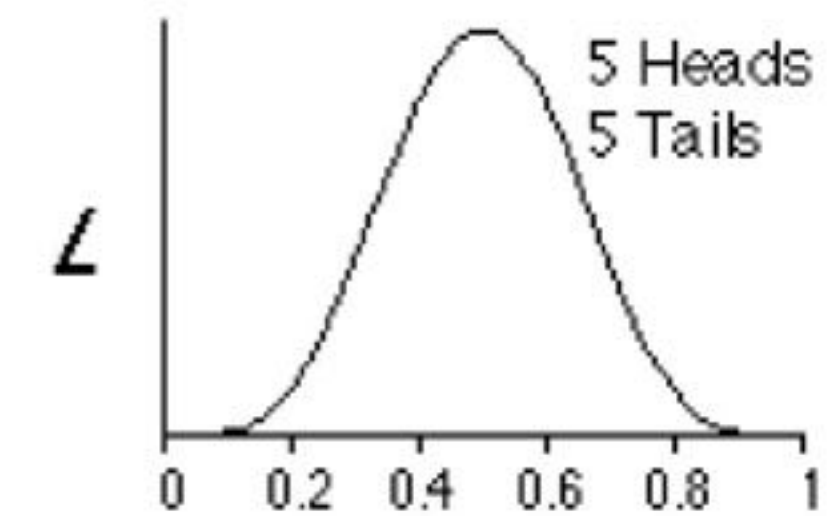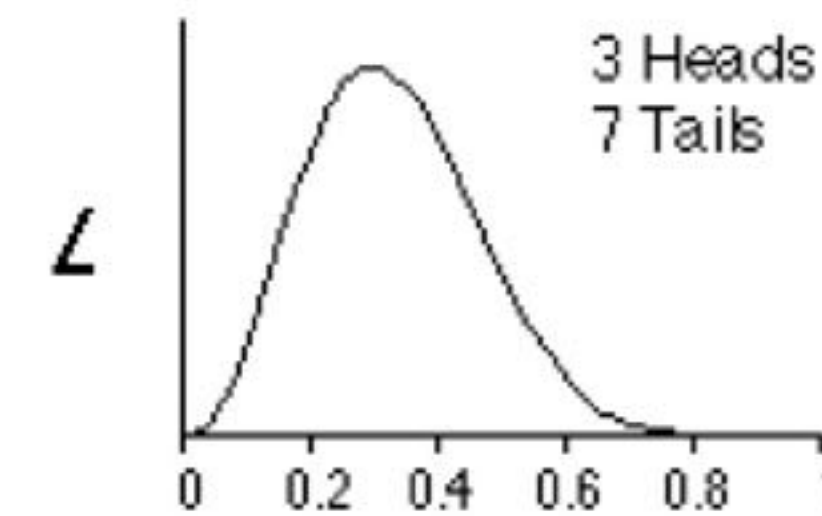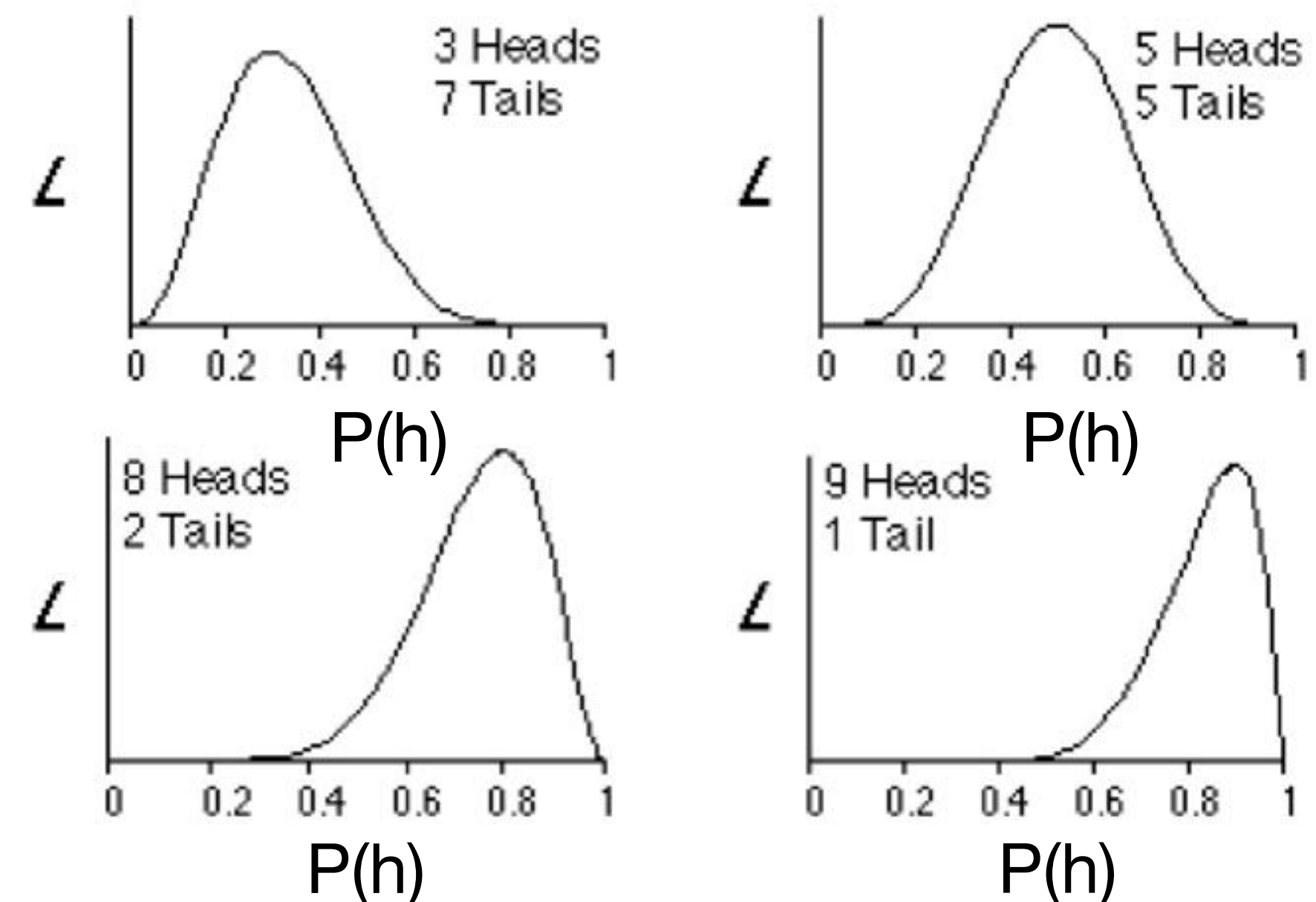
Coin Flip Model

Observed Data:

$$D = \{d_1 = h, d_2 = t, \ldots d_n = t\}$$

Model:

$$D \sim \textbf{Binomial}(n, \theta)$$

$$\theta = P(h)$$

# Log Likelihoods

Since we are usually modeling multiple data points, we need to describe the **joint likelihood** over all observations:

$$P(D \,|\, \theta) = \prod_i P(d_i \,|\, \theta)$$

This is much easier using logarithms, since we can replace multiplication with summation in log space to compute the **log likelihood**
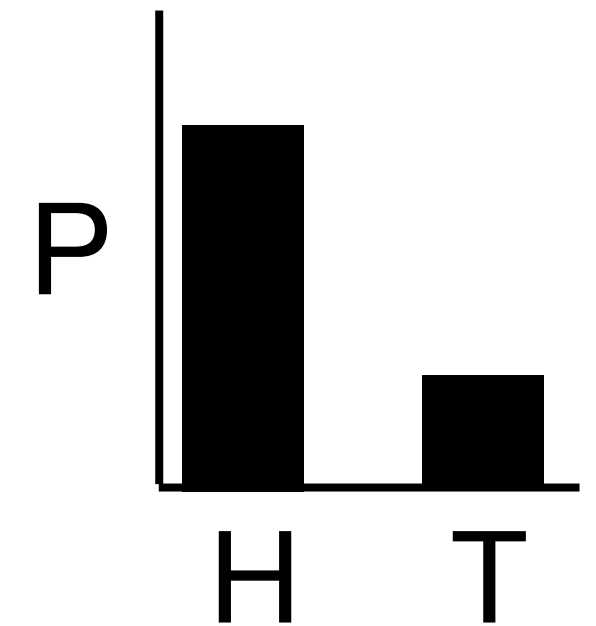
$$\log P(D \,|\, \theta) = \sum_t \log P(d_t \,|\, \theta)$$

Since probabilities are always <1, the log likelihood will always be negative. Thus, it's more convenient to express the fit of a model using the **negative log likelihood** (nLL) by inverting the sign:

$$nLL = -\log P(D \,|\, \theta)$$

The nLL expresses the amount of error or loss (aka 'log loss') and will always be greater than zero. Smaller values thus describe better model fits.
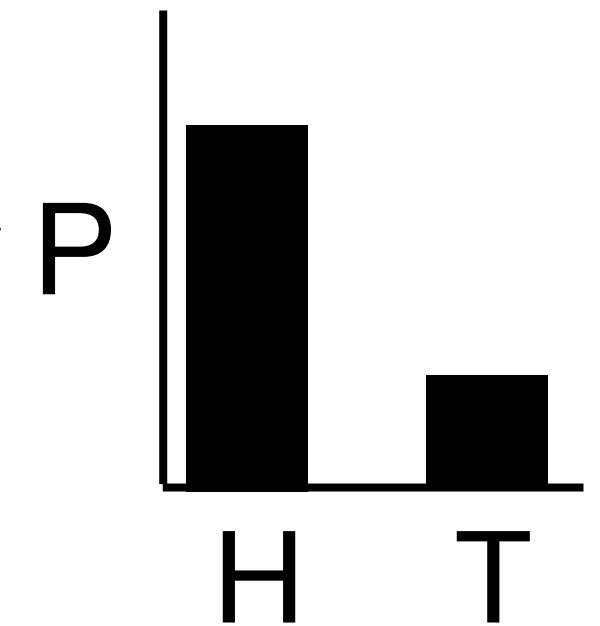
# Likelihoods as Goodness of Fit



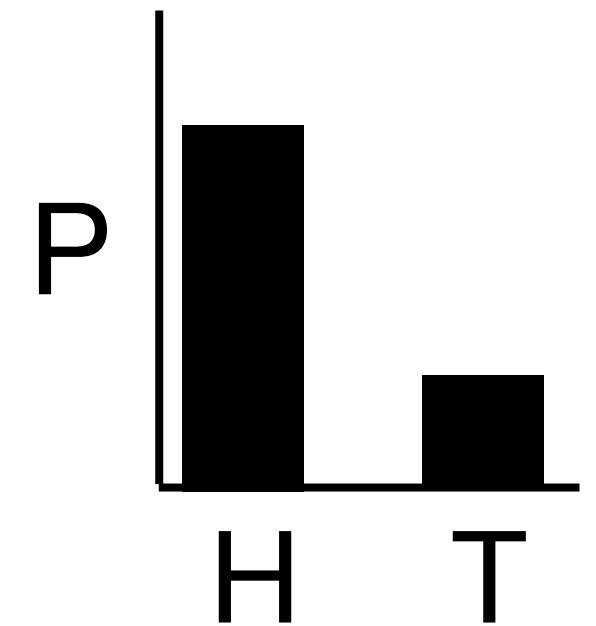| Measure | Formula | Heads | Tails |
|---|---|---|---|
| Likelihood | $P(D|\theta)$ | 80% | 20% |
| Log likelihood | $\log P(D|\theta)$ | –0.22 | –1.61 |
| Negative Log Likelihood (nLL) | $-\log P(D|\theta)$ | 0.22 | 1.61 |
| Deviance | $-2 \log P(D|\theta)$ | 0.44 | 3.22 |

# Likelihoods as Goodness of Fit

| Measure | Formula | Heads | Tails |
|---------|---------|-------|-------|
| Likelihood | $P(D\|\theta)$ | 80% | 20% |
| Log likelihood | $\log P(D\|\theta)$ | –0.22 | –1.61 |
| Negative Log Likelihood (nLL) | $-\log P(D\|\theta)$ | 0.22 | 1.61 |
| Deviance | $-2 \log P(D\|\theta)$ | 0.44 | 3.22 |

**aka Log Loss**

# Likelihoods as Goodness of Fit

| Measure | Formula | Heads | Tails |
|---|---|---|---|
| Likelihood | $P(D|\theta)$ | 80% | 20% |
| Log likelihood | $\log P(D|\theta)$ | –0.22 | –1.61 |
| Negative Log Likelihood (nLL) | $-\log P(D|\theta)$ | 0.22 | 1.61 |
| Deviance | $-2 \log P(D|\theta)$ | 0.44 | 3.22 |

**aka Log Loss**

# Likelihoods as Goodness of Fit

| Measure | Formula | Heads | Tails |
|---|---|---|---|
| Likelihood | $P(D|\theta)$ | 80% | 20% |
| Log likelihood | $\log P(D|\theta)$ | –0.22 | –1.61 |
| Negative Log Likelihood (nLL) | $-\log P(D|\theta)$ | 0.22 | 1.61 |
| Deviance | $-2 \log P(D|\theta)$ | 0.44 | 3.22 |

**aka Log Loss**

# From model simulation to likelihood functions

In practice, we can use code very simular to our model simulations to create a likelihood function

```
likelihood <- function(params, data){
    nLL <- 0 #initialize negative log likelihood
    for (d in data){ #loop through data
      predictions <- model(params) #make predictions
      observedAction <- d #define true outcome
      nLL <- nLL -log(predictions[observedAction]) #Update nLL
    }
    return(nLL)
}
```
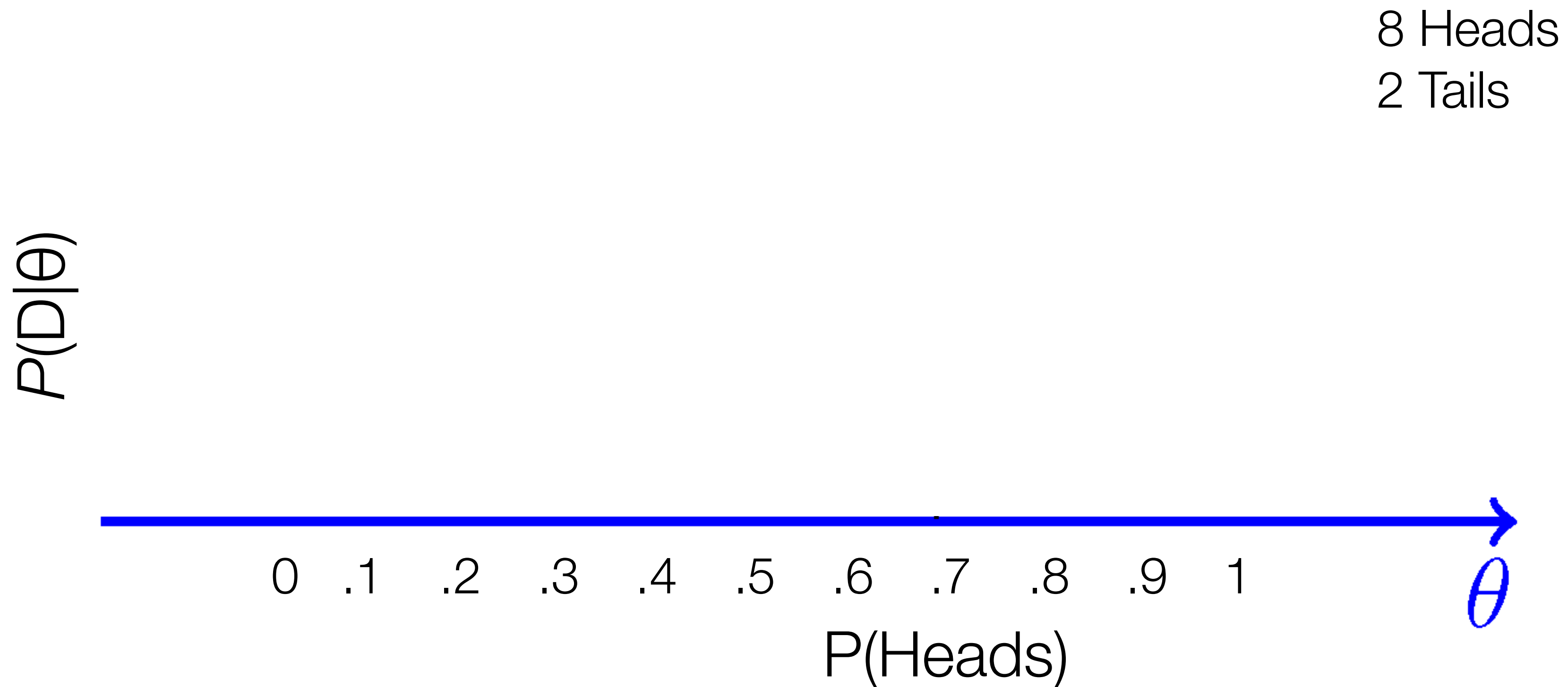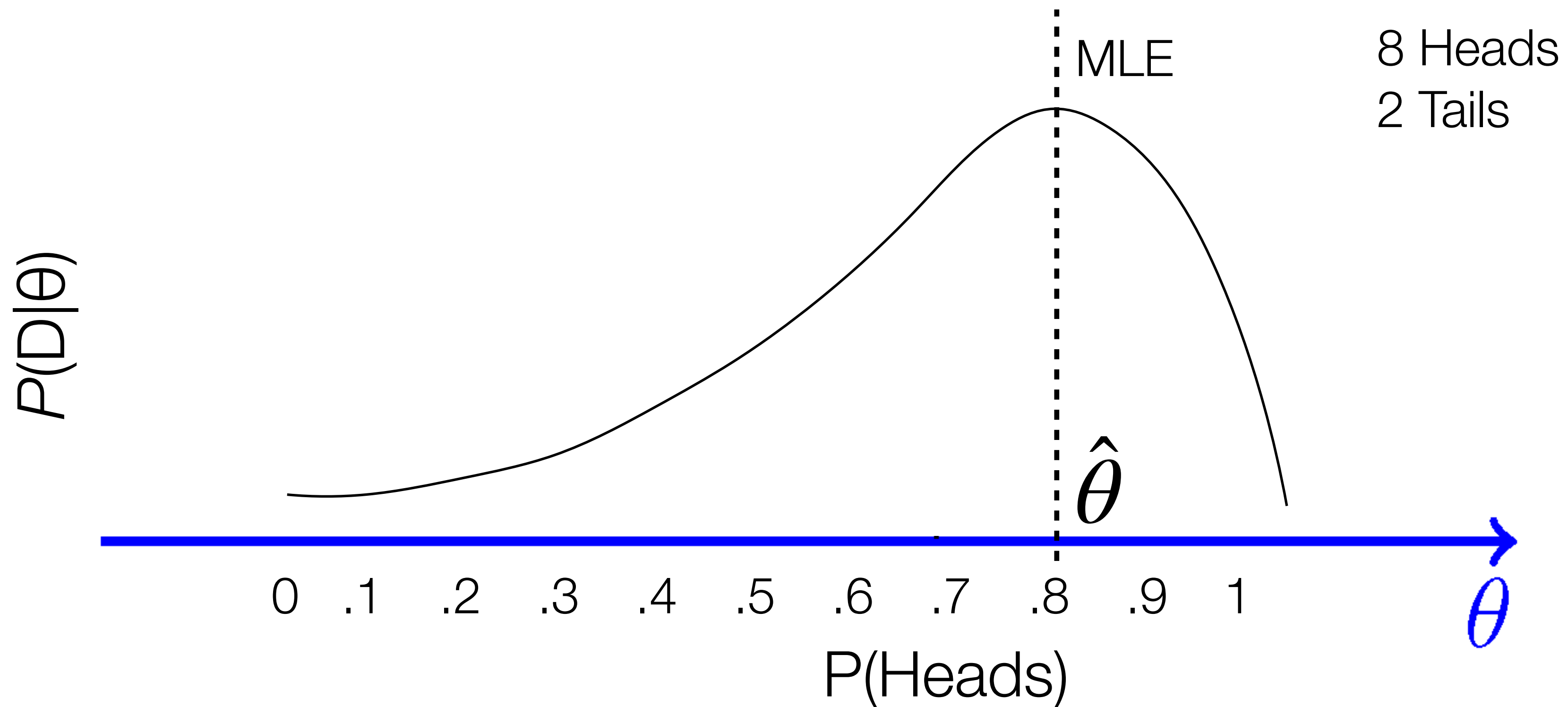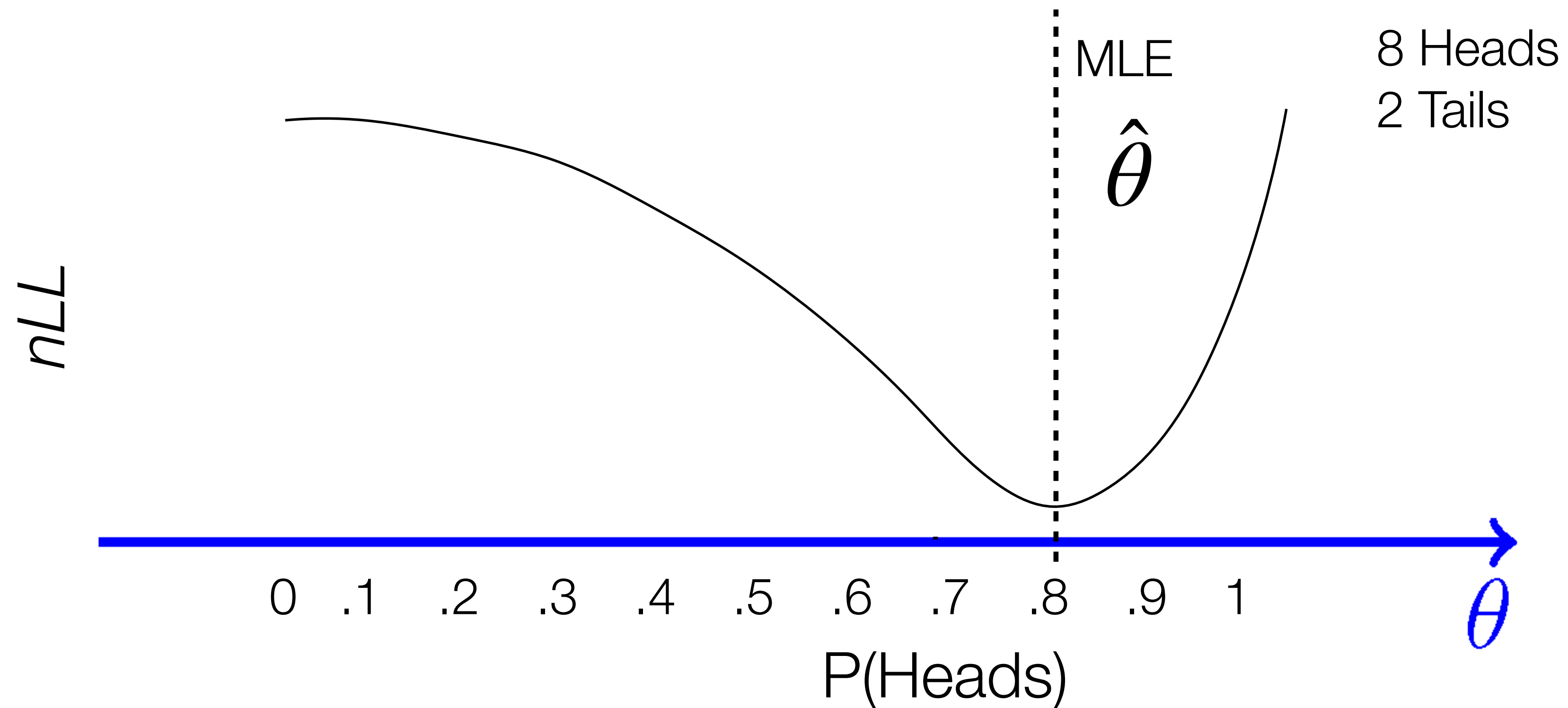
# Maximum Likelihood Estimates (MLE)

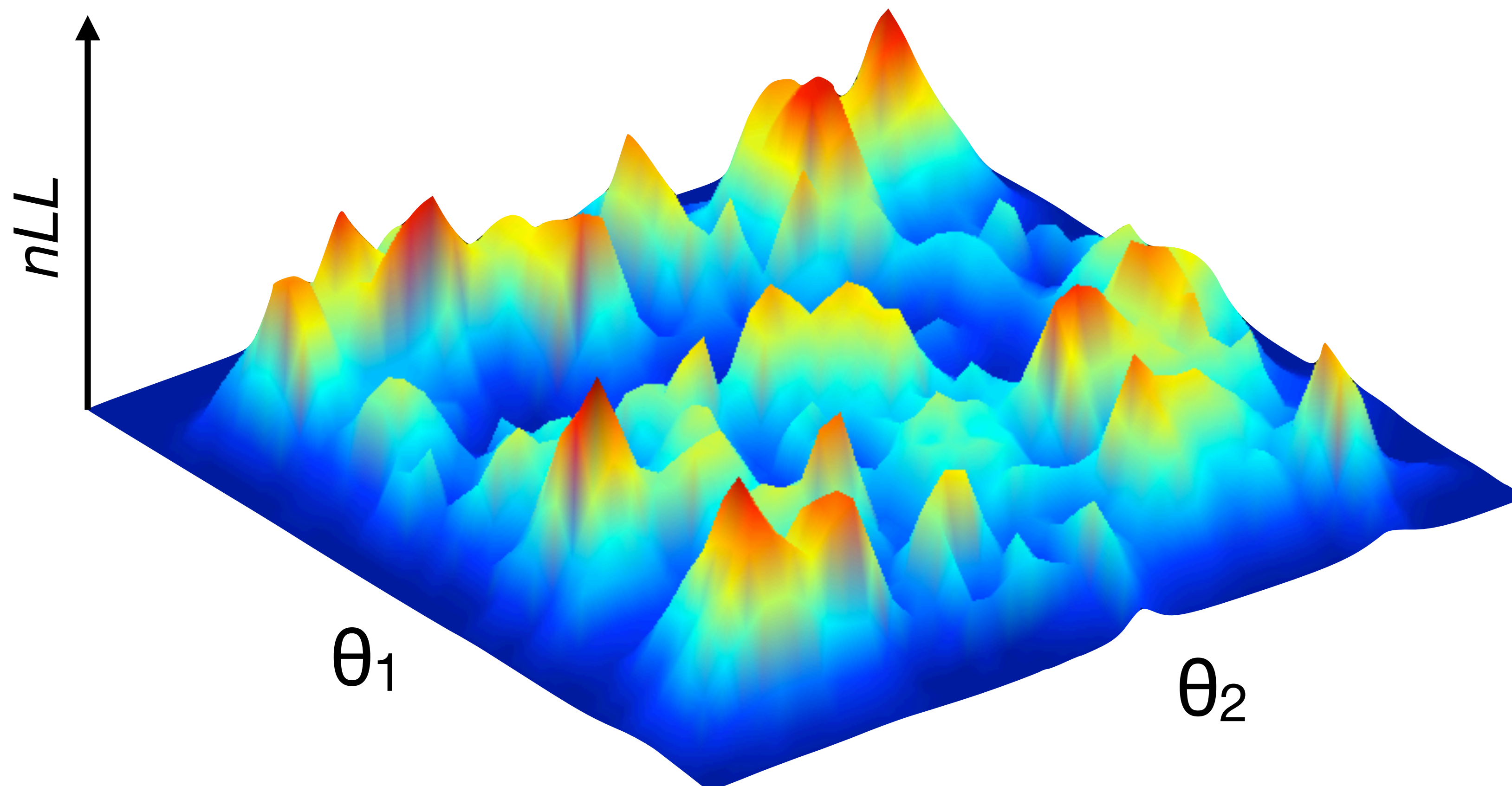Use the likelihood function to find the parameters $\hat{\theta}$ where $P(D|\theta)$ is largest

# Maximum Likelihood Estimates (MLE)

Use the likelihood function to find the parameters $\hat{\theta}$ where $P(D\,|\,\theta)$ is largest

8 Heads
2 Tails

$P(D|\theta)$

0   .1   .2   .3   .4   .5   .6   .7   .8   .9   1

P(Heads)

$\theta$

# Maximum Likelihood Estimates (MLE)

Use the likelihood function to find the parameters $\hat{\theta}$ where $P(D|\theta)$ is largest

# Maximum Likelihood Estimates (MLE)

Use the likelihood function to find the parameters $\hat{\theta}$ where $P(D|\theta)$ is largest

…. or where nLL is lowest



MLE

$\hat{\theta}$

8 Heads
2 Tails

nLL

0  .1  .2  .3  .4  .5  .6  .7  .8  .9  1

$\theta$

P(Heads)

# Computing the MLE



$nLL$

$\theta_1$

$\theta_2$

```
likelihood <-
function(params
, data)
```

Minimize nLL

Types of optimization algorithms

- Gradient descent

- Simplex methods

- Differential evolution

18

# MLE for a RL model



Simulated data

Dots: Choice (observable)

Background:
Choice probability
(unobservable latent state)

# MLE for a RL model
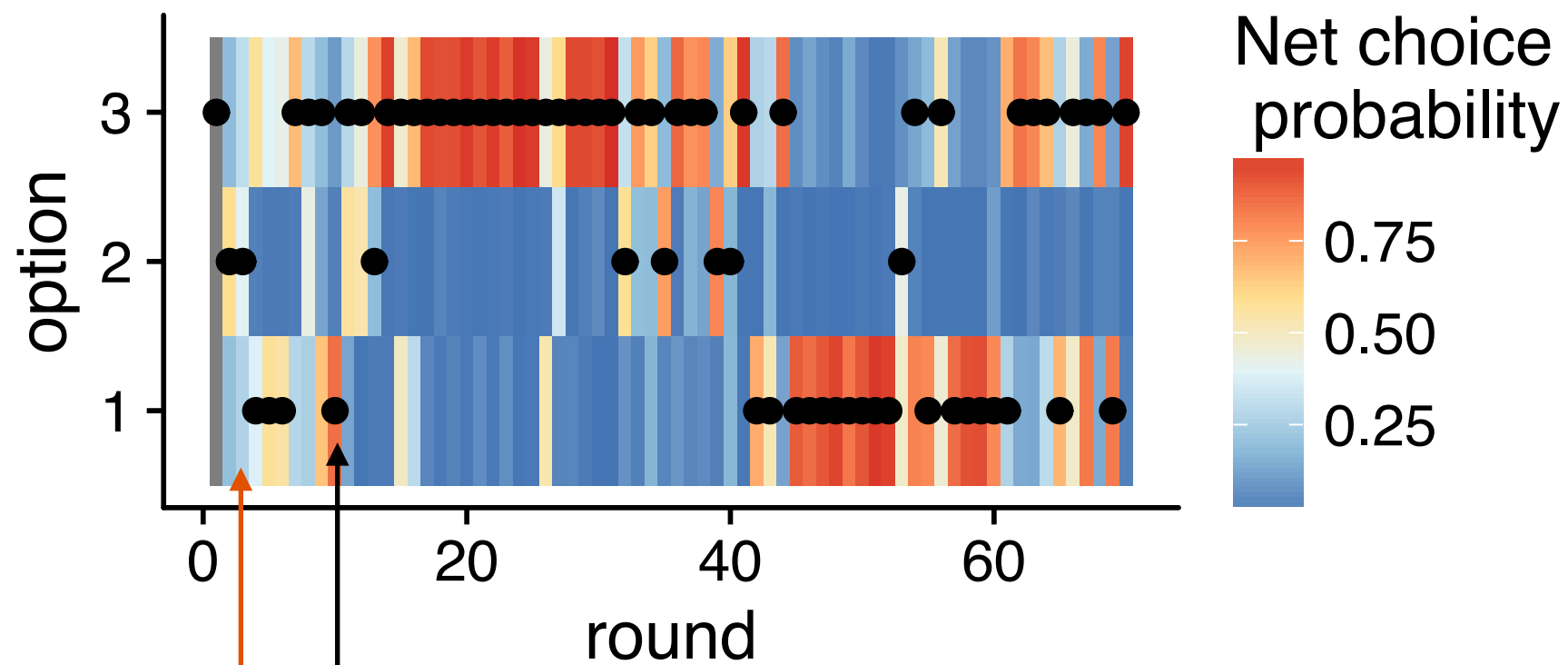
## Simulated data



Dots: Choice (observable)

Background:
Choice probability
(unobservable latent state)

Choice data from which we can calculate likelihood of the RL model

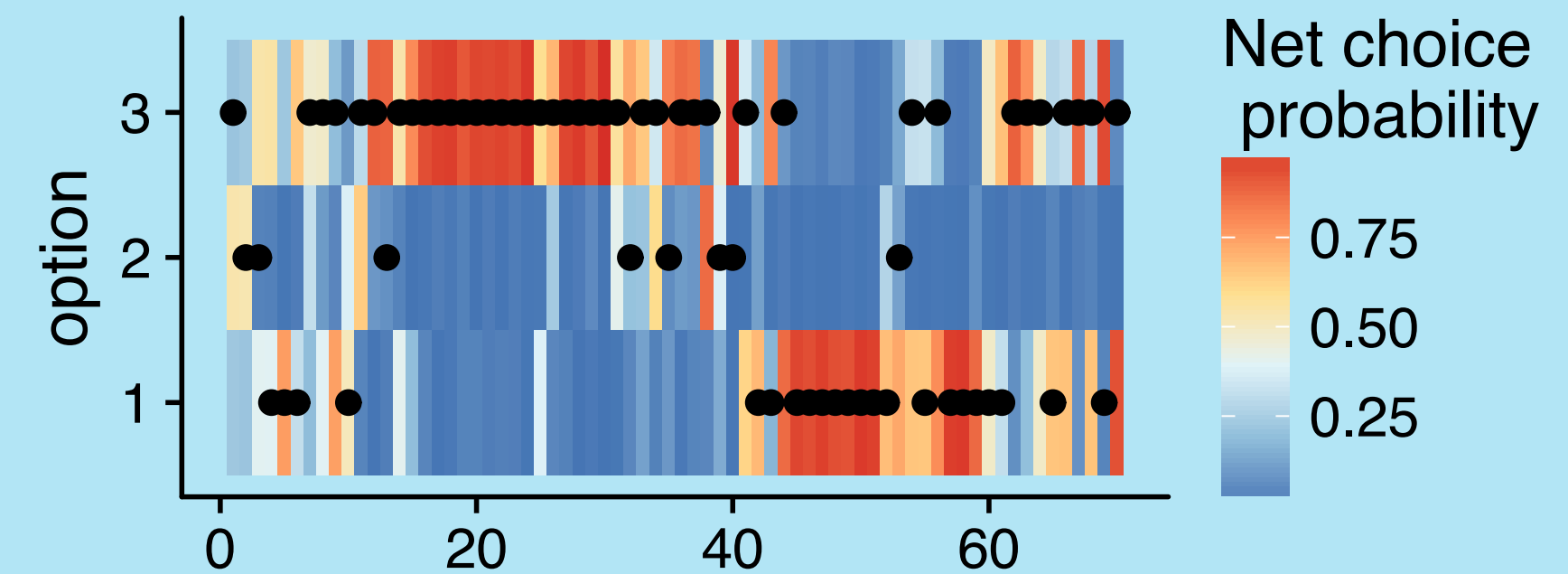# MLE for a RL model



Simulated data

Net choice probability

0.75
0.50
0.25

Dots: Choice (observable)

Background:
Choice probability
(unobservable latent state)

Search for parameters that maximise likelihood

Choice data from which we can calculate likelihood of the RL model

Estimated choice probability

Net choice probability

0.75
0.50
0.25

Maximum Likelihood Estimate (MLE)

# Q-learning agent



Loss landscape

# Parameter Space and Model Space

**Parameter Space**

**Model Space**

# Learning from social information

(5 minute break)

# Learning from social information



t = 2

t = 1

Learning from experience

23

# Learning from social information



t = 2

t = 1

Learning from experience

Social influence

# Imitating actions

Frequency-dependent copying

Probability of
choosing option a        $\propto$        frequency of other agents
                                          performing the same action

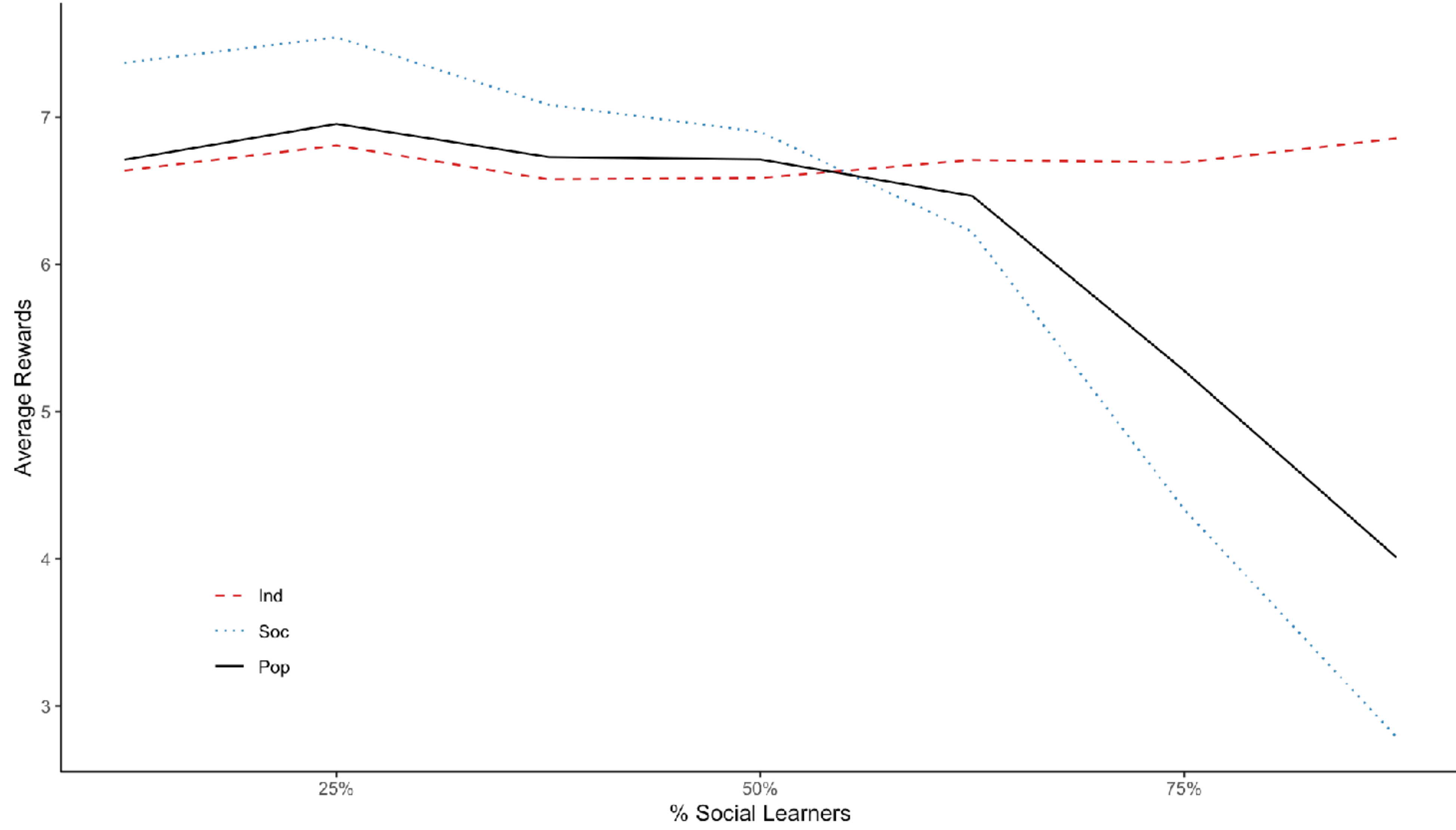$\pi_{\text{FDC}}(a)$                    $$\frac{f(a)^{\theta}}{\sum f(k)^{\theta}}$$
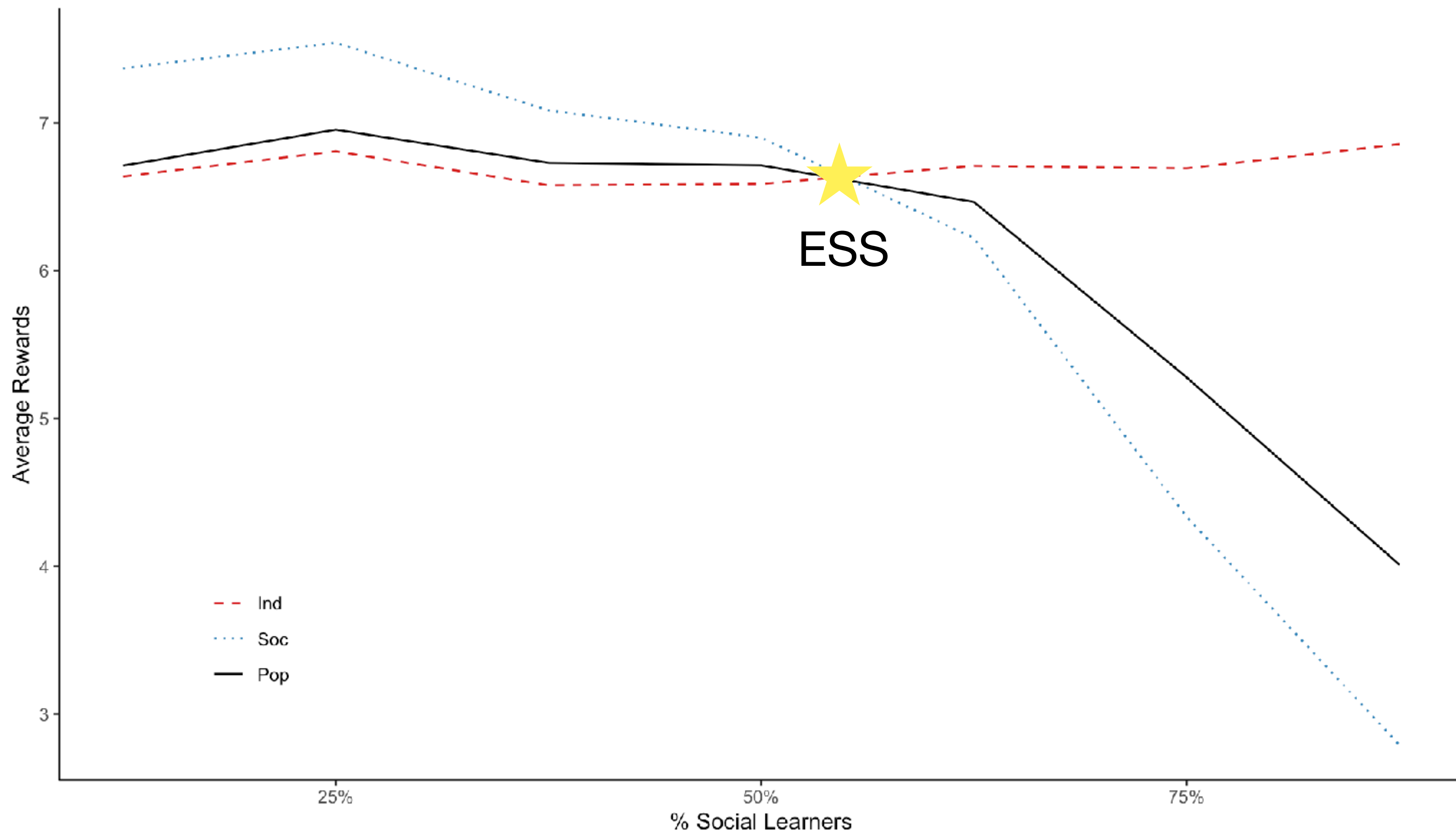
# Demo 2: Imitation and Rogers' paradox



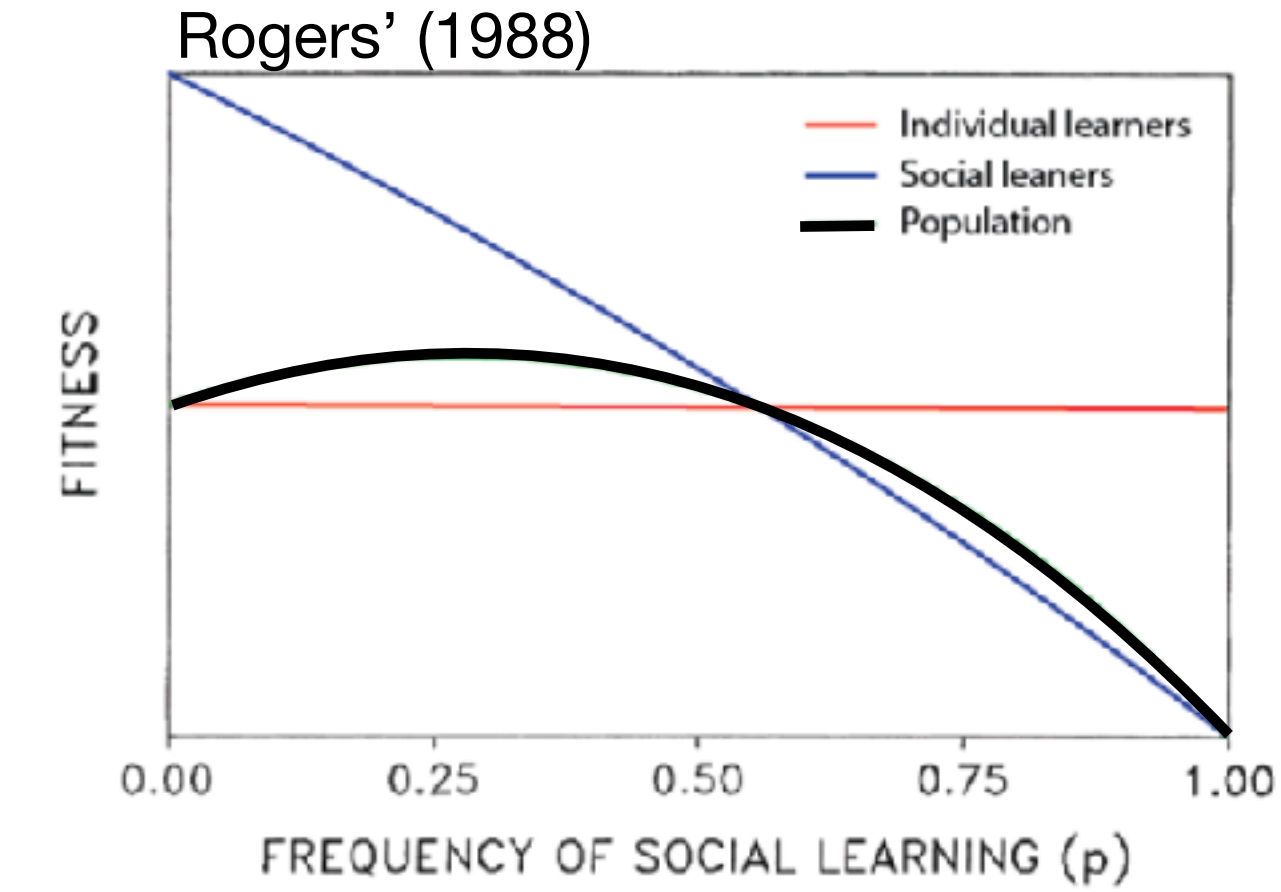How do different ratios of individual vs. social learners change the performance of each agent type?

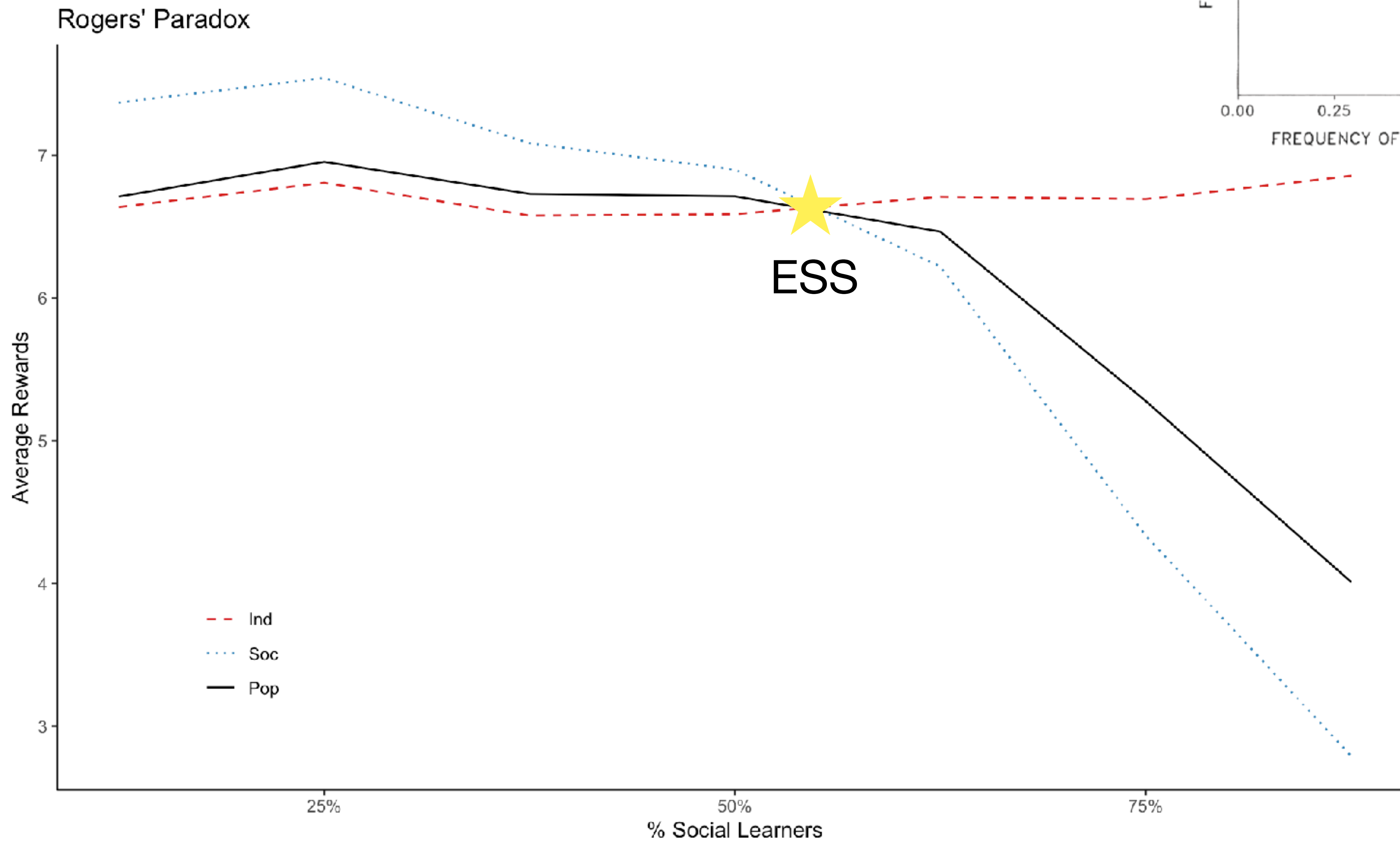Rogers' Paradox

Average Rewards (y-axis), % Social Learners (x-axis)

Legend: Ind (red dashed), Soc (blue dotted), Pop (black solid)

Rogers' Paradox

ESS

Average Rewards

% Social Learners

Ind
Soc
Pop

Rogers' Paradox

Rogers' (1988)

ESS

Average Rewards

% Social Learners

- - - Ind
....... Soc
—— Pop

# Combining imitation and value-learning

# Decision-biasing social influence

individual    social
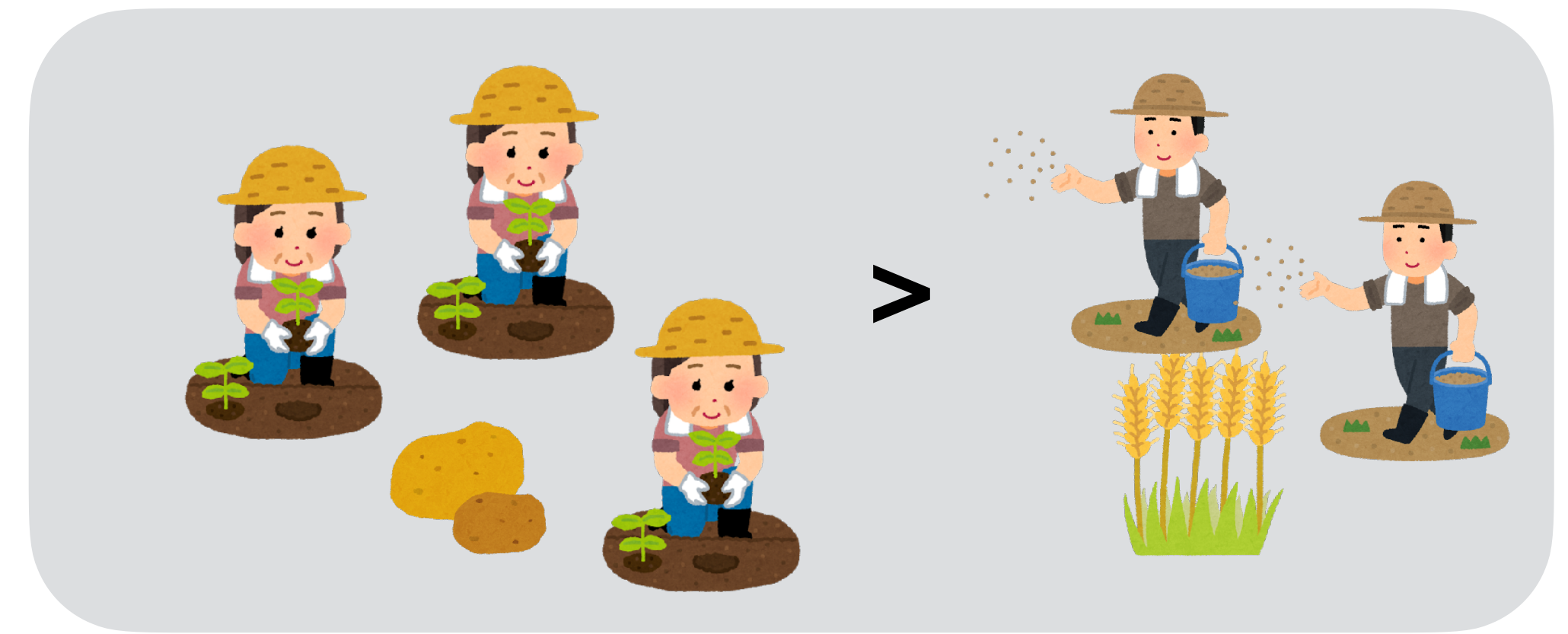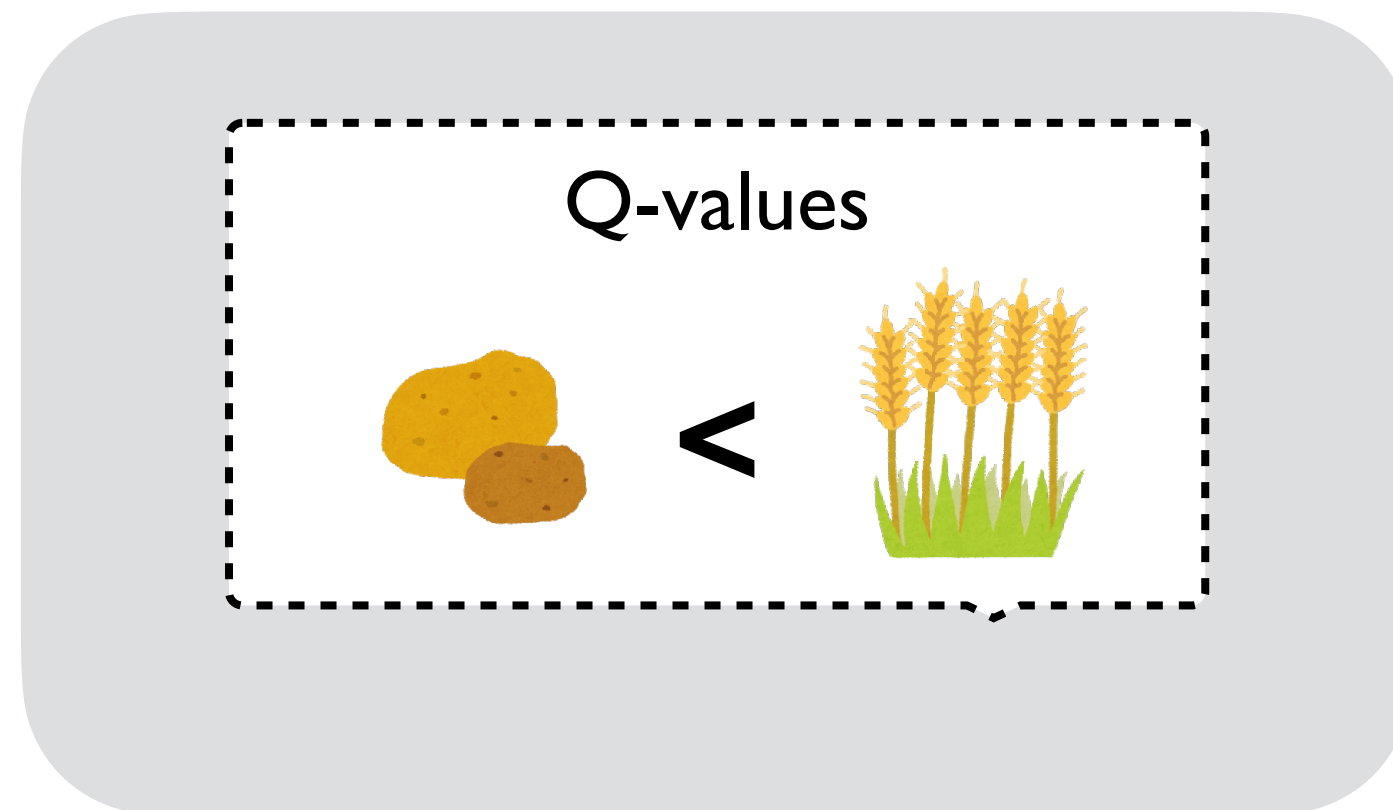
Choice probability at t = (1 - γ) Softmax + γ FDC

**Individual Q-learning**
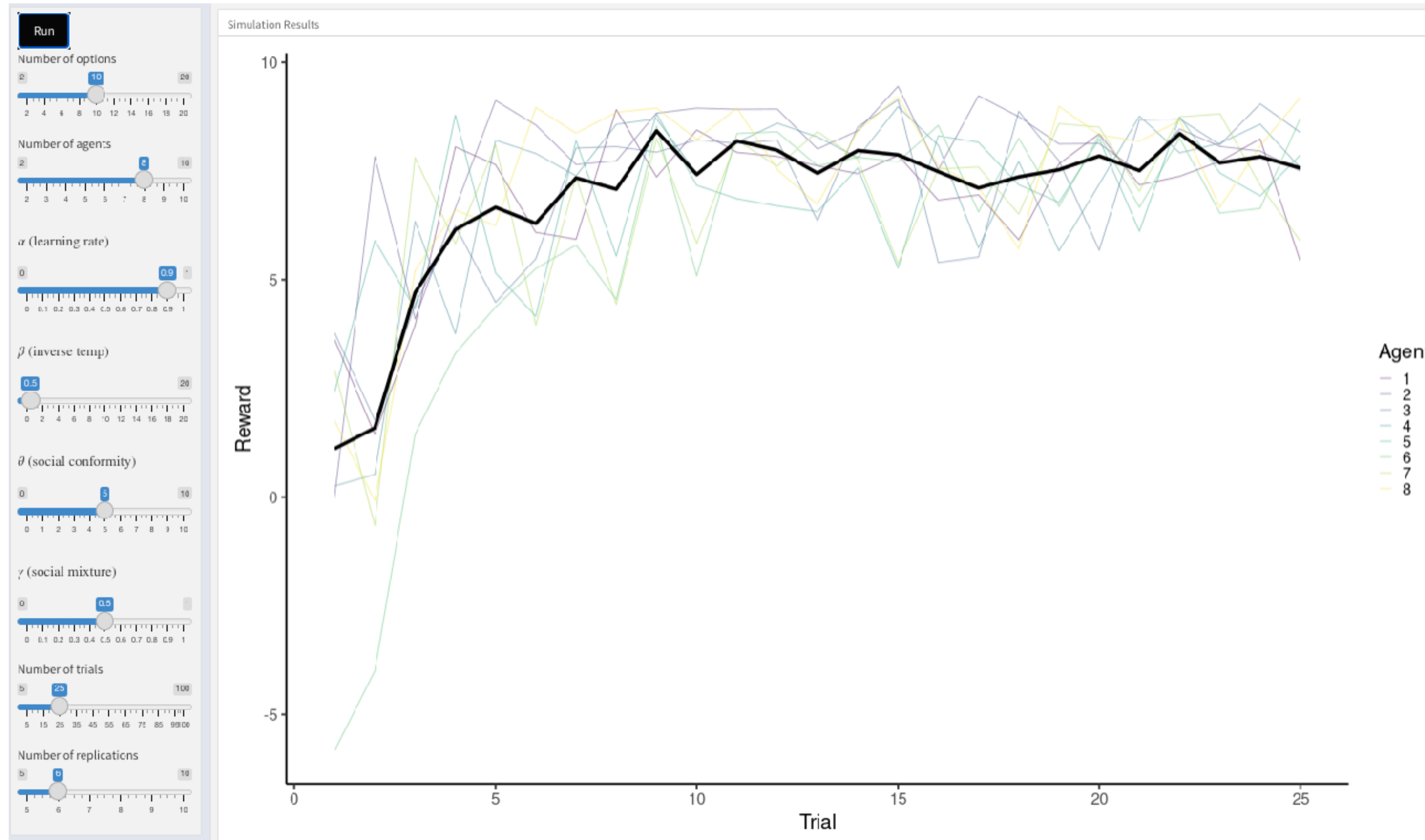
**1 - γ**

**Social frequency influence**

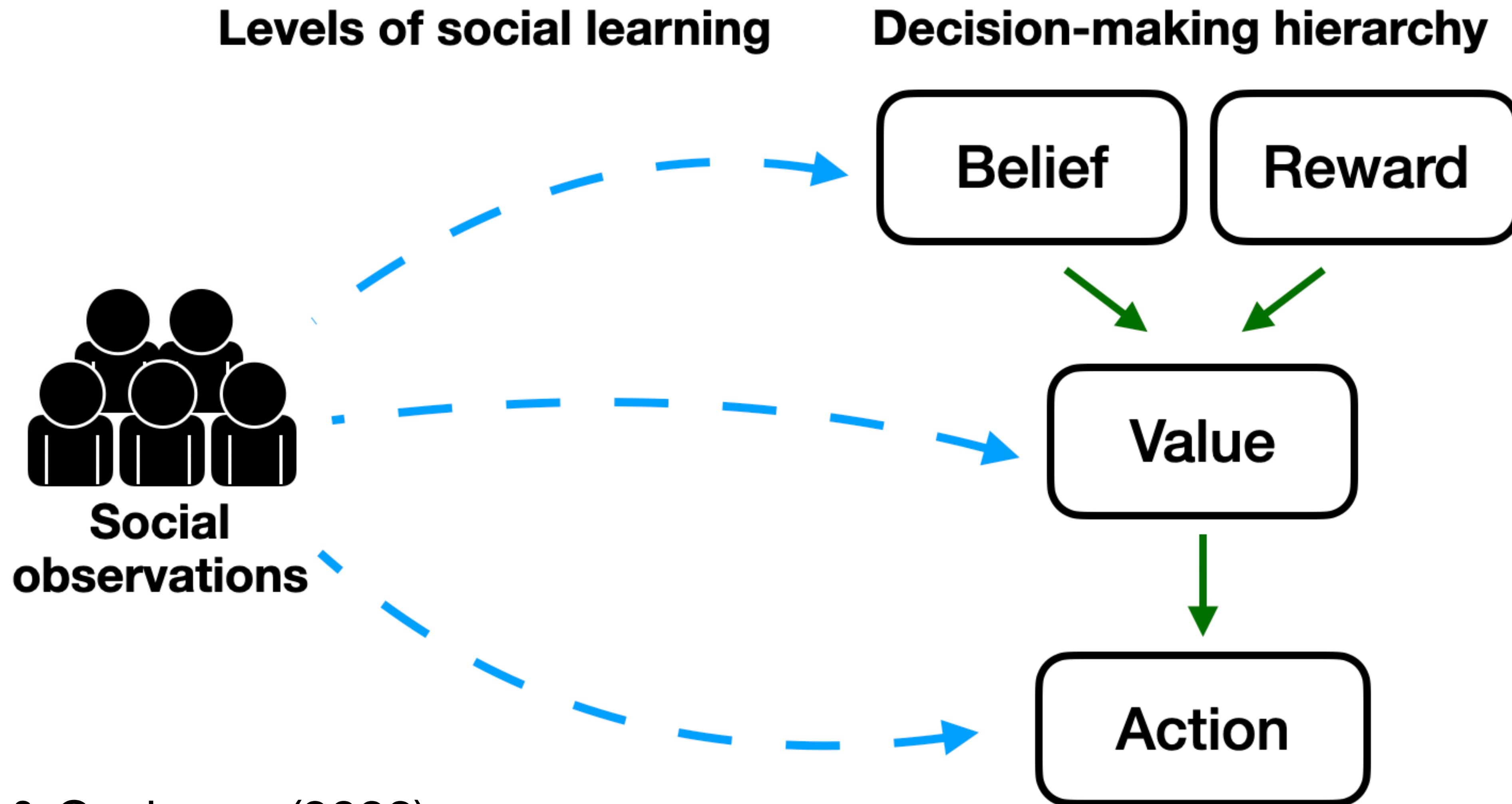**γ**

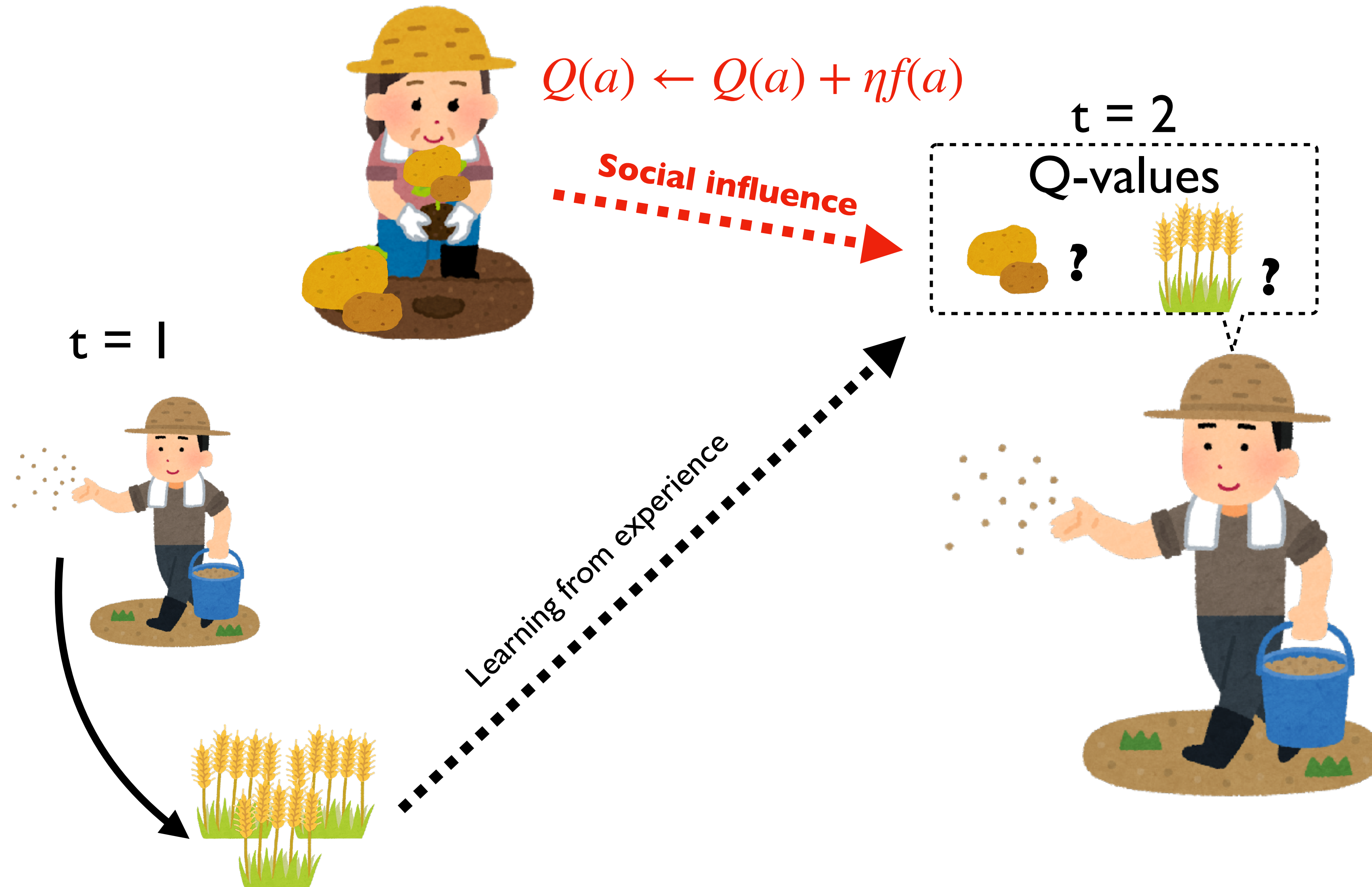$$\frac{f(i)^{\theta}}{\sum_k f(k)^{\theta}}$$

Q-values

<

>

# Demo 3: Decision-biasing



Which values of γ (social mixture) and θ (conformity exponent) typically produce the best results?

# Social influence at different levels of learning



Levels of social learning

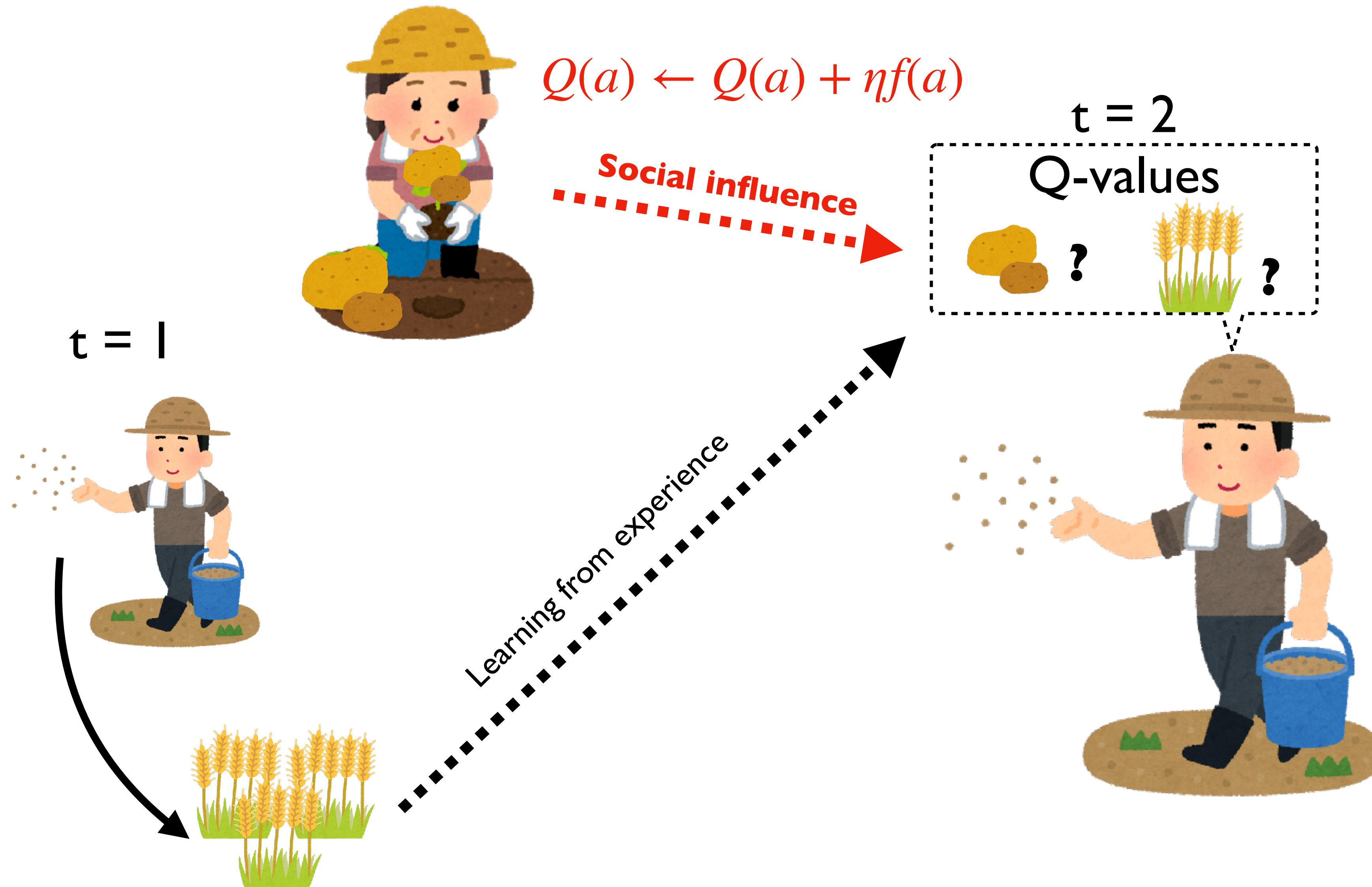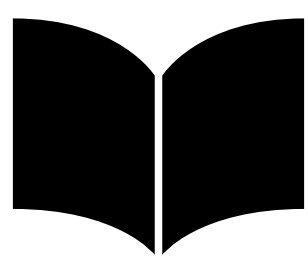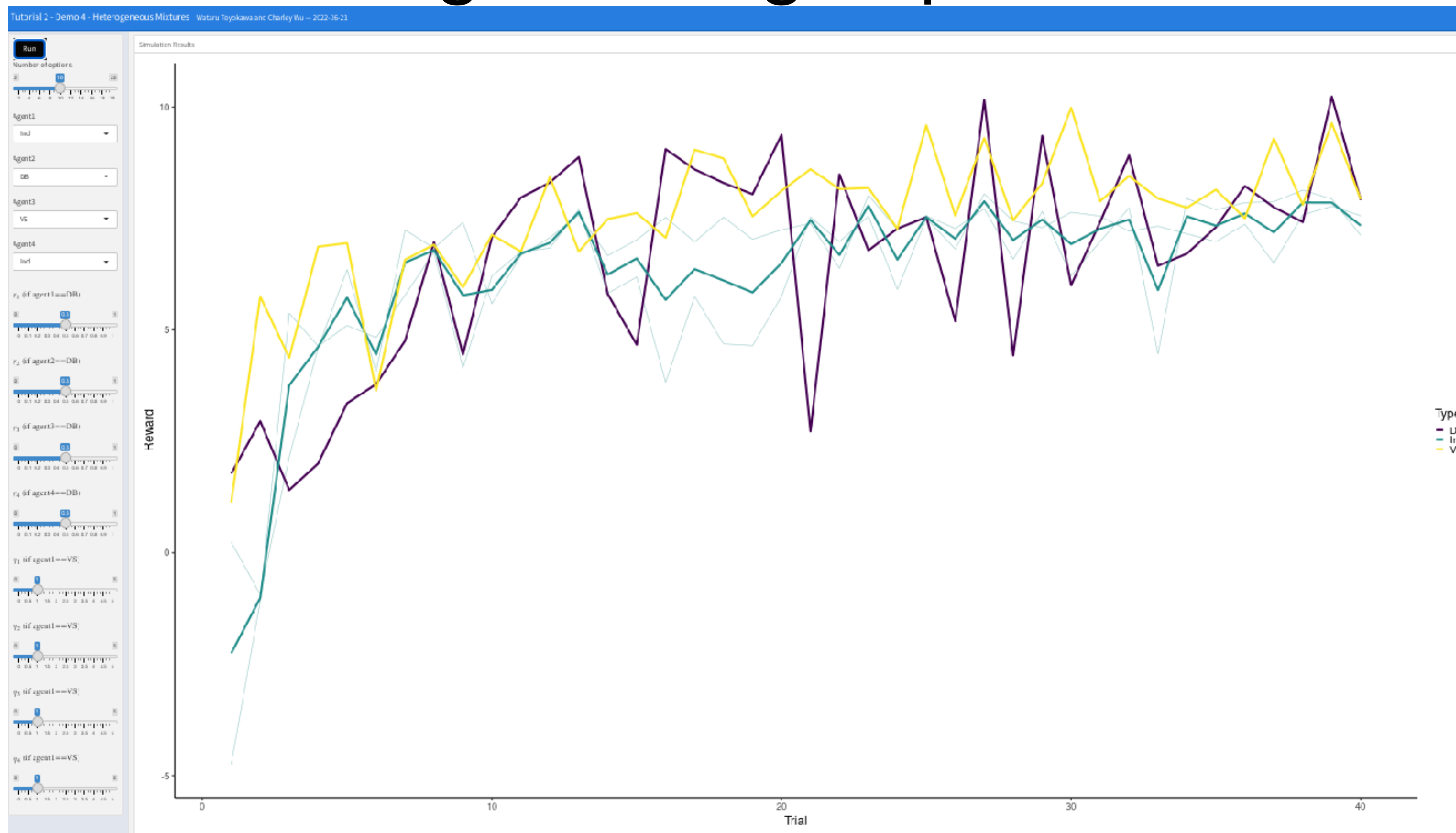Decision-making hierarchy

Wu, Vélez, & Cushman (2022)

# Value-shaping social influence

# Value-shaping social influence

value bonus

$$Q(a) \leftarrow Q(a) + \eta f(a)$$

Social influence

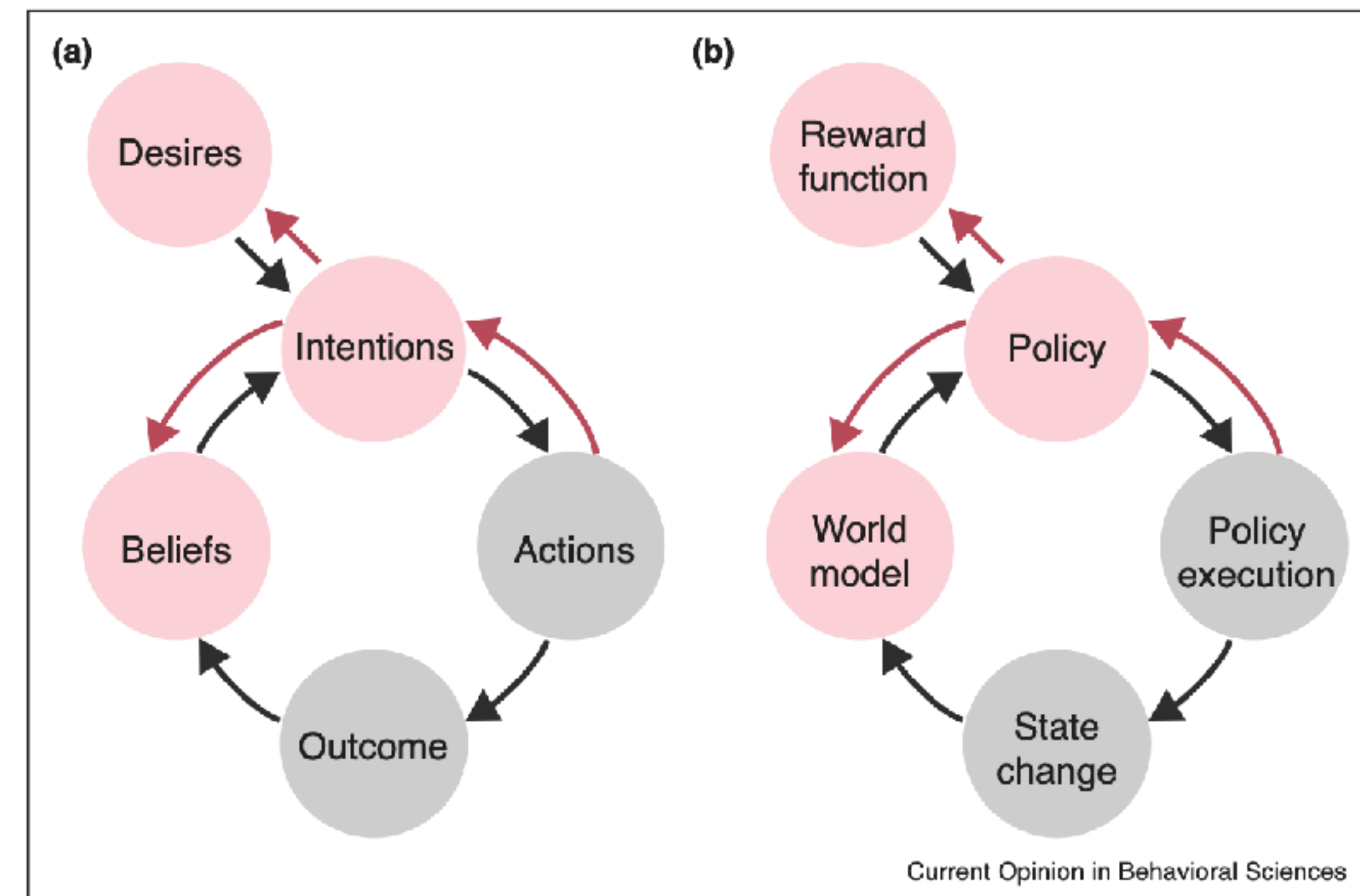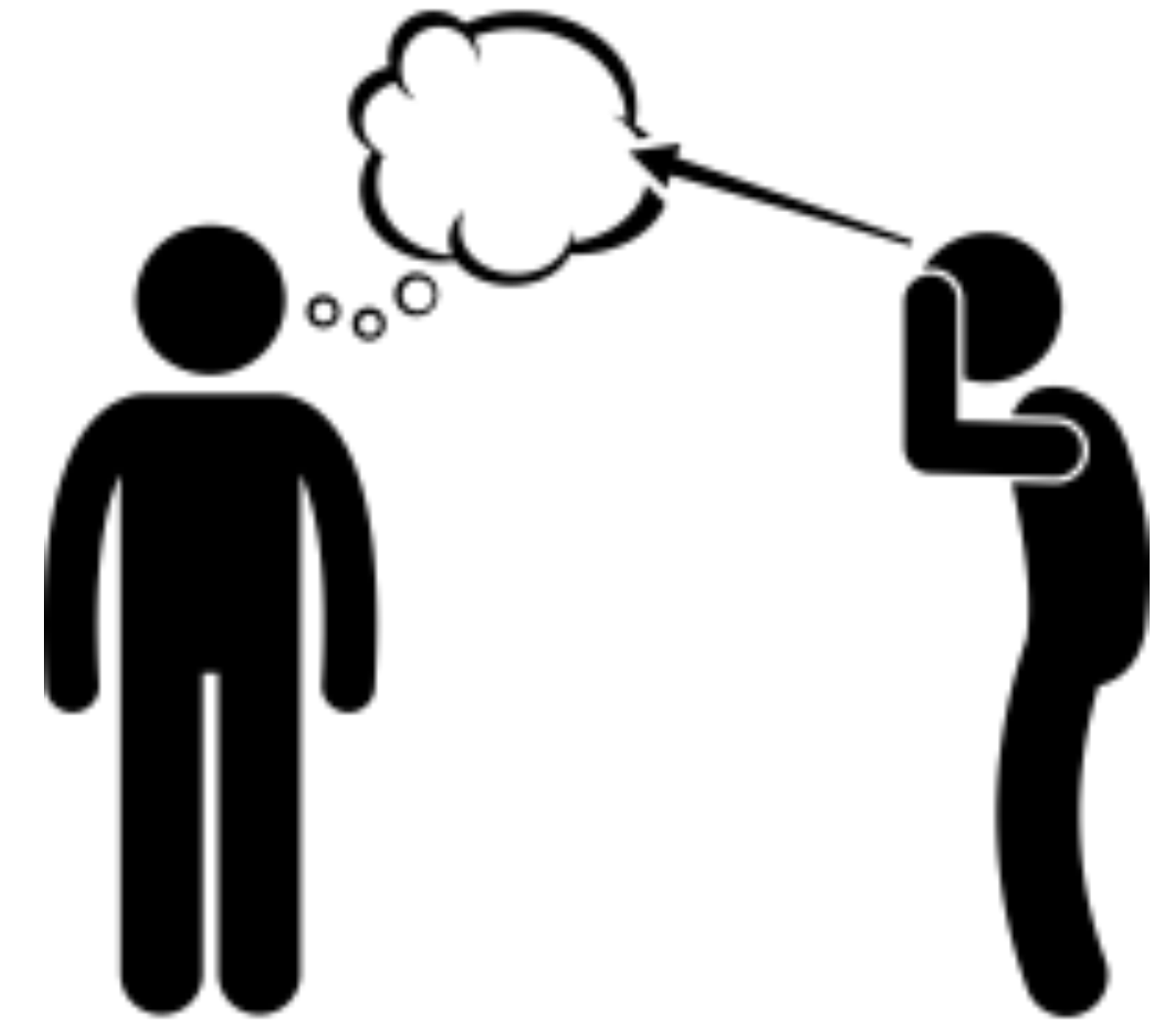t = 2

Q-values

? ?

t = 1

Learning from experience

# Demo 4: Heterogeneous groups



Which strategy perform better than others? Is it robust to different group compositions?
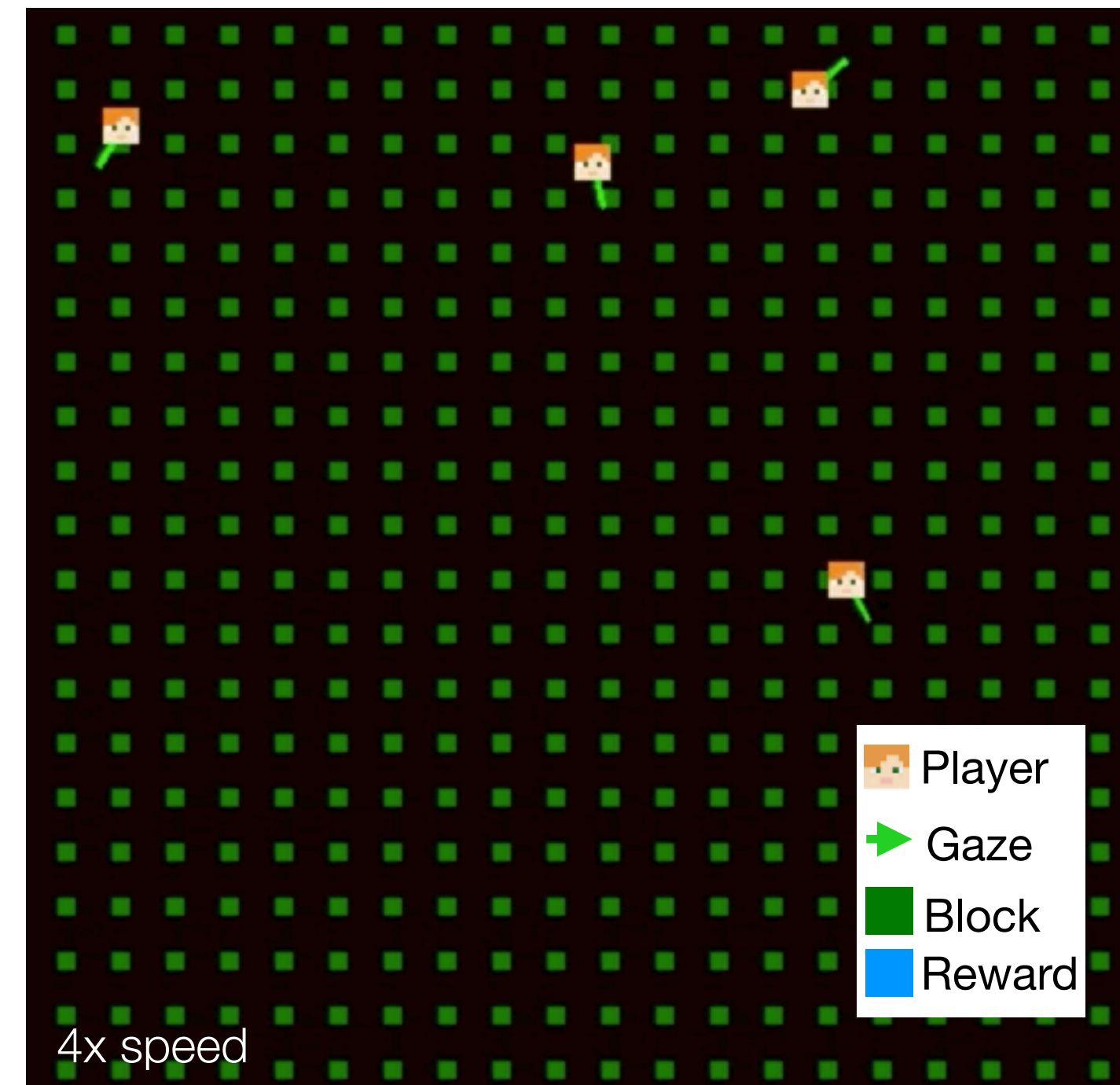
# Theory of Mind



- So far we have described very simple social learning mechanism

- Yet an important aspect of human social learning is our ability to "unpack" observed actions into imputed mental states

  - desires and intentions

  - beliefs and one's model of the world

- This is known as Theory of Mind (ToM) inference



(a) Desires → Intentions → Actions → Outcome → Beliefs

(b) Reward function → Policy → Policy execution → State change → World model

Current Opinion in Behavioral Sciences

Jara-Ettinger (2019)  33

# Scaling up to more complex tasks

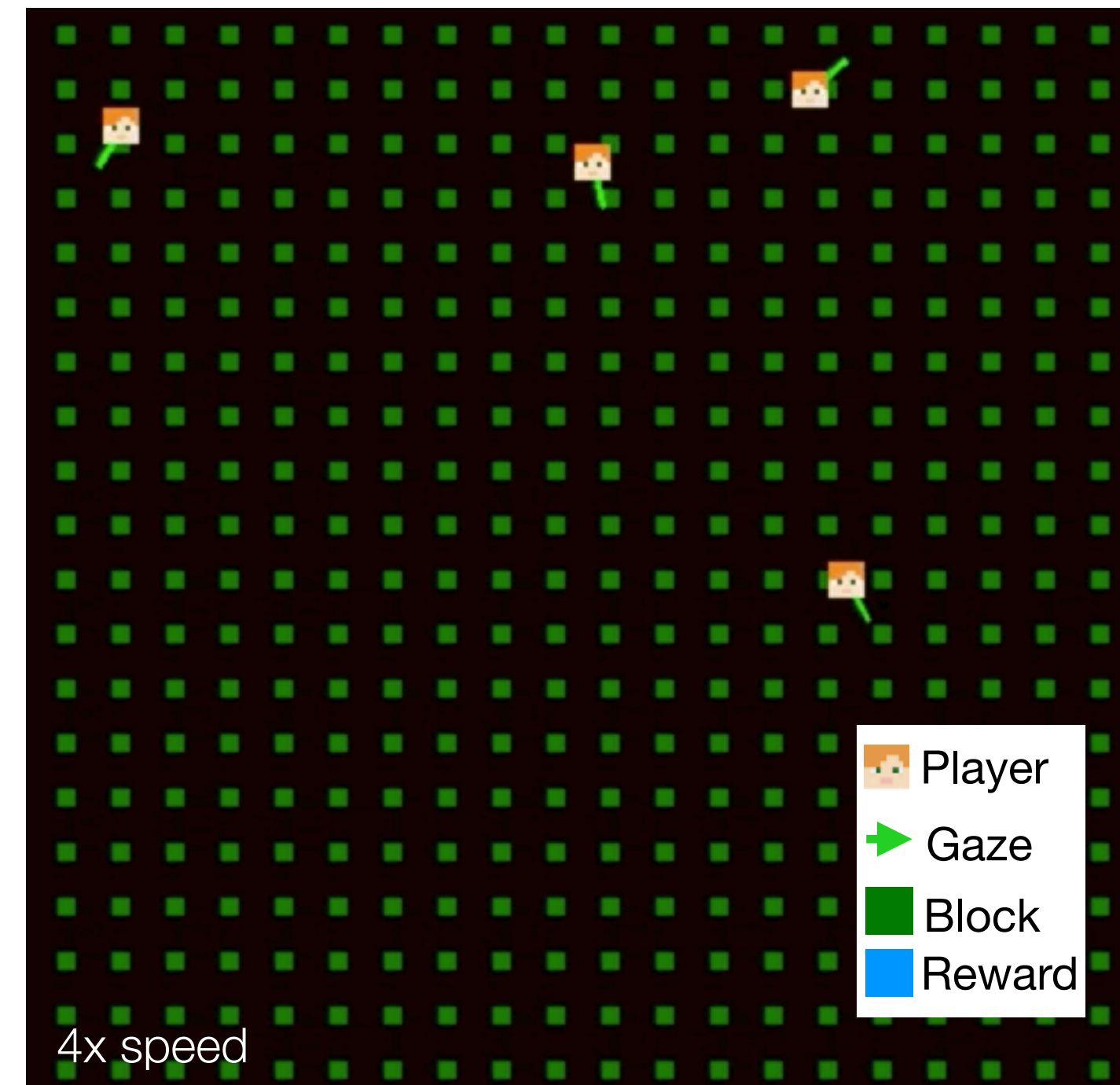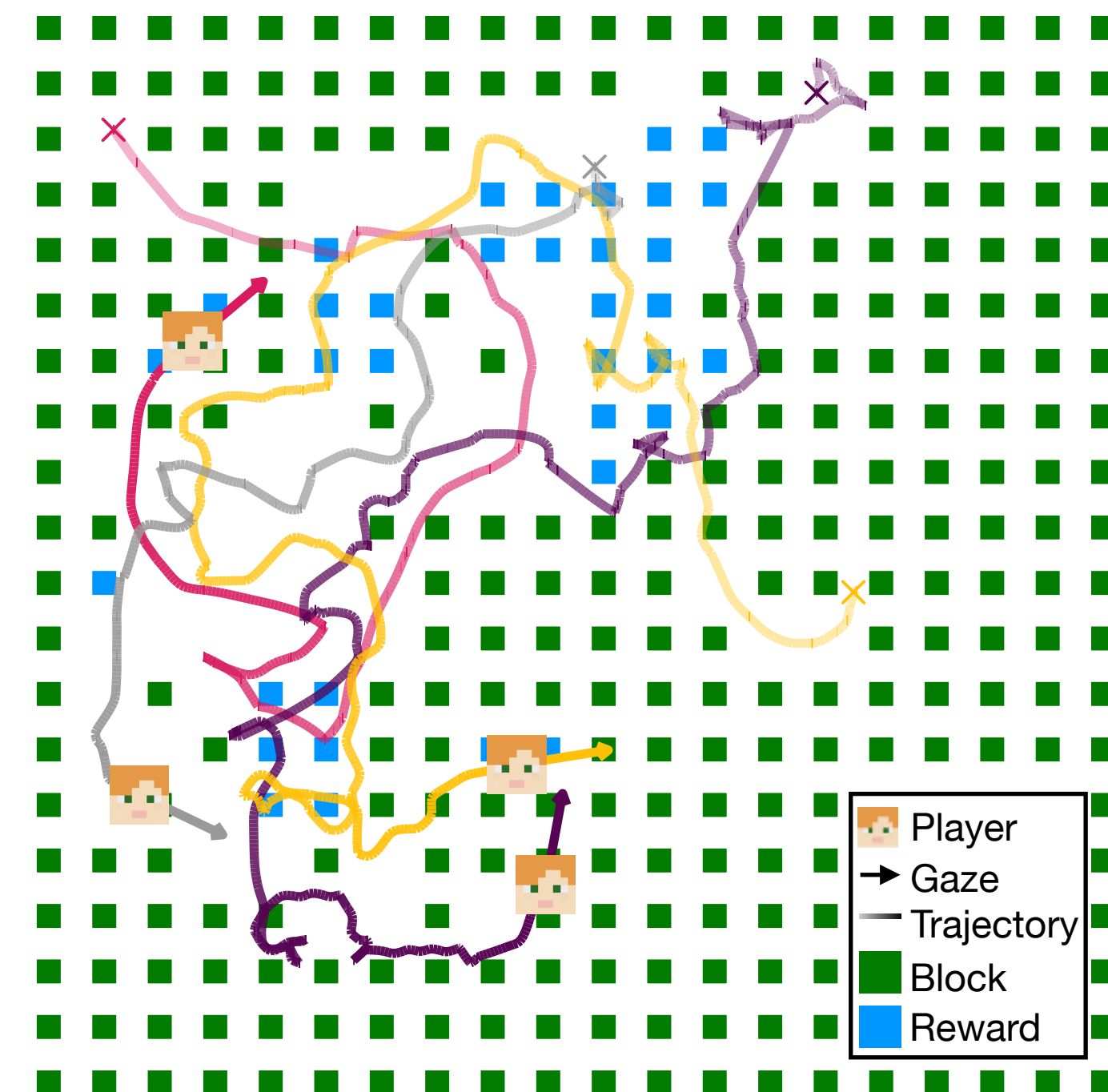Collective foraging in a dynamic and immersive virtual environment



- Participants forage for hidden rewards (blue splash) on a field of melons

- Realistic field of view creates attentional trade-offs and opportunity costs for social learning:

  - Looking at other players for social imitation comes at the cost of slower individual foraging

- Rich and dynamic social interactions through spatial position and visual gaze

# Scaling up to more complex tasks

Collective foraging in a dynamic and immersive virtual environment



Player
➤ Gaze
■ Block
■ Reward

4x speed

- Participants forage for hidden rewards (blue splash) on a field of melons

- Realistic field of view creates attentional trade-offs and opportunity costs for social learning:

  - Looking at other players for social imitation comes at the cost of slower  individual foraging

- Rich and dynamic social interactions through spatial position and visual gaze

# Scaling up to more complex tasks

Collective foraging in a dynamic and immersive virtual environment



- Participants forage for hidden rewards (blue splash) on a field of melons

- Realistic field of view creates attentional trade-offs and opportunity costs for social learning:

  - Looking at other players for social imitation comes at the cost of slower  individual foraging

- Rich and dynamic social interactions through spatial position and visual gaze

34

# Interactive Tutorial



2x speed

- Smash blocks by clicking and holding mouse (2.25 seconds)

- Some blocks contain rewards, indicated by a blue splash, visible to other players

- Other blocks have no reward

- Participants incentivized to collect as many rewards as possible

# Interactive Tutorial



2x speed

- Smash blocks by clicking and holding mouse (2.25 seconds)

- Some blocks contain rewards, indicated by a blue splash, visible to other players

- Other blocks have no reward

- Participants incentivized to collect as many rewards as possible

# Experimental Design

**Smooth**



**Random**





16 rounds with a 2x2 within-subject design

- Environment: smooth vs. random

- Condition: solo vs. groups of four

# Computational models



$$k-1 \qquad k \qquad k+1$$

Sequentially predict each of the $k$ blocks participants destroy:

$$P(\text{Choice}_{k+1}) \propto \exp(\mathbf{f}_k \cdot \mathbf{w})$$

using a softmax over a set of features $\mathbf{f}$ times weights $\mathbf{w}$

Model **f**eatures capture hypotheses about individual and social learning mechanisms (details on next slide)

Model **w**eights are estimated using hierarchical Bayesian methods in STAN with individual and group as random effects

**a** Model illustration

$P(\text{choice}_{k+1}) \propto \exp(\mathbf{f}_k \cdot \mathbf{w})$

$k-1 \qquad k \qquad k+1$

Legend:
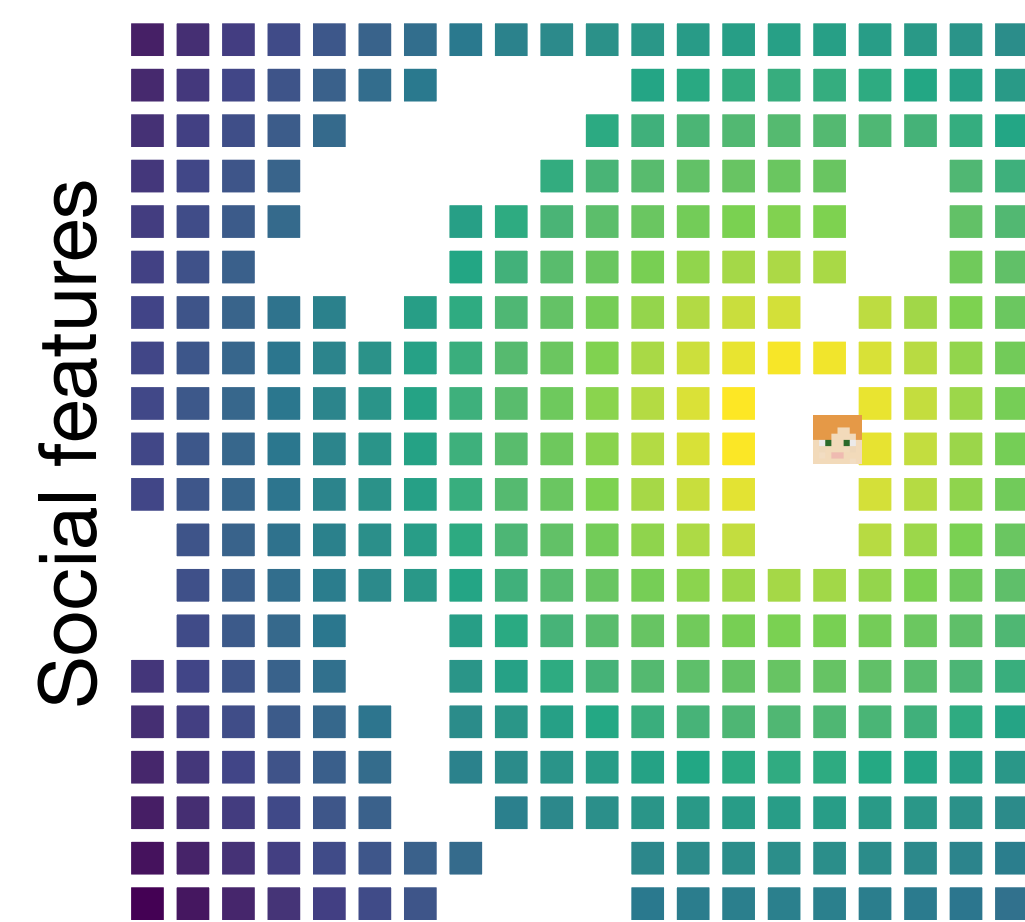- Player
- Gaze
- Trajectory
- Block
- Reward
- Target

BlockVis

Locality

GP Pred

Asocial features

$f(\text{Block})$

Ind. Obs.
- Reward
- No reward

Successful Prox.

Unsuccessful Prox.

Model Predictions

Social features

$p(\text{Choice})$
$\times$ Choice$_{k+1}$

BlockVis

Locality

GP pred

38

**a** Model illustration

BlockVis    Locality    GP Pred

Asocial features

$f(\text{Block})$

Ind. Obs.
Reward
No reward

Player
Gaze
Trajectory
Block
Reward
Target

$P(\text{choice}_{k+1})$
$\propto$
$\exp(\mathbf{f}_k \cdot \mathbf{w})$

$k-1$    $k$    $k+1$

Successful Prox.    Unsuccessful Prox.    Model Predictions

Social features

$p(\text{Choice})$
Choice$_{k+1}$
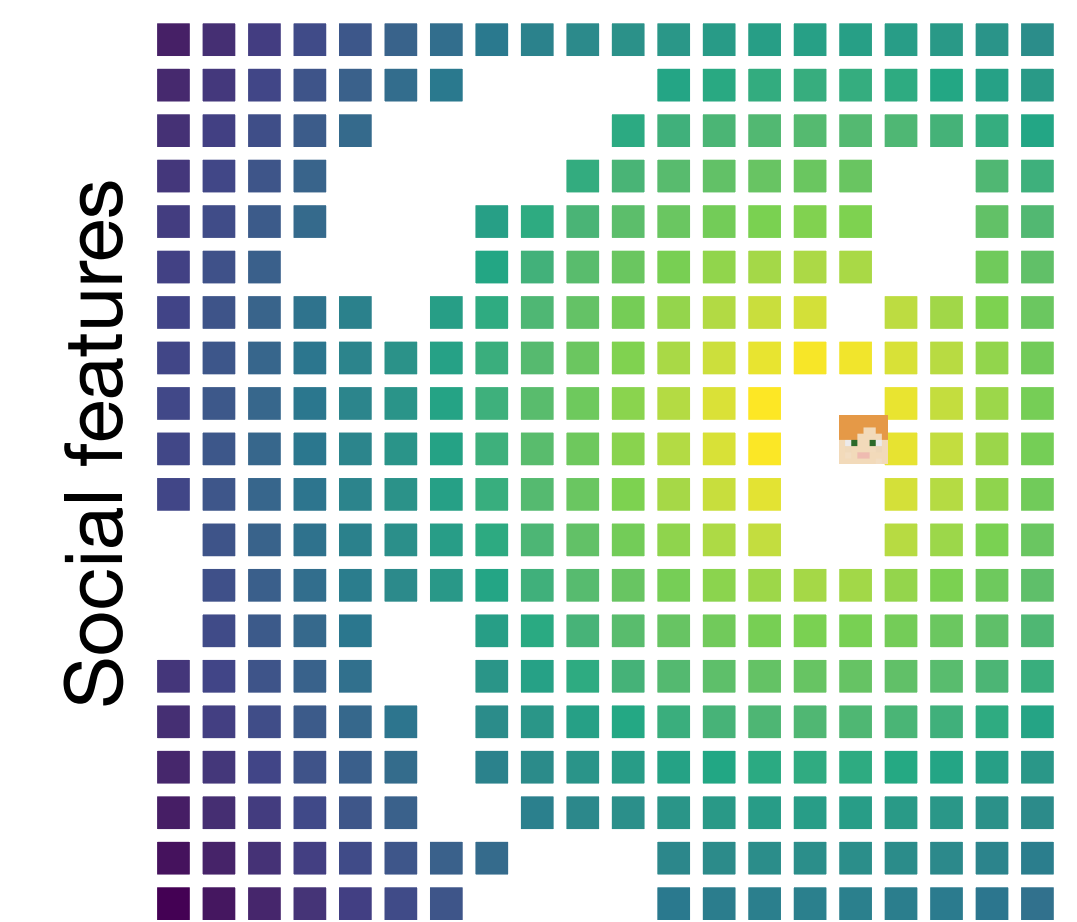
Gaussian Process (GP) asocial
RL with reward generalization

Spatial    Conceptual    Graph-structured

Stripes    Tilt

Wu, et al.,
(*NHB* 2018)    Wu, et al., (*PLOS CompBio* 2020)    Wu, Schulz & Gershman (*CBB* 2021)

Observation    Expected Reward
Hypothesis    Uncertainty

Reward

Input
(Spatial location, arbitrary features, node id)

Locality    GP pred

38

**a** Model illustration

BlockVis          Locality          GP Pred

Asocial features

$f(\text{Block})$

Ind. Obs.
Reward
No reward

Successful Prox.          Unsuccessful Prox.          Model Predictions

Social features

Player
Gaze
Trajectory
Block
Reward
Target

$P(\text{choice}_{k+1})$
$\propto$
$\exp(\mathbf{f}_k \cdot \mathbf{w})$

$k-1$          $k$          $k+1$

$p(\text{Choice})$
Choice$_{k+1}$

Gaussian Process (GP) asocial
RL with reward generalization

Spatial          Conceptual          Graph-structured

Stripes     Tilt

Wu, et al.,
(*NHB* 2018)

Wu, et al., (*PLOS
CompBio* 2020)

Wu, Schulz &
Gershman (*CBB* 2021)

$P(r=1) = \sigma(z)$

$z$

Input
(Spatial location, arbitrary features, node id)

Locality          GP pred

# a Model illustration

BlockVis   Locality   GP Pred

Asocial features

$f(\text{Block})$

Ind. Obs.
■ Reward
☐ No reward

Successful Prox.   Unsuccessful Prox.   Model Predictions

Social features

**Player**
→ Gaze
▬ Trajectory
■ Block
■ Reward
⬛ Target

$P(\text{choice}_{k+1})$
$\propto$
$\exp(\mathbf{f}_k \cdot \mathbf{w})$

$k-1$   $k$   $k+1$

$p(\text{Choice})$
✕ $\text{Choice}_{k+1}$

## Gaussian Process (GP) asocial RL with reward generalization

Spatial   Conceptual   Graph-structured

Stripes   Tilt

Wu, et al., (*NHB* 2018)   Wu, et al., (*PLOS CompBio 2020*)   Wu, Schulz & Gershman (*CBB* 2021)

$z$

$P(r=1) = \sigma(z)$

Input
(Spatial location, arbitrary features, node id)

## Social proximity

Locality
↓
Visible
↓
Successful   Unsuccessful

t-1
↓
t

GP pred

Centroid of last seen location

38

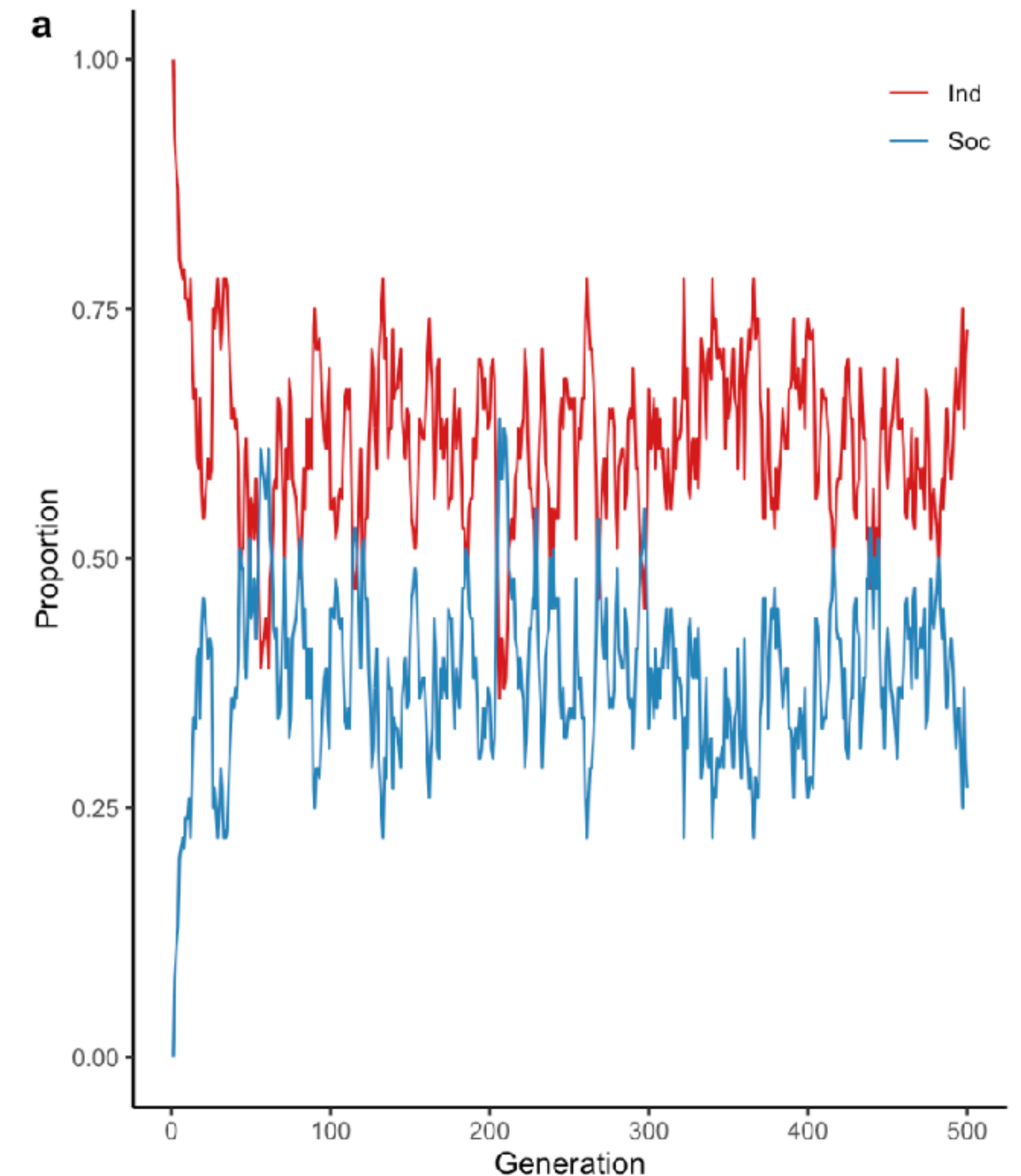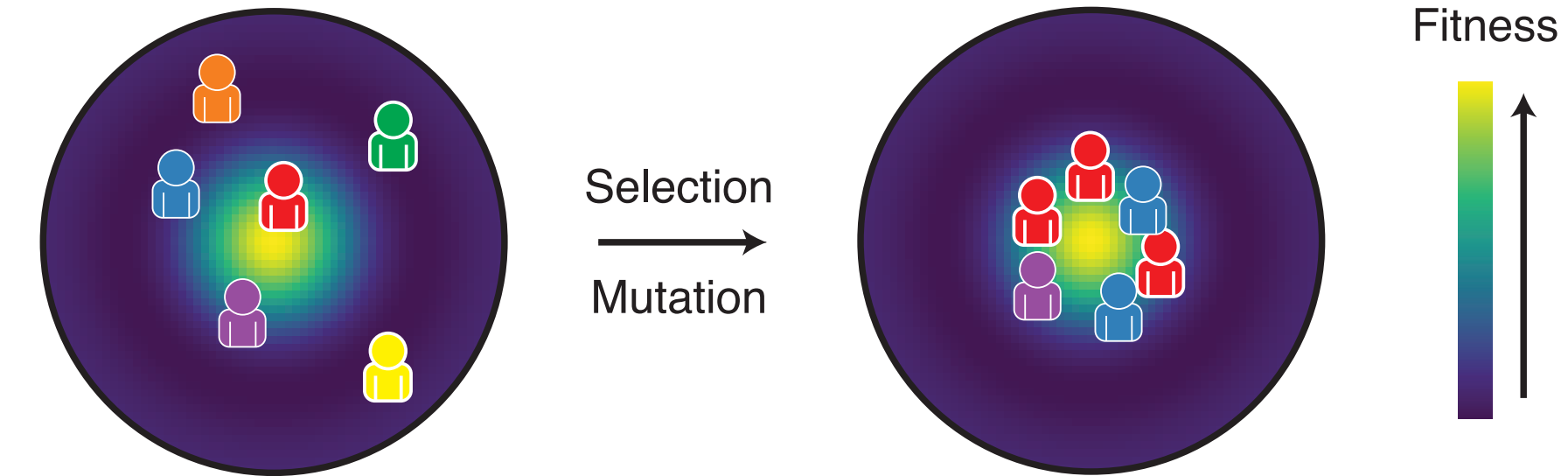# Evolutionary dynamics



- Social learning has *frequency-dependent fitness* (Rogers, 1988)

    - The best strategy to use depends on what others in the population are doing

- In order to determine the best *normative* strategy, it is often helpful to use evolutionary simulations:

    1. Initialize a population of agents

    2. Simulate performance on the task

    3. Select agents to seed the next generation (e.g., based on performance)

    4. Add mutation (change agent type, modify parameters)

    5. Repeat until convergence

# Evolutionary simulations

- Social learning despite individual differences (Witt et al., 2023)

  - People can use social information, but not verbatim

  - Exact imitation strategies might fail to account for social differences

- Decision-Biasing (**DB**)

- Value-Shaping (**VS**)

- Social Generalization (**SG**):

  - integration social info in the reward generalization process

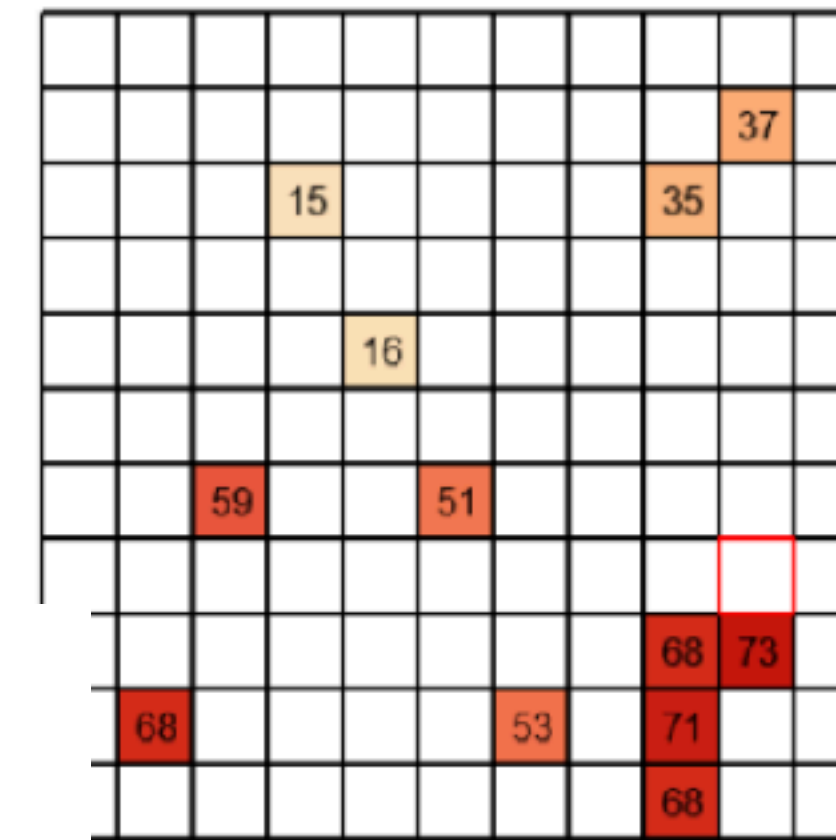  - assume social info is noisier than individual experiences



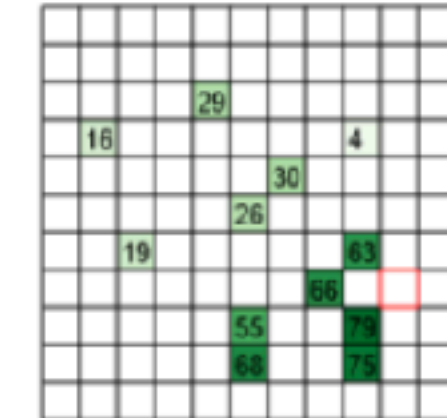Gather as much salt as possible within 14 clicks
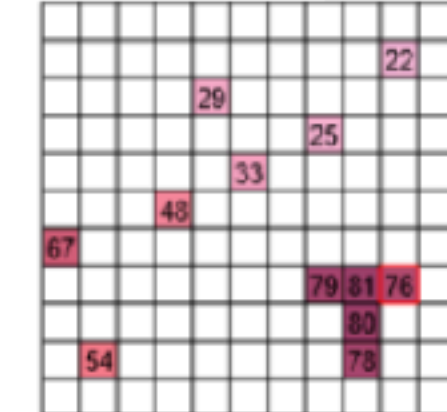
Salt concentration is correlated spatially...

.... as well as socially

# Evolutionary simulations
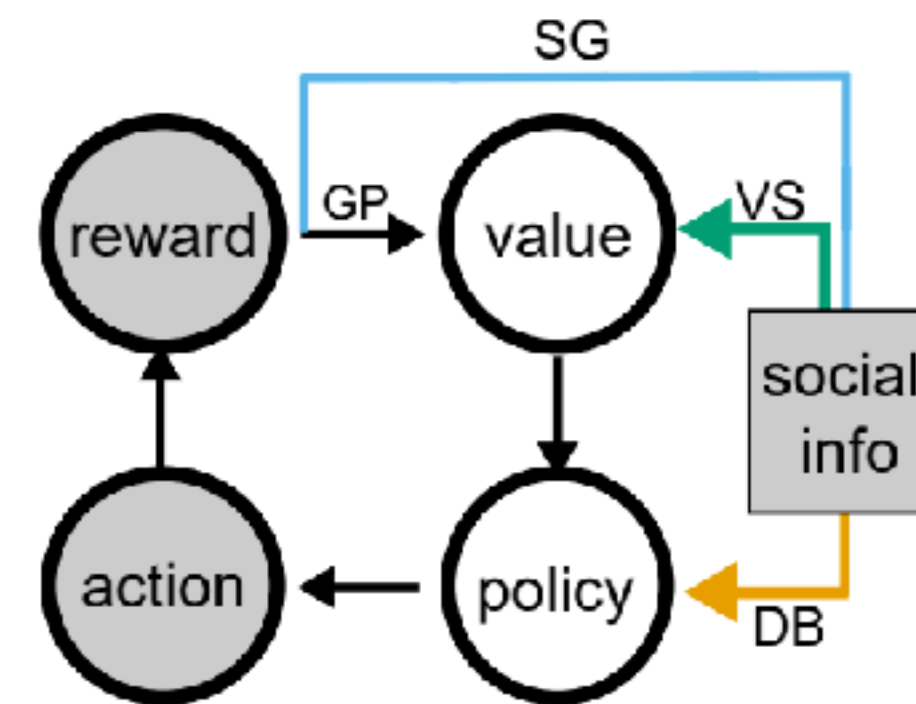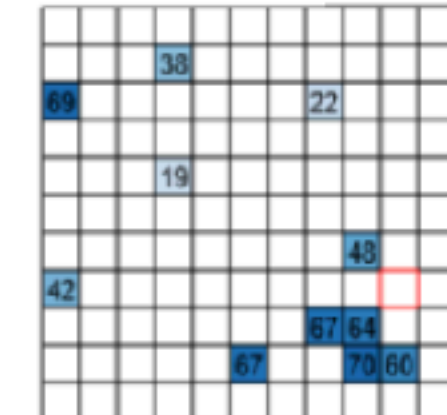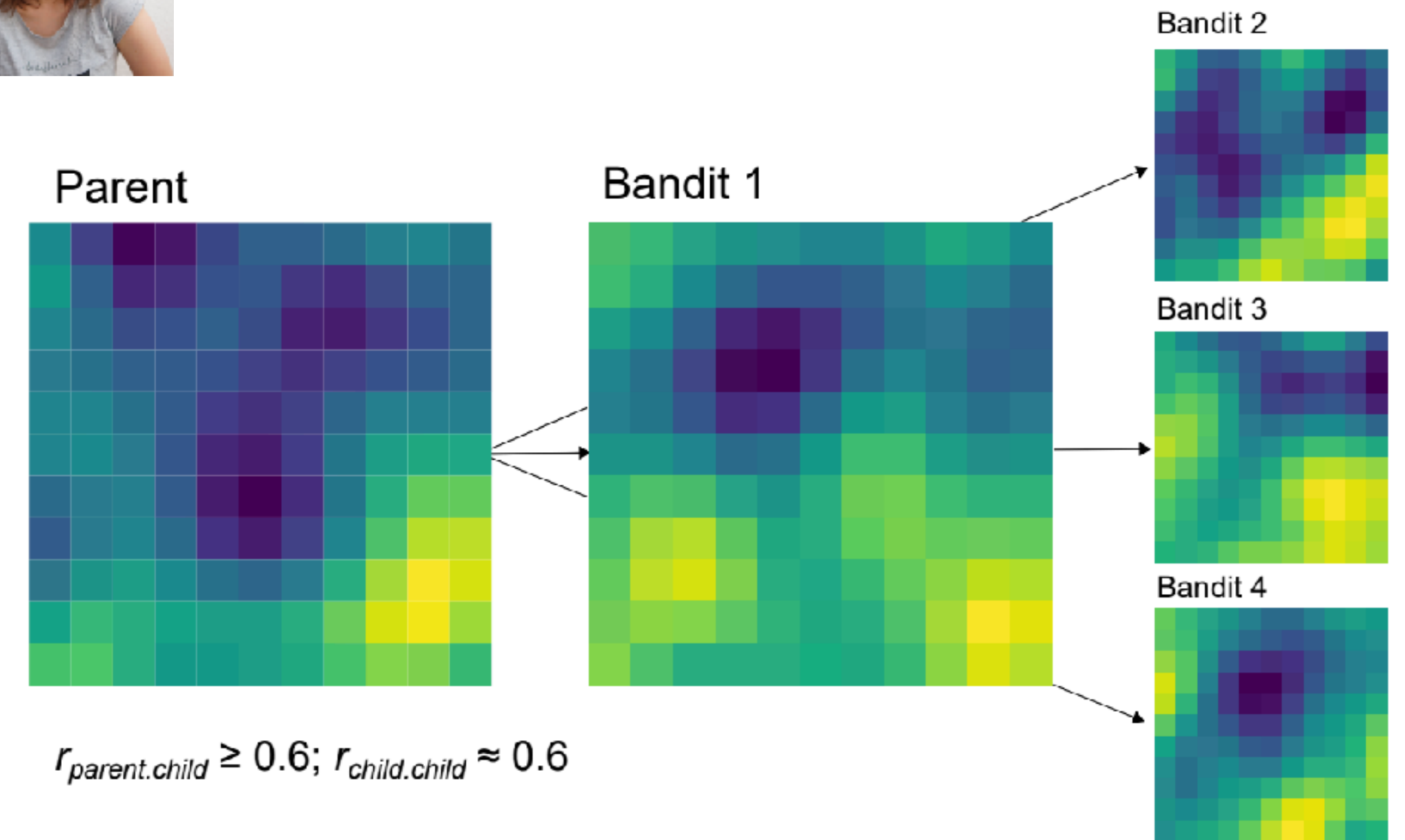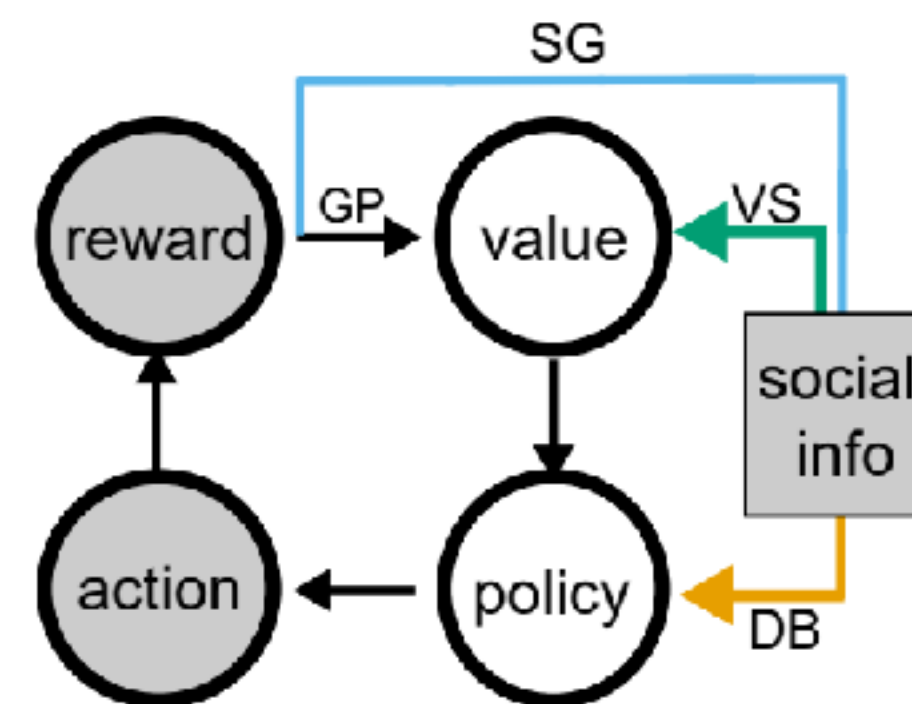
- Social learning despite individual differences (Witt et al., 2023)

  - People can use social information, but not verbatim

  - Exact imitation strategies might fail to account for social differences

- Decision-Biasing (**DB**)

- Value-Shaping (**VS**)

- Social Generalization (**SG**):

  - integration social info in the reward generalization process

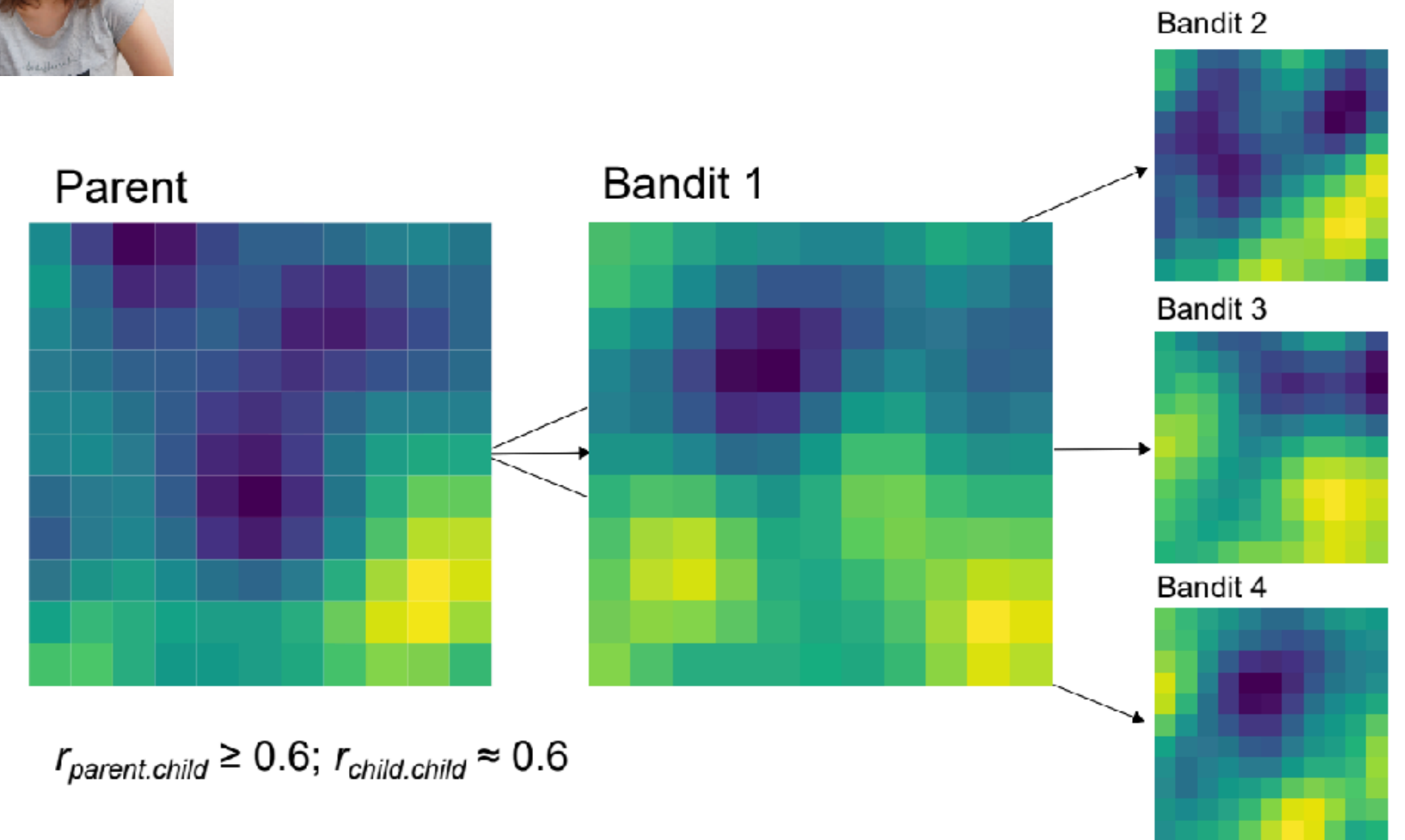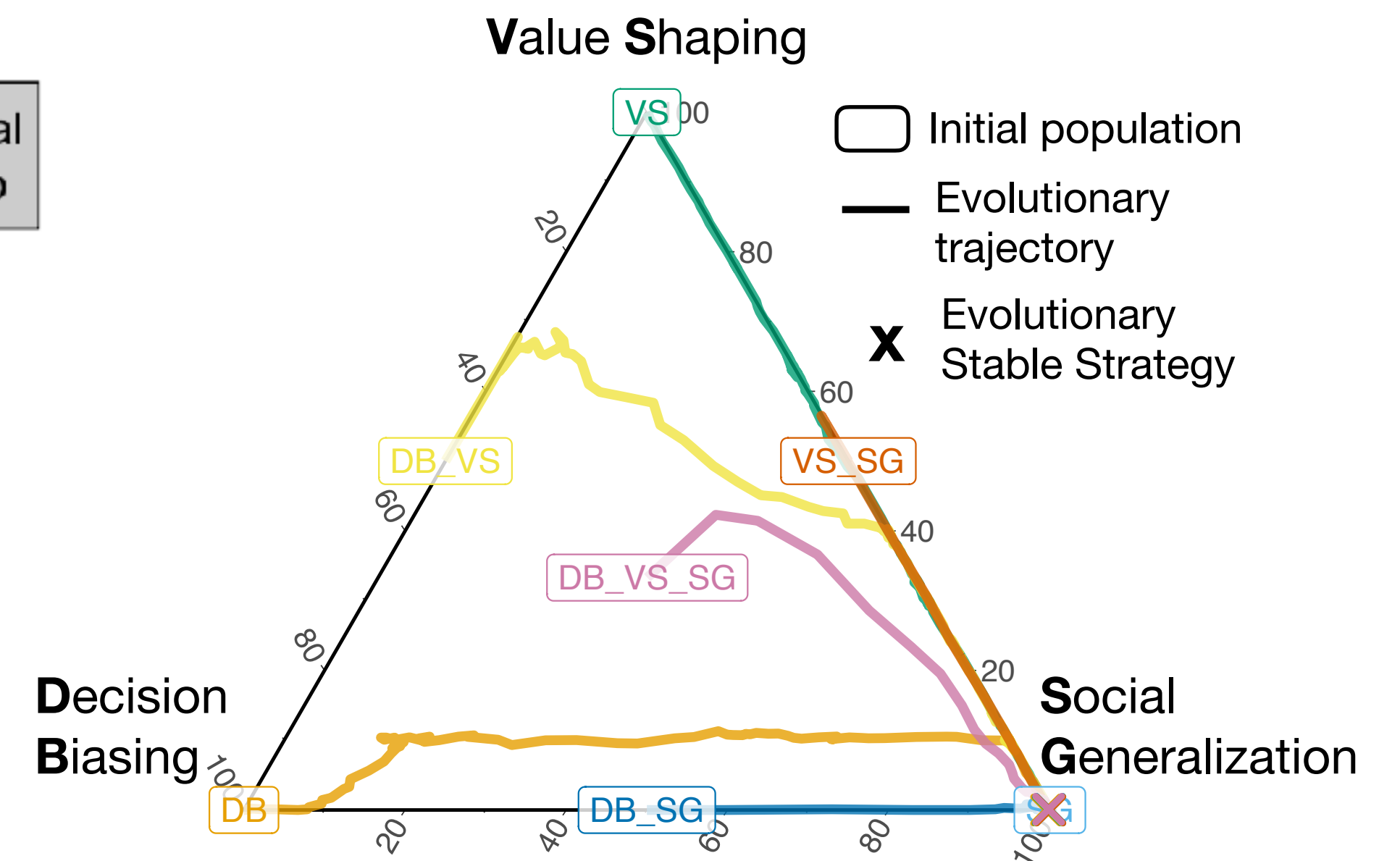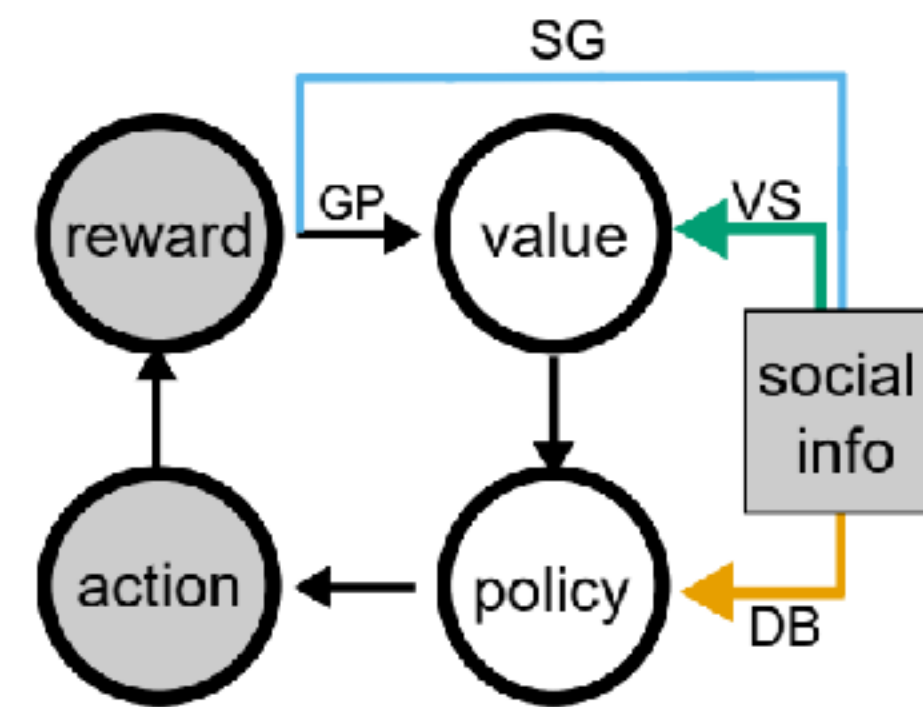  - assume social info is noisier than individual experiences



Parent     Bandit 1     Bandit 2 / Bandit 3 / Bandit 4

$r_{parent.child} \geq 0.6; r_{child.child} \approx 0.6$



SG / GP / value / VS / reward / social info / action / policy / DB

# Evolutionary simulations

# Summary and open challenges

- Social learning deploys a range of tools:

  - **imitation**: directly copy observed behaviors

  - **value-shaping:** add a heuristic bonus to observed behaviors

  - **ToM Inference**: inferring hidden value representations or hidden beliefs about the world

- However, this represents only a subset of social learning mechanisms:

  - Intelligent behavior is not only a function of each individual but also how well groups collectively solve problems

  - Over large time scales, simple innovations can cumulatively add up to produce massively complex cultural solutions

  - So far we have focused on observational learning, but social learning also involves pedagogy and explicit communication

- Yet for each mechanism we can describe verbally, we can also define a computational model that makes more precise commitments to the mechanisms of behavior

- Through experimentation and modeling, we can iteratively tweak and refine our understanding of social learning.