

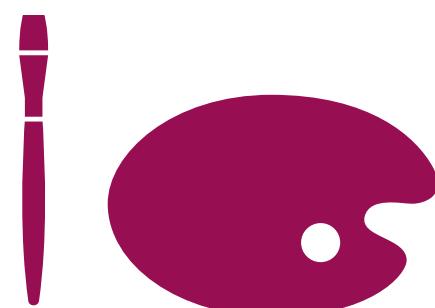
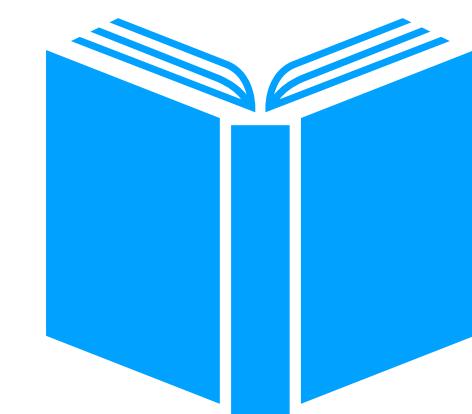
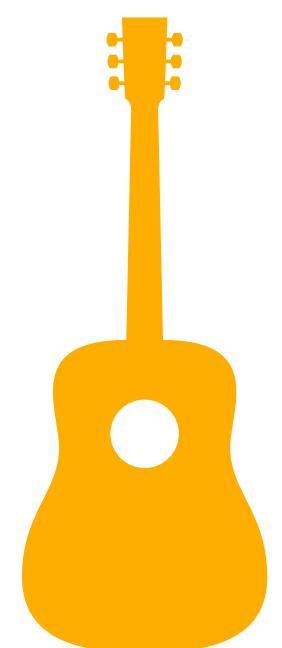
# Generalization in Reinforcement Learning

Charley Wu  
Human and Machine Cognition Lab  
[hmc-lab.com](http://hmc-lab.com)

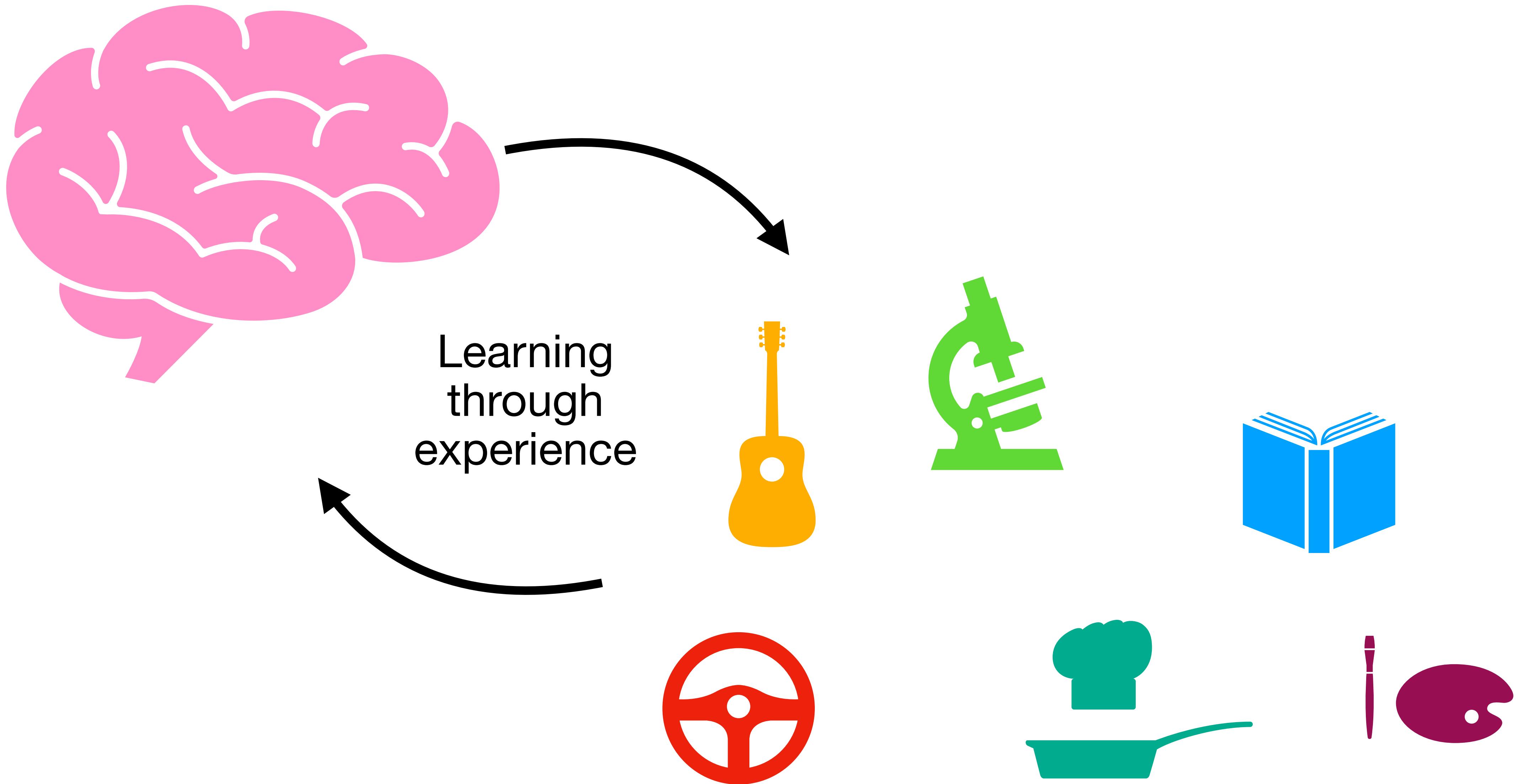
# Brains control behavior



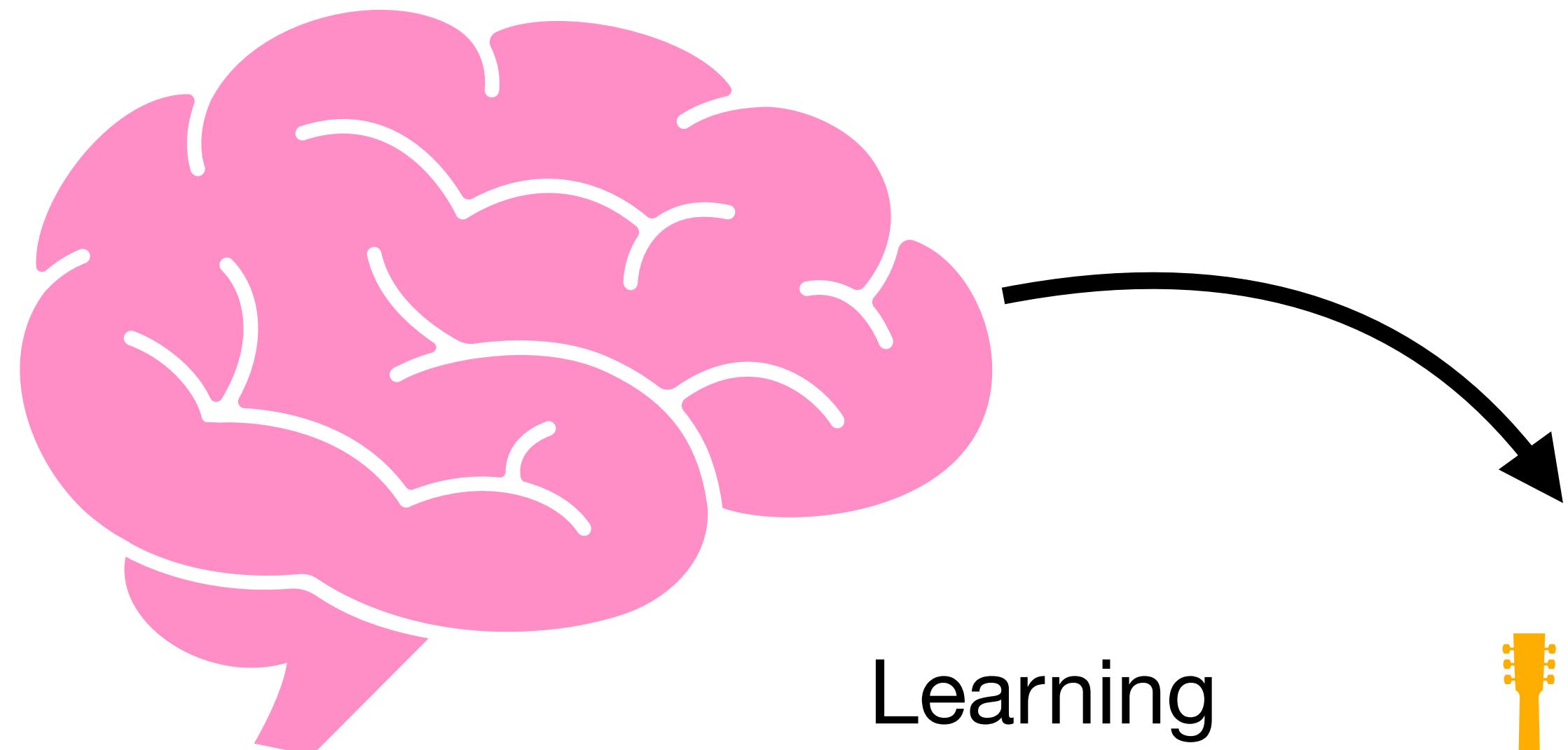
# Brains control behavior



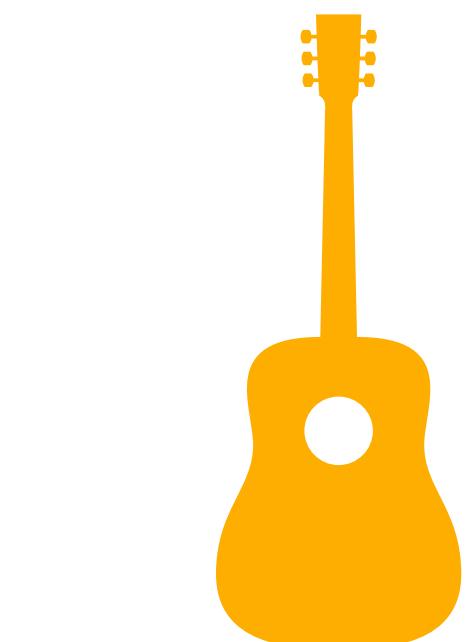
# Brains control behavior



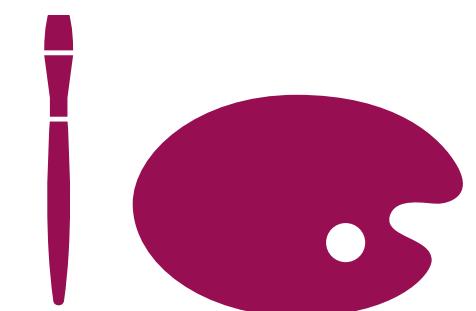
# Brains control behavior



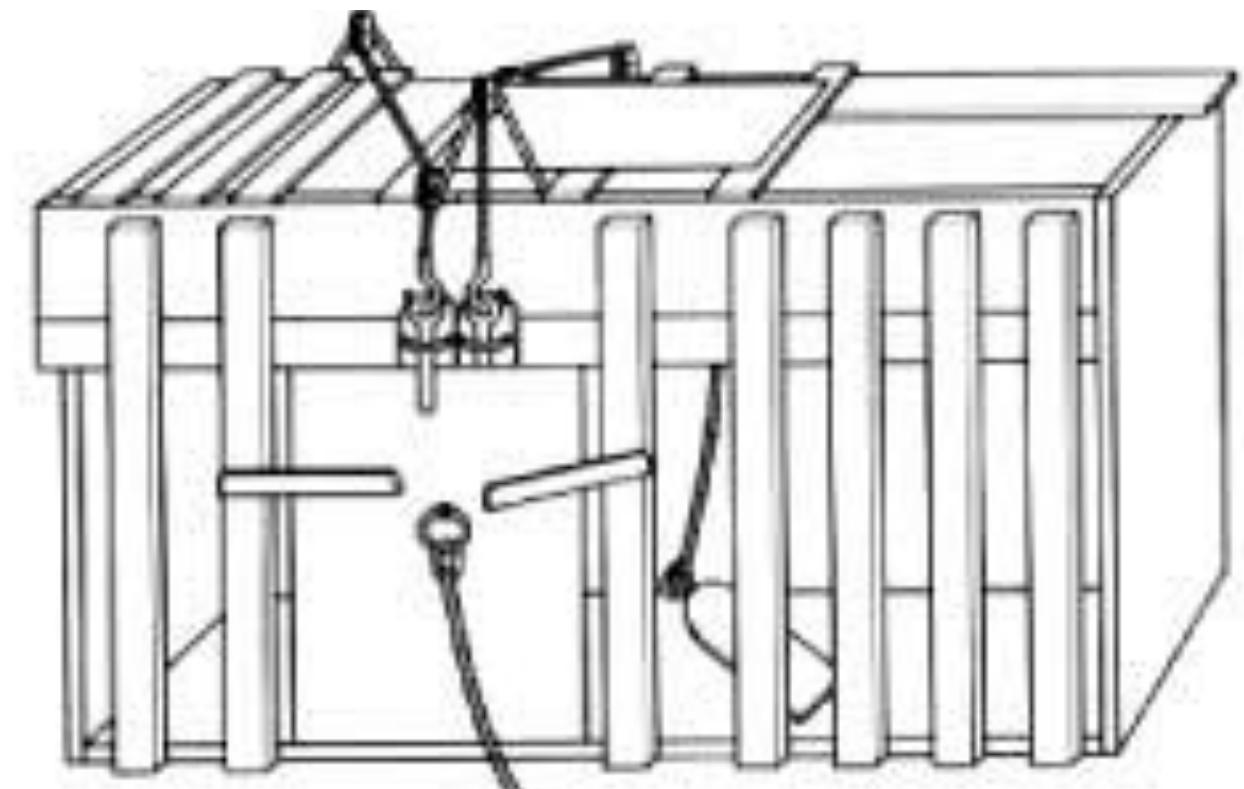
Learning  
through  
experience



***But how?***

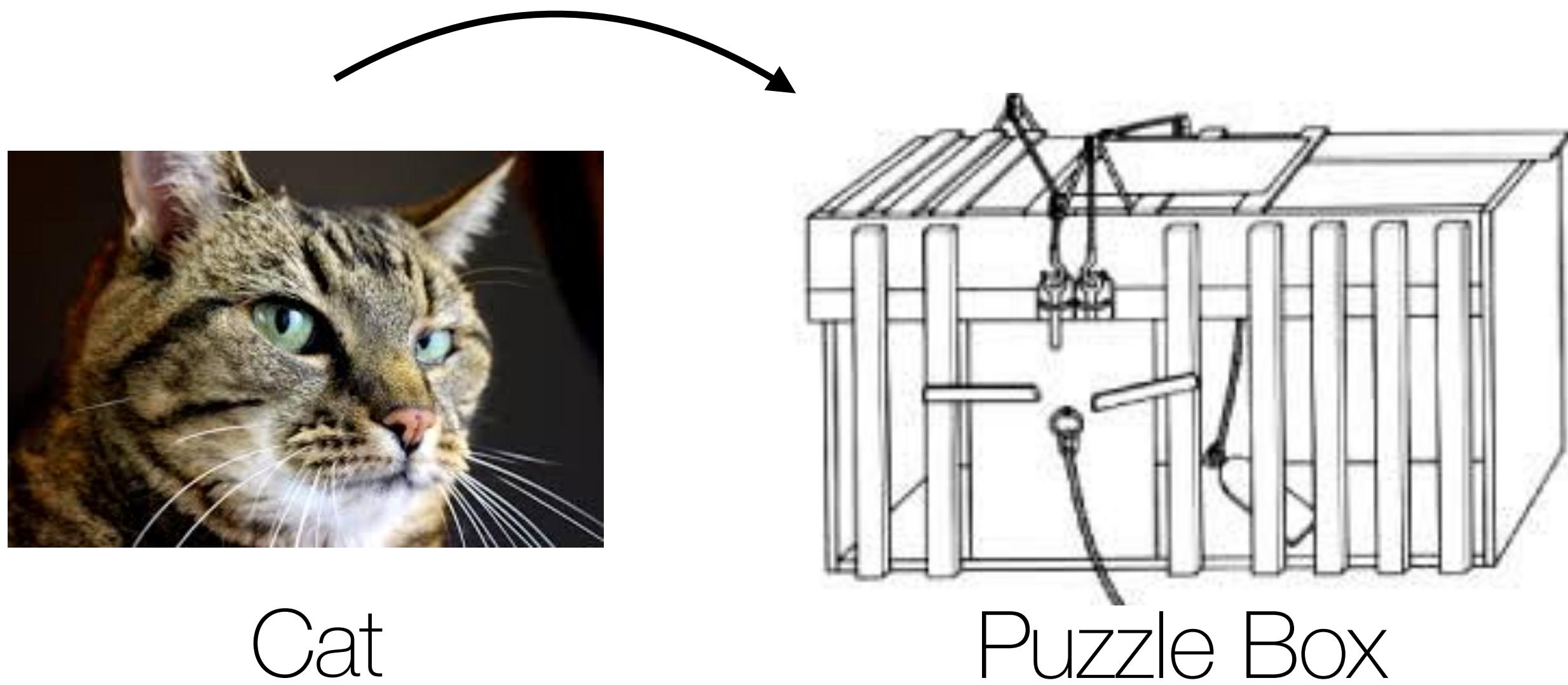


# Thorndike's (1898) Law of Effect

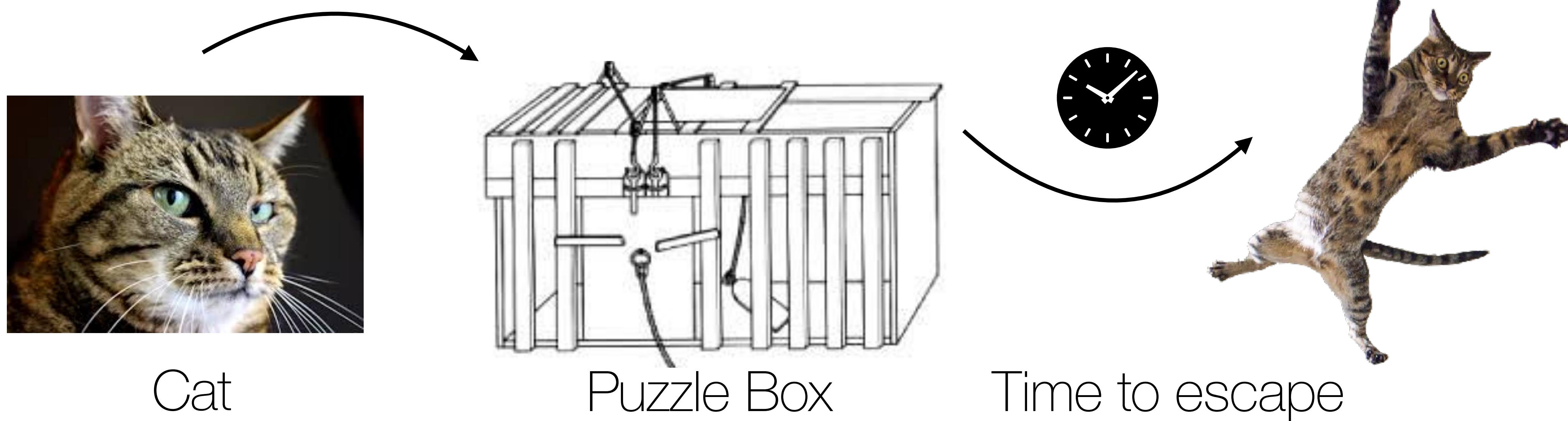


Puzzle Box

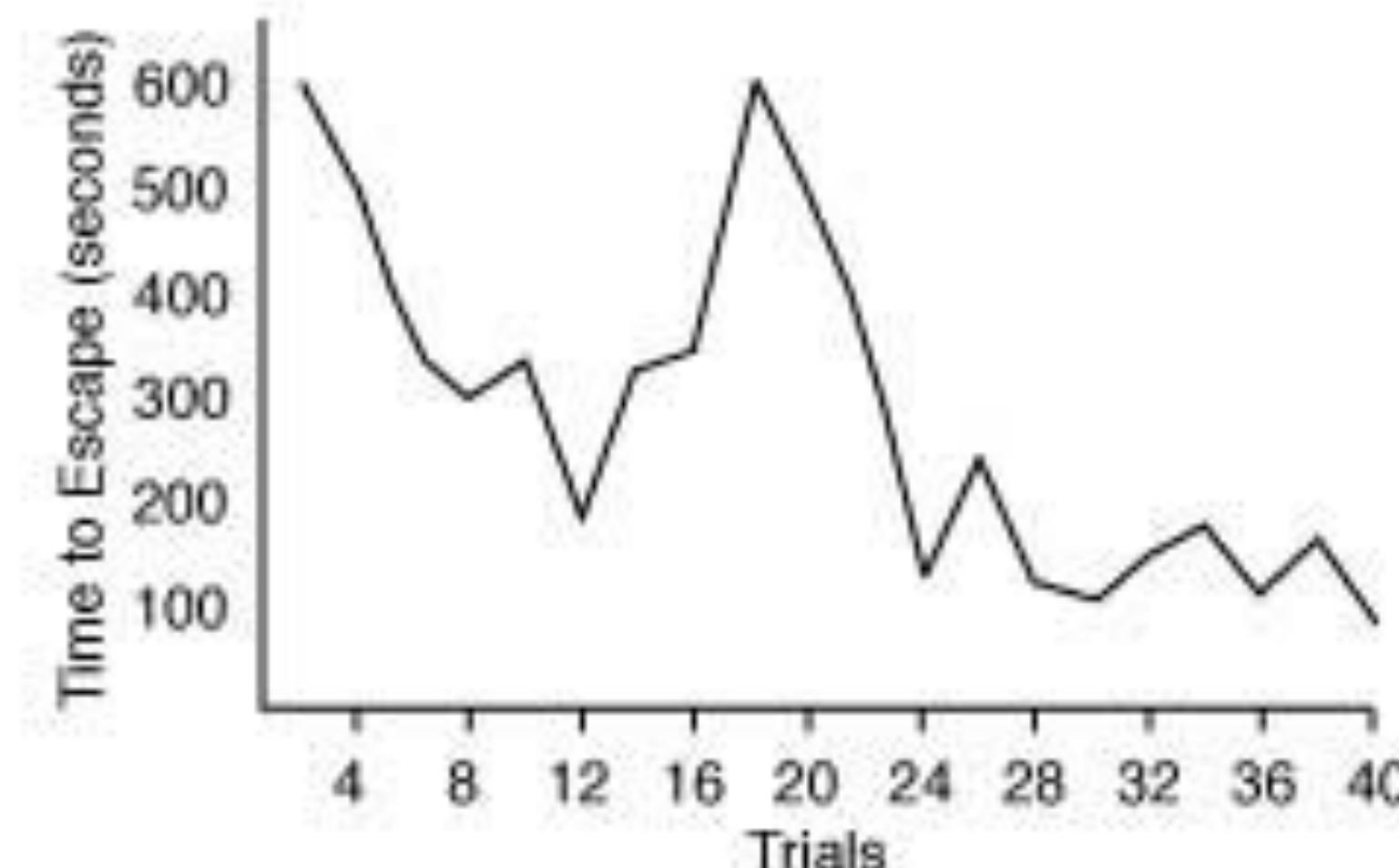
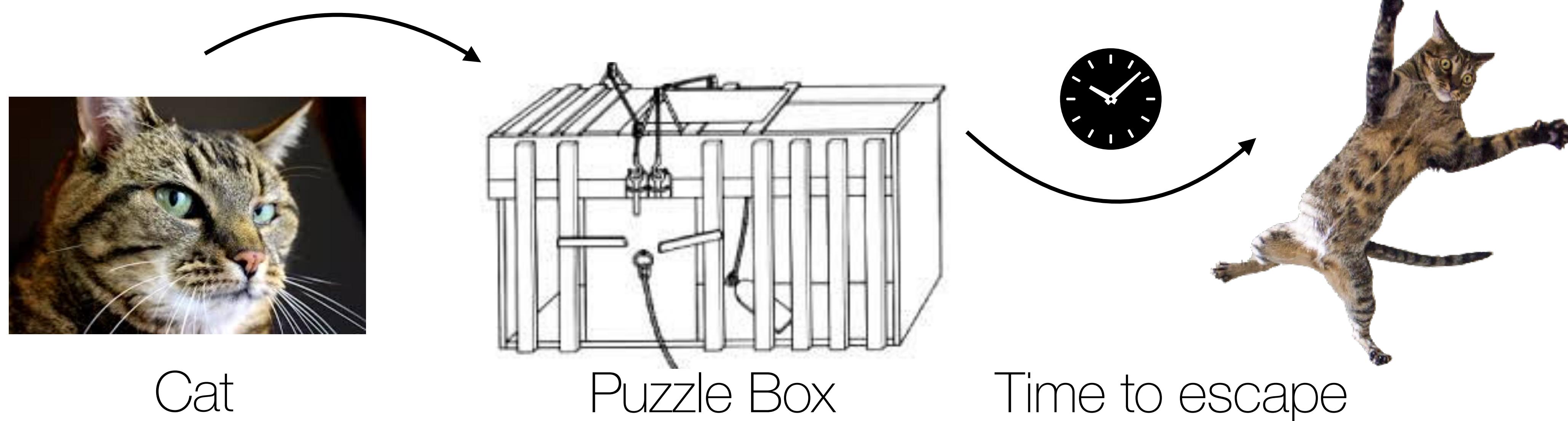
# Thorndike's (1898) Law of Effect



# Thorndike's (1898) Law of Effect



# Thorndike's (1898) Law of Effect

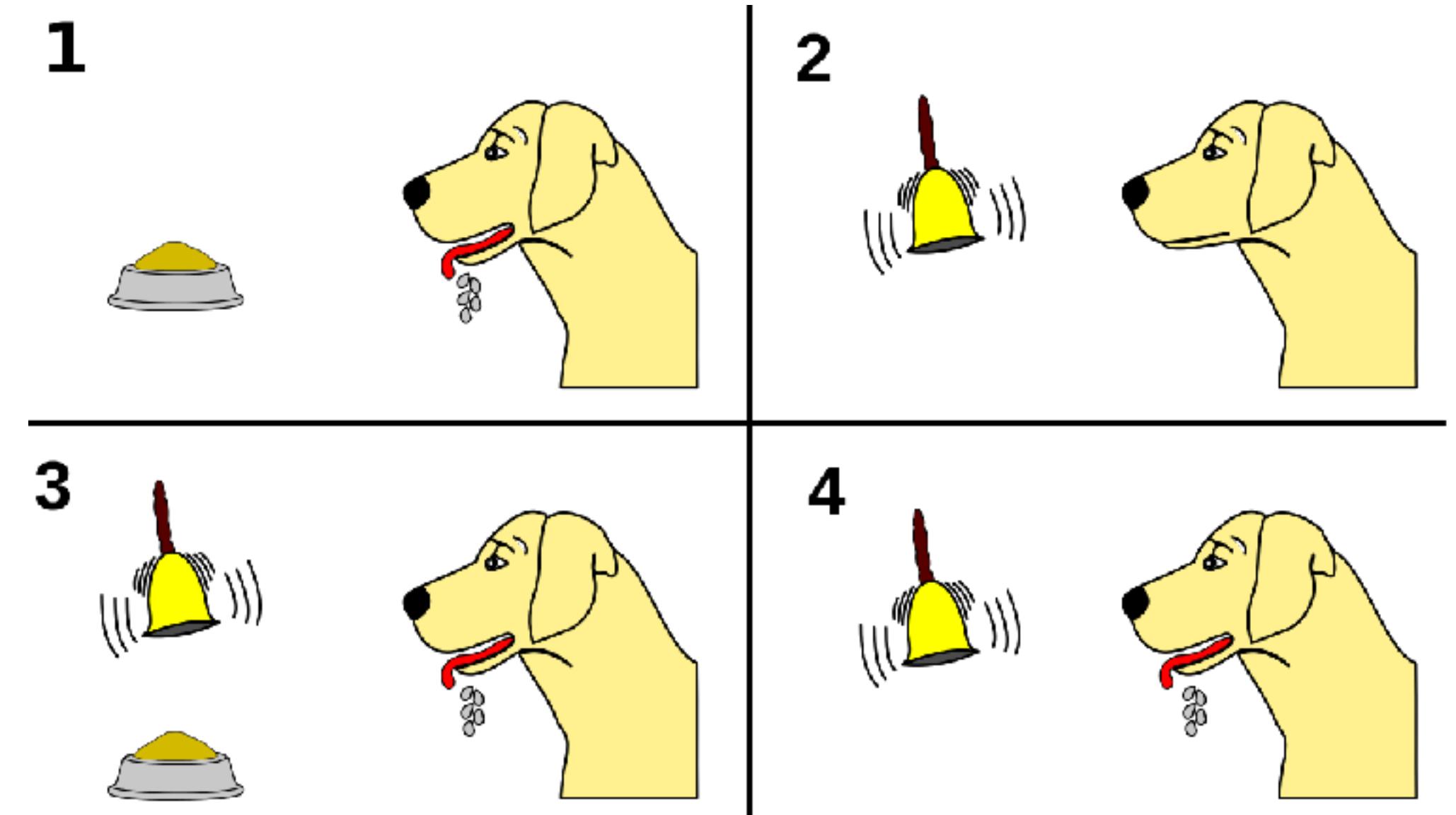


*Actions associated with satisfaction are strengthened, while those associated with discomfort become weakened.*

# Classical and Operant Conditioning

## Classical Condition (Pavlov, 1927)

Learning as the passive coupling of stimulus (bell ringing) and response (salivation), anticipating future rewards



## Operant Condition (Skinner, 1938)

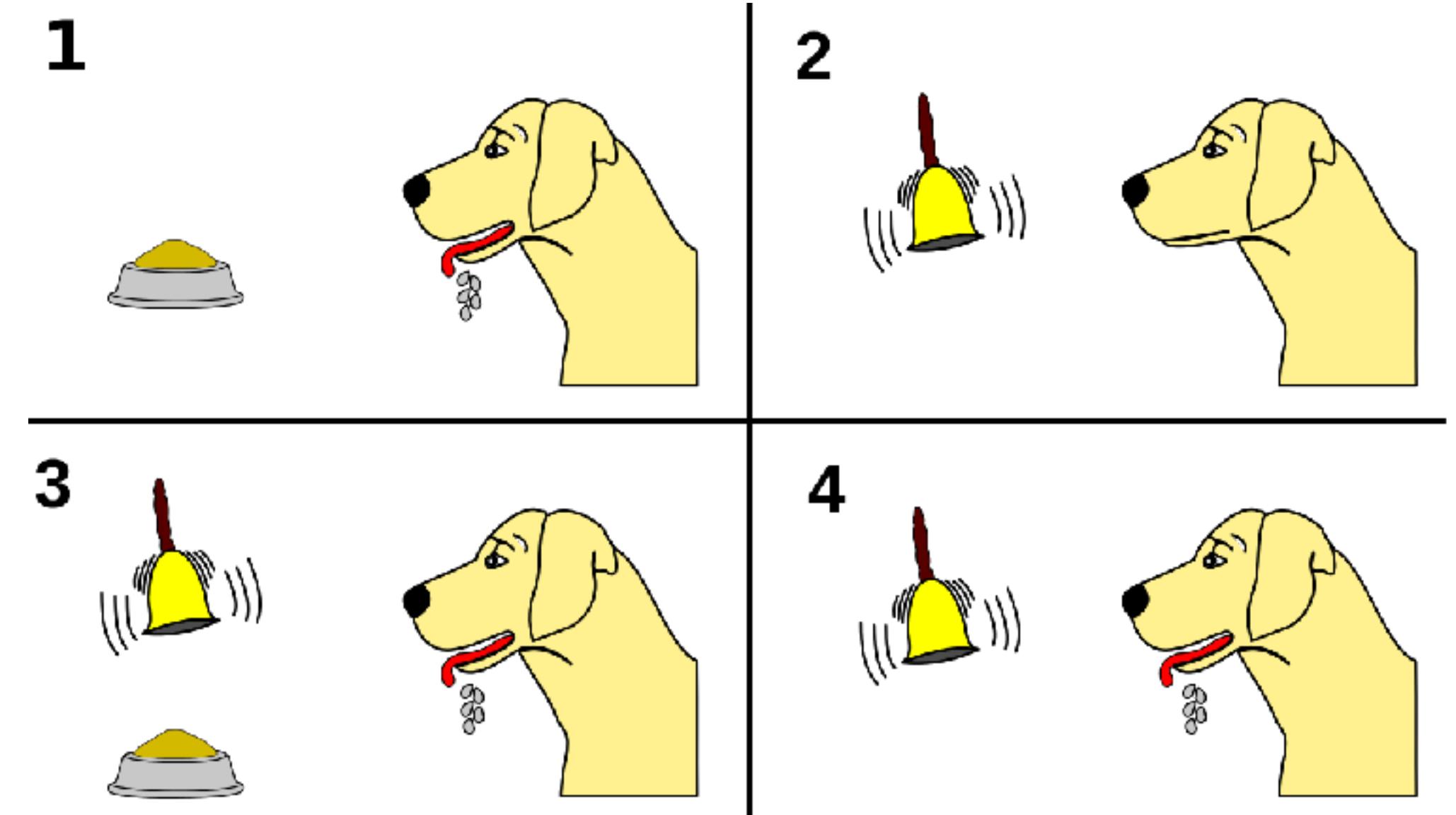
Skinner (1938): Learning as the active shaping of behavior in response to rewards or punishments



# Classical and Operant Conditioning

## Classical Condition (Pavlov, 1927)

Learning as the passive coupling of stimulus (bell ringing) and response (salivation), anticipating future rewards



## Operant Condition (Skinner, 1938)

Skinner (1938): Learning as the active shaping of behavior in response to rewards or punishments



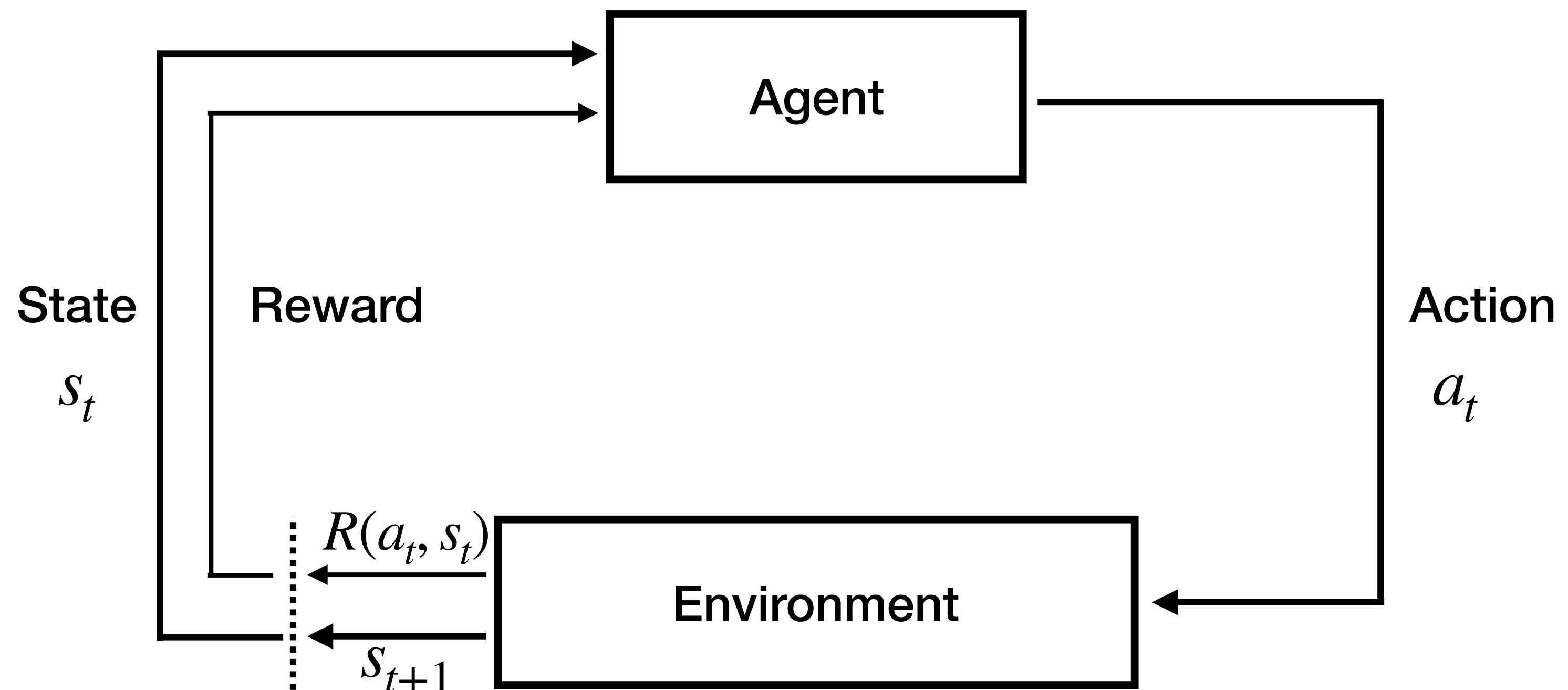
# Reinforcement Learning

## The Environment:

- governs the transition between states
- provides rewards

## The Agent:

- Learns a *value function* mapping actions onto rewards
- Implements a policy, selecting actions based on their value



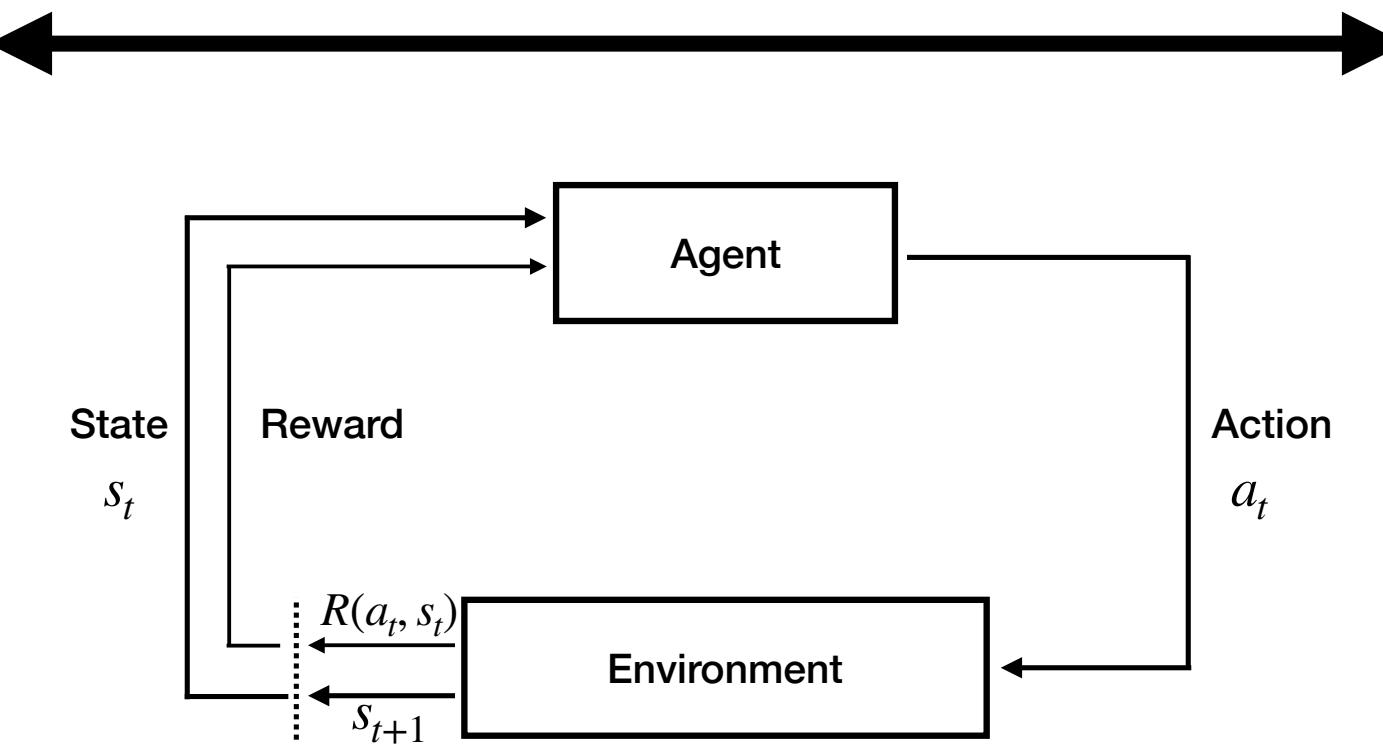
Sutton and Barto (2018 [1998])

# Neuroscience



# AI and Machine Learning

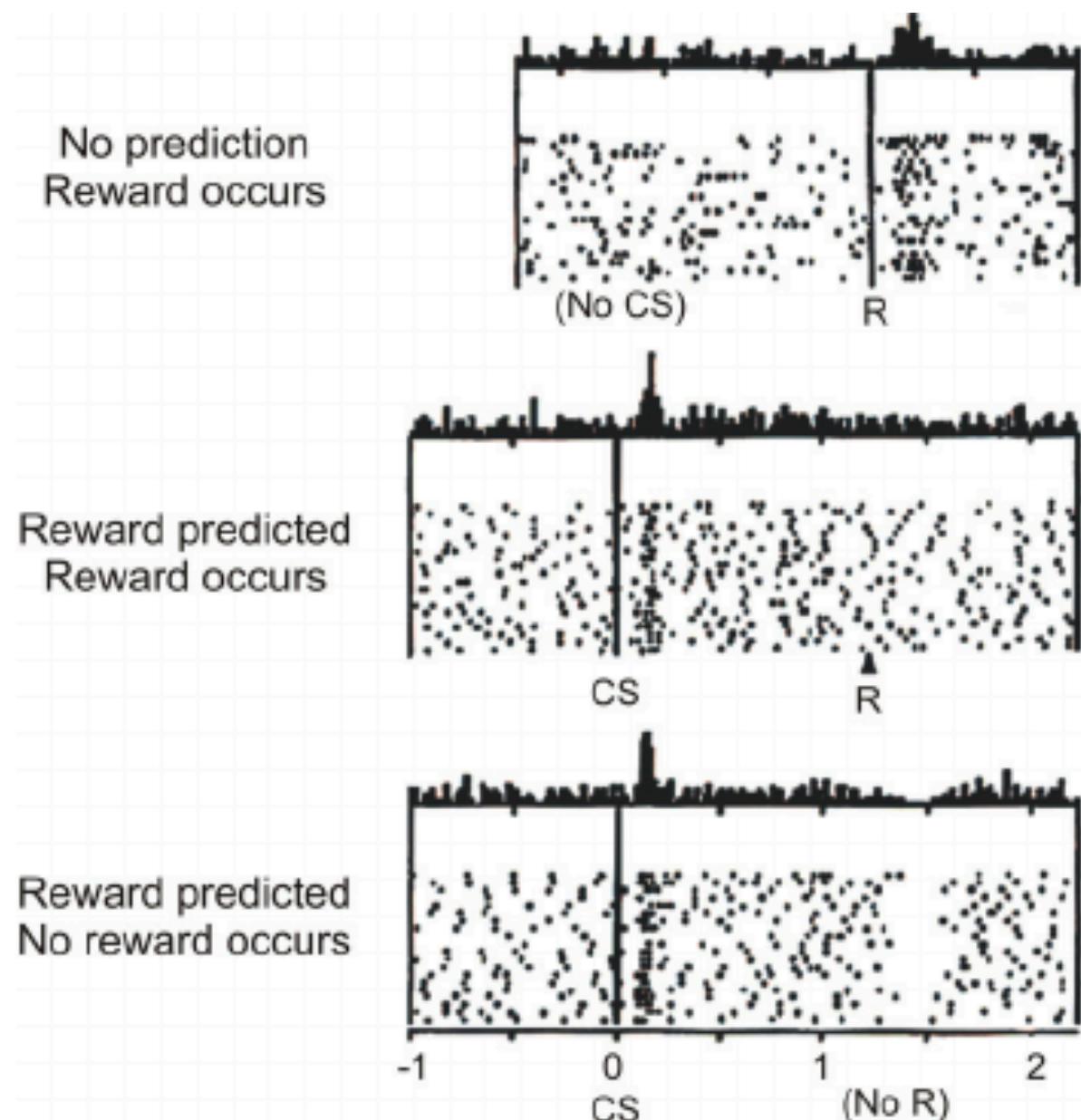
## Reinforcement Learning



# Neuroscience

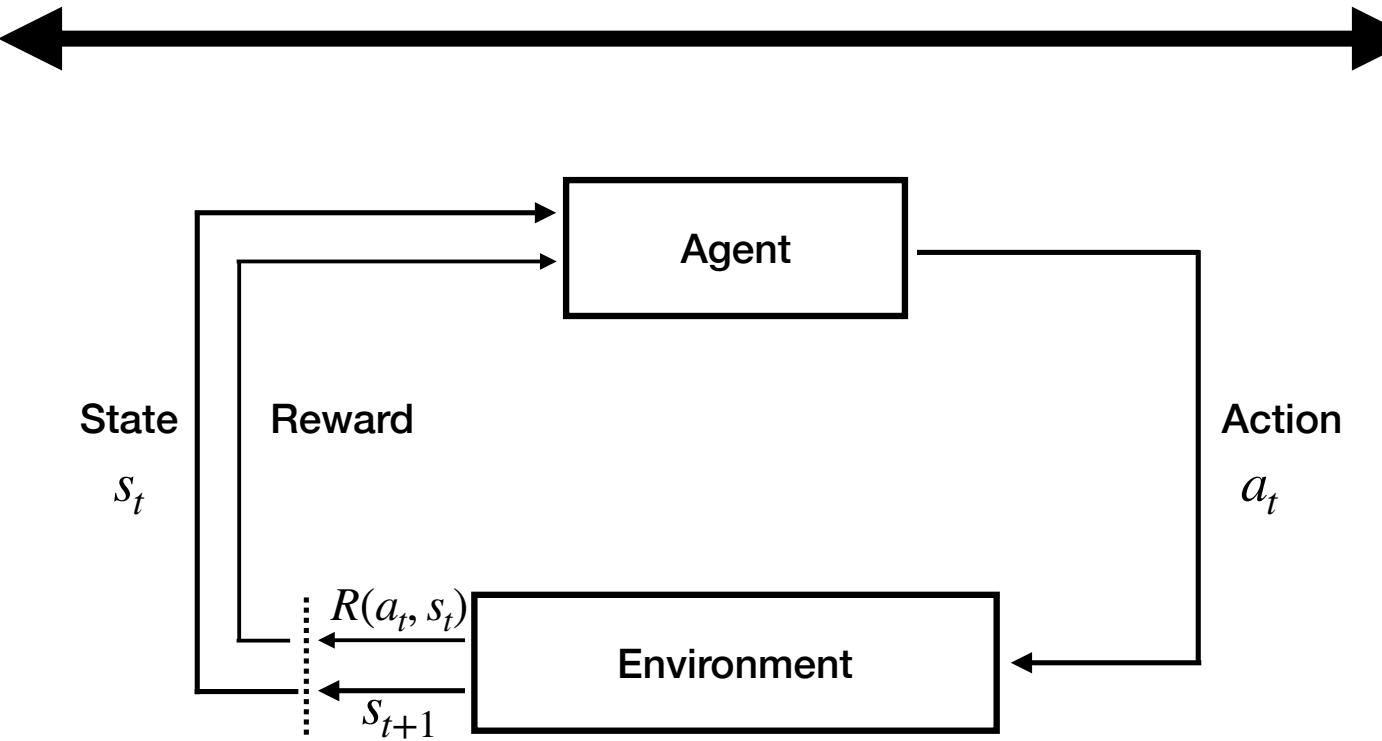


## Dopamine Reward Prediction Error



Schultz et al. (1997)

## Reinforcement Learning



# AI and Machine Learning

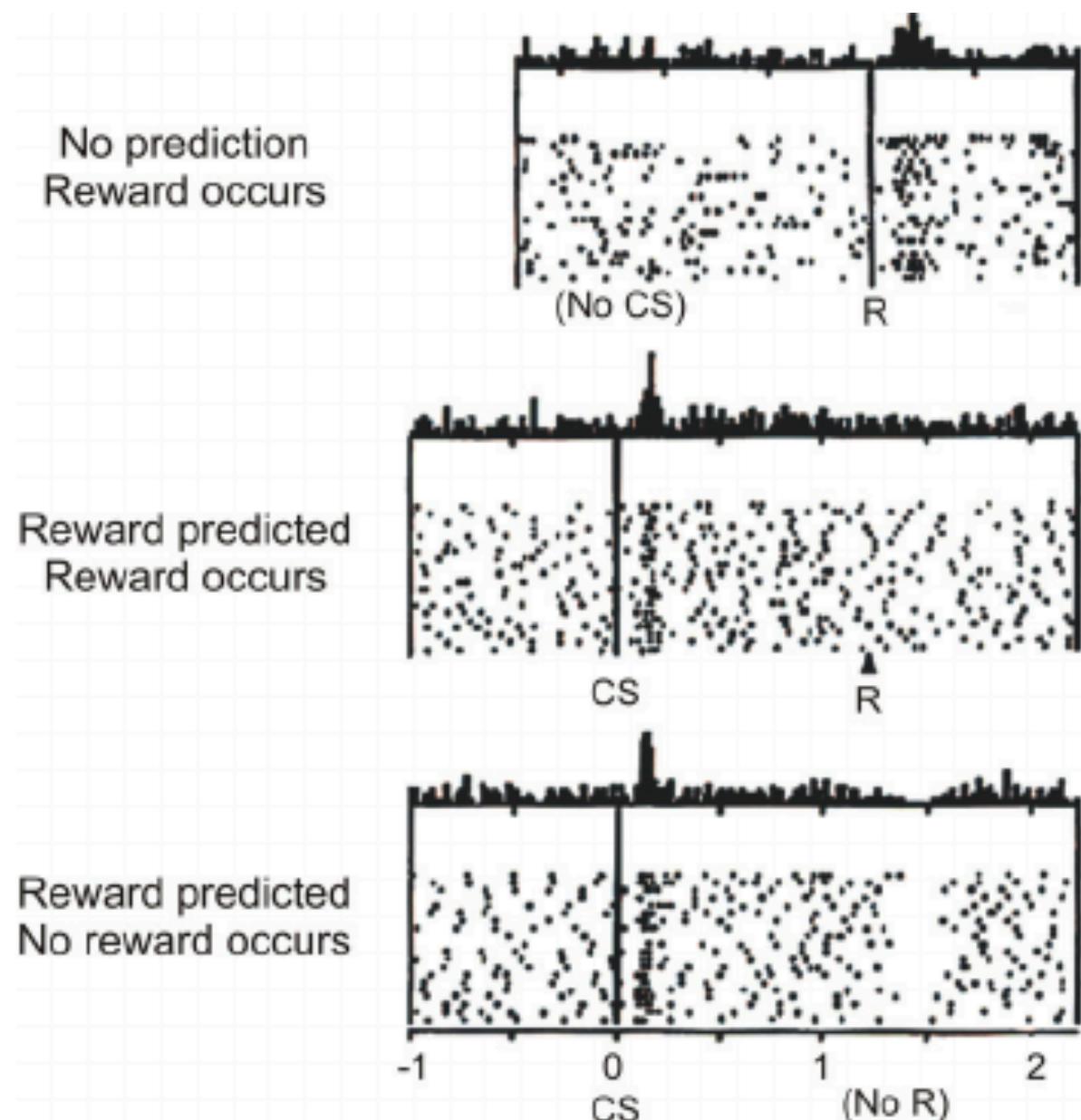


Temporal Difference  
Learning

# Neuroscience

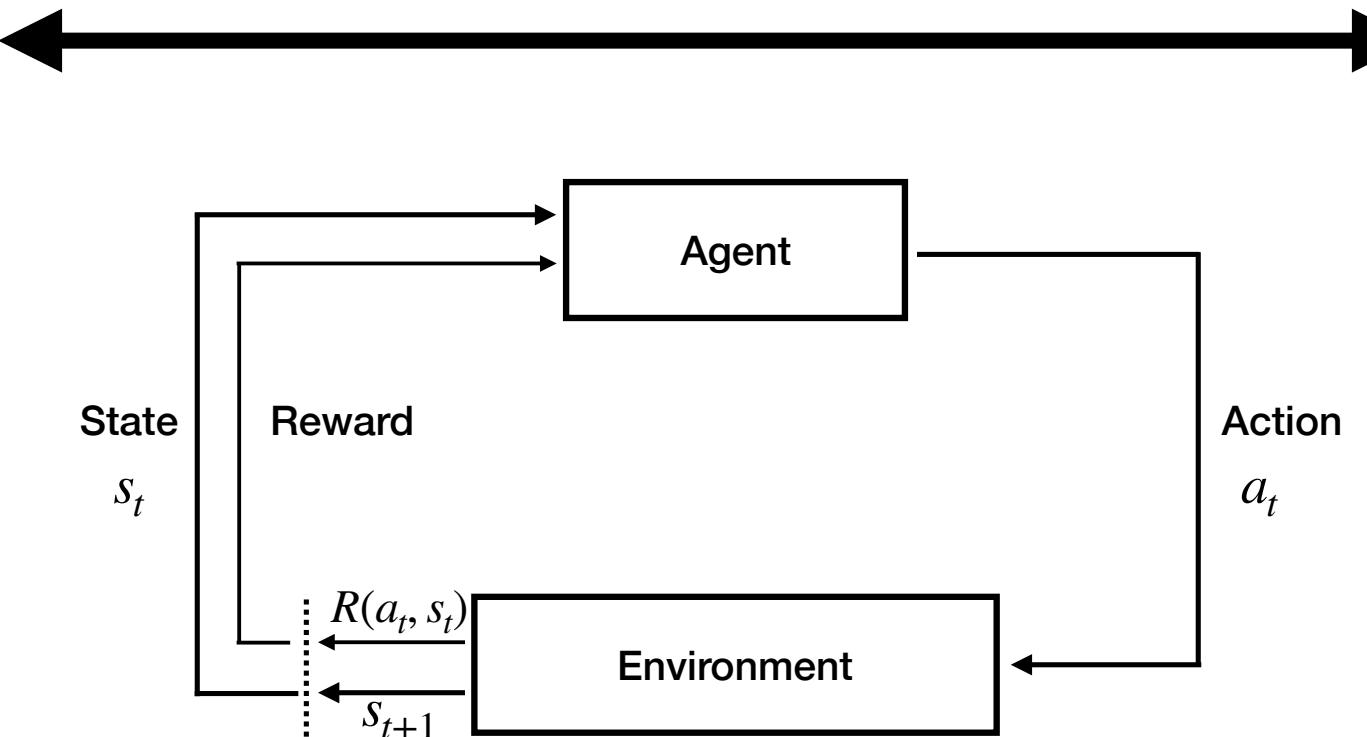


## Dopamine Reward Prediction Error



Schultz et al. (1997)

## Reinforcement Learning



# AI and Machine Learning



## AlphaGo



Silver et al. (2016)

Temporal Difference  
Learning

Monte Carlo  
Tree Search

# Outline

## **Part 1. Overview of Reinforcement Learning**

- Value functions and policies
- Tabular methods vs. value-function approximation
- Multi-armed Bandit problem
- Models of human learners

# Outline

## **Part 1. Overview of Reinforcement Learning**

- Value functions and policies
- Tabular methods vs. value-function approximation
- Multi-armed Bandit problem
- Models of human learners

~~~~~ Break ~~~~~

# Outline

## **Part 1. Overview of Reinforcement Learning**

- Value functions and policies
- Tabular methods vs. value-function approximation
- Multi-armed Bandit problem
- Models of human learners

~~~~~ Break ~~~~~

## **Part 2. Generalization guided learning**

- Search in vast spaces (Wu et al., *NHB* 2018)
- Learning like a child (Schulz et al., *PsychSci* 2019; Meder et al., *DevSci* 2021; Giron et al., *in prep*)
- Connecting spatial and conceptual search (Wu et al., *PLoS CompBio* 2020)
- Graph-structured generalization (Wu et al., *CBB* 2020)

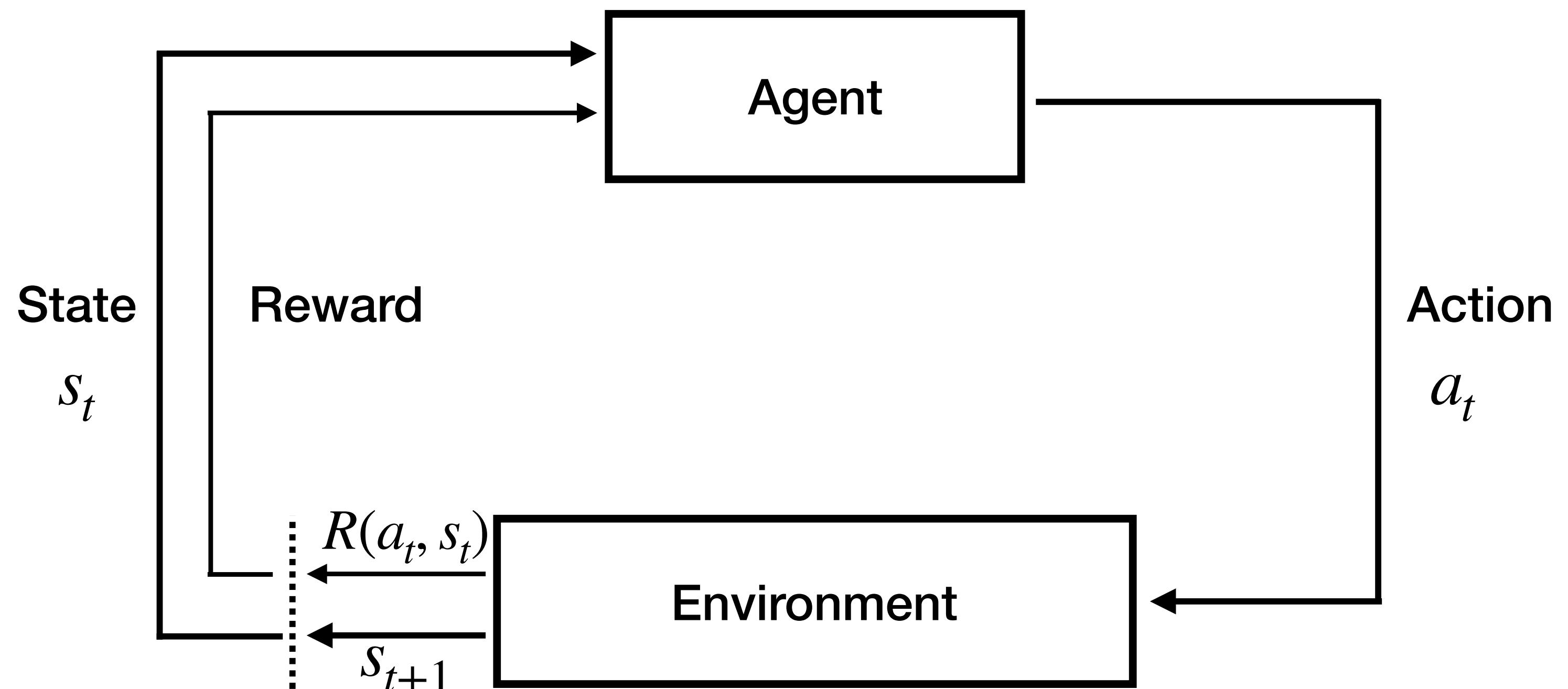
# Reinforcement Learning

## The Environment:

- governs the transition between states
- provides rewards

## The Agent:

- Learns a *value function* mapping actions onto rewards
- Implements a policy, selecting actions based on their value

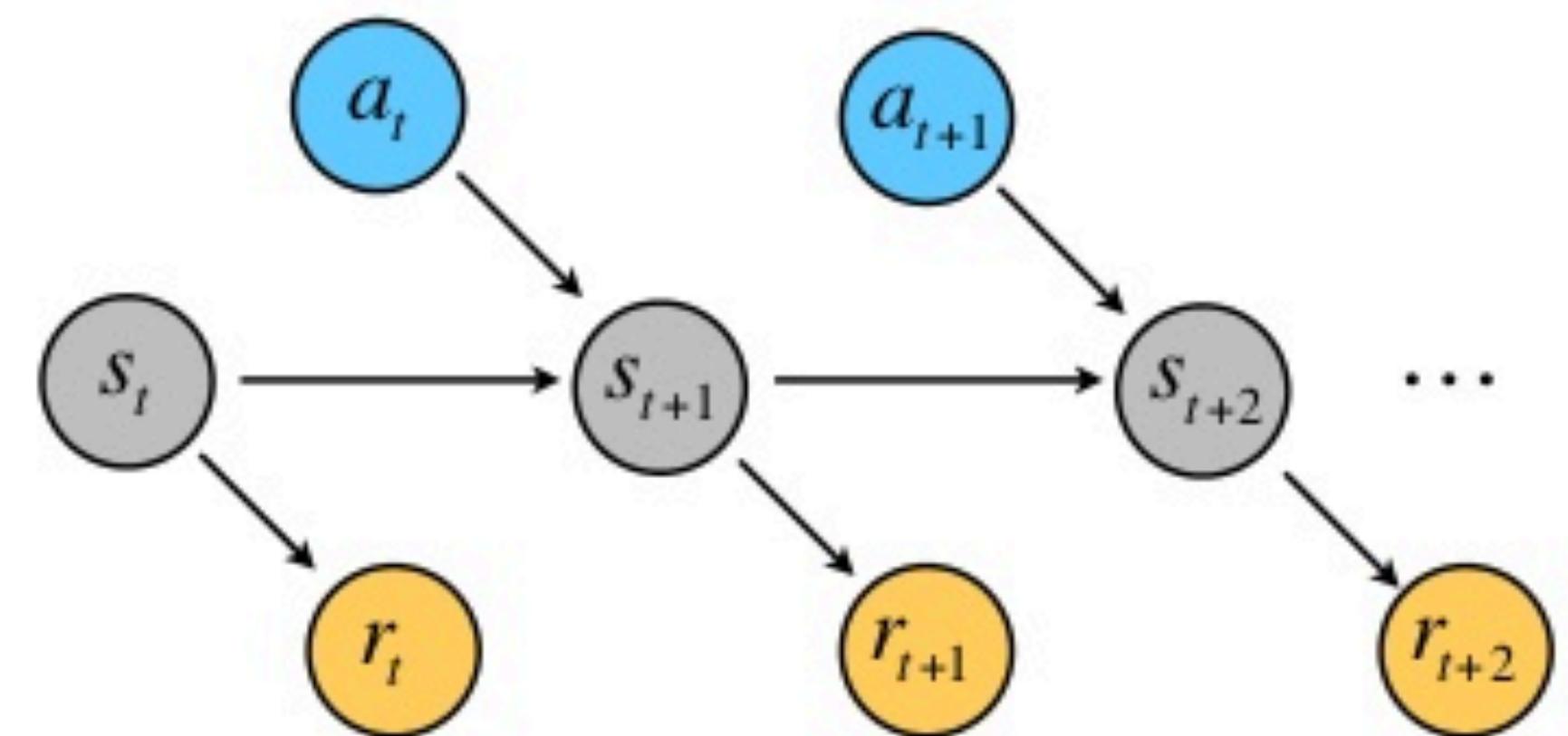


Sutton and Barto (2018 [1998])

# Environment

Markov Decision Process (MDP)

- Simplifying assumption that the system is fully defined by only the previous state (i.e., Markov Principle):  $P(s_{t+1} | s_t, a_t)$



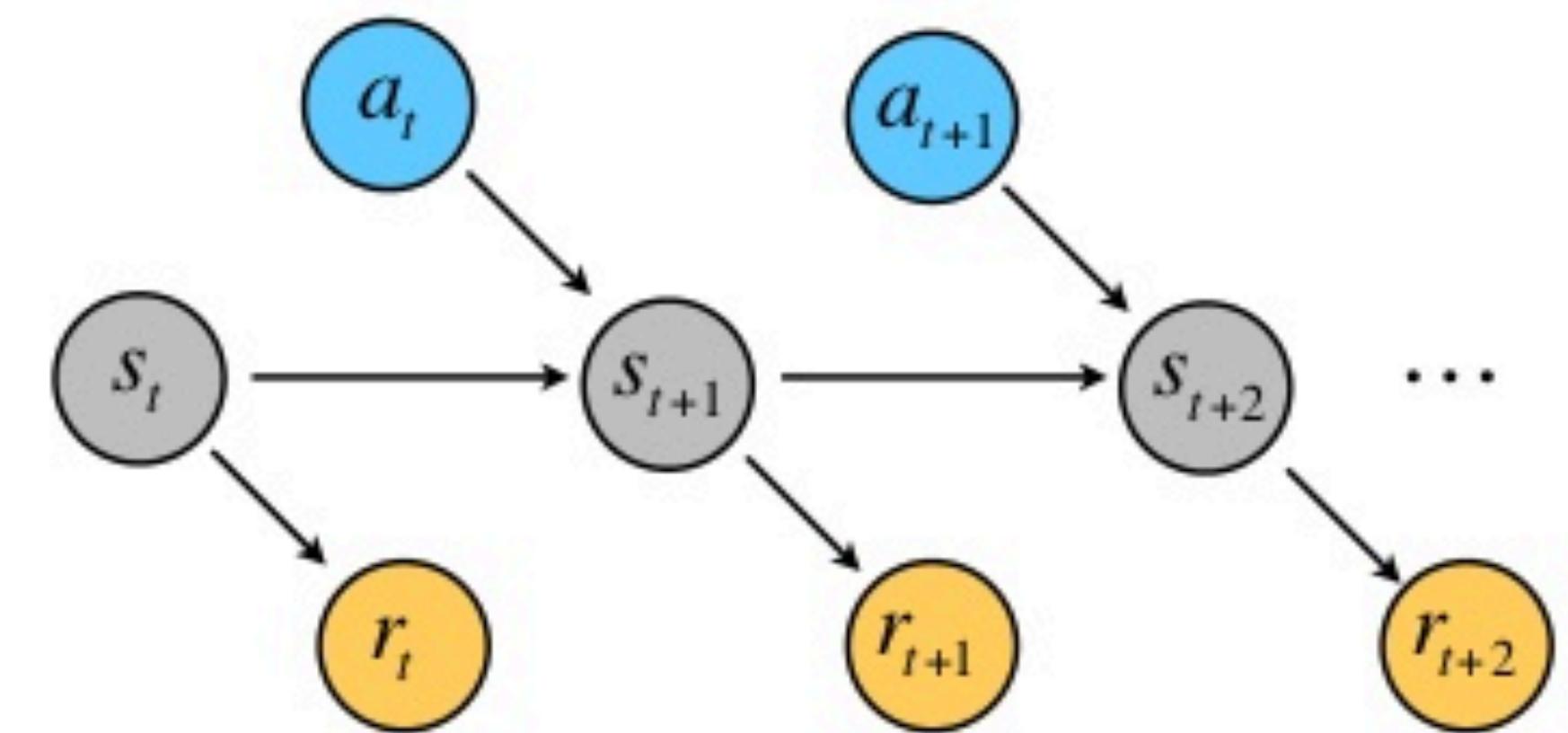
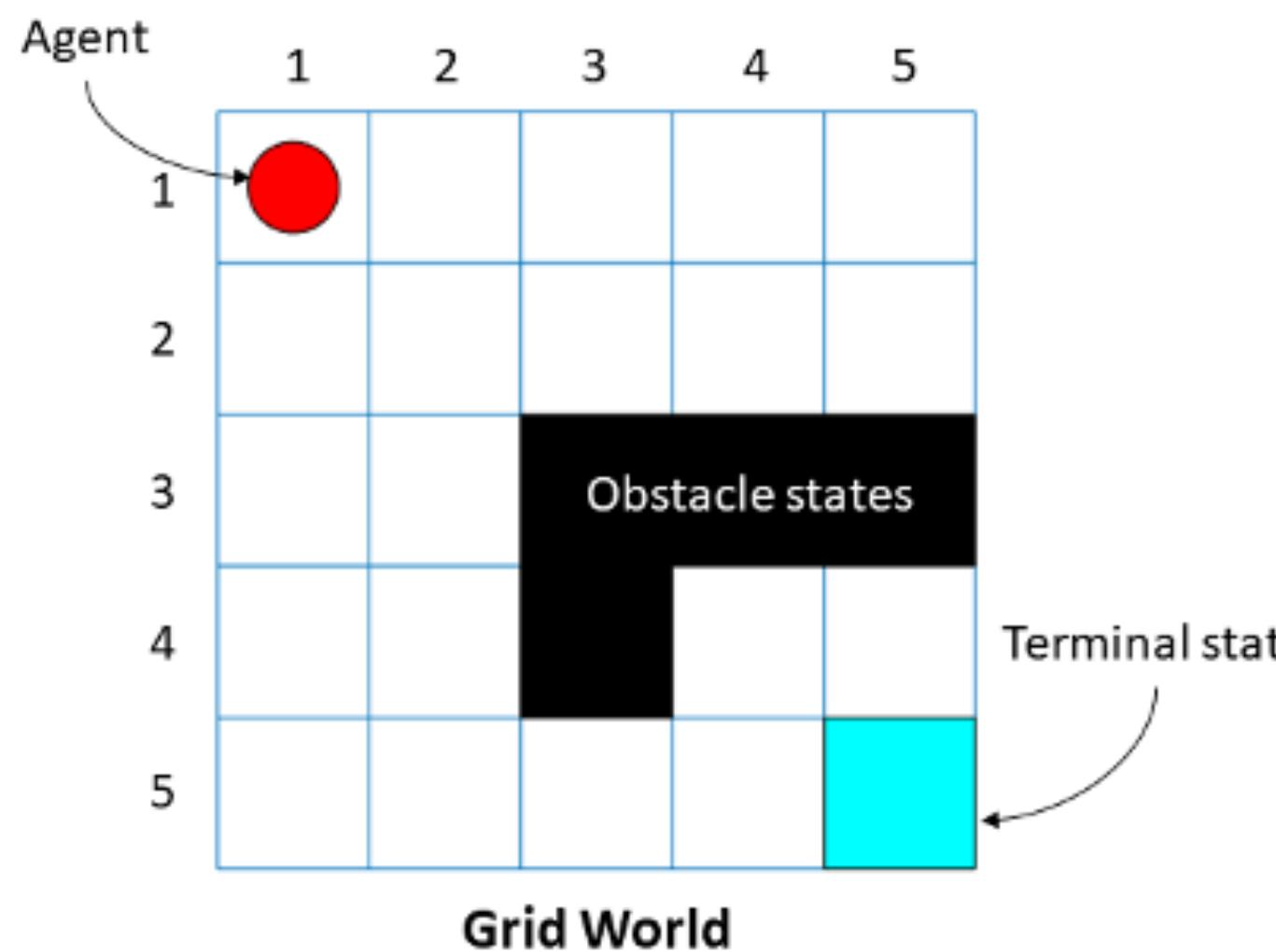
# Environment

Markov Decision Process (MDP)

- Simplifying assumption that the system is fully defined by only the previous state (i.e., Markov Principle):  $P(s_{t+1} | s_t, a_t)$

What are the states?

- Discrete locations, pixels on a screen, a set of feature values, etc...



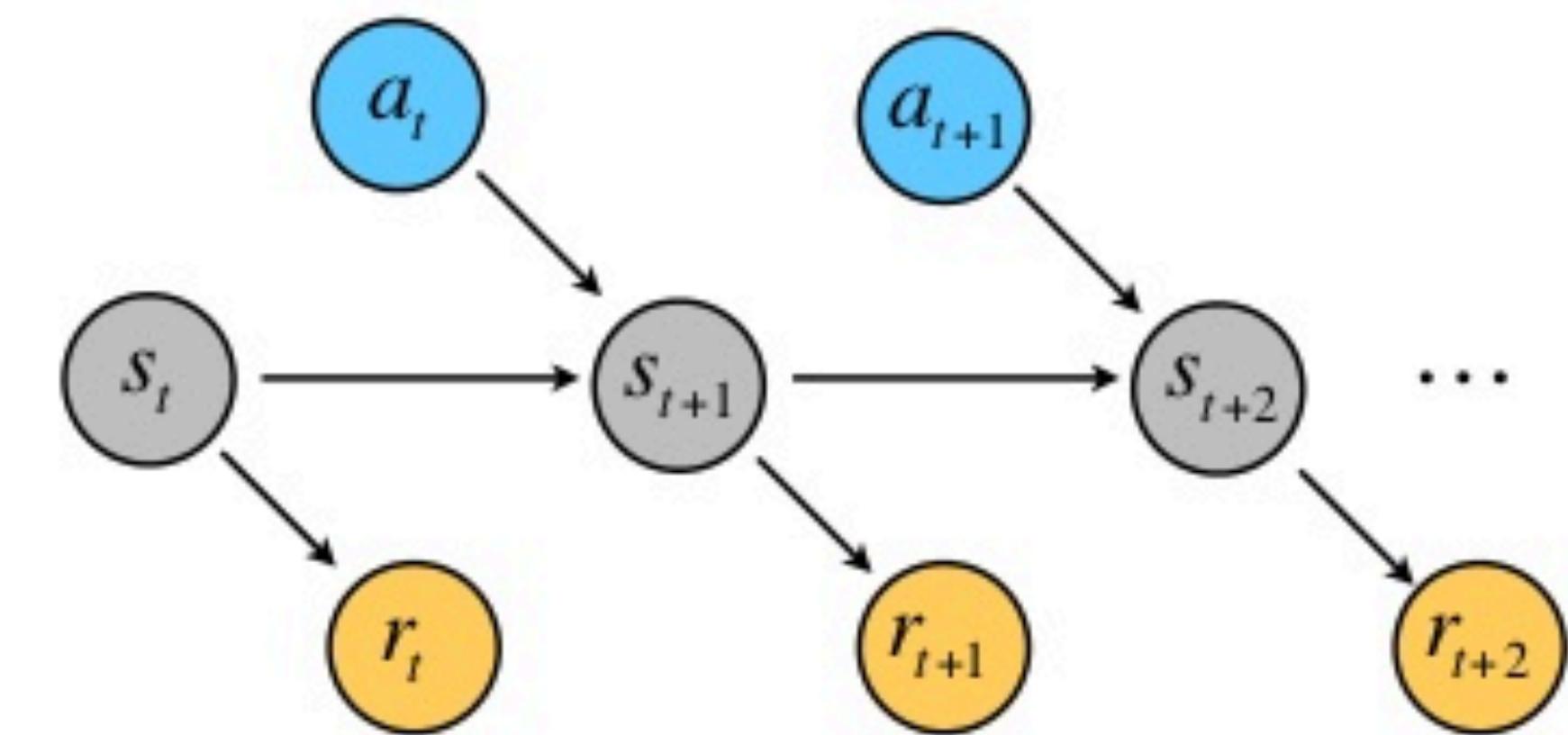
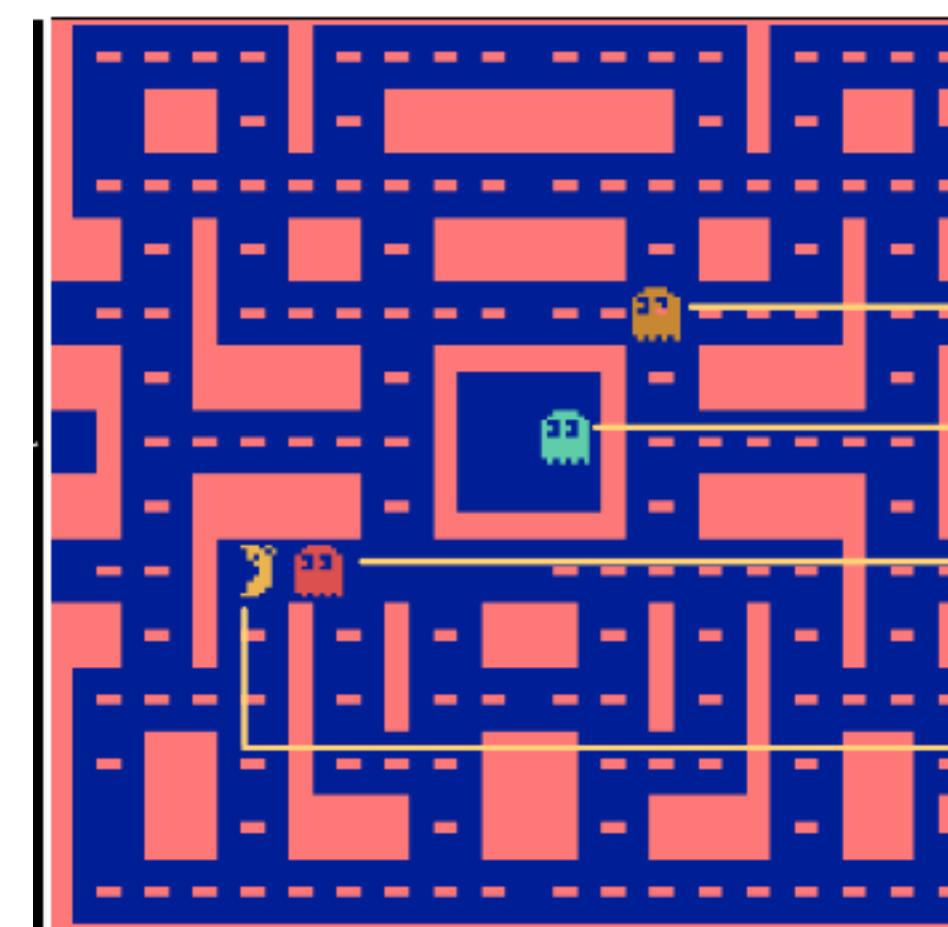
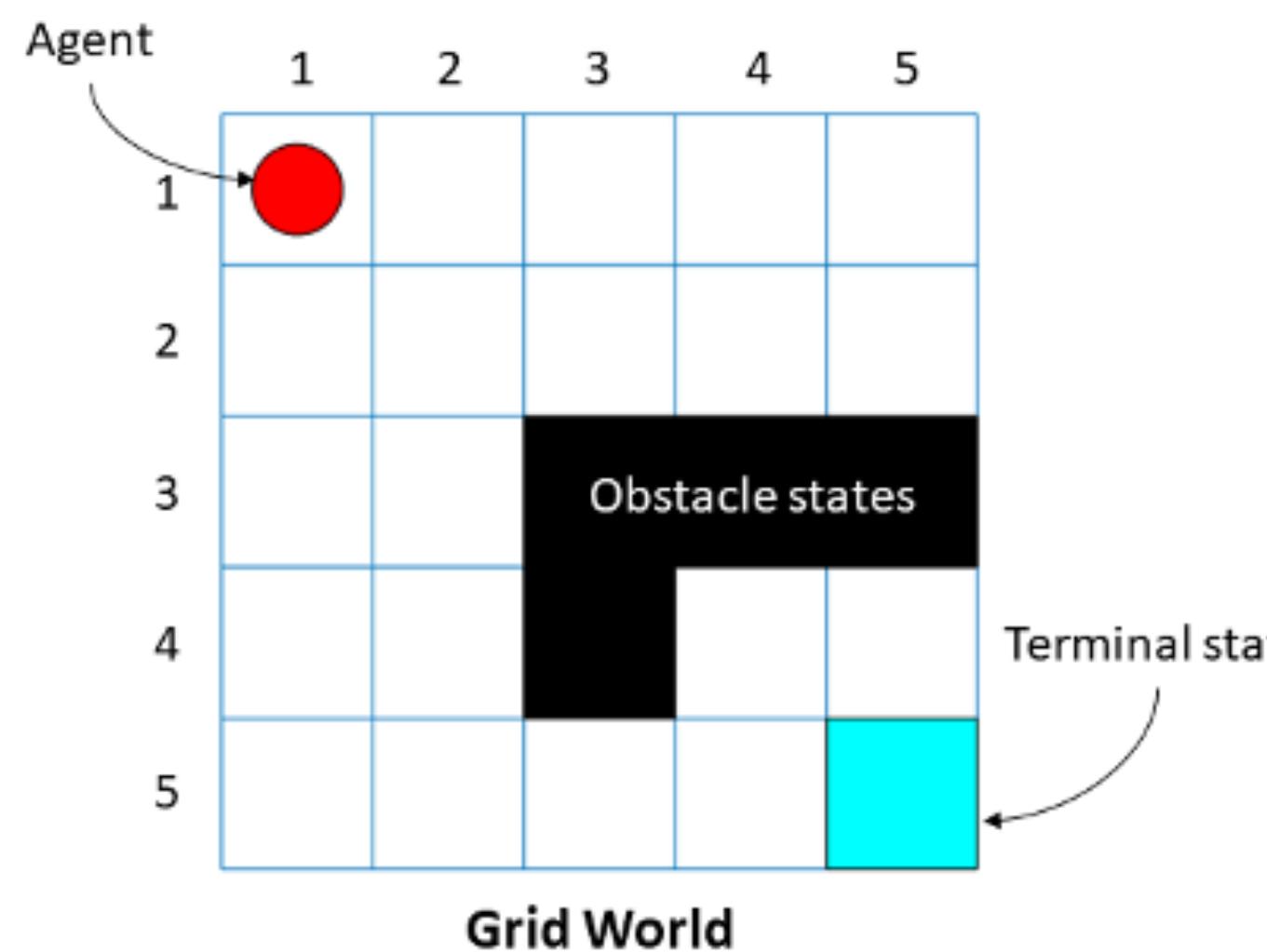
# Environment

Markov Decision Process (MDP)

- Simplifying assumption that the system is fully defined by only the previous state (i.e., Markov Principle):  $P(s_{t+1} | s_t, a_t)$

What are the states?

- Discrete locations, pixels on a screen, a set of feature values, etc...



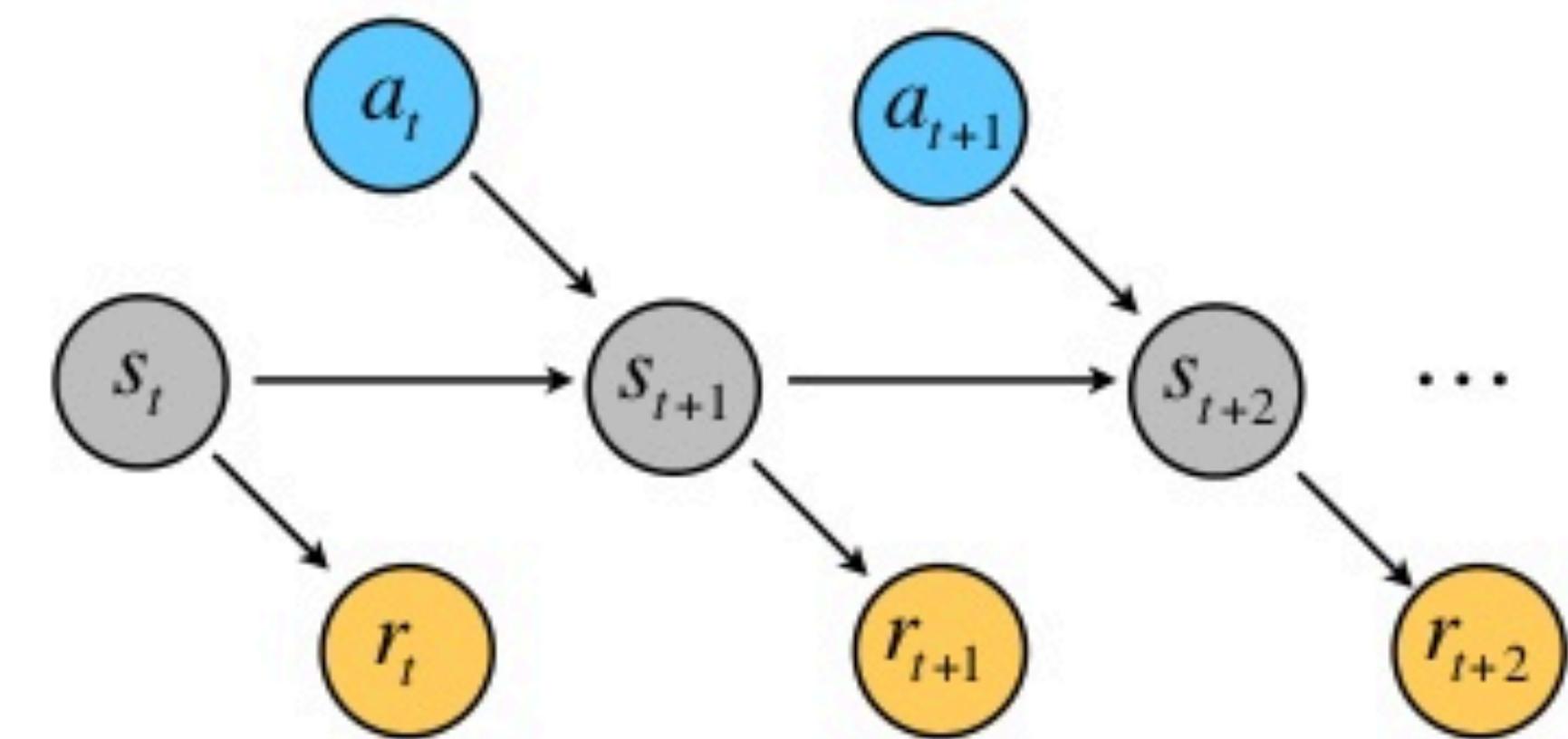
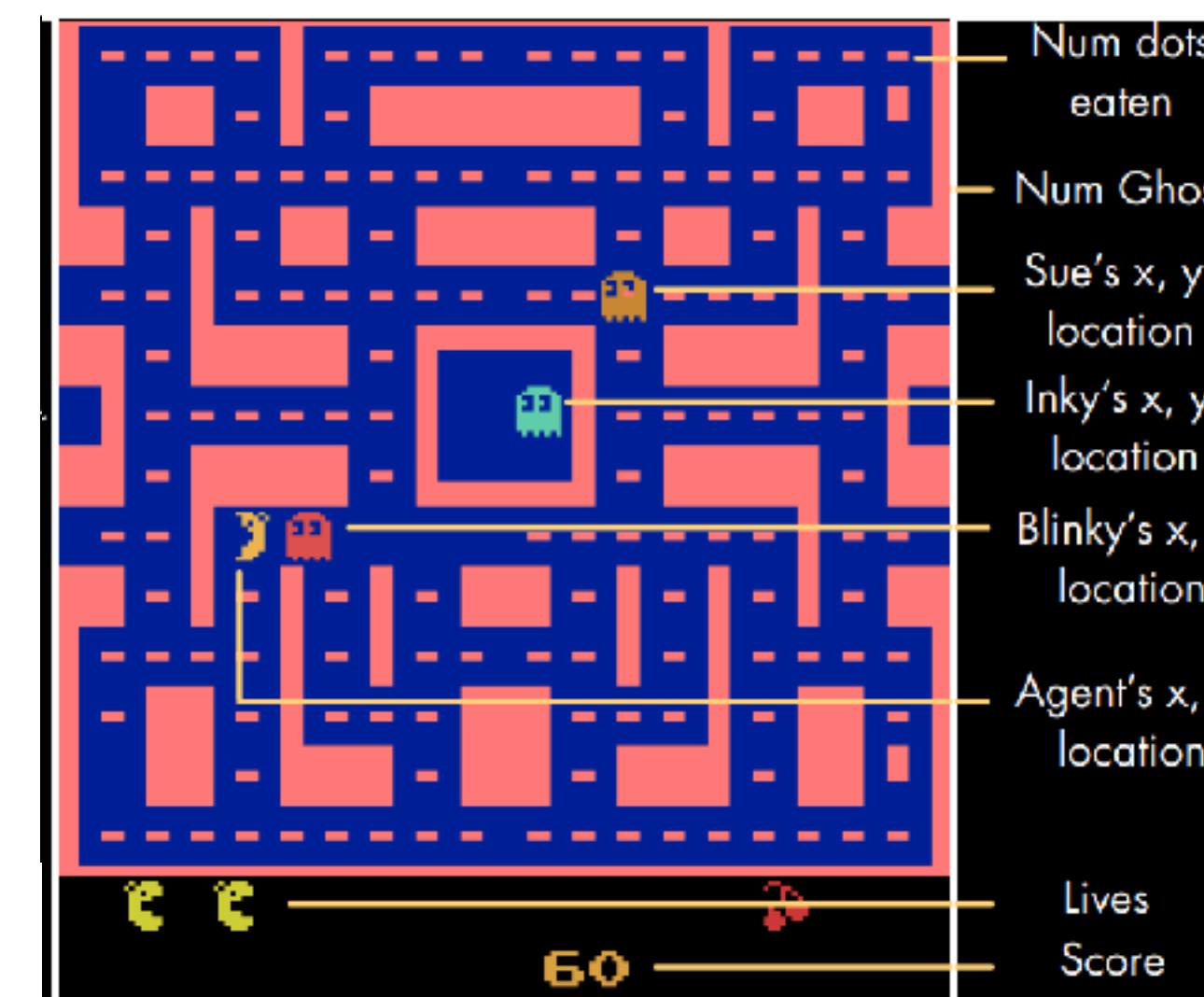
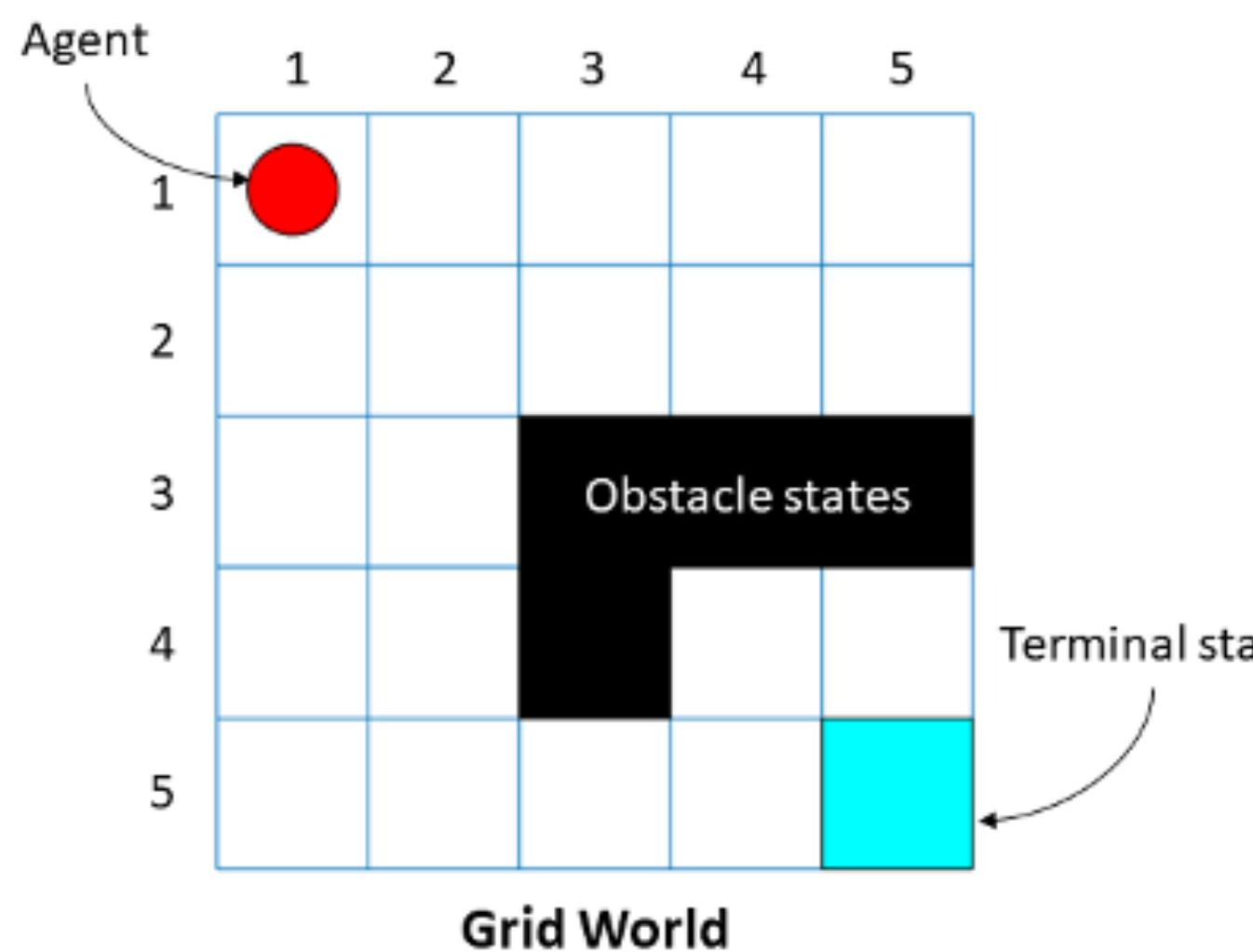
# Environment

Markov Decision Process (MDP)

- Simplifying assumption that the system is fully defined by only the previous state (i.e., Markov Principle):  $P(s_{t+1} | s_t, a_t)$

What are the states?

- Discrete locations, pixels on a screen, a set of feature values, etc...



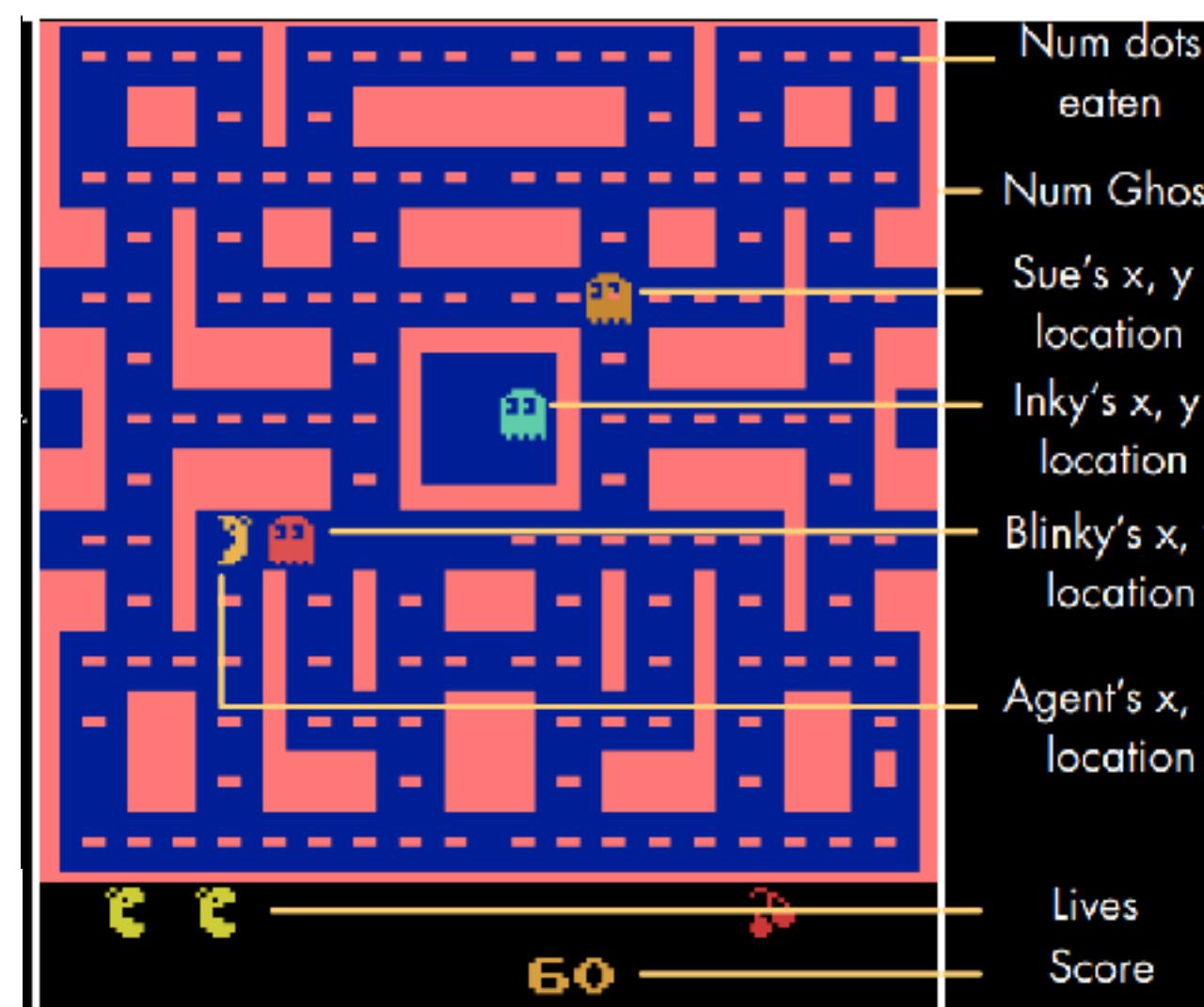
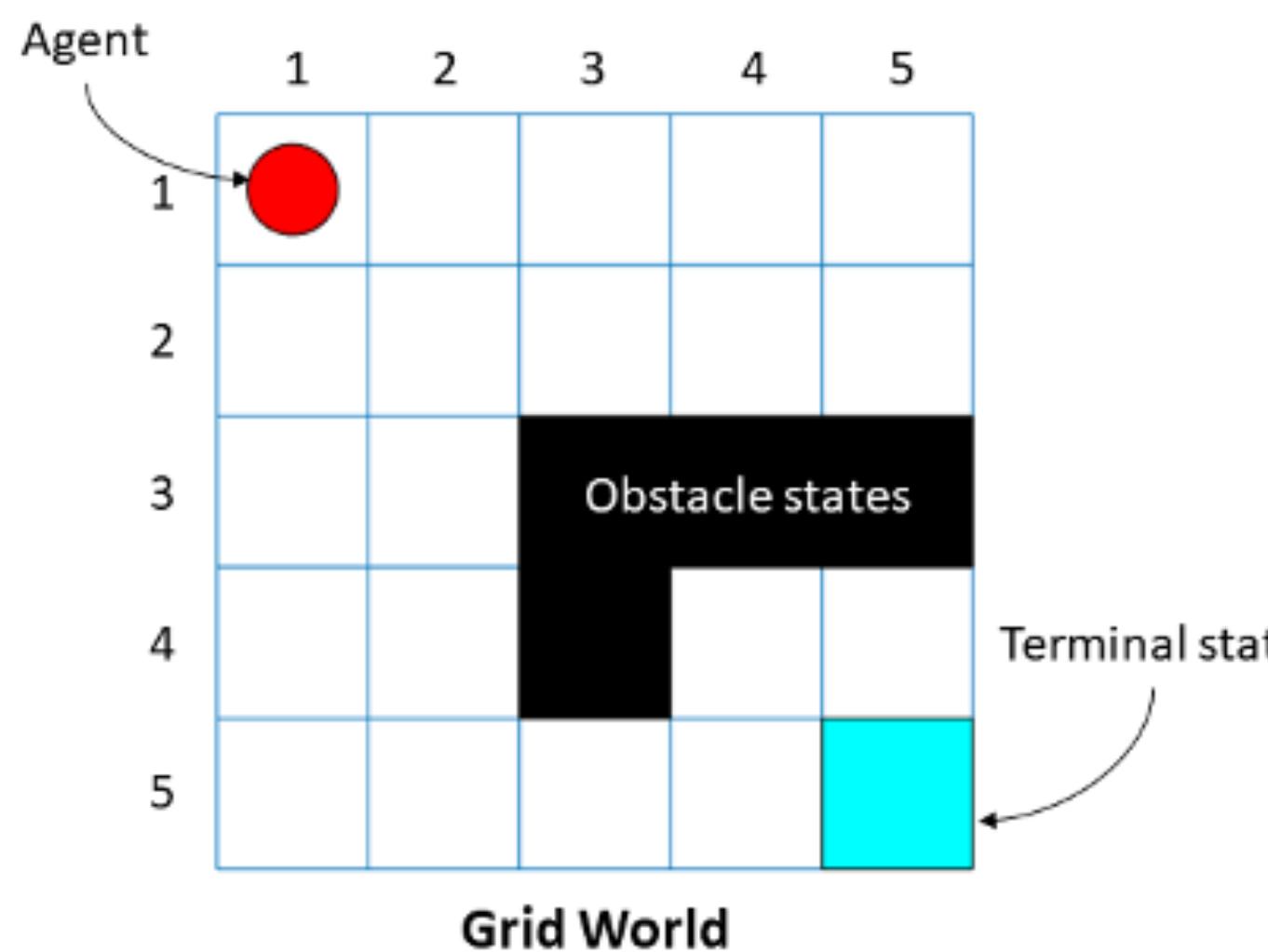
# Environment

## Markov Decision Process (MDP)

- Simplifying assumption that the system is fully defined by only the previous state (i.e., Markov Principle):  $P(s_{t+1} | s_t, a_t)$

What are the states?

- Discrete locations, pixels on a screen, a set of feature values, etc...



Partially Observable MDP (POMDP)



# Agent

# Agent

- **Experiences Rewards**
- **Learns a Policy**

# Agent

- **Experiences Rewards**

- How good is a given state?  $r_t = R(s_t)$

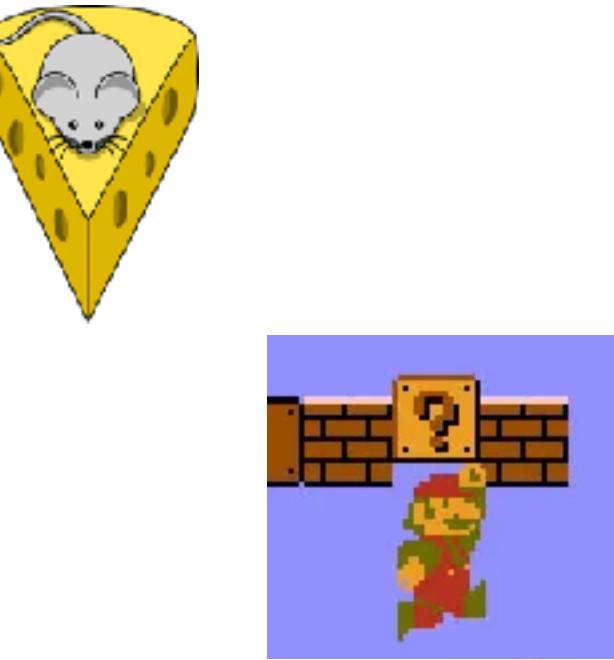


- **Learns a Policy**

# Agent

- **Experiences Rewards**

- How good is a given state?  $r_t = R(s_t)$
- How good is a state-action pair?  $r_t = R(s_t, a_t)$



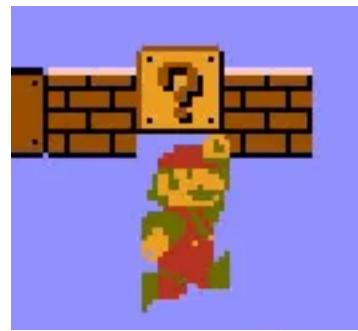
- **Learns a Policy**

# Agent

- **Experiences Rewards**

- How good is a given state?  $r_t = R(s_t)$
- How good is a state-action pair?  $r_t = R(s_t, a_t)$
- How good is a *trajectory*  $\tau = (s_0, a_0, s_1, a_1, \dots)$ :  
$$R(\tau) = \sum_{t=0}^{\infty} \gamma^k r_t$$
 where  $\gamma \in [0,1]$  is the temporal discount

- **Learns a Policy**



# Agent

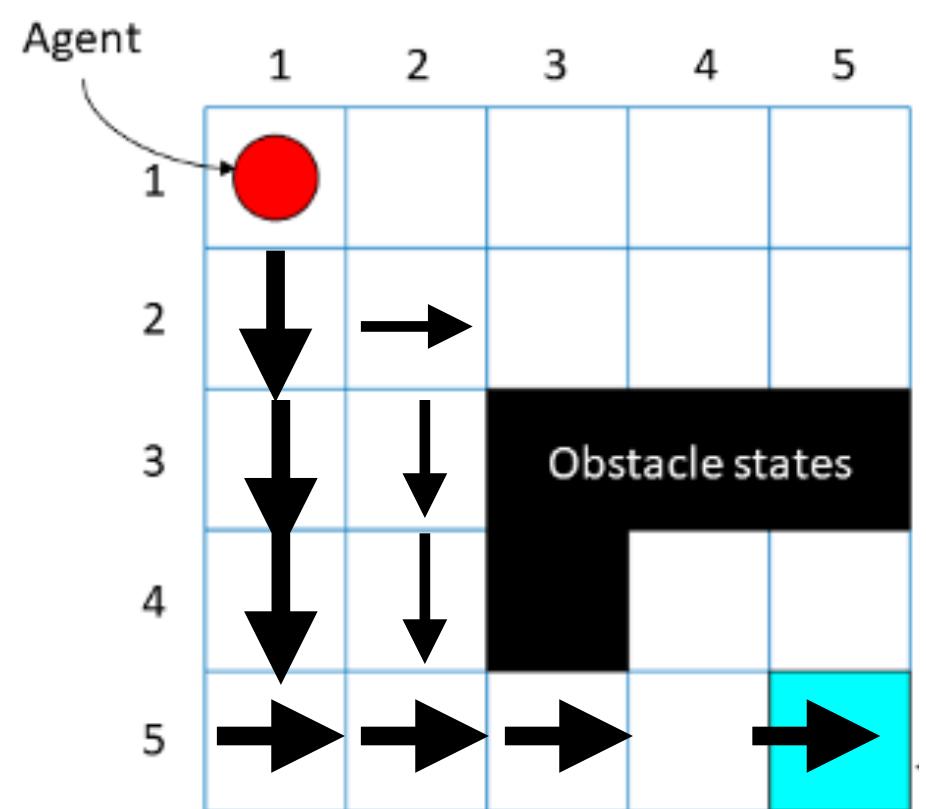
- **Experiences Rewards**

- How good is a given state?  $r_t = R(s_t)$
- How good is a state-action pair?  $r_t = R(s_t, a_t)$
- How good is a *trajectory*  $\tau = (s_0, a_0, s_1, a_1, \dots)$ :

$$R(\tau) = \sum_{t=0}^{\infty} \gamma^k r_t \text{ where } \gamma \in [0,1] \text{ is the temporal discount}$$

- **Learns a Policy**

- $\pi$  defines how to act, where  $\pi(a | s)$  is the probability of selecting action  $a$  in state  $s$
- sample trajectories from the policy  $\tau \sim \pi$



# The RL Problem

# The RL Problem

**Select a policy  $\pi^*$  that maximizes expected rewards**

# The RL Problem

**Select a policy  $\pi^*$  that maximizes expected rewards**

For this, we need to define a **value function**:

- The expected discounted returns under a policy:

$$V_\pi(s) = \mathbb{E}_{\tau \sim \pi}[R(\tau) \mid s_0 = s]$$

# The RL Problem

**Select a policy  $\pi^*$  that maximizes expected rewards**

For this, we need to define a **value function**:

- The expected discounted returns under a policy:

$$V_\pi(s) = \mathbb{E}_{\tau \sim \pi}[R(\tau) \mid s_0 = s]$$

- Let's unpack that a bit:

$$V_\pi(s) = \sum_a \pi(a \mid s) \sum_{s'} P(s' \mid s, a) [R(s', a) + \gamma V_\pi(s')]$$

- The expectation over a policy is equivalent to summing up all actions and their new state transitions, weighted by their probability
- Then we sum up the (discounted) future expected rewards

# The RL Problem

Select a policy  $\pi^*$  that maximizes expected rewards

For this, we need to define a **value function**:

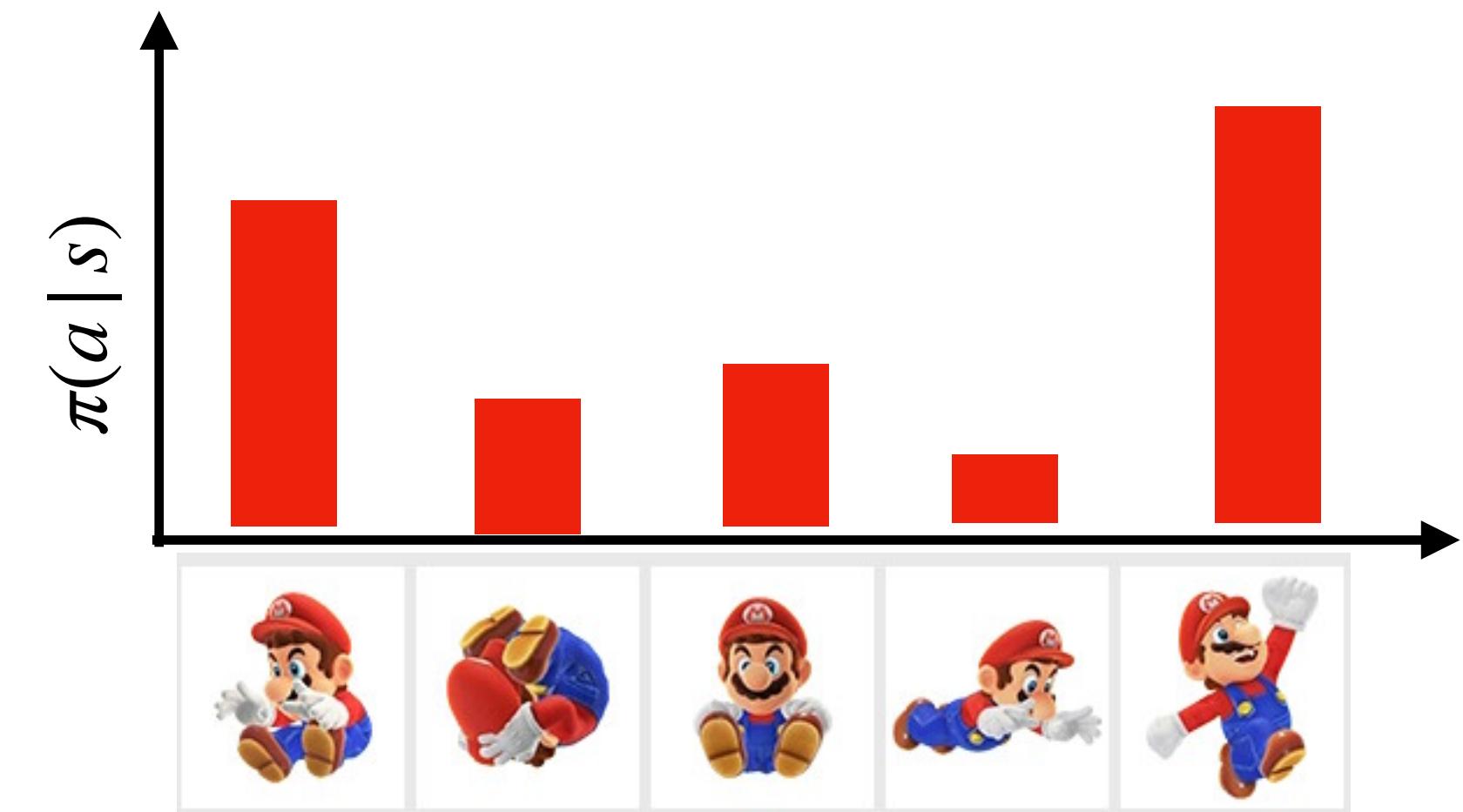
- The expected discounted returns under a policy:

$$V_\pi(s) = \mathbb{E}_{\tau \sim \pi}[R(\tau) \mid s_0 = s]$$

- Let's unpack that a bit:

$$V_\pi(s) = \sum_a \pi(a \mid s) \sum_{s'} P(s' \mid s, a) [R(s', a) + \gamma V_\pi(s')]$$

- The expectation over a policy is equivalent to summing up all actions and their new state transitions, weighted by their probability
- Then we sum up the (discounted) future expected rewards



# The RL Problem

Select a policy  $\pi^*$  that maximizes expected rewards

For this, we need to define a **value function**:

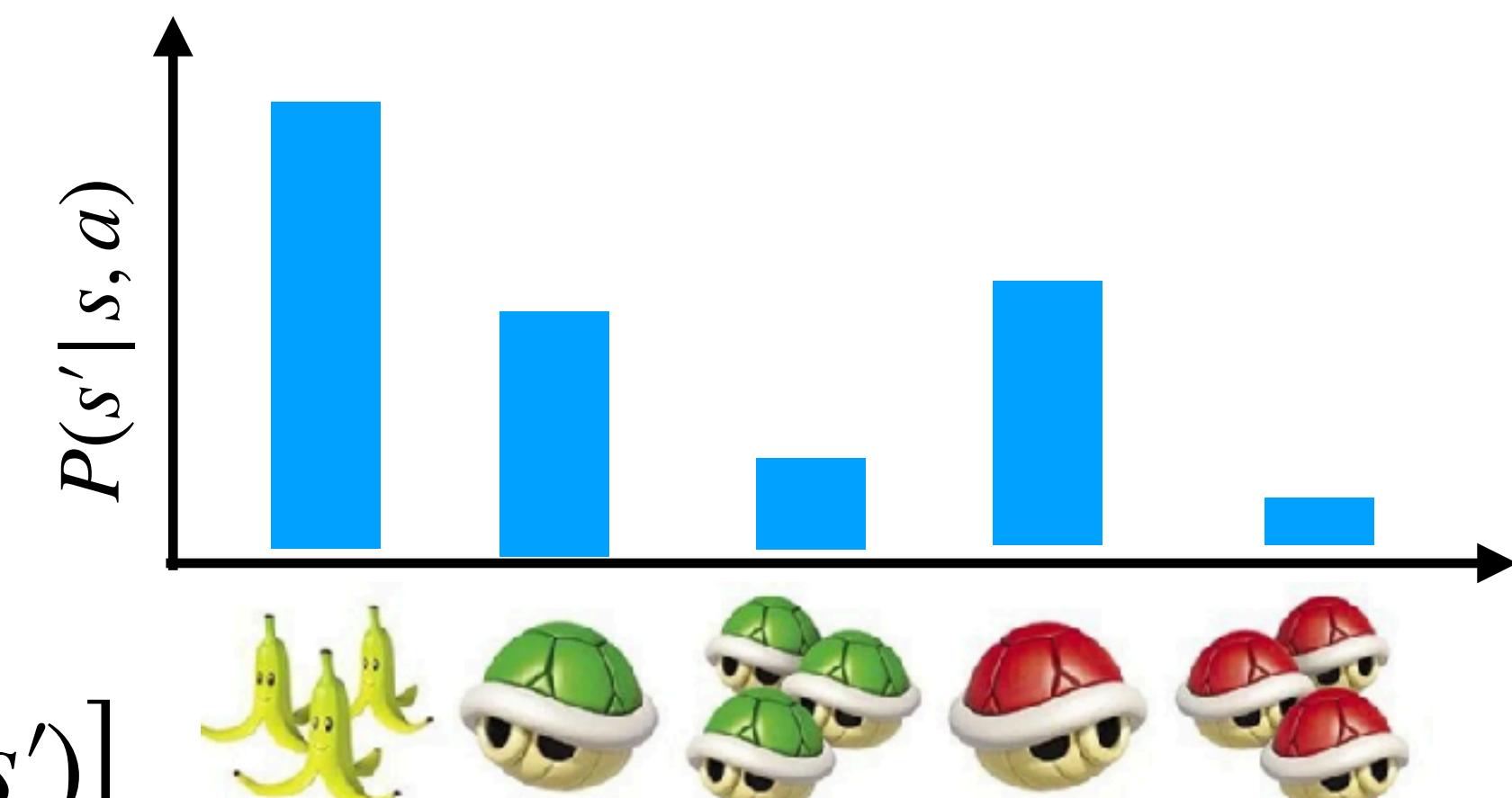
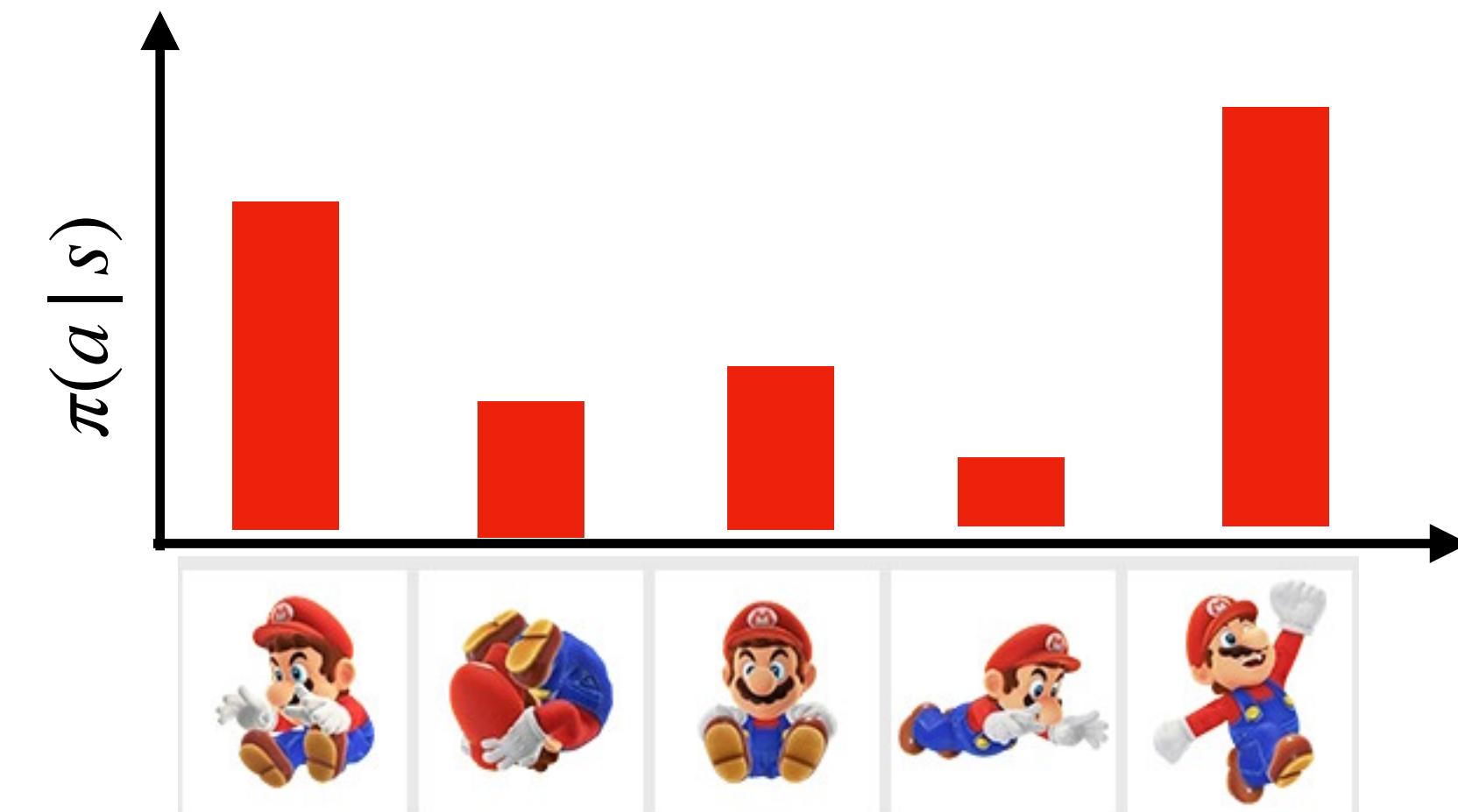
- The expected discounted returns under a policy:

$$V_\pi(s) = \mathbb{E}_{\tau \sim \pi}[R(\tau) \mid s_0 = s]$$

- Let's unpack that a bit:

$$V_\pi(s) = \sum_a \pi(a \mid s) \sum_{s'} P(s' \mid s, a) [R(s', a) + \gamma V_\pi(s')]$$

- The expectation over a policy is equivalent to summing up all actions and their new state transitions, weighted by their probability
- Then we sum up the (discounted) future expected rewards



# The RL Problem

Select a policy  $\pi^*$  that maximizes expected rewards

For this, we need to define a **value function**:

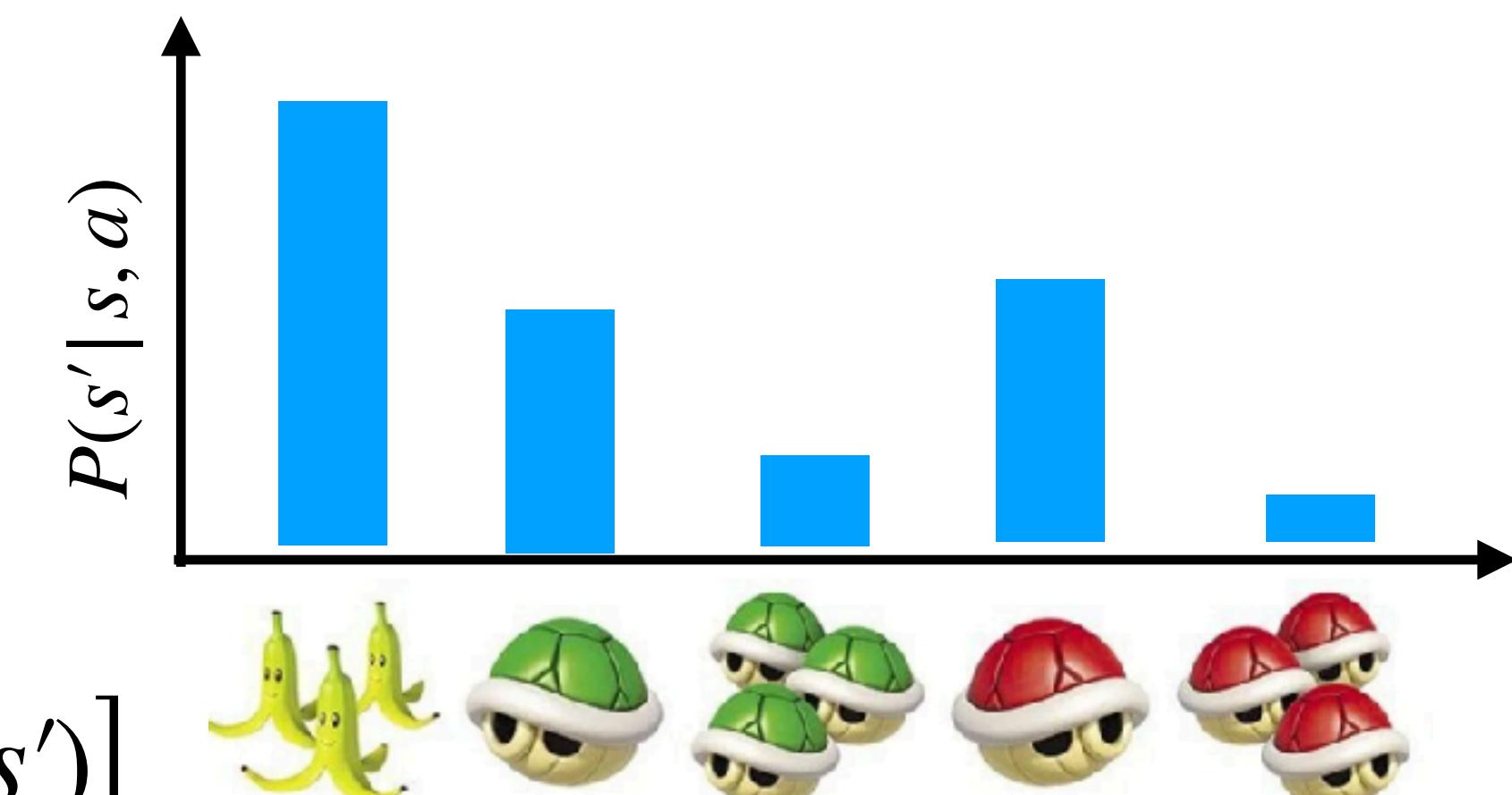
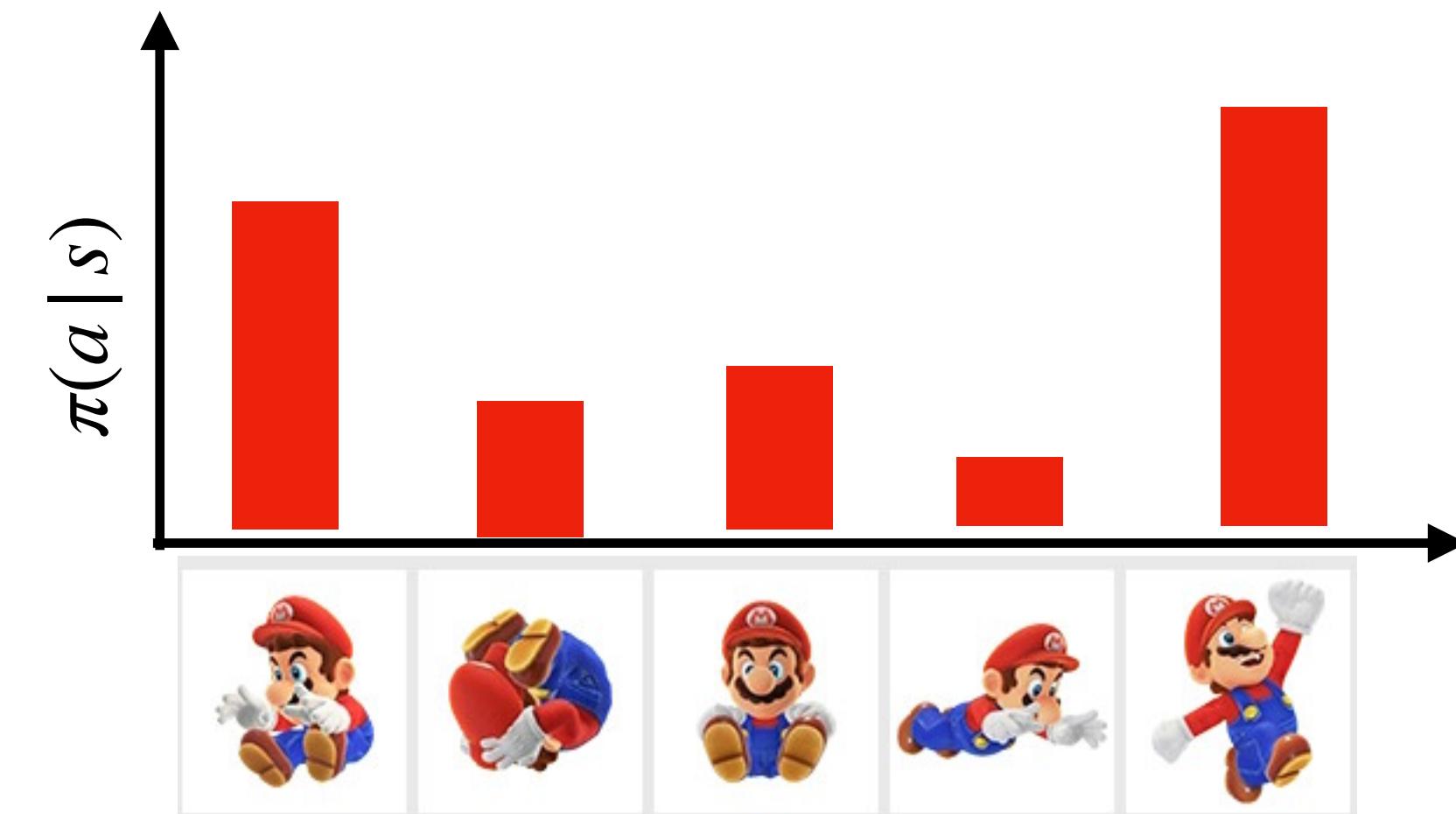
- The expected discounted returns under a policy:

$$V_\pi(s) = \mathbb{E}_{\tau \sim \pi}[R(\tau) \mid s_0 = s]$$

- Let's unpack that a bit:

$$V_\pi(s) = \sum_a \pi(a \mid s) \sum_{s'} P(s' \mid s, a) [R(s', a) + \gamma V_\pi(s')]$$

- The expectation over a policy is equivalent to summing up all actions and their new state transitions, weighted by their probability
- Then we sum up the (discounted) future expected rewards



# Finding an optimal policy via Bellman Equations

- Bellman equations are a concept from dynamic programming that provide a recursive method for optimization:

$$V_{\pi}(s) = \sum_a \pi(a | s) \sum_{s'} P(s' | s, a) [R(s', a) + \gamma V_{\pi}(s')]$$

- Theoretically, we can define an optimal value function:

$$V_*(s) = \max_a \sum_{s'} P(s' | s, a) [R(s, a) + \gamma V_*(s')]$$

- Optimal policy:

$$\pi_* = \arg \max_a V_*(s)$$

# Finding an optimal policy via Bellman Equations

- Bellman equations are a concept from dynamic programming that provide a recursive method for optimization:

$$V_{\pi}(s) = \sum_a \pi(a | s) \sum_{s'} P(s' | s, a) [R(s', a) + \gamma V_{\pi}(s')]$$

recursive

- Theoretically, we can define an optimal value function:

$$V_*(s) = \max_a \sum_{s'} P(s' | s, a) [R(s, a) + \gamma V_*(s')]$$

- Optimal policy:

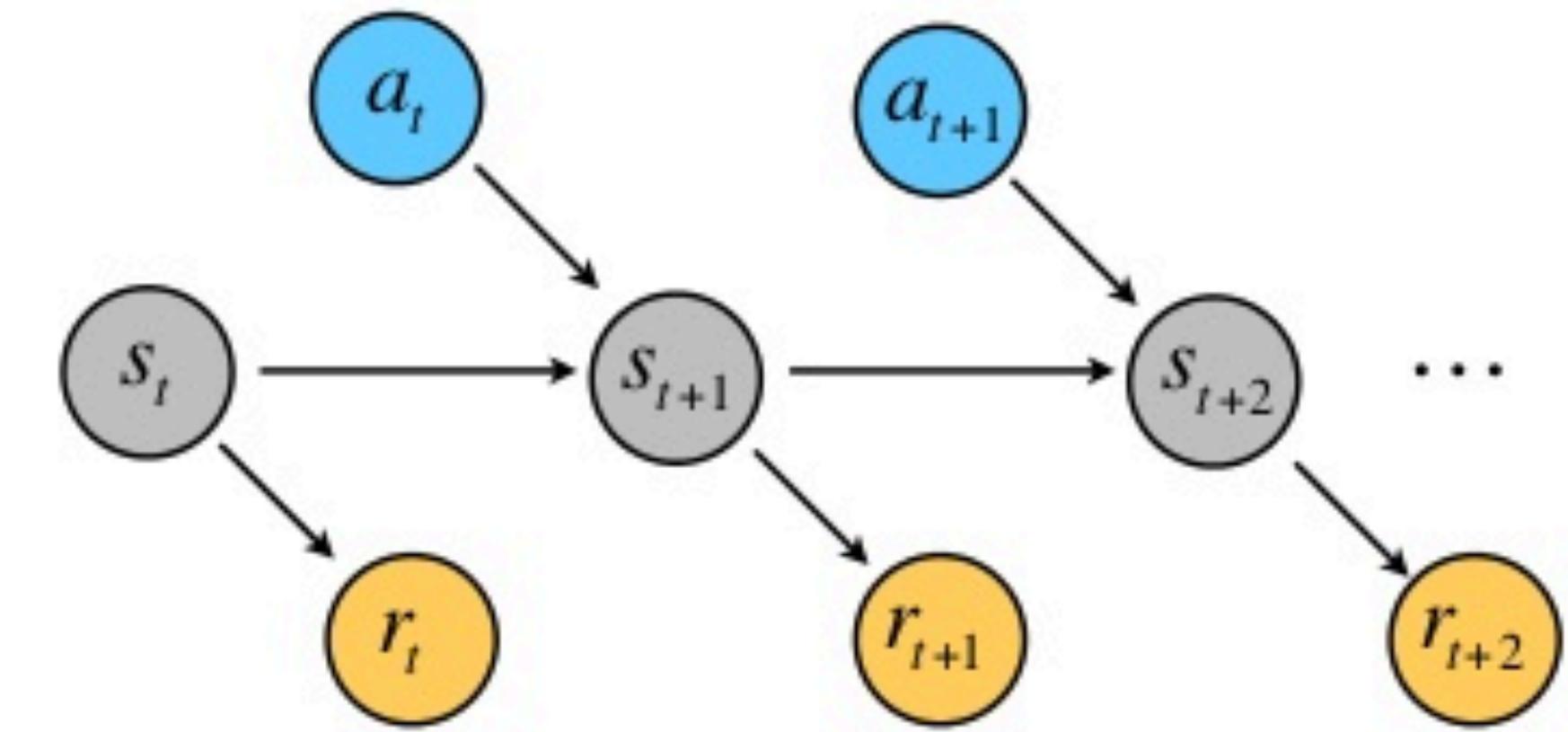
$$\pi_* = \arg \max_a V_*(s)$$

# Finding an optimal policy via Bellman Equations

- Bellman equations are a concept from dynamic programming that provide a recursive method for optimization:

$$V_{\pi}(s) = \sum_a \pi(a | s) \sum_{s'} P(s' | s, a) [R(s', a) + \gamma V_{\pi}(s')]$$

recursive



- Theoretically, we can define an optimal value function:

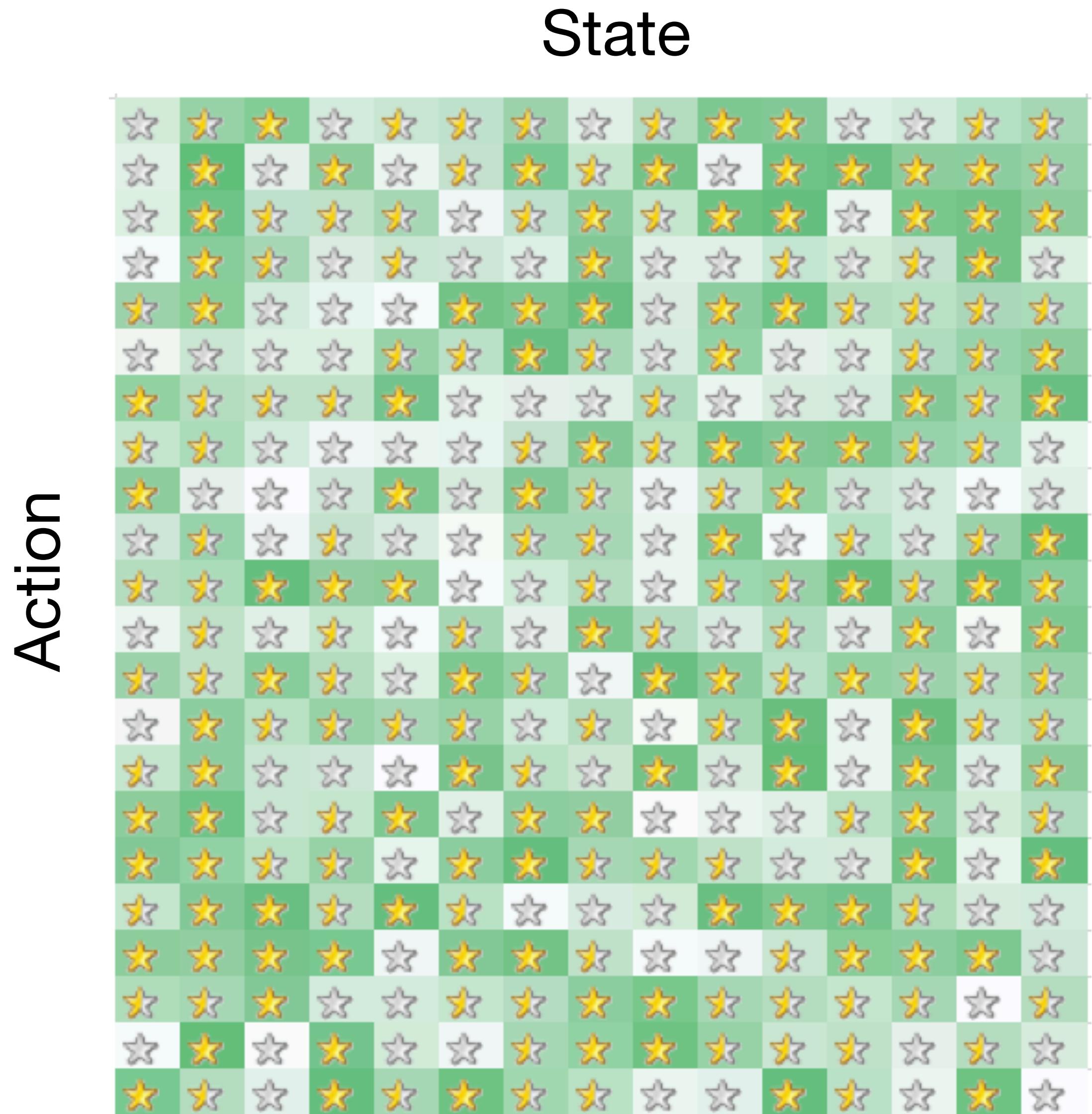
$$V_*(s) = \max_a \sum_{s'} P(s' | s, a) [R(s, a) + \gamma V_*(s')]$$

- Optimal policy:

$$\pi_* = \arg \max_a V_*(s)$$

# Tabular methods

- Based on methods from Dynamic programming (Bellman, 1957), Tabular methods were first proposed as solutions for RL problems by Minsky (1961)
- Think of a giant lookup table, where we store a value representation
- Value iteration and policy iteration are examples
- Caveat: solutions require repeat visits to each state, which is infeasible in most real-world problems



# Tabular methods

- Based on methods from Dynamic programming (Bellman, 1957), Tabular methods were first proposed as solutions for RL problems by Minsky (1961)
- Think of a giant lookup table, where we store a value representation
- Value iteration and policy iteration are examples
- Caveat: solutions require repeat visits to each state, which is infeasible in most real-world problems

| Action |   |   |   |   |   |   |   |   |   | State |   |   |   |   |   |   |   |   |   |
|--------|---|---|---|---|---|---|---|---|---|-------|---|---|---|---|---|---|---|---|---|
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |
| ?      | ? | ? | ? | ? | ? | ? | ? | ? | ? | ?     | ? | ? | ? | ? | ? | ? | ? | ? | ? |

# Value iteration

Iteratively visit all states and update the value function until a “good enough” solution has been reached.

1. Initialize the value function as  $V_0(s) = 0$  for all states

2. For all  $s$  in  $\mathcal{S}$ :

$$V_{k+1}(s) = \max_{a \in A} \sum_{s'} P(s' | s, a)[R(s, a) + \gamma V_k(s')]$$

until  $\max_{s \in \mathcal{S}} |V_k(s) - V_{k-1}(s)| < \theta$  Bellman residual

# Value iteration

Iteratively visit all states and update the value function until a “good enough” solution has been reached.

1. Initialize the value function as  $V_0(s) = 0$  for all states

2. For all  $s$  in  $\mathcal{S}$ :

$$V_{k+1}(s) = \max_{a \in A} \sum_{s'} P(s' | s, a)[R(s, a) + \gamma V_k(s')]$$

until  $\max_{s \in \mathcal{S}} |V_k(s) - V_{k-1}(s)| < \theta$  Bellman residual

$V_k$  converges on  $V_*$  as  $k \rightarrow \infty$ , and perhaps sooner, but with many costly sweeps through the state space

# Policy Iteration

Alternate between evaluating a policy and then improving the policy.

Start with  $\pi_0$  (typically a random policy), and then iterate for all  $s \in \mathcal{S}$  in each step

- **Policy Evaluation**

$$V_{\pi_k}(s) = \mathbb{E}_{\pi_k} \left[ R(s', a) + \gamma V_{\pi_k}(s') \right]$$

- **Policy Improvement**

$$\pi_{k+1} = \arg \max_a \sum_{s'} P(s' | s, a) \left[ R(s, a) + \gamma V_{\pi_k} \right]$$

# Policy Iteration

Alternate between evaluating a policy and then improving the policy.

Start with  $\pi_0$  (typically a random policy), and then iterate for all  $s \in \mathcal{S}$  in each step

- **Policy Evaluation**

$$V_{\pi_k}(s) = \mathbb{E}_{\pi_k} \left[ R(s', a) + \gamma V_{\pi_k}(s') \right]$$

- **Policy Improvement**

$$\pi_{k+1} = \arg \max_a \sum_{s'} P(s' | s, a) \left[ R(s, a) + \gamma V_{\pi_k} \right]$$

Policy can converge faster than value function, but still requires visiting all states  $2n$  times and lacks convergence guarantees



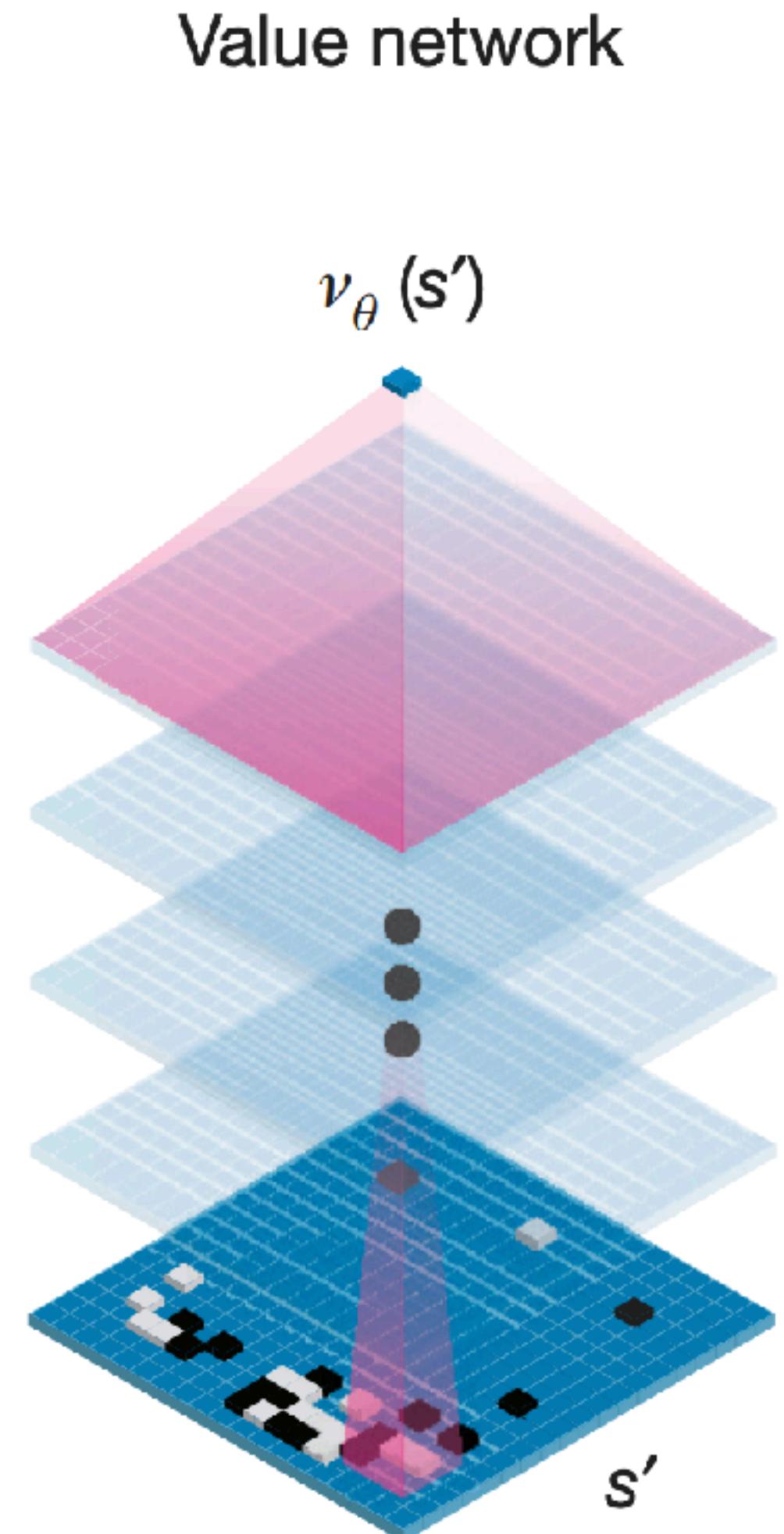
# Value function approximation

Instead of learning the value for each state independently, learn a value function mapping each state to some value:

$$f: s \in \mathcal{S} \rightarrow V_{\theta}(s)$$

... a variety of methods are available including:

- linear function approximation (e.g., regression)
- Neural networks
- Gaussian process regression (non-parametric)



# Value function approximation

Instead of learning the value for each state independently, learn a value function mapping each state to some value:

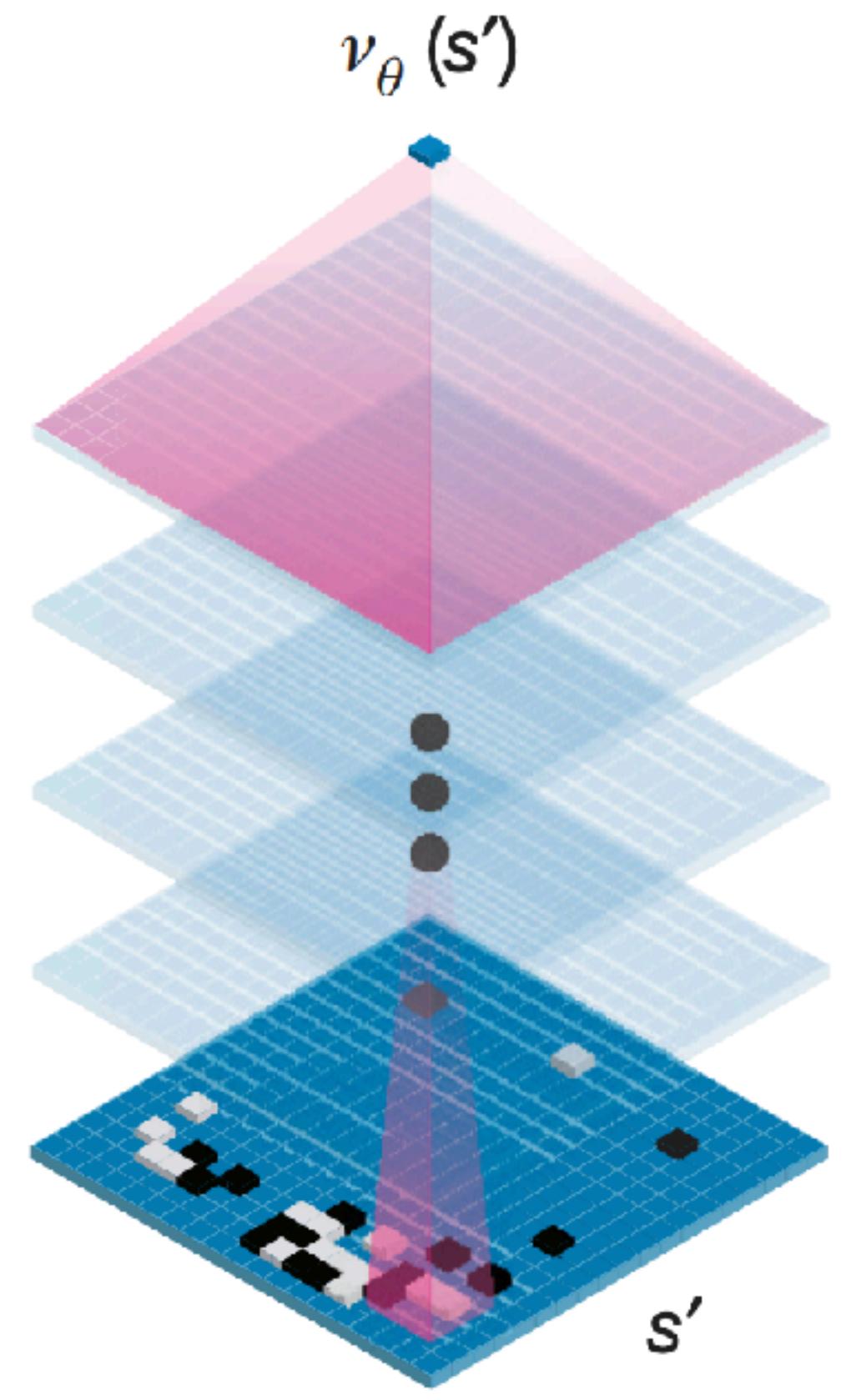
$$f: s \in \mathcal{S} \rightarrow V_{\theta}(s)$$

... a variety of methods are available including:

- linear function approximation (e.g., regression)
- Neural networks
- Gaussian process regression (non-parametric)



Value network



Silver et al. (2016)

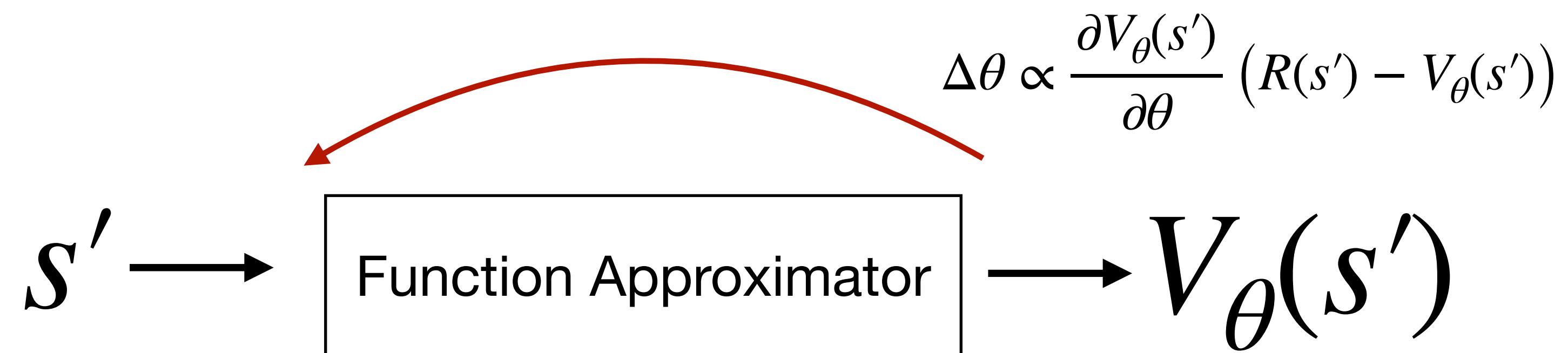
# Value function approximation

Instead of learning the value for each state independently, learn a value function mapping each state to some value:

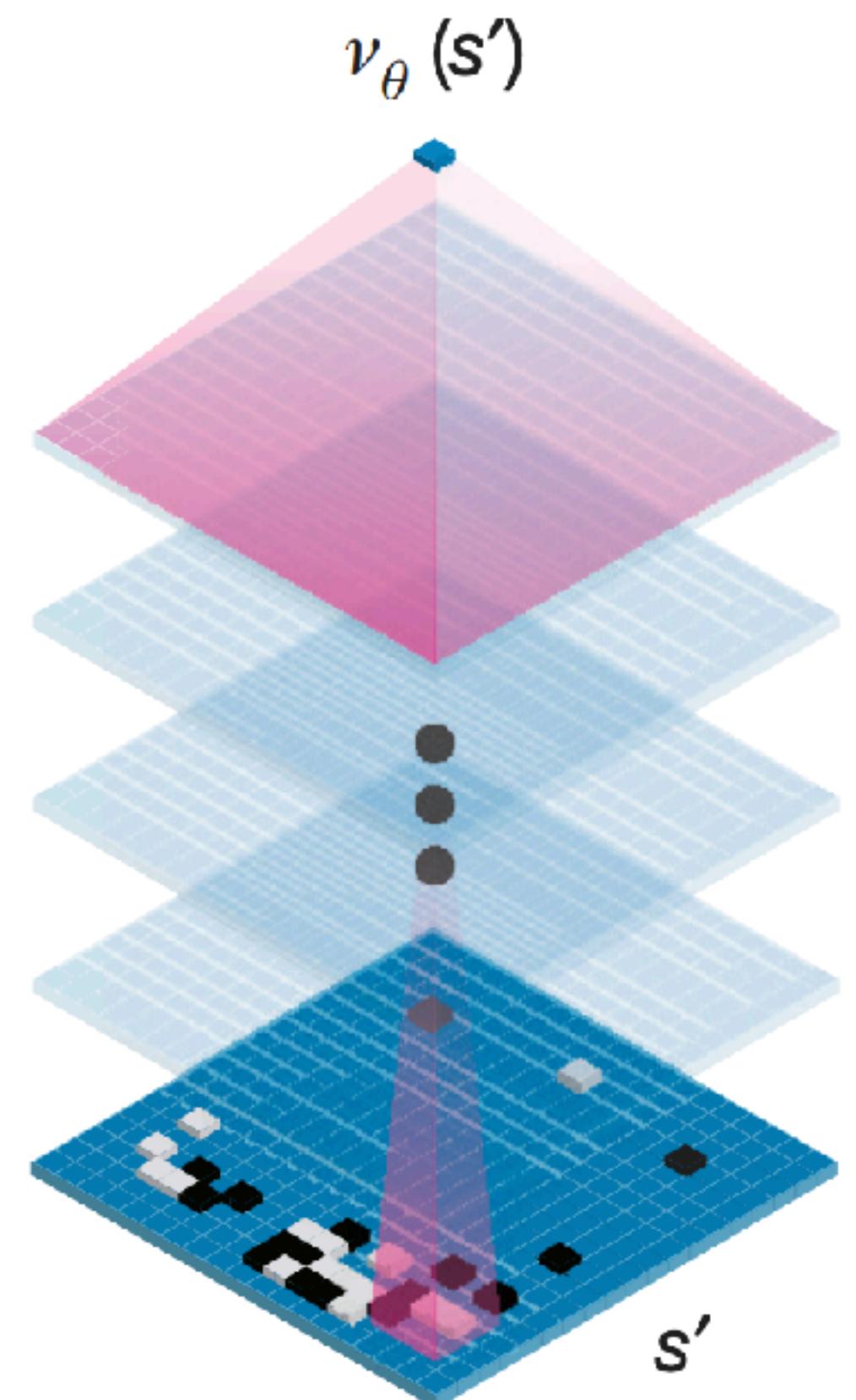
$$f: s \in \mathcal{S} \rightarrow V_{\theta}(s)$$

... a variety of methods are available including:

- linear function approximation (e.g., regression)
- Neural networks
- Gaussian process regression (non-parametric)



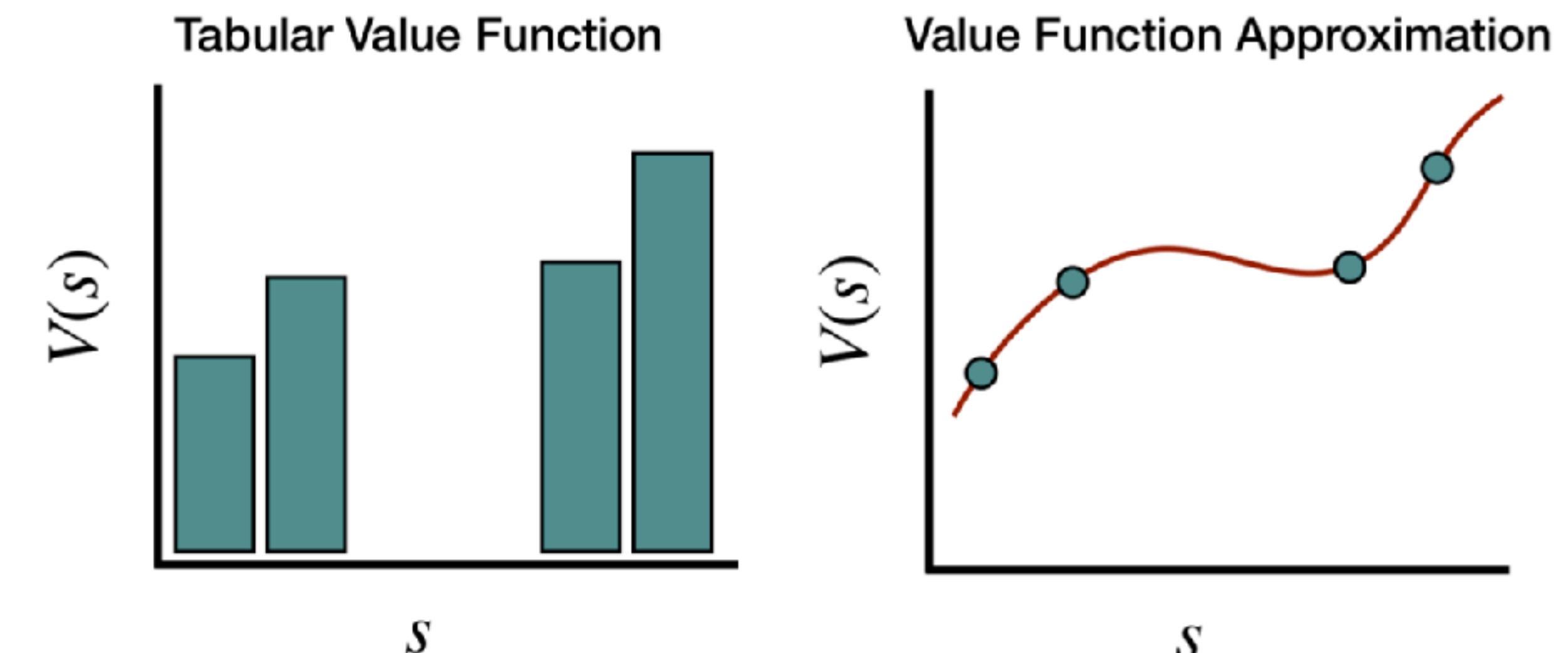
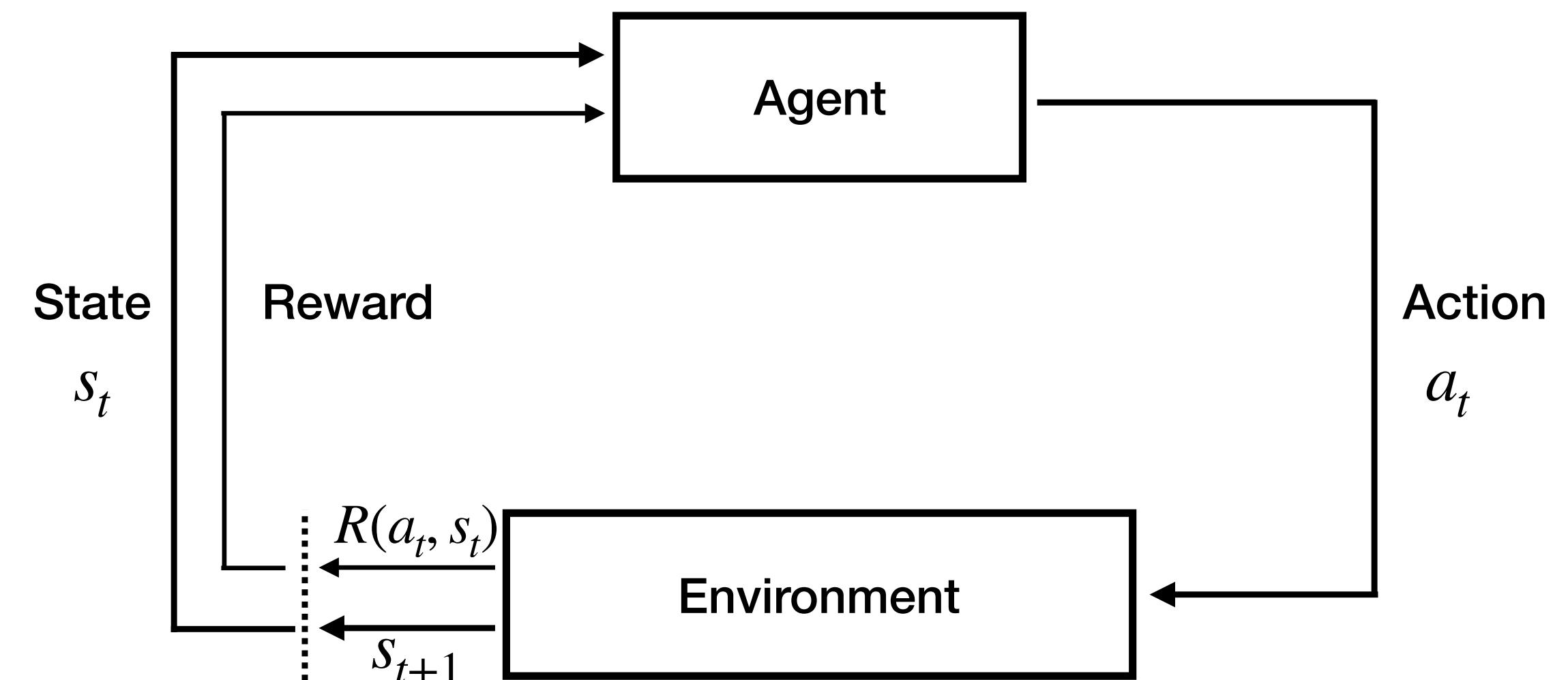
Value network



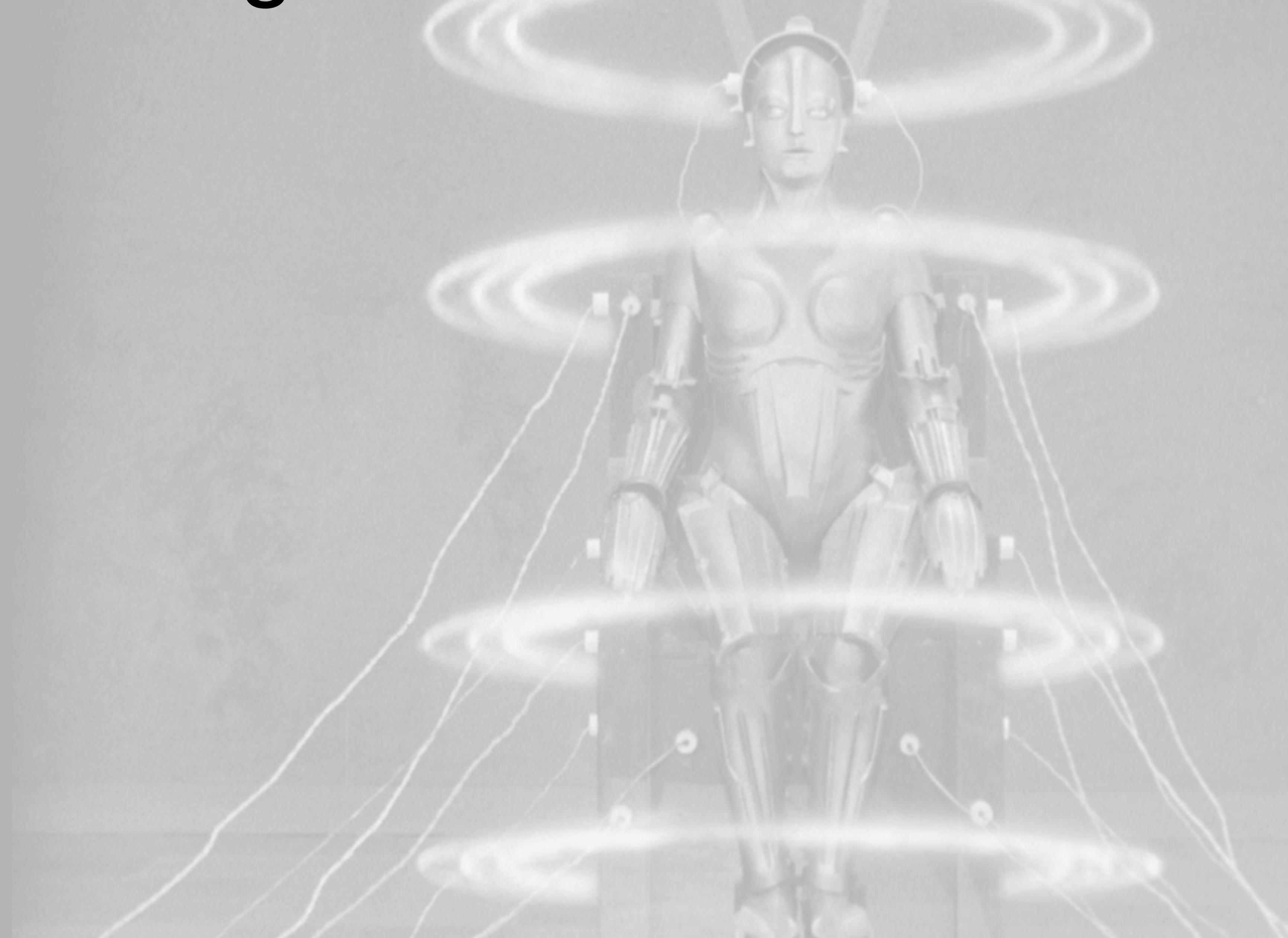
Silver et al. (2016)

# Quick recap

- RL framework defines interactions between an **agent** and the **environment**
  - The environment defines the transitions between states and provides rewards
  - The agent learns a value function and then turns this into a policy
- Traditional solutions to RL problems can be broadly classified into either **tabular methods** or **value function approximation**



# Modeling Human Learners



# Multi-armed bandit

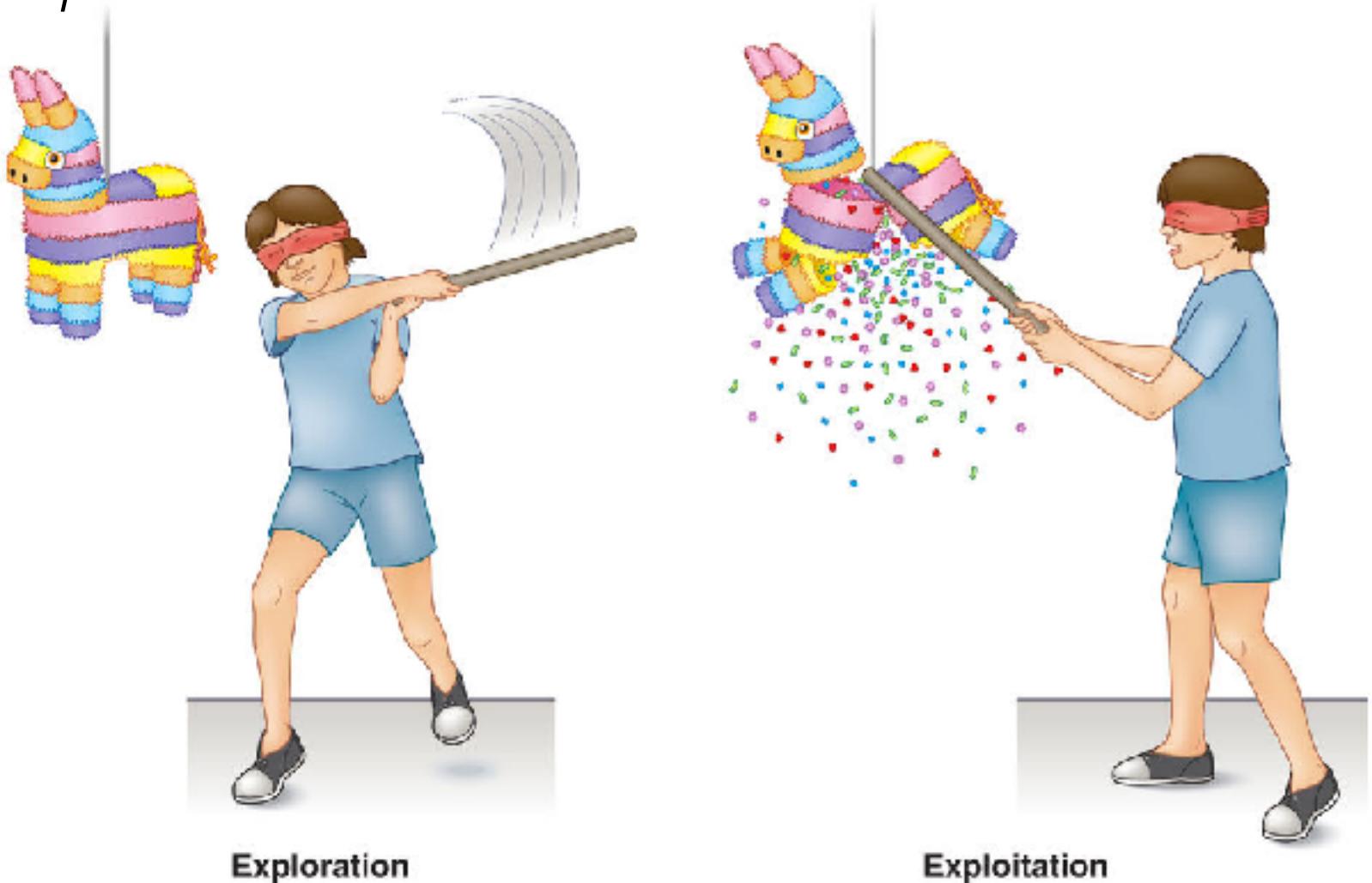
# Multi-armed bandit

- Colorful metaphor for a row of slot machines



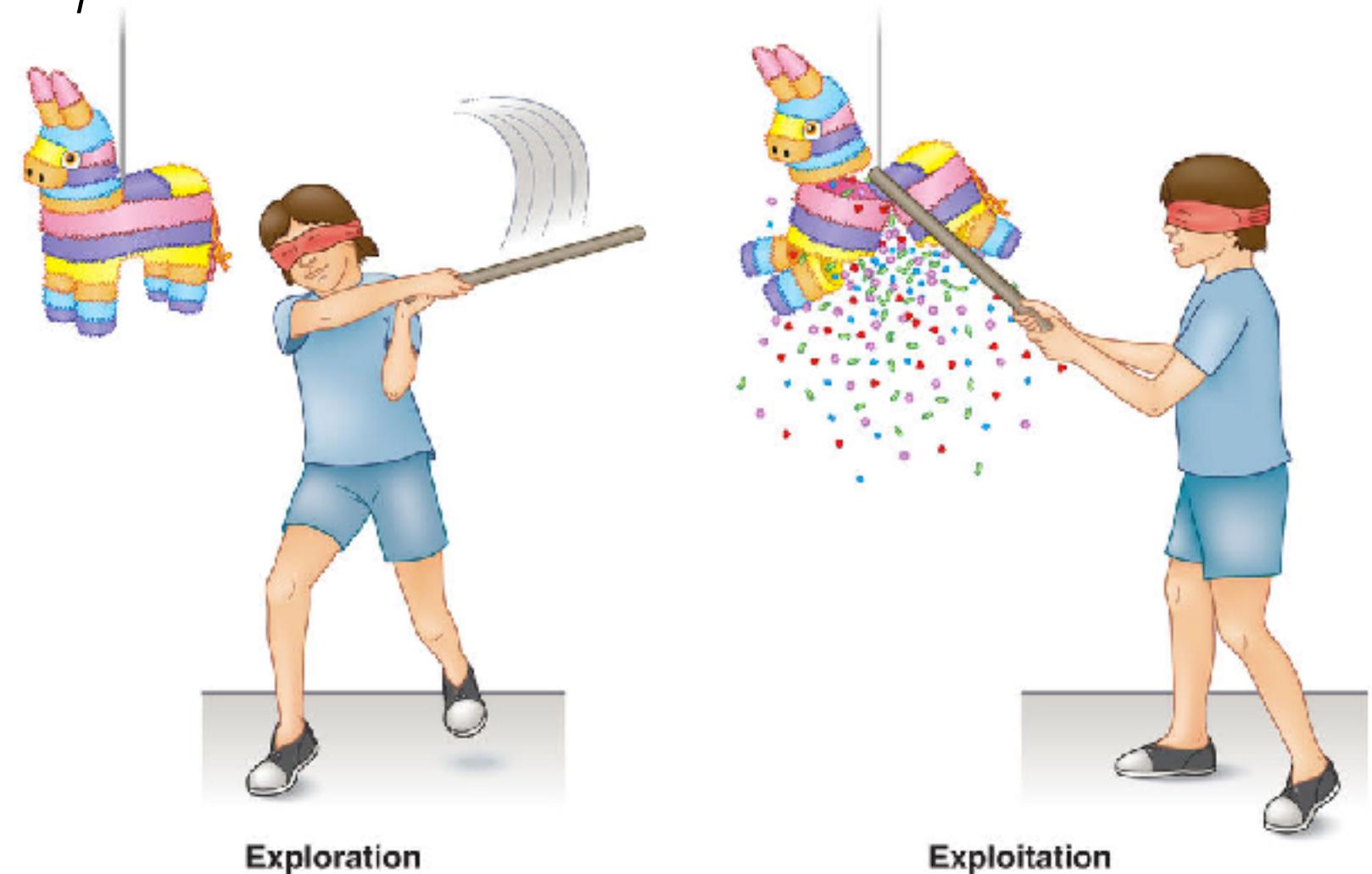
# Multi-armed bandit

- Colorful metaphor for a row of slot machines
- One of the simplest RL problems for studying the *exploration-exploitation dilemma*:
  - **exploring** untried options to acquire information
  - **exploiting** known options to acquire immediate rewards



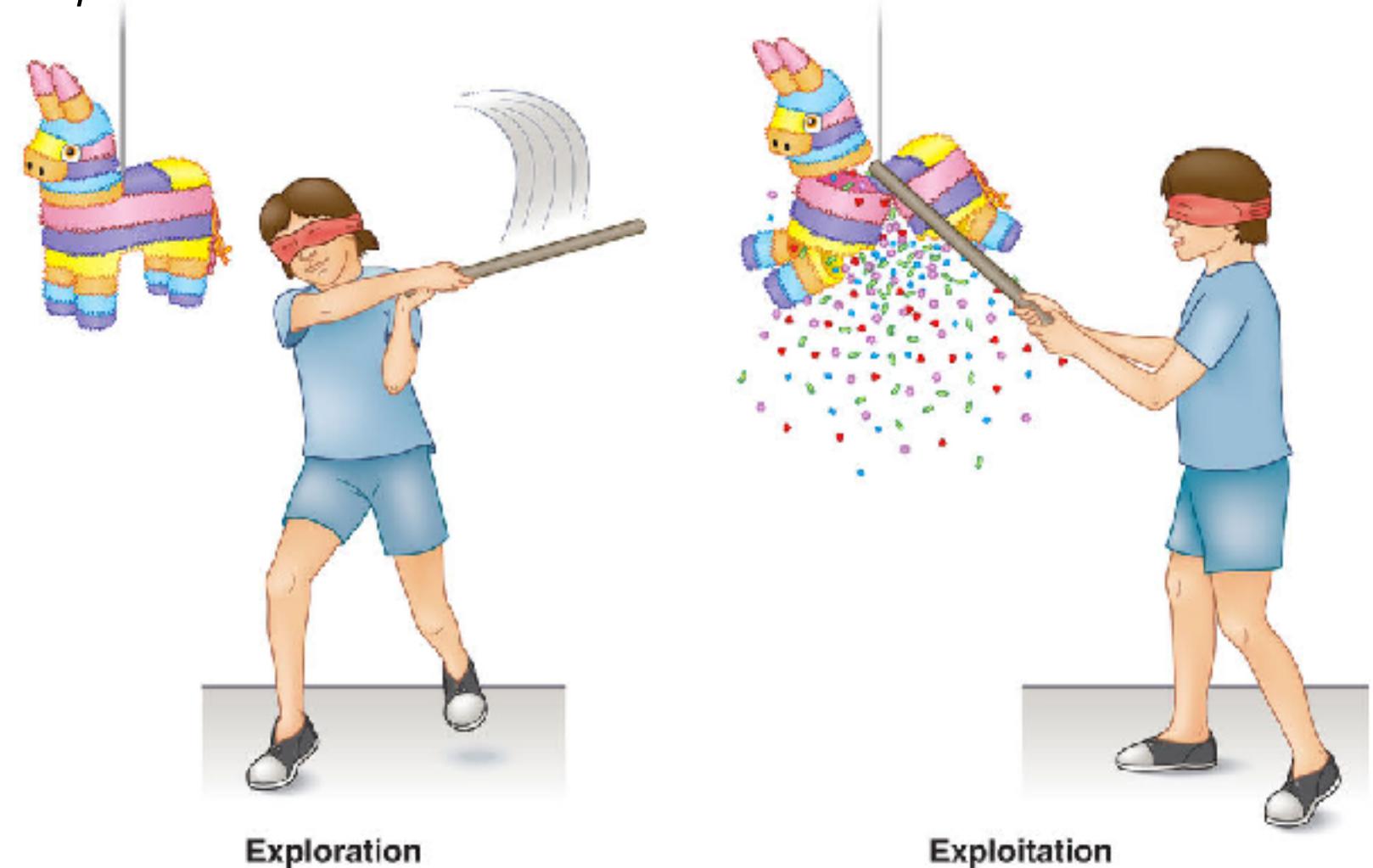
# Multi-armed bandit

- Colorful metaphor for a row of slot machines
- One of the simplest RL problems for studying the *exploration-exploitation dilemma*:
  - **exploring** untried options to acquire information
  - **exploiting** known options to acquire immediate rewards
- We can think of it as an MDP with a single state:
  - No need to discount future states
  - Instead of state-value function  $V(s)$  we can represent the value of an action  $V(a)$  or more properly, using a action-value function  $Q(s, a)$



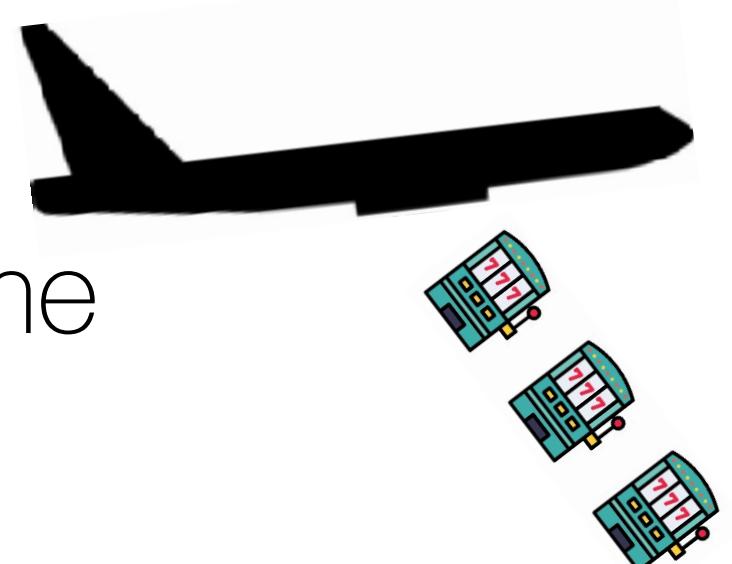
# Multi-armed bandit

- Colorful metaphor for a row of slot machines
- One of the simplest RL problems for studying the *exploration-exploitation dilemma*:
  - **exploring** untried options to acquire information
  - **exploiting** known options to acquire immediate rewards
- We can think of it as an MDP with a single state:
  - No need to discount future states
  - Instead of state-value function  $V(s)$  we can represent the value of an action  $V(a)$  or more properly, using a action-value function  $Q(s, a)$
  - Despite its simplicity, optimal solutions are generally only available with infinite time



# Multi-armed bandit

- Colorful metaphor for a row of slot machines
- One of the simplest RL problems for studying the *exploration-exploitation dilemma*:
  - **exploring** untried options to acquire information
  - **exploiting** known options to acquire immediate rewards
- We can think of it as an MDP with a single state:
  - No need to discount future states
  - Instead of state-value function  $V(s)$  we can represent the value of an action  $V(a)$  or more properly, using a action-value function  $Q(s, a)$
- Despite its simplicity, optimal solutions are generally only available with infinite time



# Brief aside: State-Value function vs. Action-value function

So far, I've focused on describing state-value functions:

$$V_{\pi}(s) = \mathbb{E}_{\tau \sim \pi}[R(\tau) \mid s_0 = s]$$

However, you can also describe a RL model using the action-value function:

$$Q_{\pi}(s, a) = \mathbb{E}_{\tau \sim \pi}[R(\tau) \mid s_0 = s, a_0 = a]$$

Both are equivalent under:

$$V_{\pi}(s) = \sum_{a \in A} \pi(a \mid s) * Q_{\pi}(s, a)$$

# Rescorla-Wagner and Q-learning: the Delta Rule

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972) model**

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$

  
 $\delta_t$

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner** (1972) model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$

**Q-Learning** (Watkins, 1989) uses the same update rule

$$\delta_t$$

$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$

# Rescorla-Wagner and Q-learning: the Delta Rule

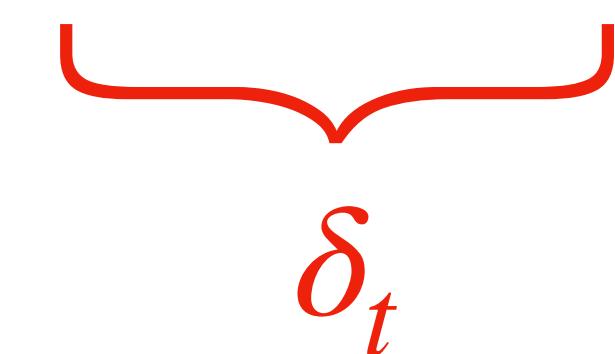
**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

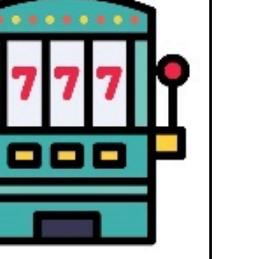
$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$


$$\delta_t$$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |   |   |
|---|---|---|---|
|   |  |  |  |
| Q | 0   | 0   | 0   |

$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

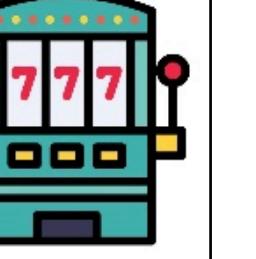
- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$



$\delta_t$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |   |   |
|---|---|---|---|
|   |  |  |  |
| Q | 0   | 0   | 0   |

$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$

0

0

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

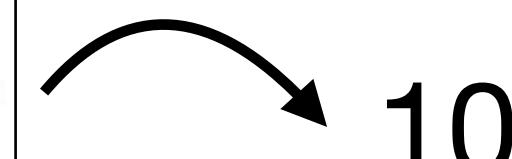
$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |   |   |
|---|---|---|---|
|   |   |   |   |
| Q | 0 | 0 | 0 |



$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$

0

0

$\delta_t$

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

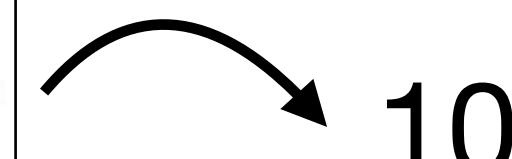
$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |   |   |
|---|---|---|---|
|   |   |   |   |
| Q | 0 | 0 | 0 |



$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$

0                          10                          0

$\delta_t$

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

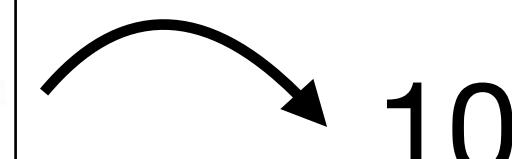
$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |   |   |
|---|---|---|---|
|   |   |   |   |
| Q | 0 | 0 | 0 |



$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$

0      1      10      0

$\delta_t$

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

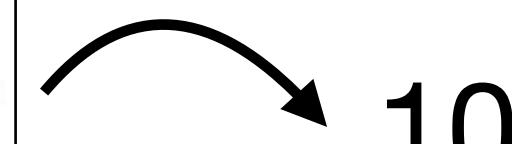
$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |   |    |
|---|---|---|----|
|   |   |   |    |
| Q | 0 | 0 | 10 |



$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$

  
 $\delta_t$

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

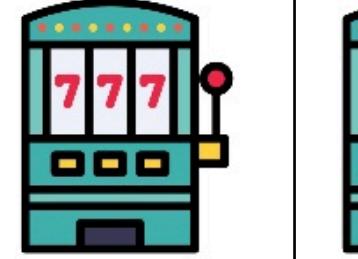
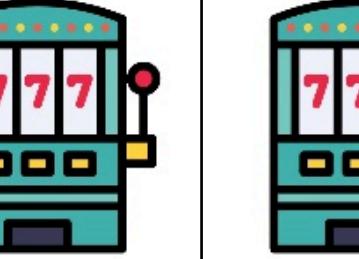
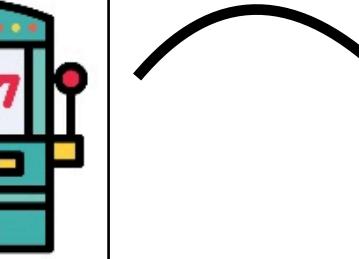
$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

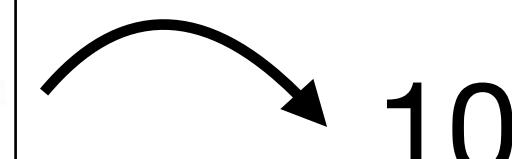
- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$


$$\delta_t$$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |  |   |
|---|---|--|---|
|   |  |  |  |
| Q | 0   | 0  | 10  |



$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$


# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |   |    |   |
|---|---|---|----|---|
|   |   |   |    | 7 |
| Q | 0 | 0 | 10 |   |

$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$

10      1      7      10

$\delta_t$

# Rescorla-Wagner and Q-learning: the Delta Rule

**Rescorla-Wagner (1972)** model

- Learning as an active process of making predictions about the world
- Predict the value of stimulus  $\mathbf{x}_t$  with a linear combination of weights  $\mathbf{w}_t$ :

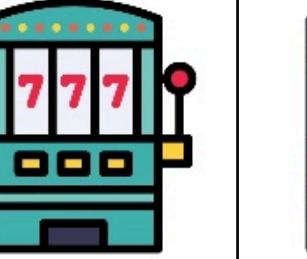
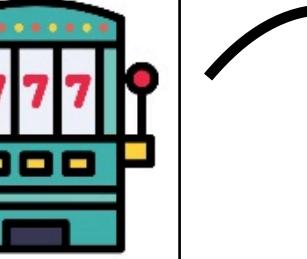
$$V(\mathbf{x}_t) = \mathbf{w}_t^\top \mathbf{x}_t$$

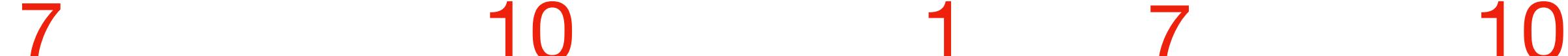
- Weights are updated based on prediction error (aka Delta rule)

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta [r_t - V(\mathbf{x}_t)]$$


$$\delta_t$$

**Q-Learning** (Watkins, 1989) uses the same update rule

|   |   |  |   |   |
|---|---|--|---|---|
|   |  |  |  | 7 |
| Q | 0   | 0  | 7   |   |

$$Q(s_t, a_t)_{new} \leftarrow Q(s_t, a_t)_{old} + \eta [R(s, a) - Q(s, a)_{old}]$$


# Temporal Difference Learning

For simplicity, I'm omitting the temporal difference (TD) error, which looks like:

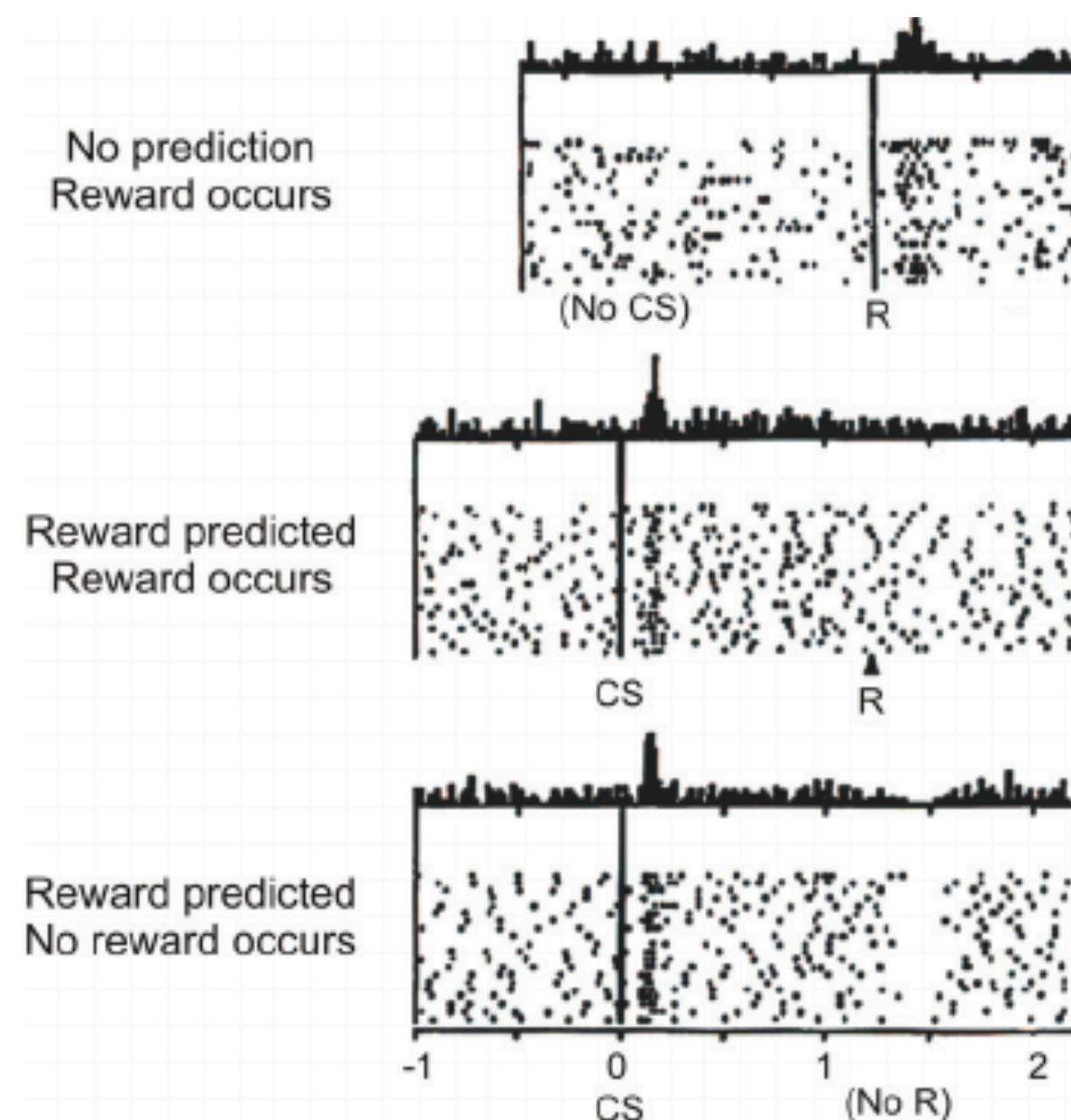
$$V(s) \leftarrow V(s) + \eta \left( r + \underline{\gamma V(s')} - V(s) \right)$$

# Temporal Difference Learning

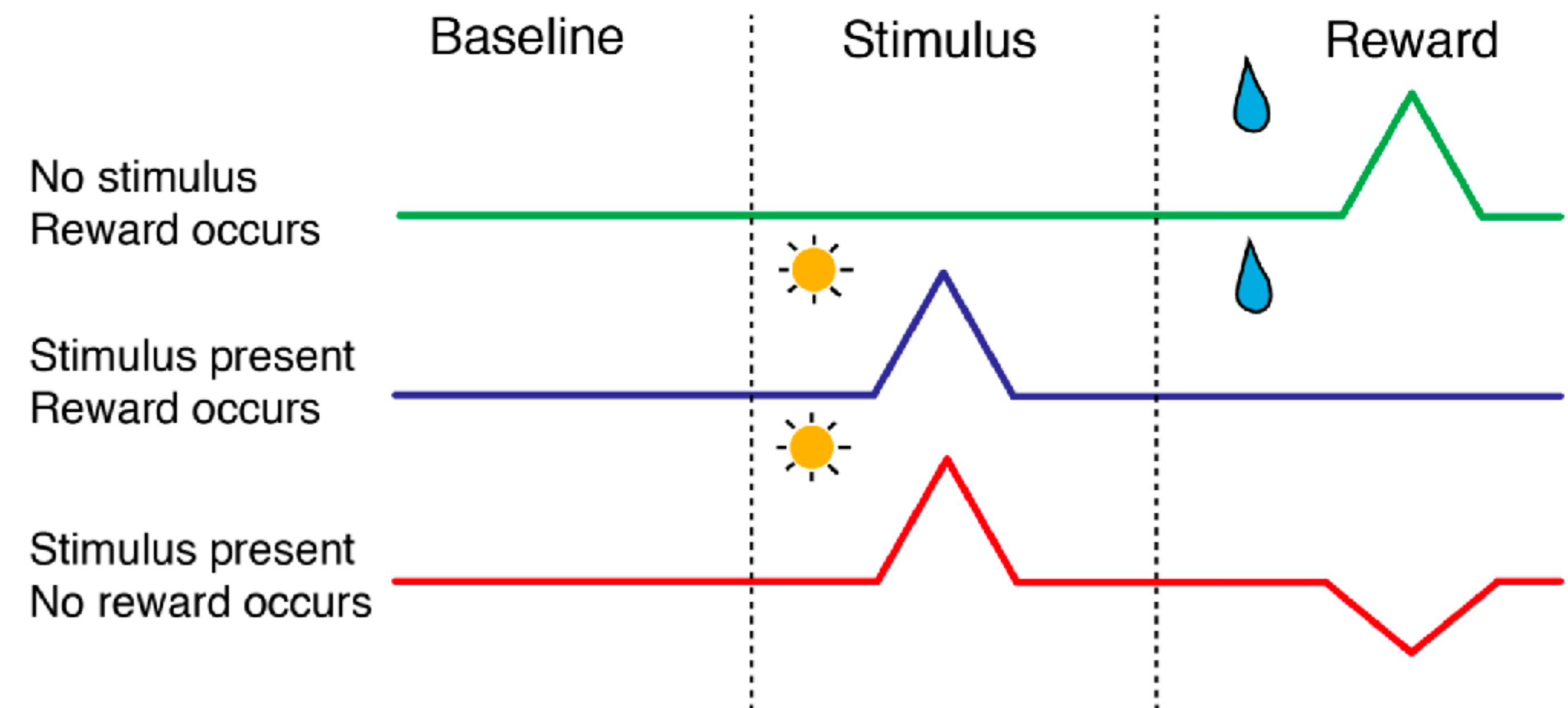
For simplicity, I'm omitting the temporal difference (TD) error, which looks like:

$$V(s) \leftarrow V(s) + \eta (r + \underline{\gamma V(s')} - V(s))$$

## Dopamine Reward Prediction Error



Schultz et al. (1997)



# Policy aka choice function

Goal of Human RL?

# Policy aka choice function

## Goal of Human RL?

Predict behavior using a model. This tells us we have understood some aspect of human learning.

# Policy aka choice function

## Goal of Human RL?

Predict behavior using a model. This tells us we have understood some aspect of human learning.

We want to make probabilistic predictions, since people don't behave deterministically.

# Policy aka choice function

## Goal of Human RL?

Predict behavior using a model. This tells us we have understood some aspect of human learning.

We want to make probabilistic predictions, since people don't behave deterministically.

## Epsilon greedy:

$$\text{action} = \begin{cases} \max_a Q(a) & \text{with } p(1 - \epsilon) \\ \text{random action} & \text{with } p(\epsilon) \end{cases}$$

# Policy aka choice function

## Goal of Human RL?

Predict behavior using a model. This tells us we have understood some aspect of human learning.

We want to make probabilistic predictions, since people don't behave deterministically.

### Epsilon greedy:

$$\text{action} = \begin{cases} \max_a Q(a) & \text{with } p(1 - \epsilon) \\ \text{random action} & \text{with } p(\epsilon) \end{cases}$$

### Softmax:

$$P(a_i) = \frac{\exp(Q(a_i)/\tau)}{\sum_j^n \exp(Q(a_j)/\tau)}$$

# Policy aka choice function

## Goal of Human RL?

Predict behavior using a model. This tells us we have understood some aspect of human learning.

We want to make probabilistic predictions, since people don't behave deterministically.

### Epsilon greedy:

$$\text{action} = \begin{cases} \max_a Q(a) & \text{with } p(1 - \epsilon) \\ \text{random action} & \text{with } p(\epsilon) \end{cases}$$

### Softmax:

$$\begin{array}{c} \text{Q-values} \\ \begin{bmatrix} 1.3 \\ 5.1 \\ 0.7 \\ 1.1 \end{bmatrix} \end{array} \rightarrow \frac{\exp(Q(a_i)/\tau)}{\sum_j^n \exp(Q(a_j)/\tau)} \rightarrow \begin{array}{c} \text{Probabilities} \\ \begin{bmatrix} 0.002 \\ 0.90 \\ 0.05 \\ 0.02 \end{bmatrix} \end{array}$$

# Policy aka choice function

## Goal of Human RL?

Predict behavior using a model. This tells us we have understood some aspect of human learning.

We want to make probabilistic predictions, since people don't behave deterministically.

### Epsilon greedy:

$$\text{action} = \begin{cases} \max_a Q(a) & \text{with } p(1 - \epsilon) \\ \text{random action} & \text{with } p(\epsilon) \end{cases}$$

### Softmax:

$$\begin{array}{c} \text{Q-values} \\ \begin{bmatrix} 1.3 \\ 5.1 \\ 0.7 \\ 1.1 \end{bmatrix} \end{array} \rightarrow \frac{\exp(Q(a_i)/\tau)}{\sum_j^n \exp(Q(a_j)/\tau)} \rightarrow \begin{array}{c} \text{Probabilities} \\ \begin{bmatrix} 0.002 \\ 0.90 \\ 0.05 \\ 0.02 \end{bmatrix} \end{array}$$

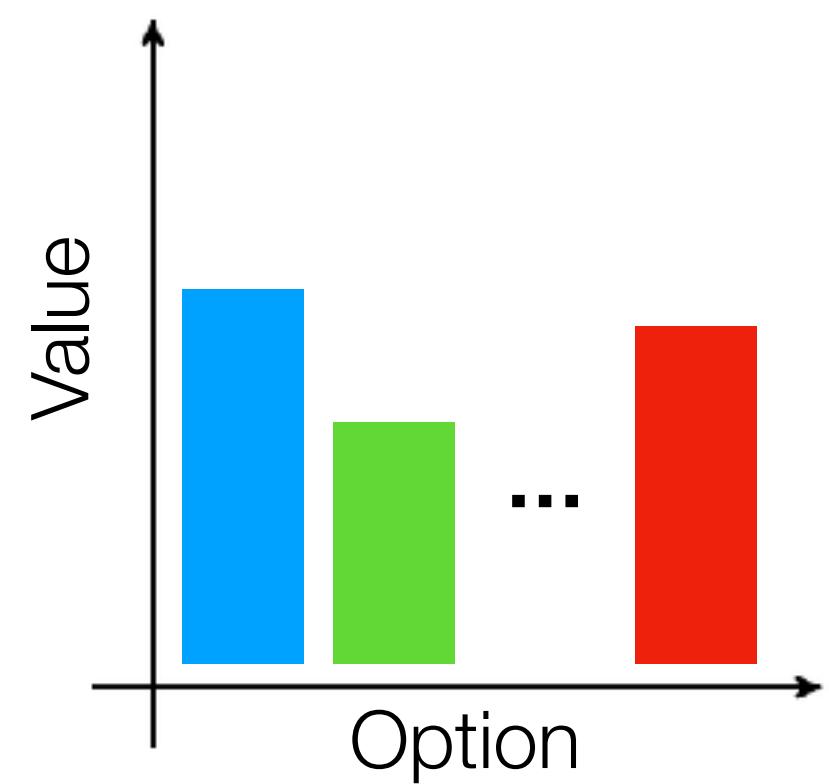
**$\epsilon$  and  $\tau$  model forms of exploration, but don't distinguish between experienced vs. inexperienced options**

# Bayesian reinforcement learning

We can describe a Bayesian variant of the RW model using a Kalman filter:

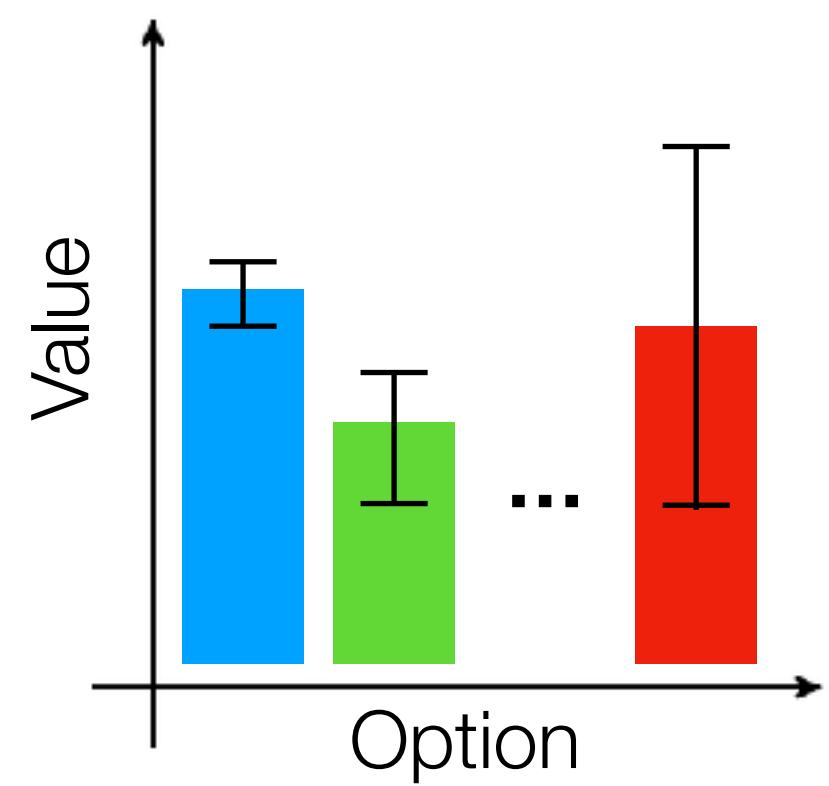
# Bayesian reinforcement learning

We can describe a Bayesian variant of the RW model using a Kalman filter:



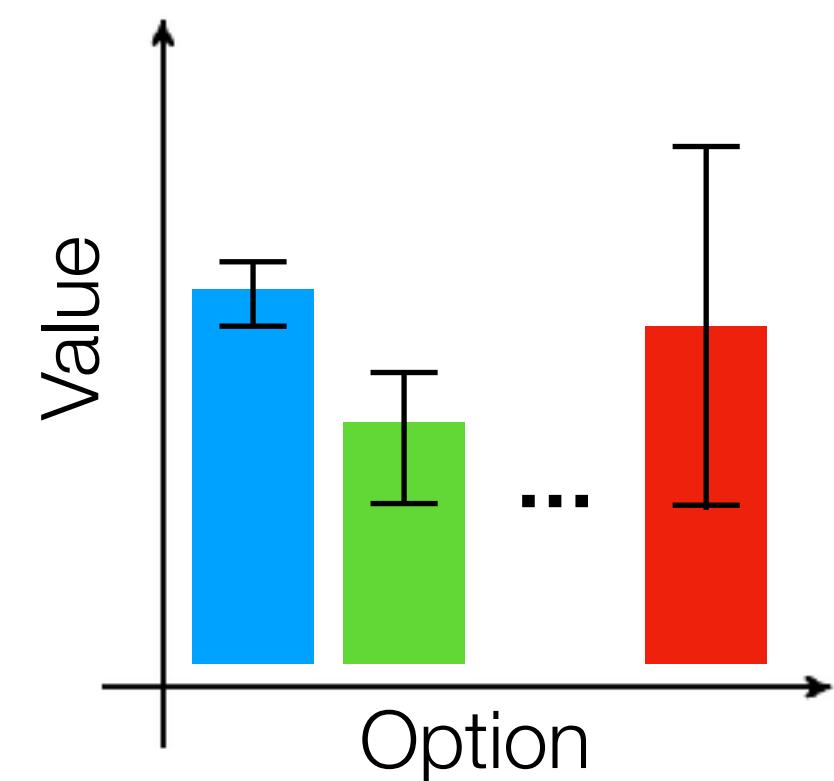
# Bayesian reinforcement learning

We can describe a Bayesian variant of the RW model using a Kalman filter:



# Bayesian reinforcement learning

We can describe a Bayesian variant of the RW model using a Kalman filter:



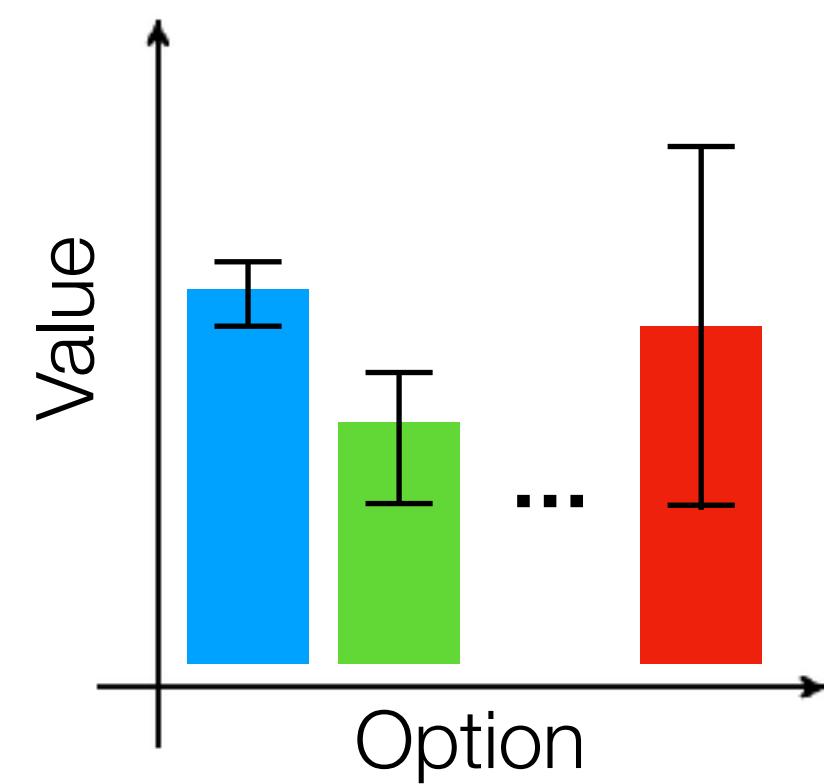
Normally distributed posterior over rewards, where the mean is  $Q_{i,t}$ :

$$p(r_{i,t} | \mathcal{D}_{t-1}) = \mathcal{N}(Q_{i,t}, \sigma_{i,t}^2)$$

where  $\mathcal{D}_t = [a_0, r_0, a_1, r_1, \dots]$  collect previous choices and rewards

# Bayesian reinforcement learning

We can describe a Bayesian variant of the RW model using a Kalman filter:



Normally distributed posterior over rewards, where the mean is  $Q_{i,t}$ :

$$p(r_{i,t} | \mathcal{D}_{t-1}) = \mathcal{N}(Q_{i,t}, \sigma_{i,t}^2)$$

where  $\mathcal{D}_t = [a_0, r_0, a_1, r_1, \dots]$  collect previous choices and rewards

Bayesian updates:

Mean:  $Q_{i,t+1} = Q_{i,t} + k_{i,t} [r_{i,t} - Q_{i,t}]$

Variance:  $\sigma_{i,t+1}^2 = [1 - k_{i,t}] \sigma_{i,t}^2$

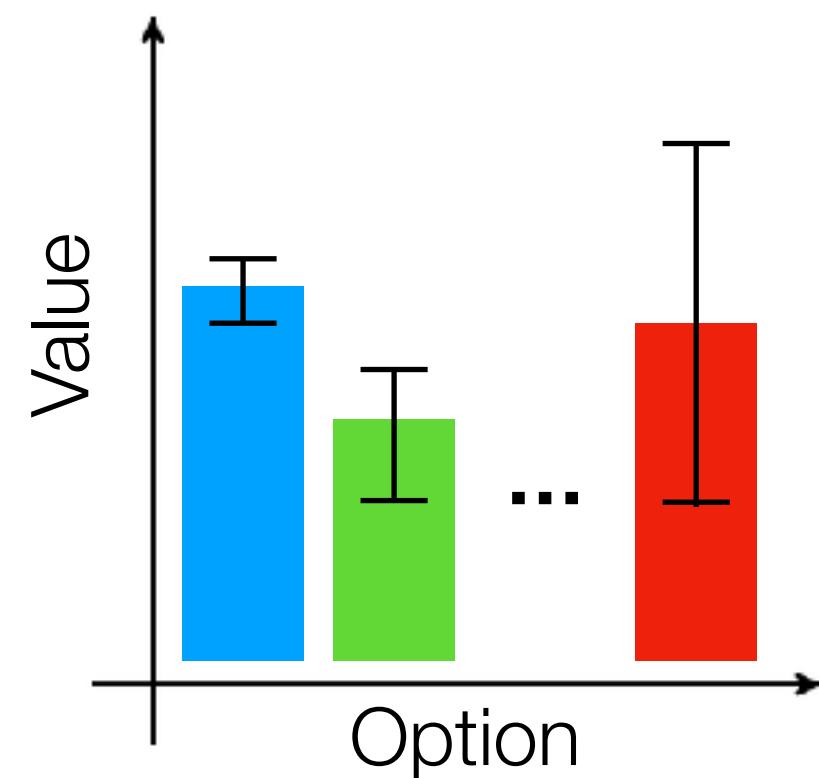
Kalman Gain (learning rate):

$$k_{i,t} = \begin{cases} \frac{\sigma_{i,t}^2}{\sigma_{i,t}^2 + \sigma_\epsilon^2} & \text{if } a_t = i \\ 0 & \text{otherwise} \end{cases}$$

Error variance  $\sigma_\epsilon^2$  is a free parameter

# Bayesian reinforcement learning

We can describe a Bayesian variant of the RW model using a Kalman filter:



Normally distributed posterior over rewards, where the mean is  $Q_{i,t}$ :

$$p(r_{i,t} | \mathcal{D}_{t-1}) = \mathcal{N}(Q_{i,t}, \sigma_{i,t}^2)$$

where  $\mathcal{D}_t = [a_0, r_0, a_1, r_1, \dots]$  collect previous choices and rewards

Bayesian updates:

Mean:  $Q_{i,t+1} = Q_{i,t} + k_{i,t} [r_{i,t} - Q_{i,t}]$

Variance:  $\sigma_{i,t+1}^2 = [1 - k_{i,t}] \sigma_{i,t}^2$

Kalman Gain (learning rate):

$$k_{i,t} = \begin{cases} \frac{\sigma_{i,t}^2}{\sigma_{i,t}^2 + \sigma_\epsilon^2} & \text{if } a_t = i \\ 0 & \text{otherwise} \end{cases}$$

Strictly, this is a KF variant known as a Bayesian mean tracker (BMT), assuming stationary rewards

Error variance  $\sigma_\epsilon^2$  is a free parameter

# Policies for Bayesian models

We can now use uncertainty estimates to inform our policy and explore more efficiently.

# Policies for Bayesian models

We can now use uncertainty estimates to inform our policy and explore more efficiently.

**Thompson Sampling:**

$$P(a_i) = P(r_i > r_{j \neq i})$$

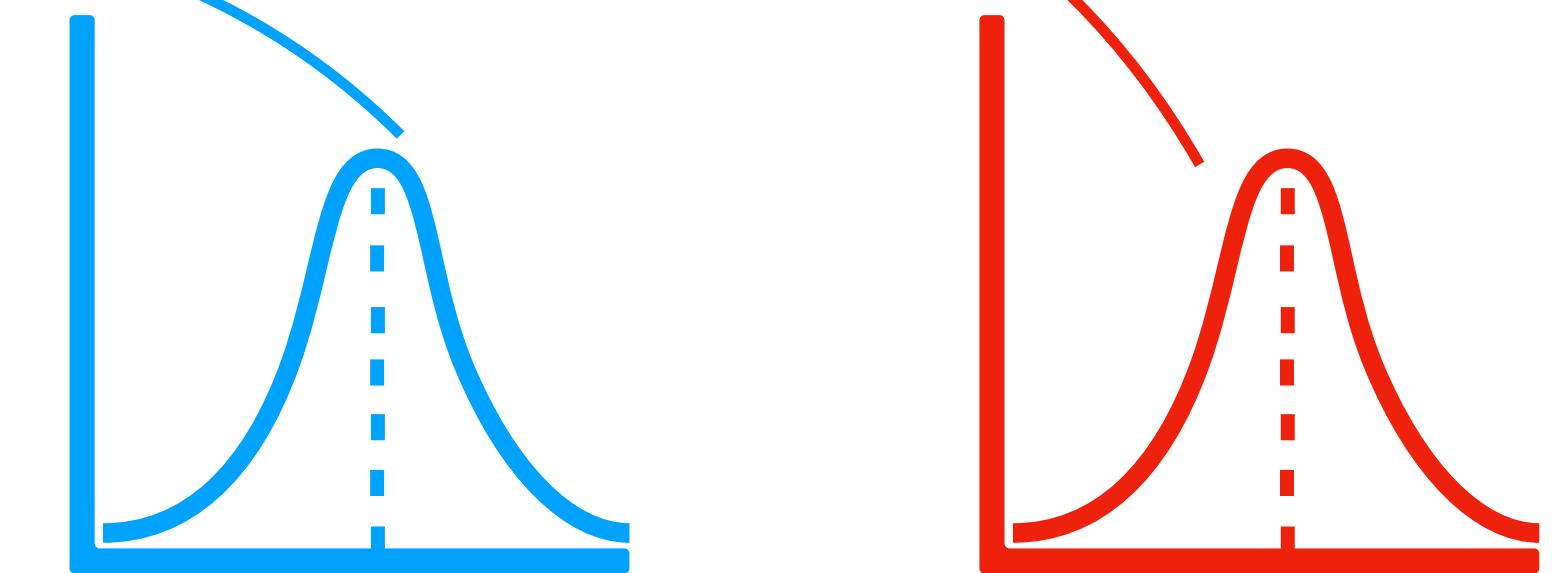
# Policies for Bayesian models

We can now use uncertainty estimates to inform our policy and explore more efficiently.

**Thompson Sampling:**

$$P(a_i) = P(r_i > r_{j \neq i})$$

Monte Carlo Samples



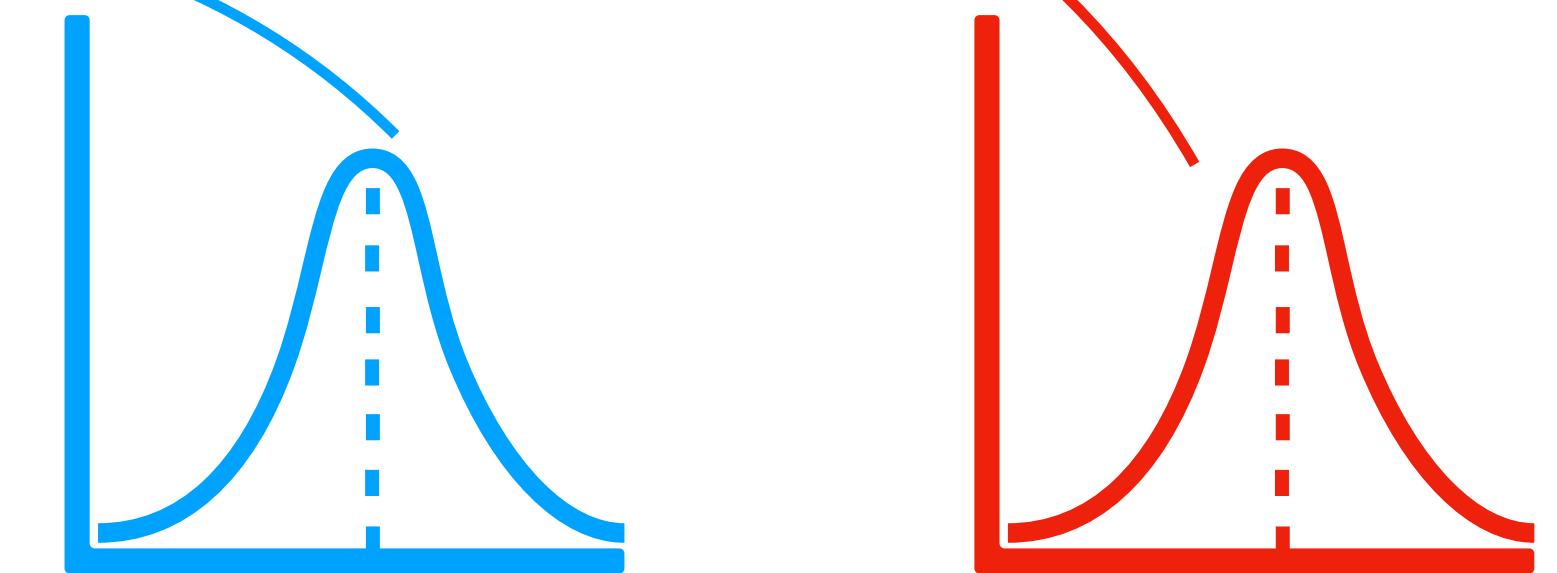
# Policies for Bayesian models

We can now use uncertainty estimates to inform our policy and explore more efficiently.

**Thompson Sampling:**

$$P(a_i) = P(r_i > r_{j \neq i})$$

Monte Carlo Samples



**Upper Confidence Bound Sampling:**

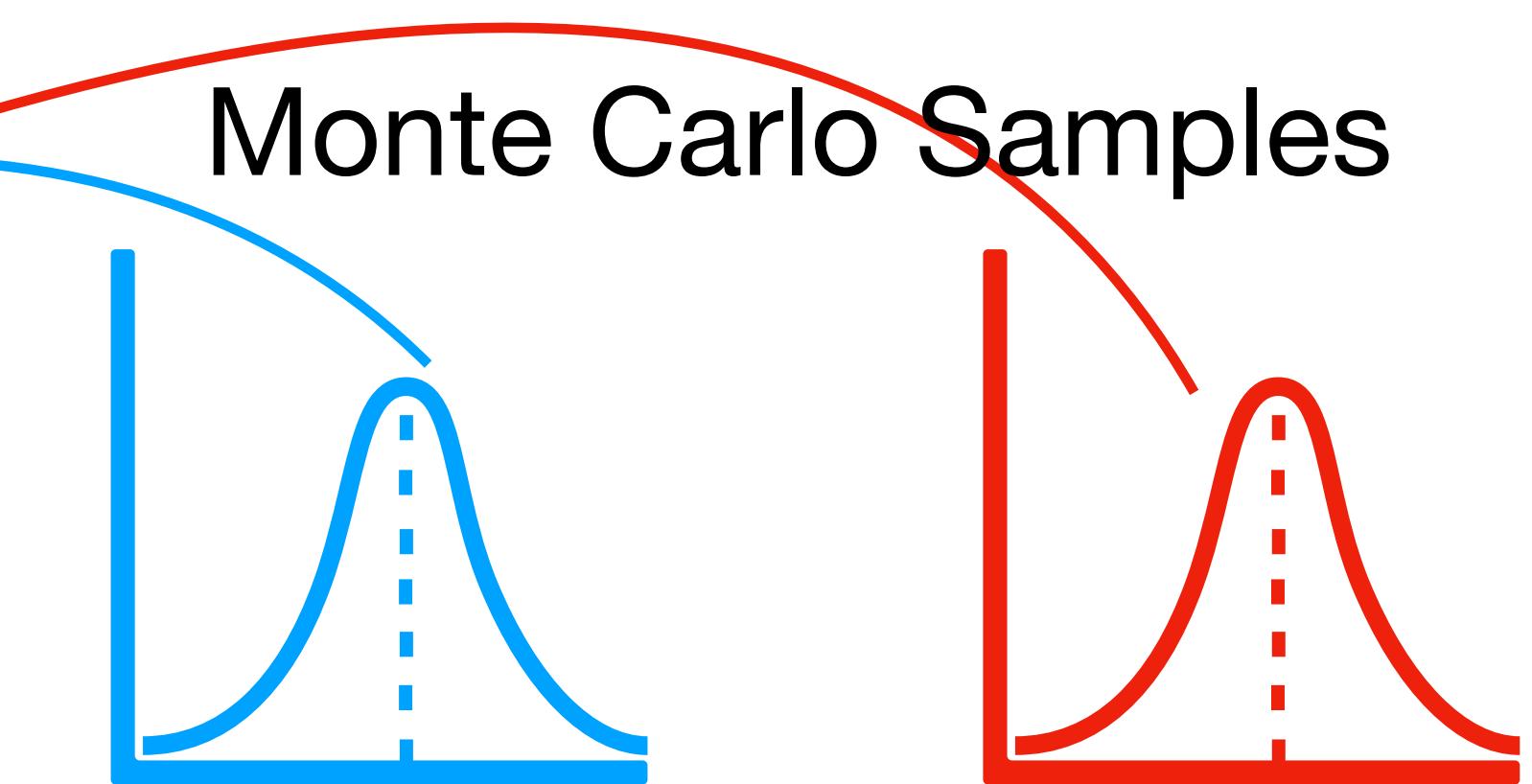
$$UCB(a_i) = Q_i + \beta \sqrt{\sigma_i^2}$$

# Policies for Bayesian models

We can now use uncertainty estimates to inform our policy and explore more efficiently.

**Thompson Sampling:**

$$P(a_i) = P(r_i > r_{j \neq i})$$

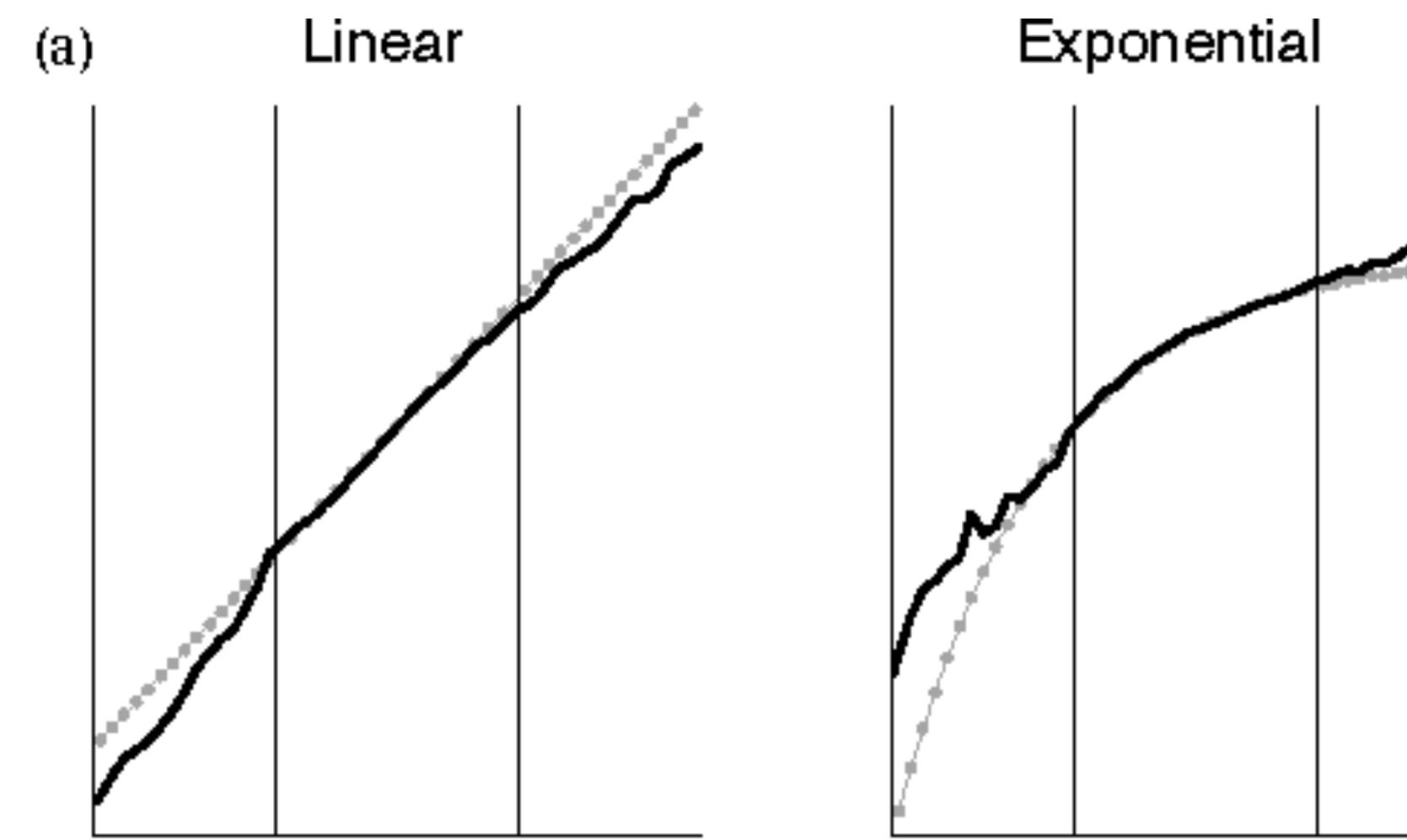
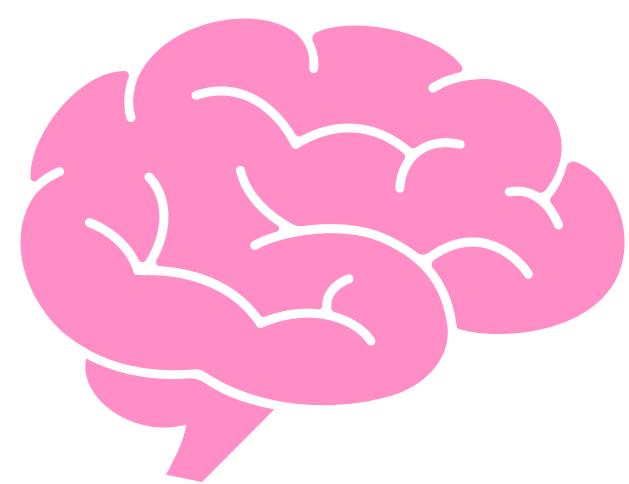


**Upper Confidence Bound Sampling:**

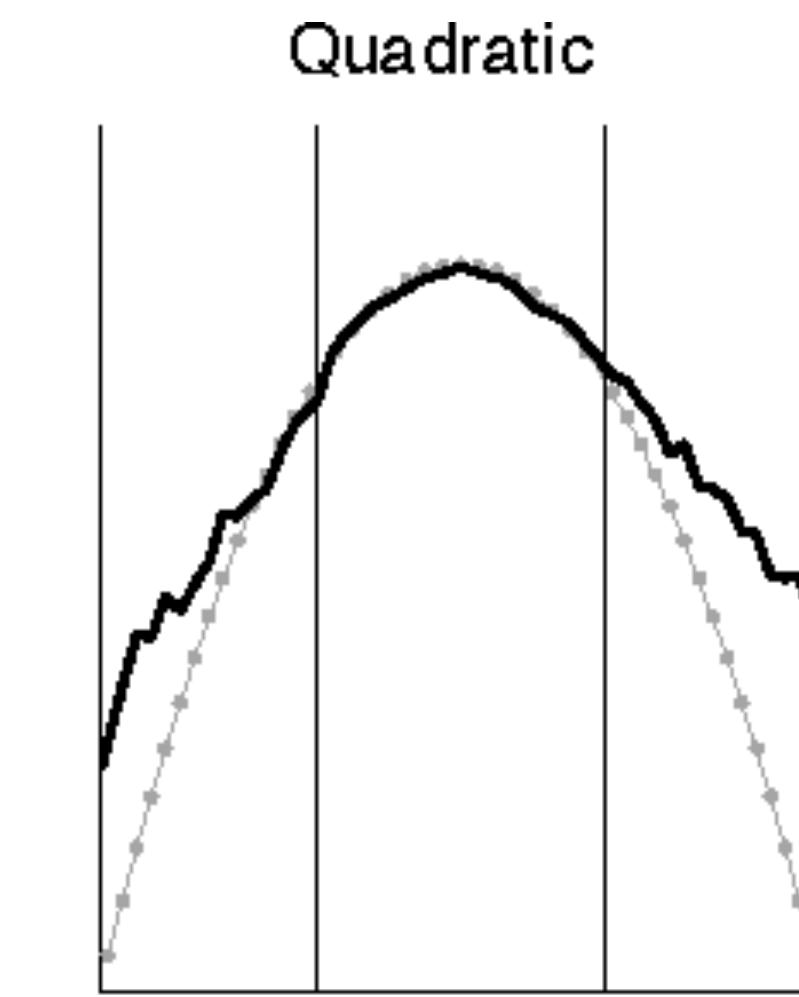
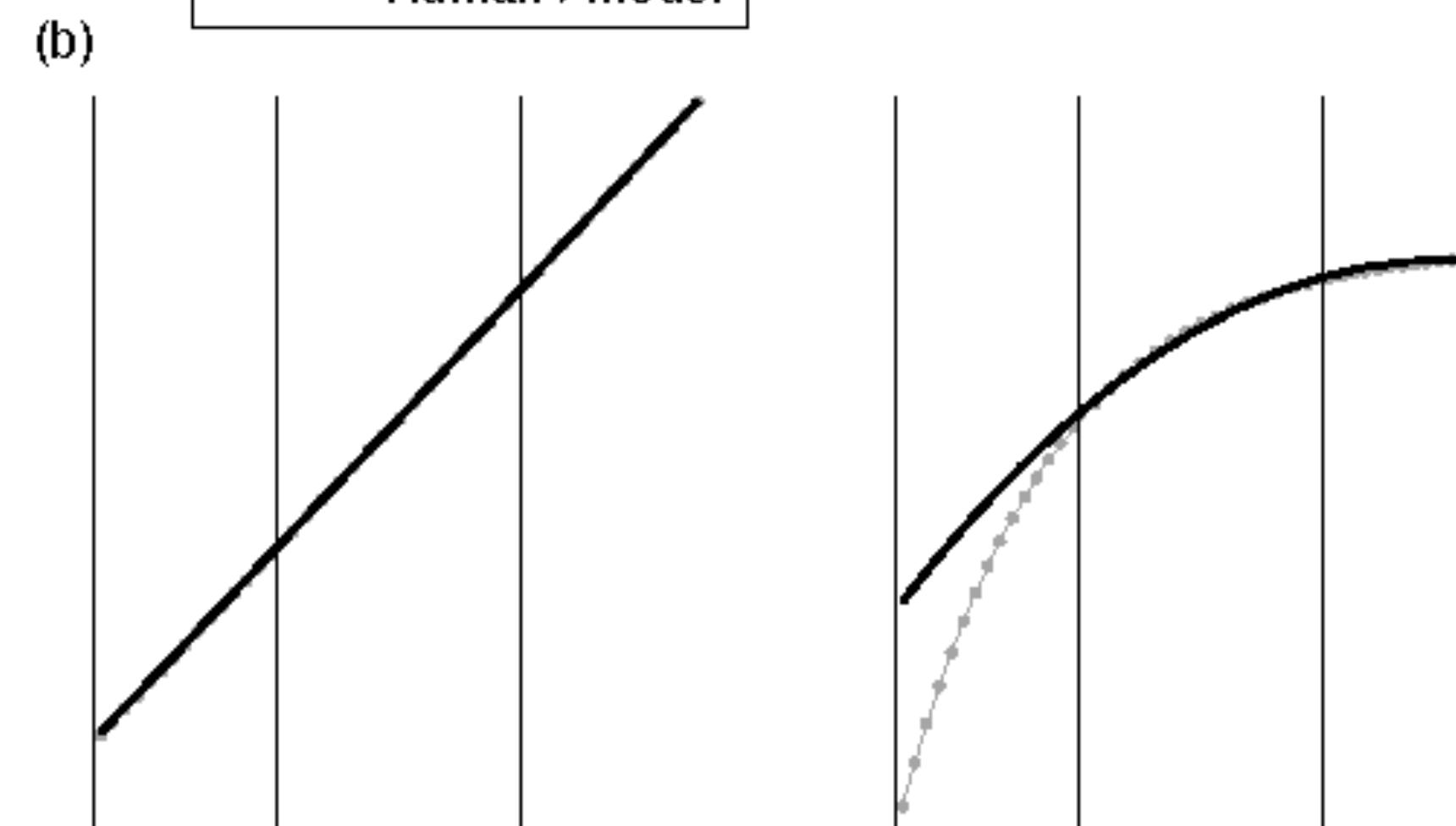
$$UCB(a_i) = Q_i + \beta \sqrt{\sigma_i^2}$$

$\beta$  models exploration directed towards uncertain choices

# Function approximation?



— Function  
— Human / Model



(c)

| Model  | Linear | Exponential | Quadratic |
|--------|--------|-------------|-----------|
| EXAM   | .999   | .997        | .961      |
| Linear | .999   | .989        | .470      |
| Quad   | .997   | .997        | .901      |
| RBF    | .999   | .997        | .882      |
| LQ     | .999   | .997        | .886      |
| LR     | .999   | .997        | .892      |
| RQ     | .998   | .994        | .878      |
| LRQ    | .999   | .995        | .877      |

# Summary

# Summary

RL framework for learning value functions and policies

# Summary

RL framework for learning value functions and policies

- Tabular methods vs. value function approximation

# Summary

RL framework for learning value functions and policies

- Tabular methods vs. value function approximation
- Exhaustive iteration over states vs. predictions about novel states

# Summary

RL framework for learning value functions and policies

- Tabular methods vs. value function approximation
- Exhaustive iteration over states vs. predictions about novel states

Modeling human learners with RL models

# Summary

RL framework for learning value functions and policies

- Tabular methods vs. value function approximation
- Exhaustive iteration over states vs. predictions about novel states

Modeling human learners with RL models

- Can we understand something about the efficiency of human learning?

# Summary

RL framework for learning value functions and policies

- Tabular methods vs. value function approximation
- Exhaustive iteration over states vs. predictions about novel states

Modeling human learners with RL models

- Can we understand something about the efficiency of human learning?
- How humans navigate the exploration-exploitation dilemma

# Summary

RL framework for learning value functions and policies

- Tabular methods vs. value function approximation
- Exhaustive iteration over states vs. predictions about novel states

Modeling human learners with RL models

- Can we understand something about the efficiency of human learning?
  - How humans navigate the exploration-exploitation dilemma
- RW and Q-learning resemble tabular methods that have access to random or noisy exploration

# Summary

RL framework for learning value functions and policies

- Tabular methods vs. value function approximation
- Exhaustive iteration over states vs. predictions about novel states

Modeling human learners with RL models

- Can we understand something about the efficiency of human learning?
  - How humans navigate the exploration-exploitation dilemma
- RW and Q-learning resemble tabular methods that have access to random or noisy exploration
- Kalman filter is a Bayesian model, which can use uncertainty-informed exploration

# Summary

RL framework for learning value functions and policies

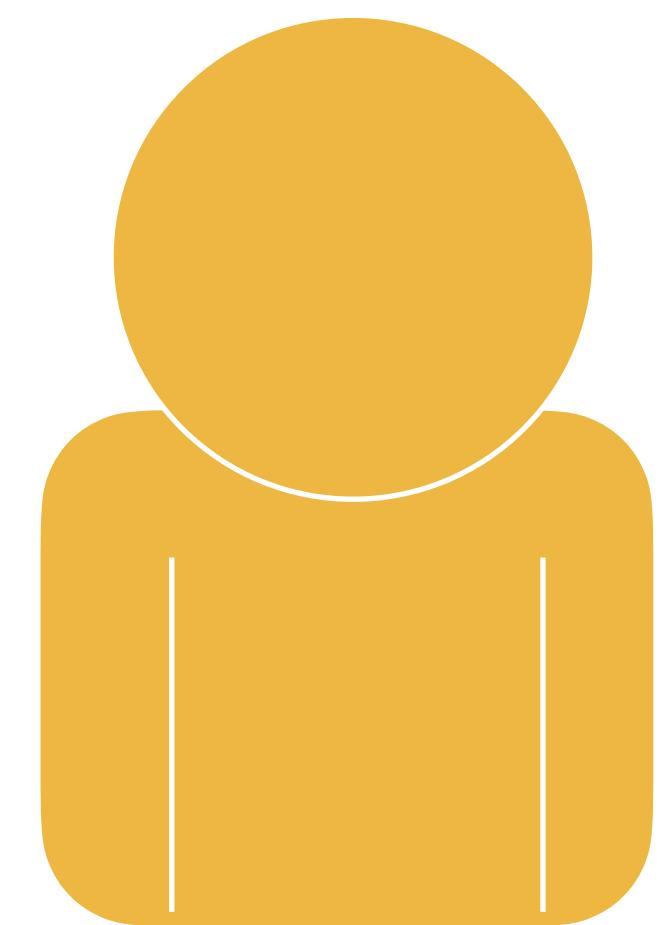
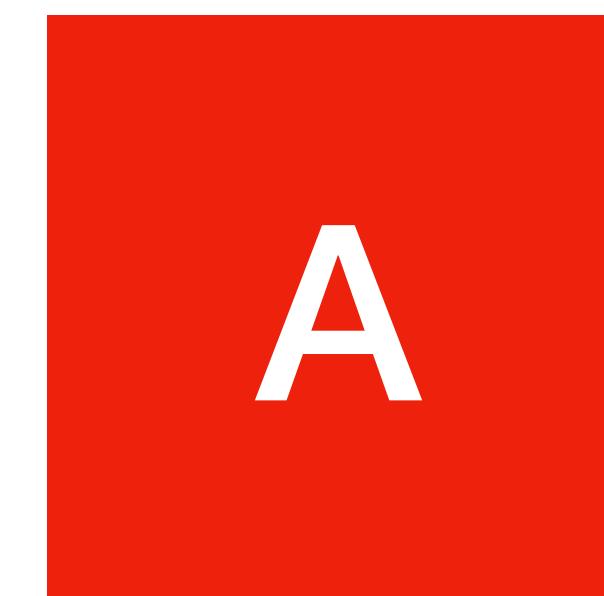
- Tabular methods vs. value function approximation
- Exhaustive iteration over states vs. predictions about novel states

Modeling human learners with RL models

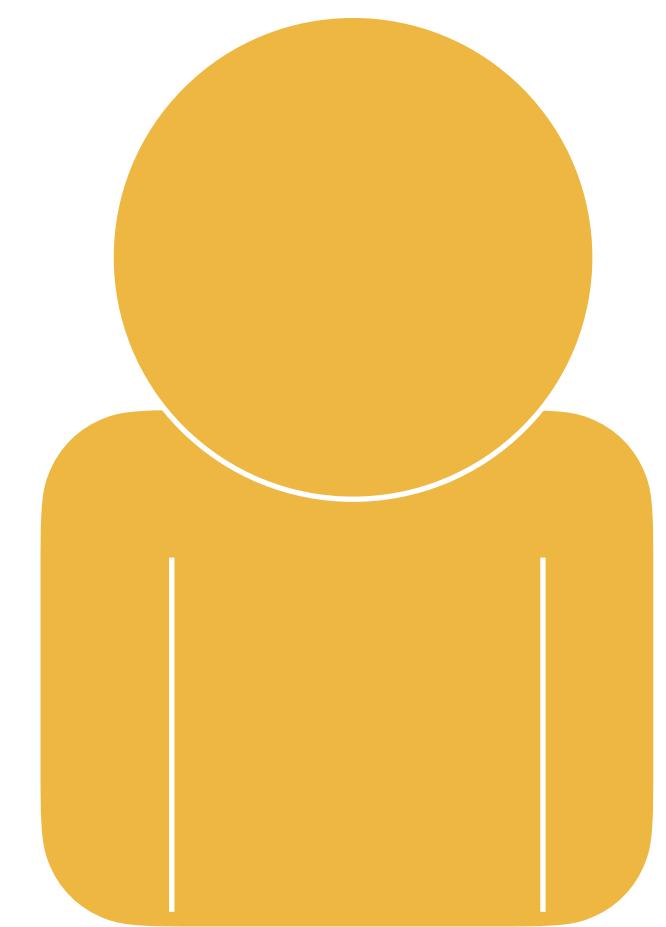
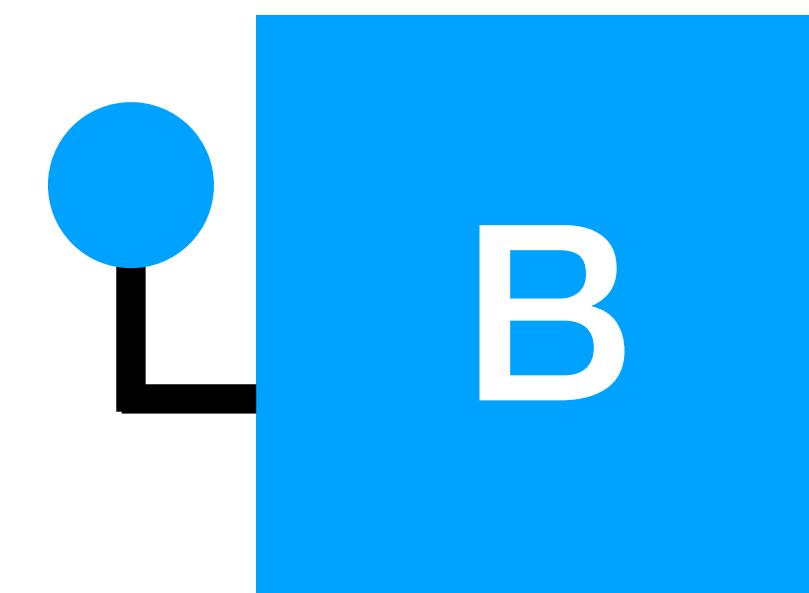
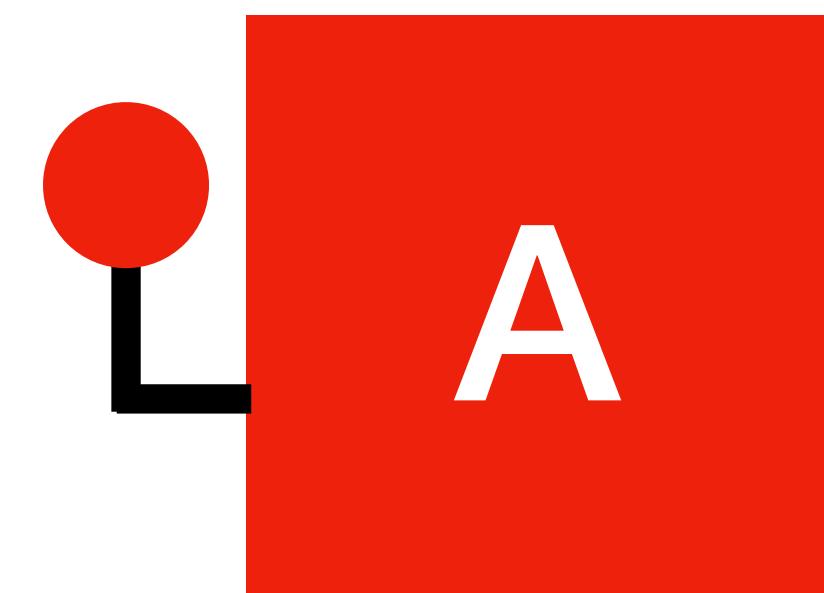
- Can we understand something about the efficiency of human learning?
  - How humans navigate the exploration-exploitation dilemma
- RW and Q-learning resemble tabular methods that have access to random or noisy exploration
- Kalman filter is a Bayesian model, which can use uncertainty-informed exploration
- But what about function approximation?

# Break

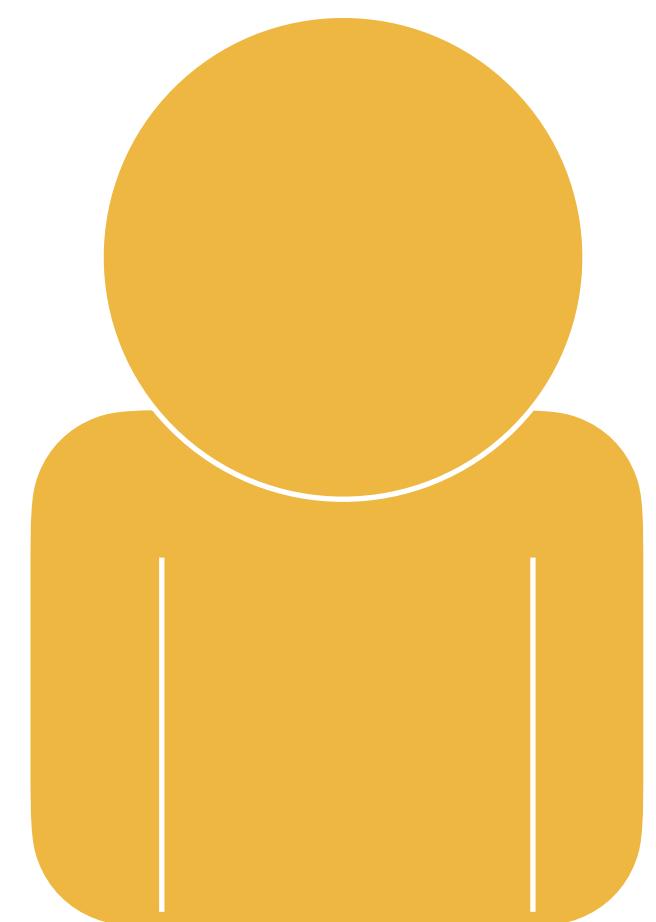
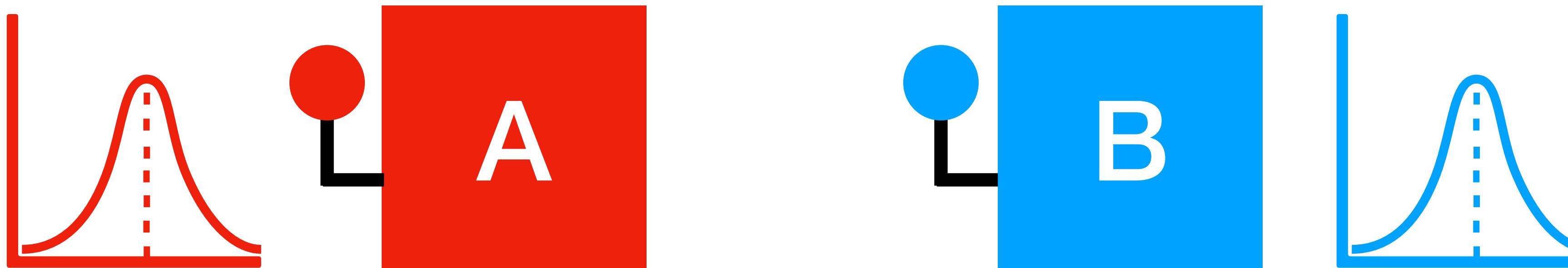
# Human learning in the lab



# Human learning in the lab

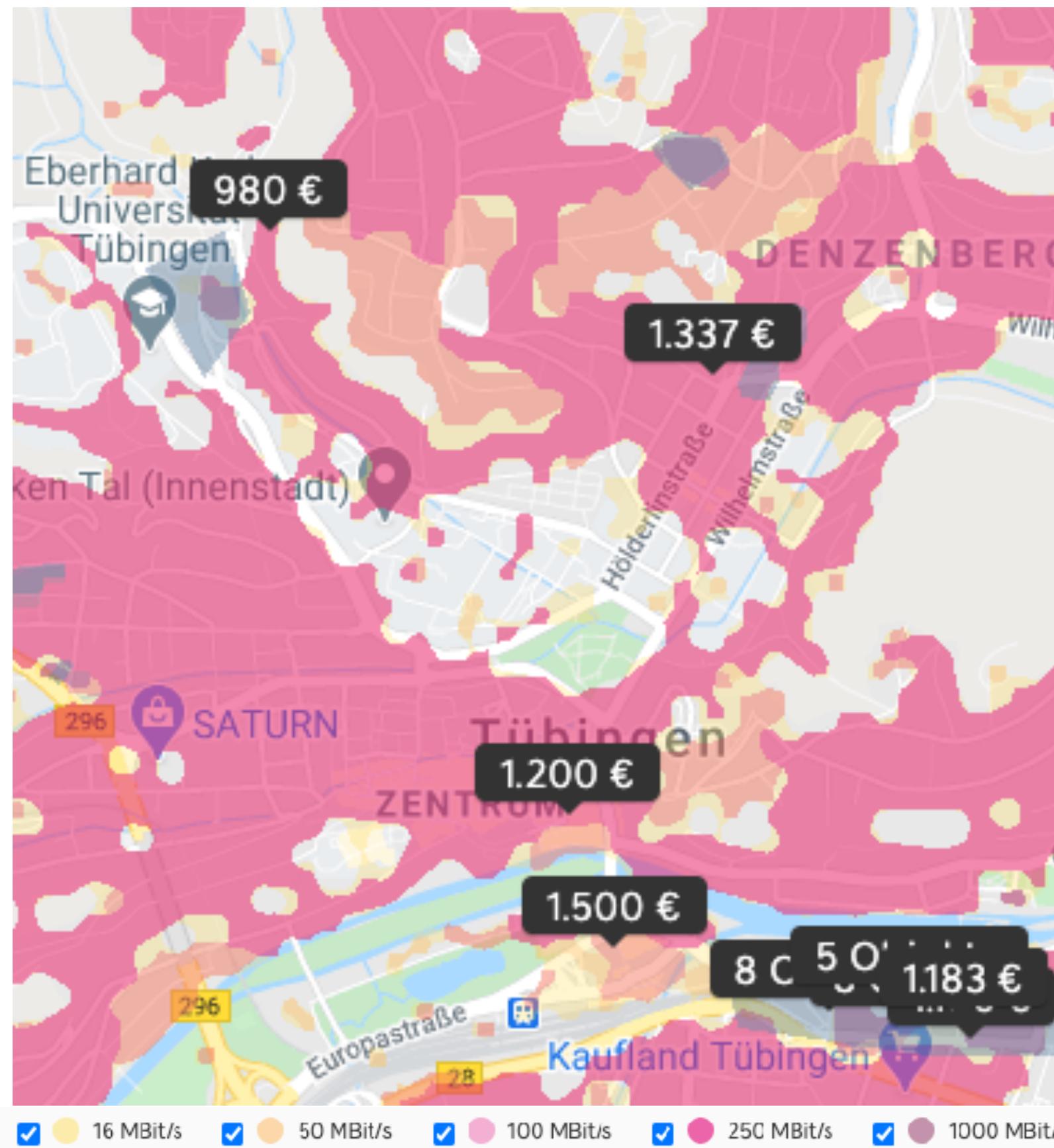


# Human learning in the lab



# Real life problems

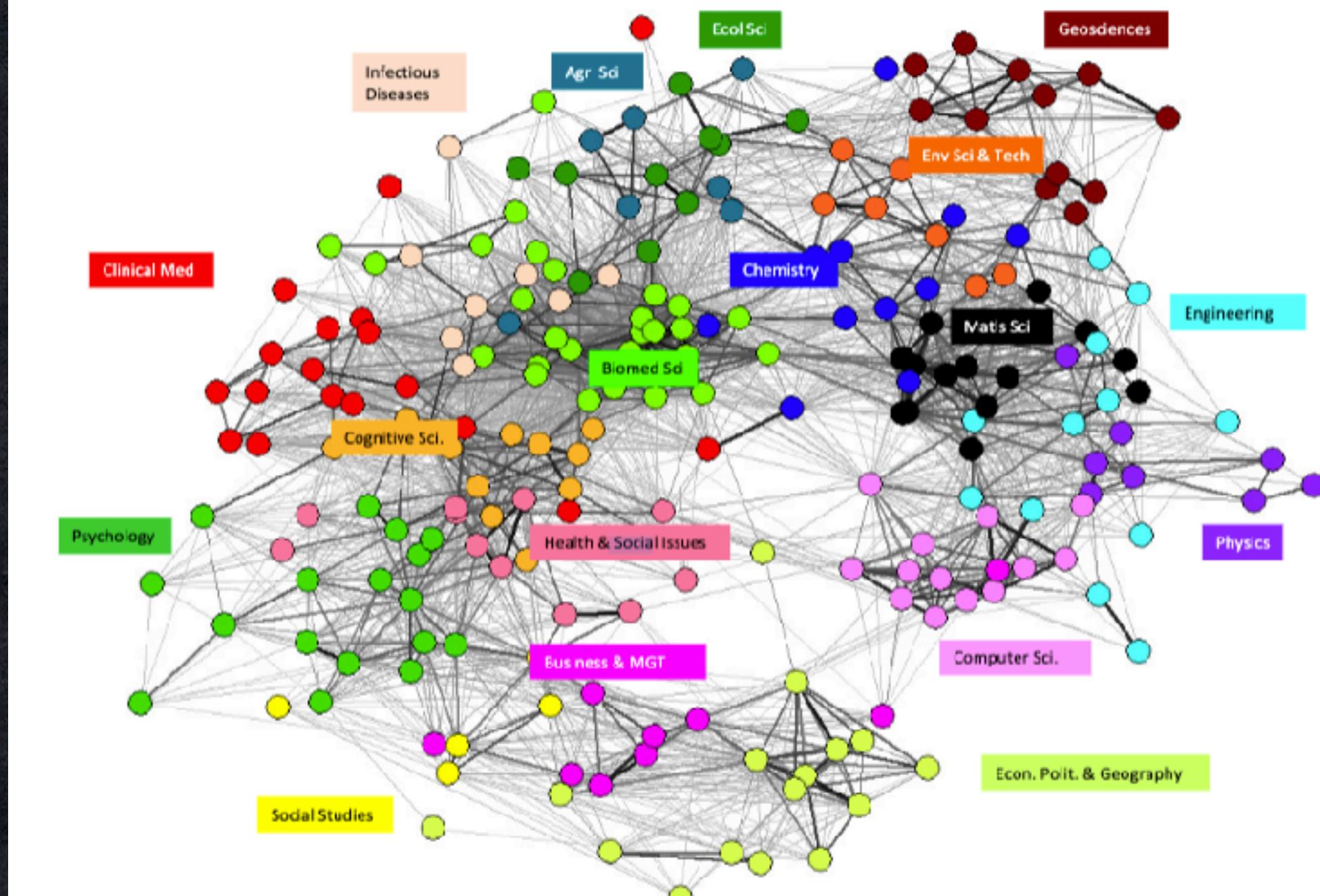
Finding a place to live



Picking what to eat



Choosing a research topic



# Exploration-Exploitation Dilemma



**Exploration**



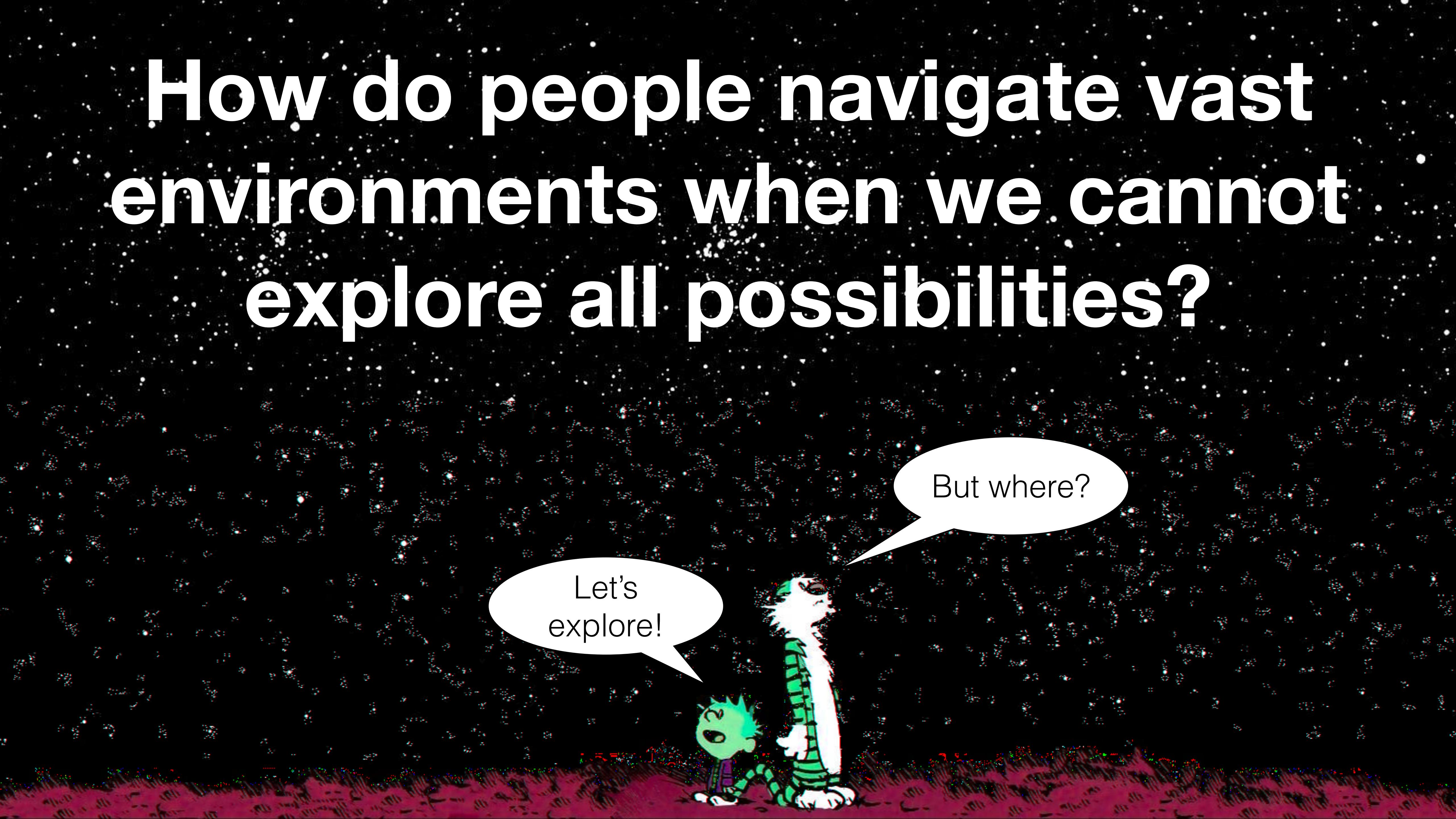
**Exploitation**



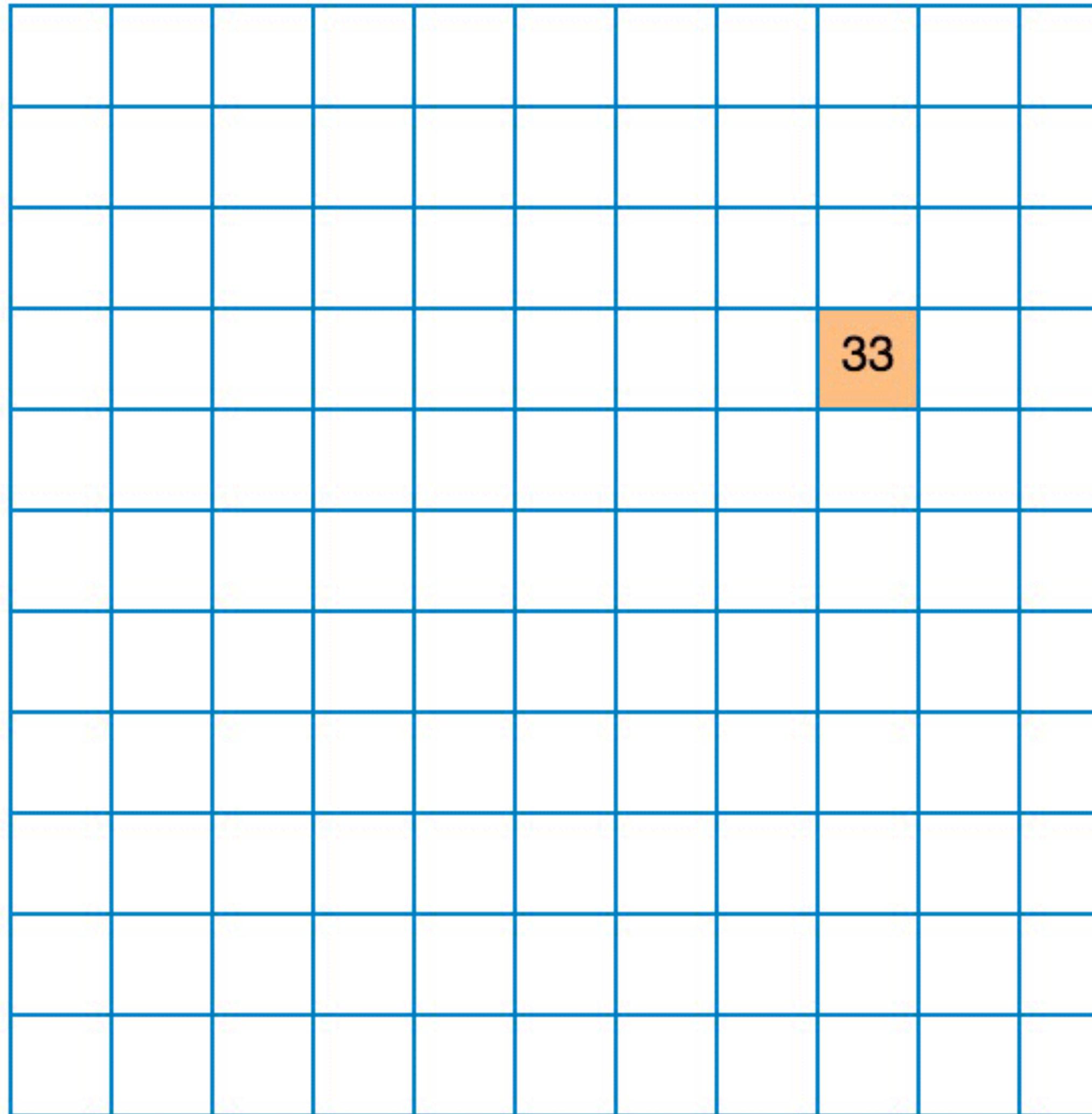
Let's  
explore!

But where?

# How do people navigate vast environments when we cannot explore all possibilities?



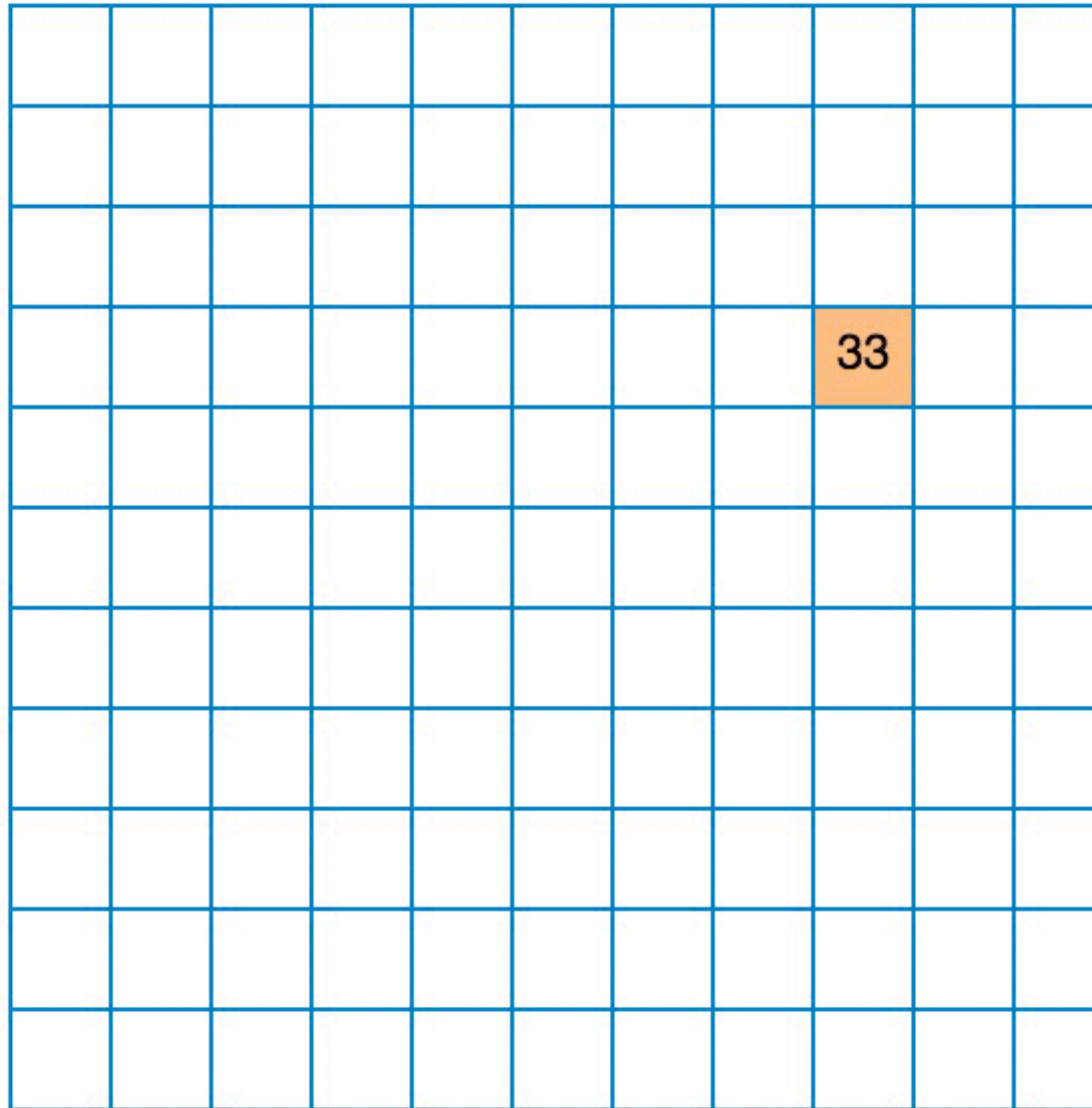
# Spatially Correlated Bandit



Wu et al., (*Nature Human Behaviour* 2018)

- click tiles on the grid
- maximize reward
- each tile has normally distributed rewards
- nearby tiles have similar rewards
- limited search horizon privileges good generalization & efficient exploration

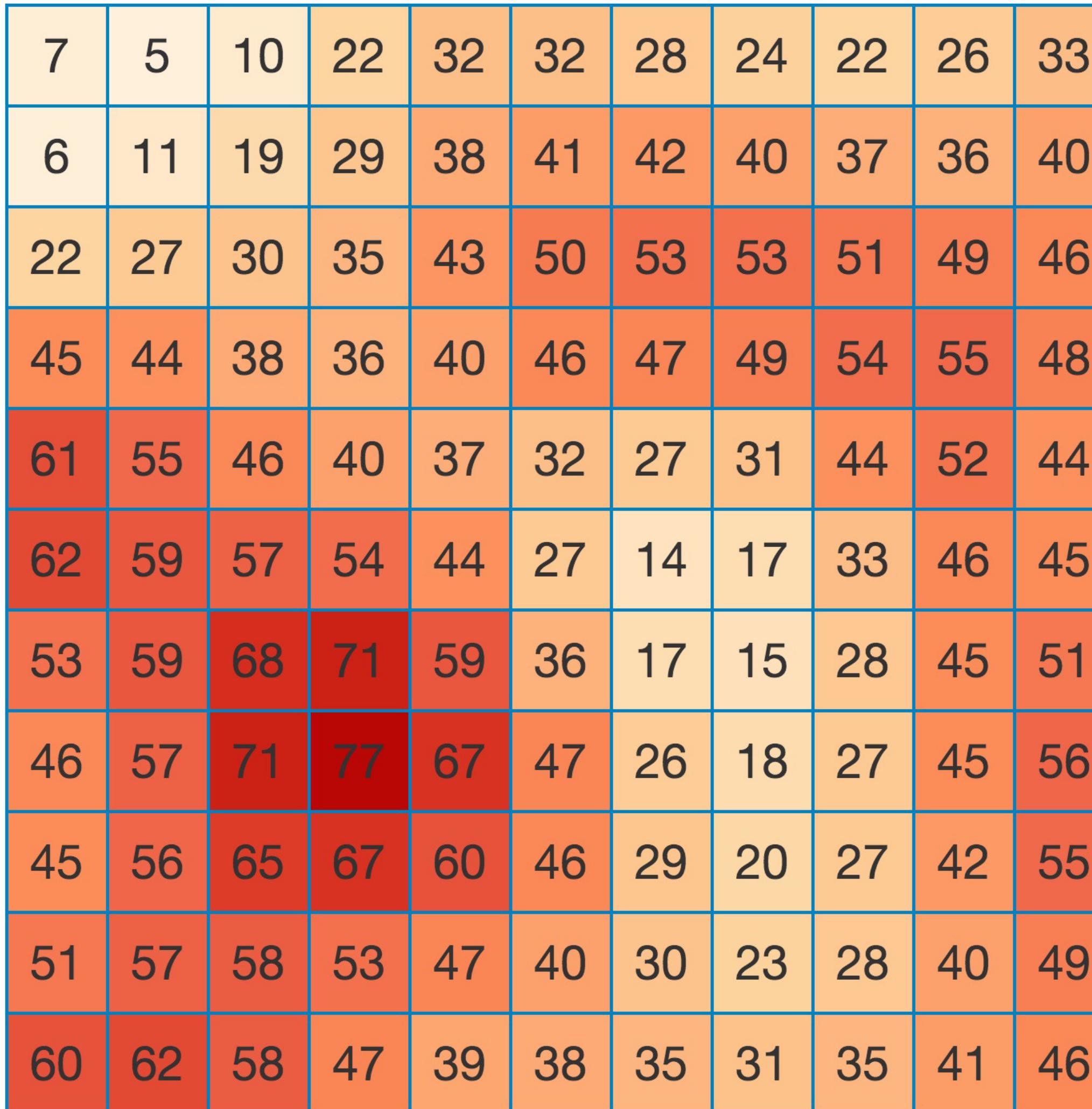
# Spatially Correlated Bandit



Wu et al., (*Nature Human Behaviour* 2018)

- click tiles on the grid
- maximize reward
- each tile has normally distributed rewards
- nearby tiles have similar rewards
- limited search horizon privileges good generalization & efficient exploration

# Spatially Correlated Bandit



Wu et al., (*Nature Human Behaviour* 2018)

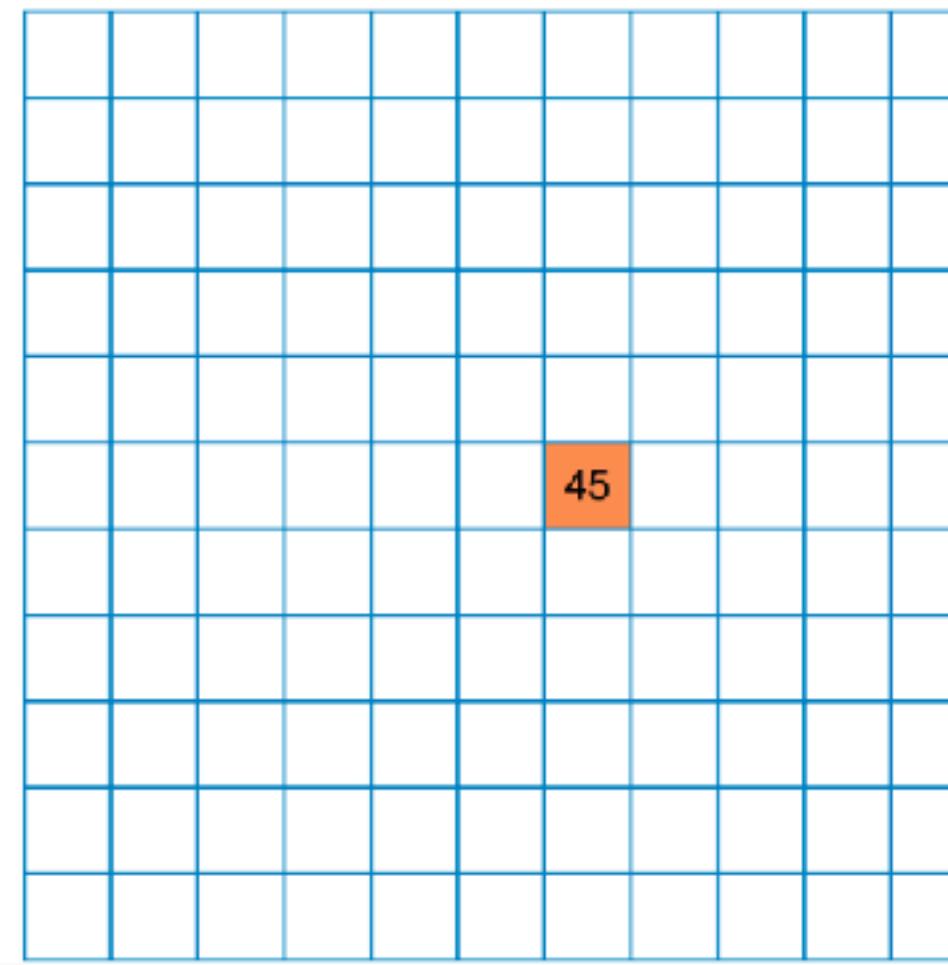
- click tiles on the grid
- maximize reward
- each tile has normally distributed rewards
- nearby tiles have similar rewards
- limited search horizon privileges good generalization & efficient exploration

# Experiment 1

30-Armed Bandit (Univariate)

| Rough  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
|--------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 23     | 20 | 32 | 45 | 58 | 63 | 45 | 52 | 54 | 53 | 45 | 49 | 45 | 28 | 21 | 56 | 51 | 54 | 25 | 5  | 7  | 42 | 48 | 35 | 30 | 54 | 59 | 42 | 24 | 35 |
| Smooth |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 46     | 44 | 40 | 39 | 37 | 27 | 21 | 30 | 47 | 59 | 65 | 69 | 65 | 51 | 41 | 47 | 57 | 55 | 44 | 39 | 43 | 46 | 36 | 17 | 5  | 11 | 28 | 43 | 47 | 40 |

Number of grids left: **6**  
Number of clicks left: **40**



Participants acquired rewards by clicking  
on new or previously revealed tiles



# Experiment 2

121-Armed Bandit (Bivariate)

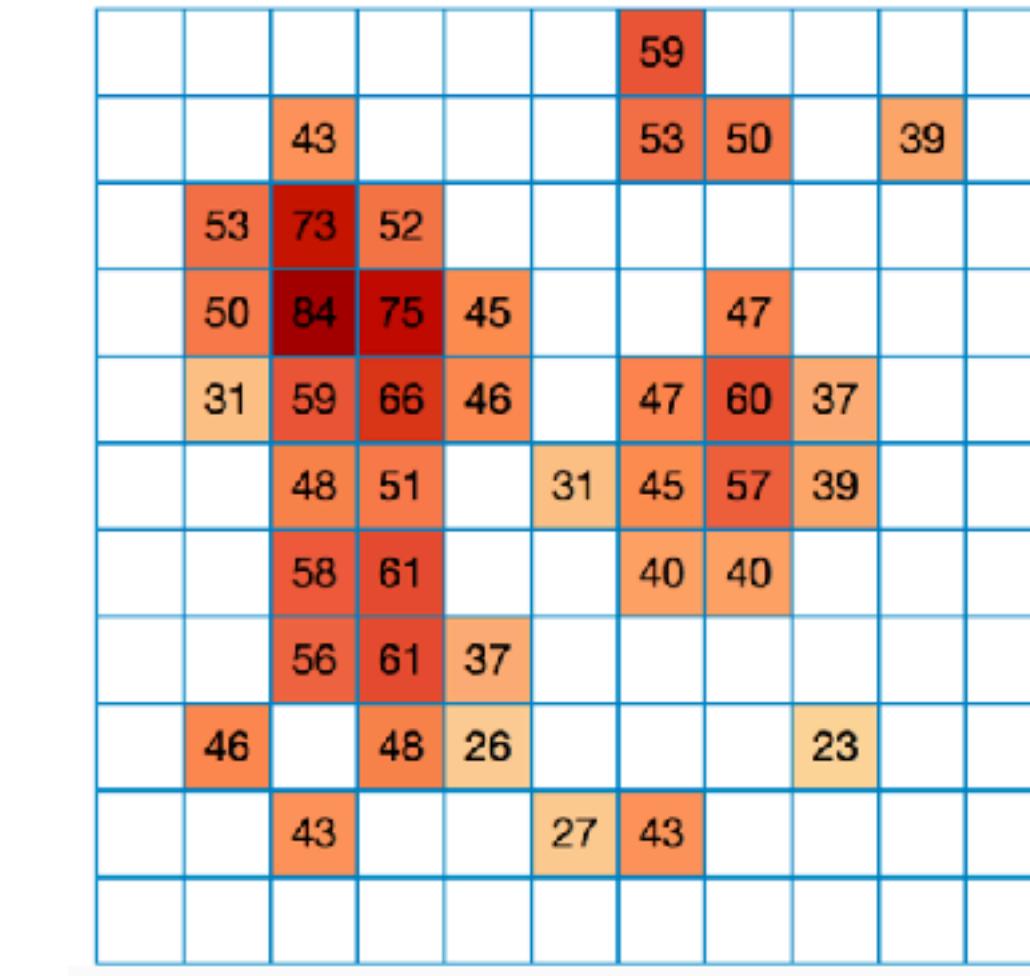
| Rough  |    |    |    |    |    |    |    |    |    |    |  |
|--------|----|----|----|----|----|----|----|----|----|----|--|
| 39     | 42 | 54 | 50 | 44 | 70 | 72 | 72 | 57 | 35 | 12 |  |
| 31     | 21 | 28 | 41 | 44 | 44 | 47 | 64 | 54 | 40 | 27 |  |
| 51     | 37 | 30 | 29 | 26 | 24 | 31 | 48 | 45 | 33 | 25 |  |
| 64     | 55 | 59 | 43 | 36 | 38 | 32 | 26 | 30 | 27 | 28 |  |
| 58     | 62 | 41 | 19 | 23 | 26 | 29 | 37 | 44 | 33 | 27 |  |
| 65     | 76 | 33 | 5  | 7  | 13 | 27 | 50 | 48 | 17 | 25 |  |
| 38     | 44 | 35 | 24 | 17 | 17 | 16 | 33 | 48 | 33 | 45 |  |
| 45     | 29 | 21 | 31 | 50 | 34 | 23 | 29 | 41 | 42 | 37 |  |
| 49     | 24 | 13 | 23 | 38 | 38 | 38 | 34 | 47 | 63 | 40 |  |
| 55     | 31 | 24 | 24 | 26 | 23 | 22 | 26 | 40 | 64 | 66 |  |
| 64     | 53 | 41 | 50 | 59 | 42 | 19 | 29 | 28 | 29 | 49 |  |
| Smooth |    |    |    |    |    |    |    |    |    |    |  |
| 7      | 5  | 10 | 22 | 32 | 32 | 28 | 24 | 22 | 26 | 33 |  |
| 6      | 11 | 19 | 29 | 38 | 41 | 42 | 40 | 37 | 36 | 40 |  |
| 22     | 27 | 30 | 35 | 43 | 50 | 53 | 53 | 51 | 49 | 46 |  |
| 45     | 44 | 38 | 36 | 40 | 46 | 47 | 49 | 54 | 55 | 48 |  |
| 61     | 55 | 46 | 40 | 37 | 32 | 27 | 31 | 44 | 52 | 44 |  |
| 62     | 59 | 57 | 54 | 44 | 27 | 14 | 17 | 33 | 46 | 45 |  |
| 53     | 59 | 68 | 71 | 59 | 36 | 17 | 15 | 28 | 45 | 51 |  |
| 46     | 57 | 71 | 77 | 67 | 47 | 26 | 18 | 27 | 45 | 56 |  |
| 45     | 56 | 65 | 67 | 60 | 46 | 29 | 20 | 27 | 42 | 55 |  |
| 51     | 57 | 58 | 53 | 47 | 40 | 30 | 23 | 28 | 40 | 49 |  |
| 60     | 62 | 58 | 47 | 39 | 38 | 35 | 31 | 35 | 41 | 46 |  |

# Experiment 3

121-Armed Bandit (Natural)

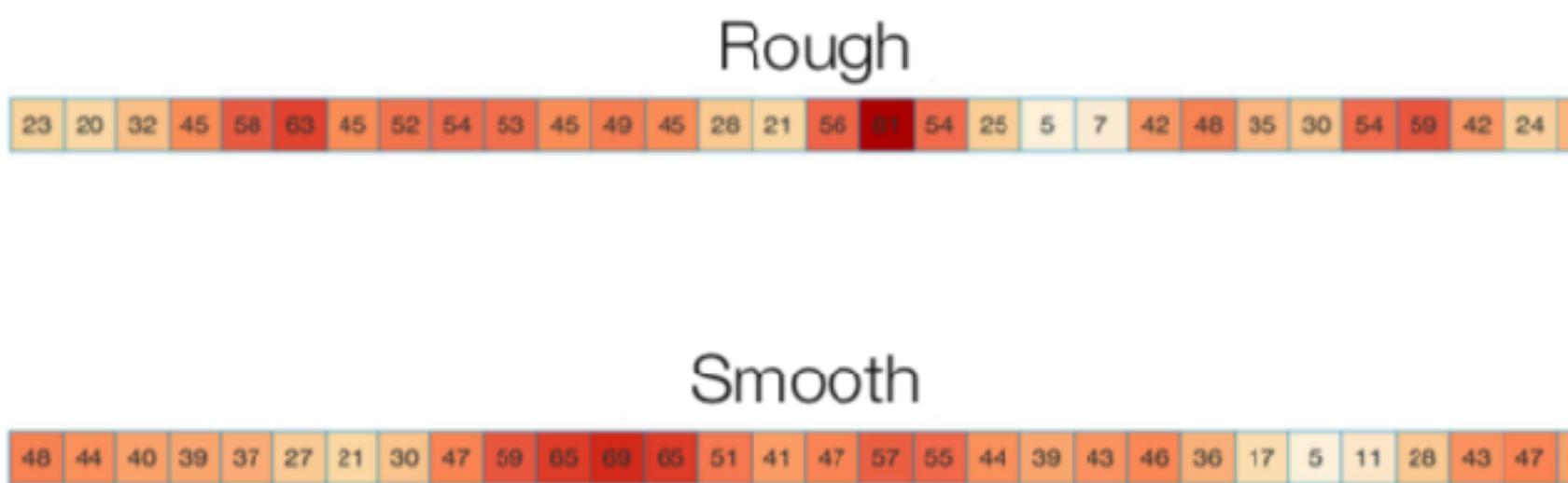
| Crop Yield |    |    |    |    |    |    |    |    |    |    |  |
|------------|----|----|----|----|----|----|----|----|----|----|--|
| 38         | 38 | 37 | 41 | 38 | 35 | 39 | 5  | 35 | 41 | 54 |  |
| 39         | 19 | 18 | 29 | 46 | 33 | 32 | 55 | 35 | 37 | 57 |  |
| 27         | 19 | 21 | 21 | 47 | 46 | 19 | 72 | 38 | 42 | 40 |  |
| 22         | 23 | 20 | 18 | 44 | 35 | 36 | 53 | 27 | 45 | 47 |  |
| 23         | 14 | 15 | 24 | 35 | 55 | 51 | 53 | 46 | 38 | 63 |  |
| 19         | 17 | 12 | 25 | 38 | 74 | 77 | 71 | 46 | 55 | 63 |  |
| 20         | 19 | 21 | 23 | 6  | 57 | 60 | 63 | 58 | 60 | 69 |  |
| 19         | 19 | 14 | 28 | 31 | 33 | 47 | 58 | 68 | 70 | 61 |  |
| 16         | 19 | 19 | 22 | 28 | 31 | 54 | 69 | 74 | 74 | 65 |  |
| 21         | 11 | 31 | 30 | 21 | 43 | 50 | 60 | 70 | 72 | 57 |  |
| 27         | 21 | 27 | 31 | 45 | 45 | 66 | 56 | 53 | 50 | 42 |  |

Number of grids left: **6**  
Number of clicks left: **1**

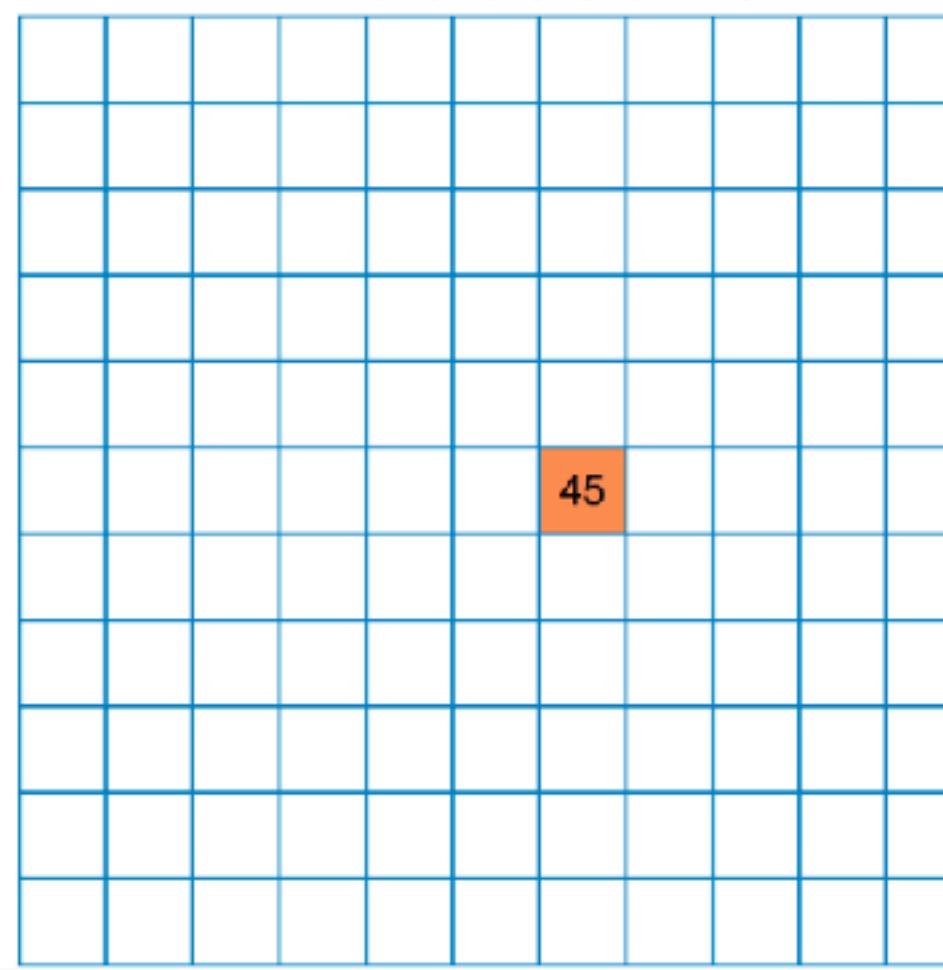


# Experiment 1

30-Armed Bandit (Univariate)



Number of grids left: **6**  
Number of clicks left: **40**



Participants acquired rewards by clicking  
on new or previously revealed tiles



# Experiment 2

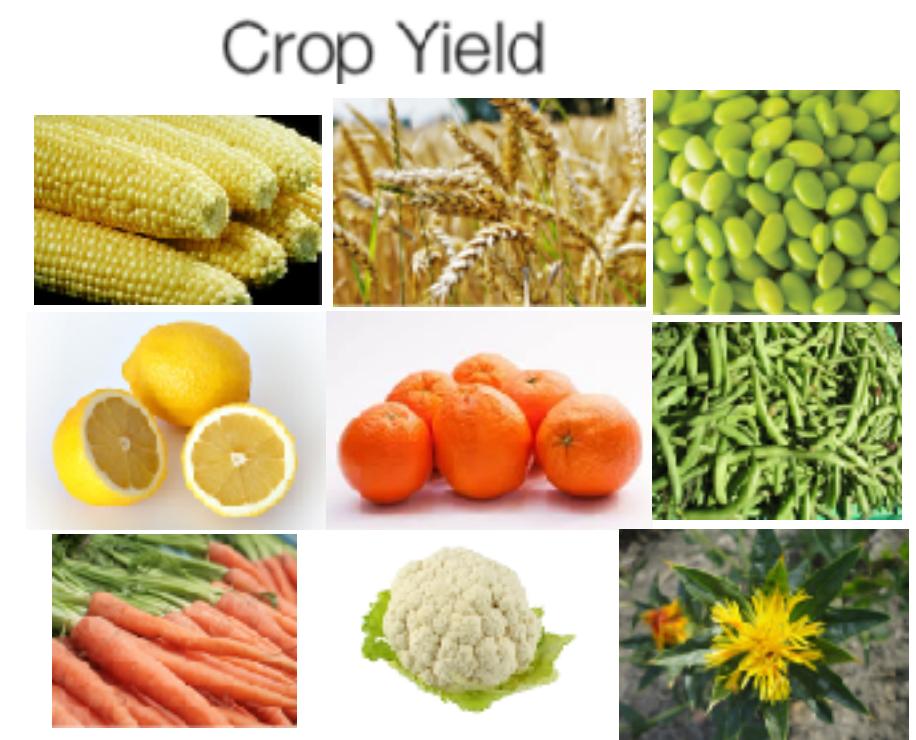
121-Armed Bandit (Bivariate)

|    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|
| 39 | 42 | 54 | 50 | 44 | 70 | 72 | 72 | 57 | 35 | 12 |
| 31 | 21 | 28 | 41 | 44 | 44 | 47 | 64 | 54 | 40 | 27 |
| 51 | 37 | 30 | 29 | 26 | 24 | 31 | 48 | 45 | 33 | 25 |
| 64 | 55 | 58 | 43 | 36 | 38 | 32 | 26 | 30 | 27 | 28 |
| 58 | 62 | 41 | 19 | 23 | 26 | 29 | 37 | 44 | 33 | 27 |
| 65 | 76 | 33 | 5  | 7  | 13 | 27 | 50 | 48 | 17 | 25 |
| 38 | 44 | 35 | 24 | 17 | 17 | 16 | 33 | 48 | 33 | 45 |
| 45 | 29 | 21 | 31 | 50 | 34 | 23 | 29 | 41 | 42 | 37 |
| 49 | 24 | 13 | 23 | 38 | 38 | 38 | 34 | 47 | 63 | 40 |
| 55 | 31 | 24 | 24 | 26 | 23 | 22 | 26 | 40 | 64 | 66 |
| 64 | 53 | 41 | 50 | 59 | 42 | 19 | 29 | 28 | 29 | 49 |

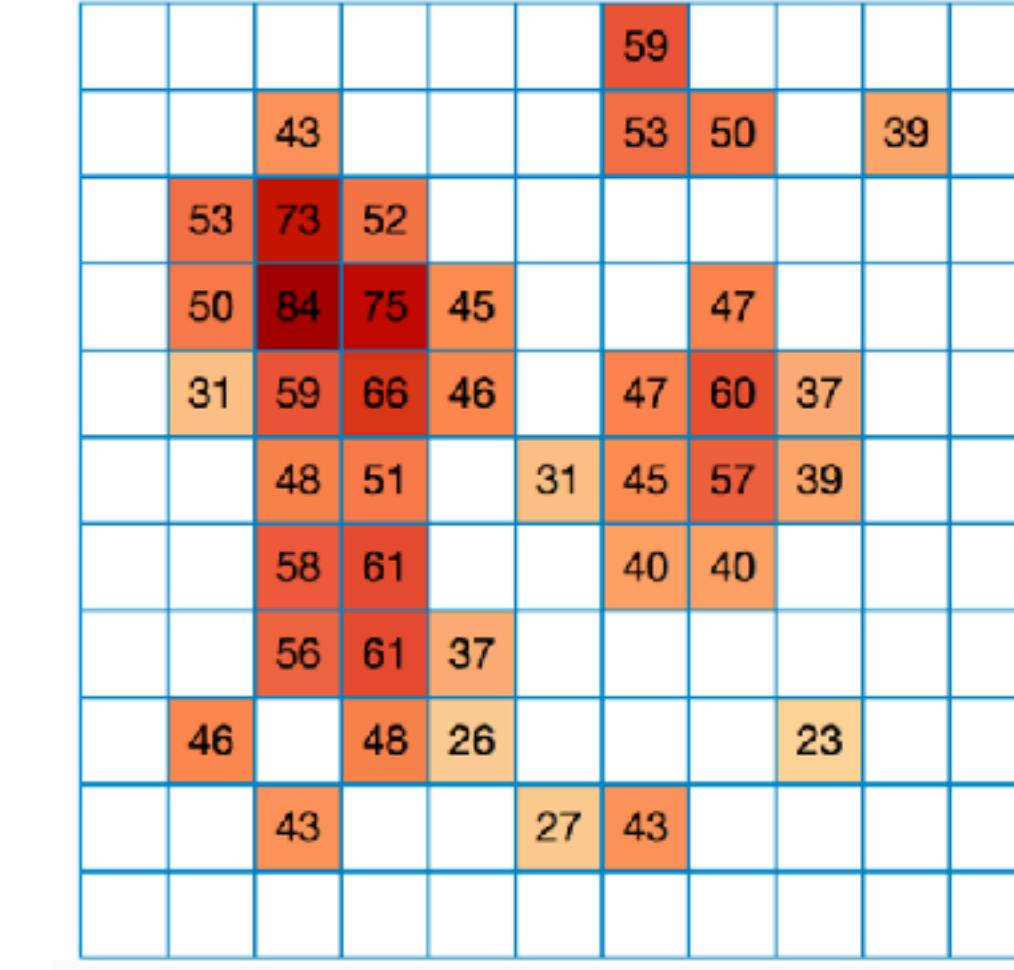
|    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|
| 7  | 5  | 10 | 22 | 32 | 32 | 28 | 24 | 22 | 26 | 33 |
| 6  | 11 | 19 | 29 | 38 | 41 | 42 | 40 | 37 | 36 | 40 |
| 22 | 27 | 30 | 35 | 43 | 50 | 53 | 53 | 51 | 49 | 46 |
| 45 | 44 | 38 | 36 | 40 | 46 | 47 | 49 | 54 | 55 | 48 |
| 61 | 55 | 46 | 40 | 37 | 32 | 27 | 31 | 44 | 52 | 44 |
| 62 | 59 | 57 | 54 | 44 | 27 | 14 | 17 | 33 | 46 | 45 |
| 53 | 59 | 68 | 71 | 59 | 36 | 17 | 15 | 28 | 45 | 51 |
| 46 | 57 | 71 | 77 | 67 | 47 | 26 | 18 | 27 | 45 | 56 |
| 45 | 56 | 65 | 67 | 60 | 46 | 29 | 20 | 27 | 42 | 55 |
| 51 | 57 | 58 | 53 | 47 | 40 | 30 | 23 | 28 | 40 | 49 |
| 60 | 62 | 58 | 47 | 39 | 38 | 35 | 31 | 35 | 41 | 46 |

# Experiment 3

121-Armed Bandit (Natural)



Number of grids left: **6**  
Number of clicks left: **1**



# Modeling Human Search

# Modeling Human Search

- Exploration is not performed blindly



# Modeling Human Search

- Exploration is not performed blindly
- Search is guided by generalization

# Modeling Human Search

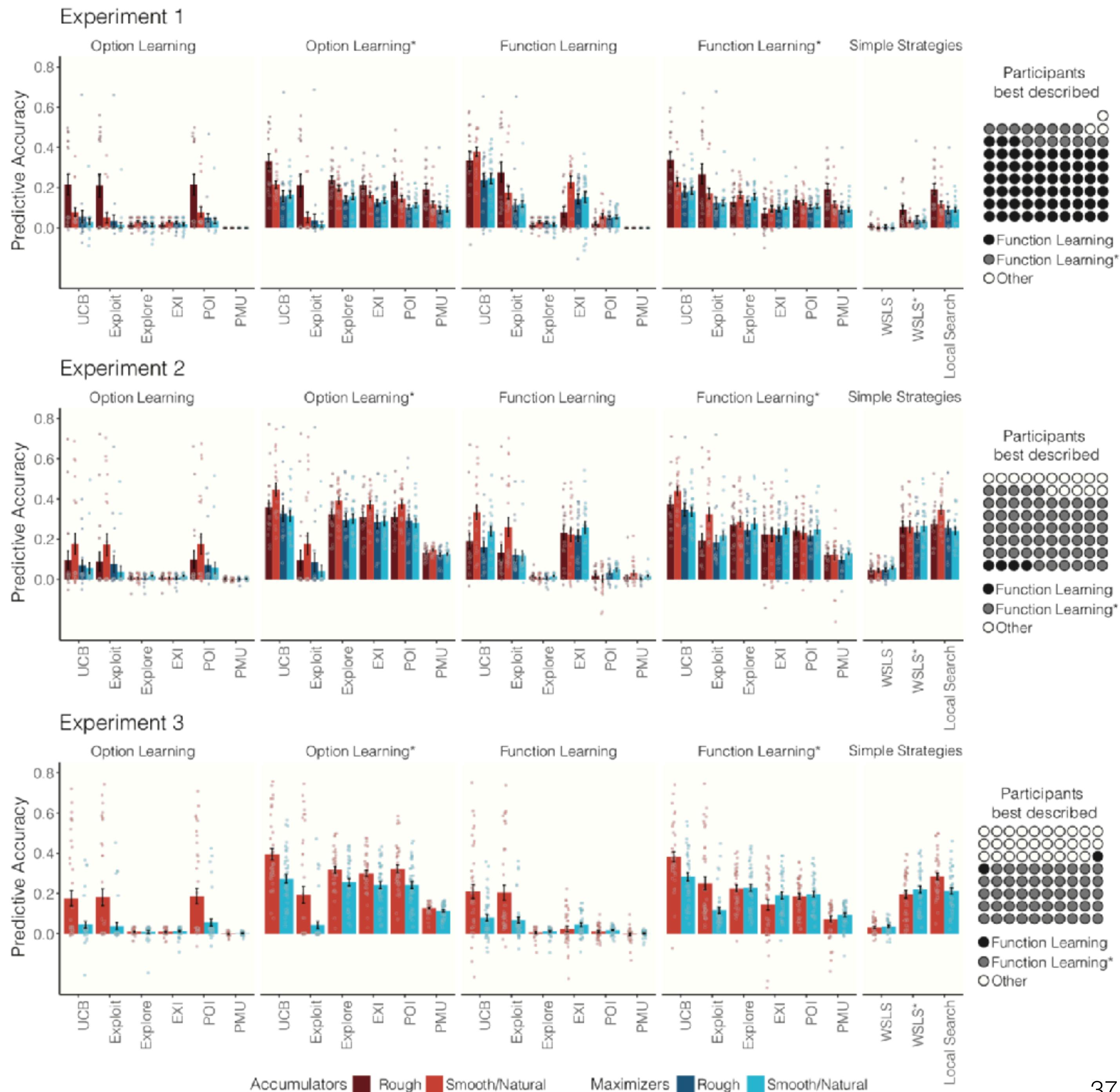
- Exploration is not performed blindly
- Search is guided by generalization
- Generalization as Bayesian inference about novel options:
  - Expected reward
  - Uncertainty

# Modeling Human Search

- Exploration is not performed blindly
- Search is guided by generalization
- Generalization as Bayesian inference about novel options:
  - Expected reward
  - Uncertainty
- Human search is directed towards both ingredients

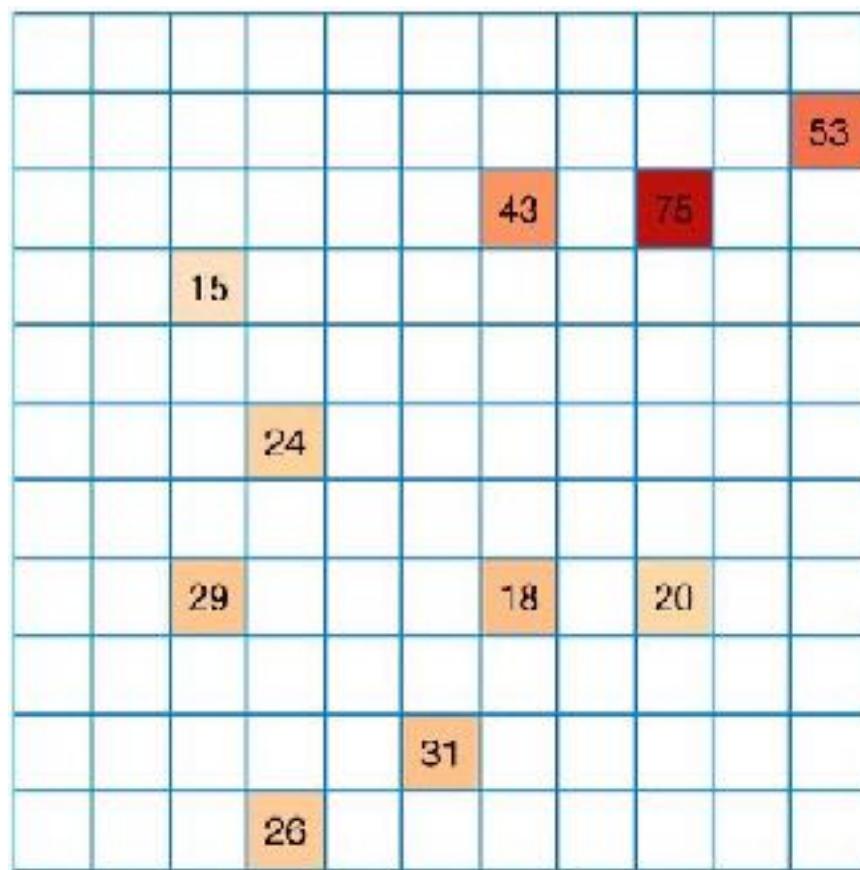
# Model comparison

- We performed a large-scale comparison of 27 different models using cross-validated (out-of-sample) predictions
- Some heuristic models but mostly reinforcement learning models \* sampling strategies
- ... here, I focus on the best model, which consistently outperformed all others across a variety of manipulation checks

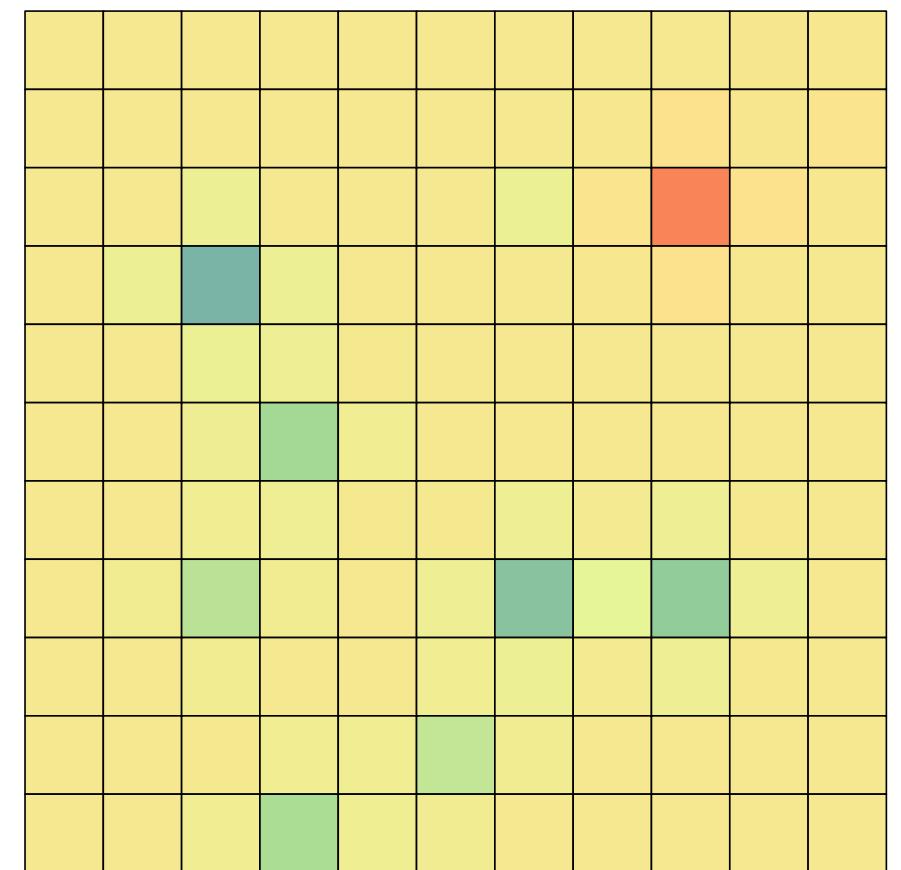


# GP-UCB Model

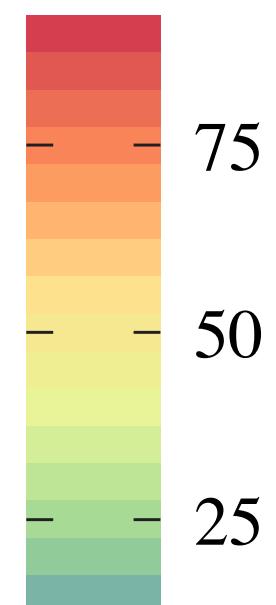
Observations



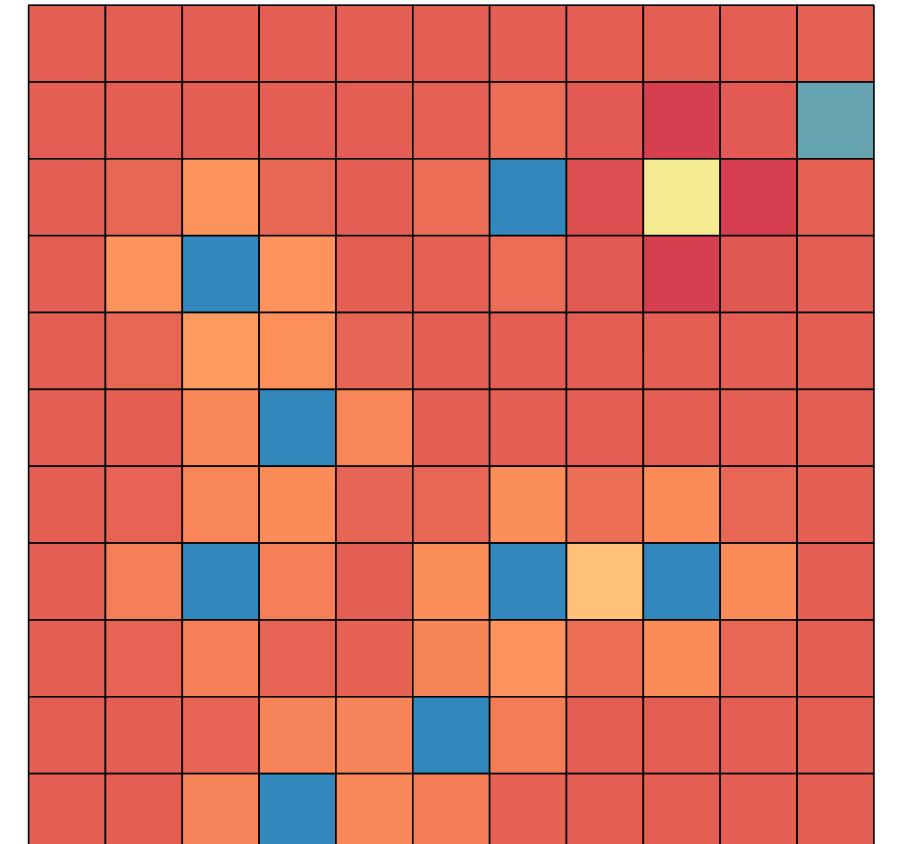
Gaussian Process (GP)



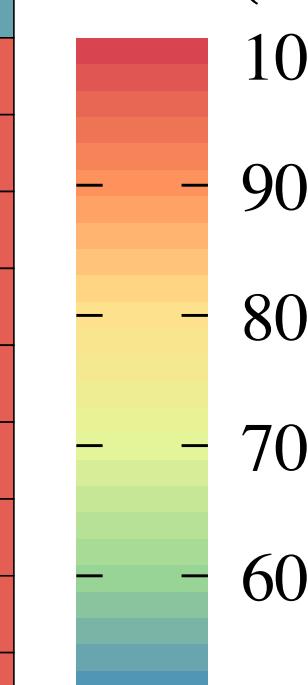
$\mu(\mathbf{x})$



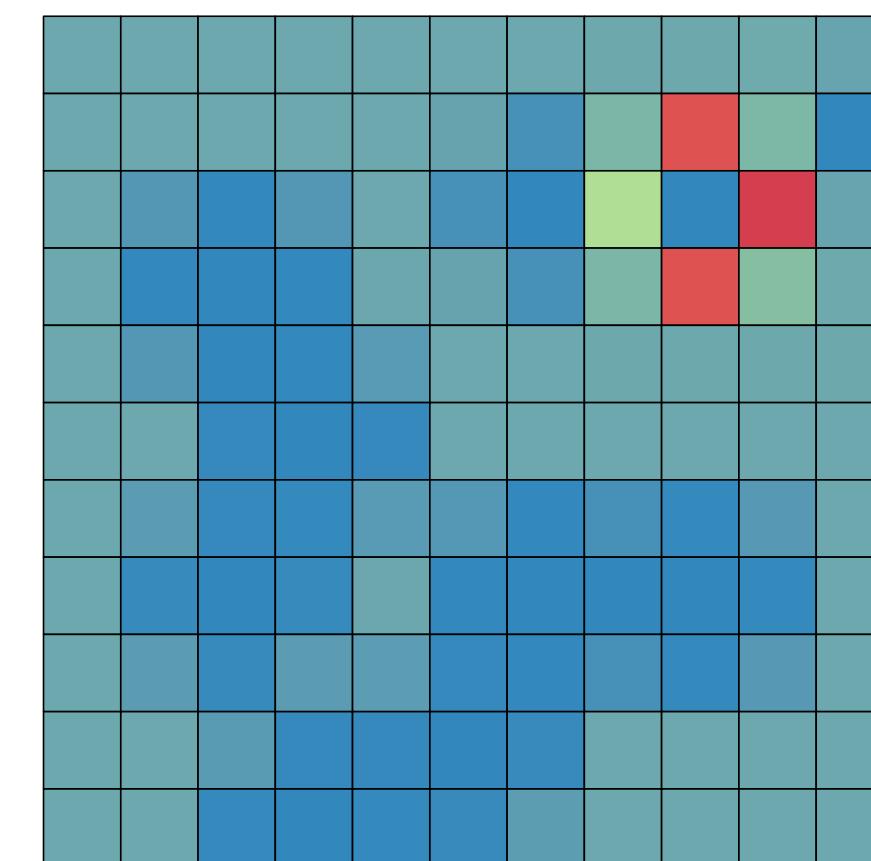
Upper Confidence Bound  
(UCB) Sampling



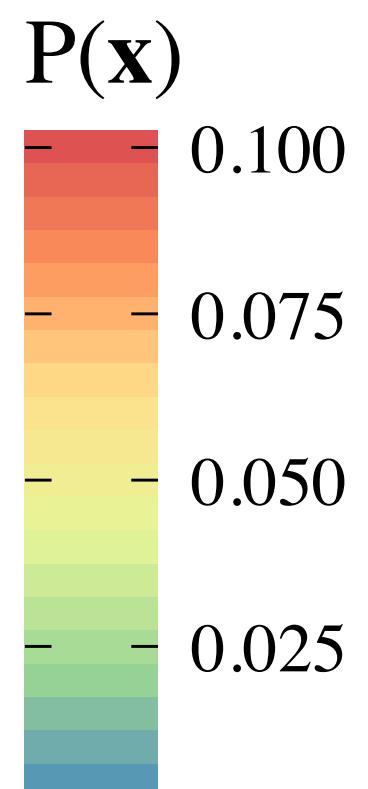
UCB( $\mathbf{x}$ )



Softmax Choice Rule

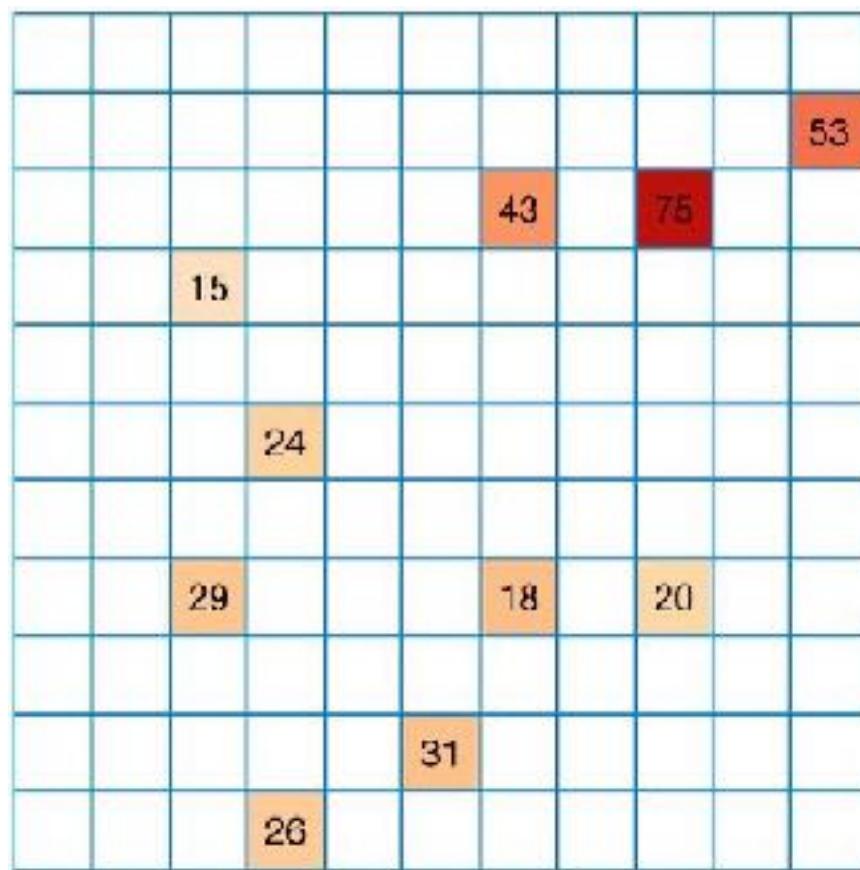


P( $\mathbf{x}$ )

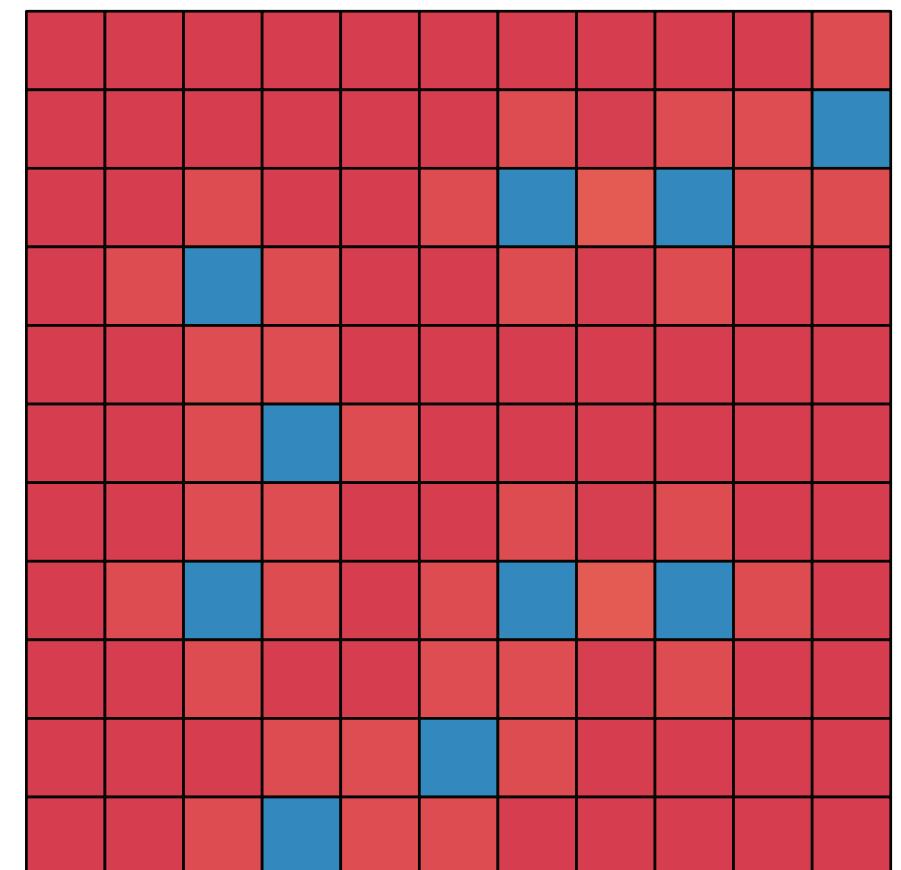


# GP-UCB Model

Observations



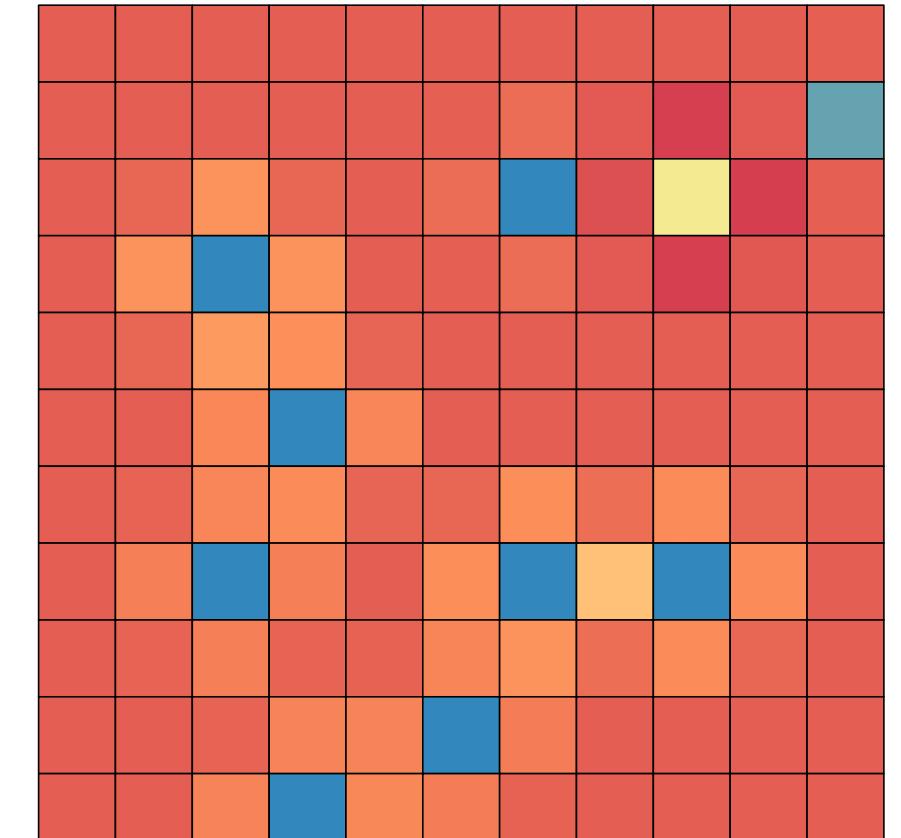
Gaussian Process (GP)



$\sigma(\mathbf{x})$

100  
75  
50  
25

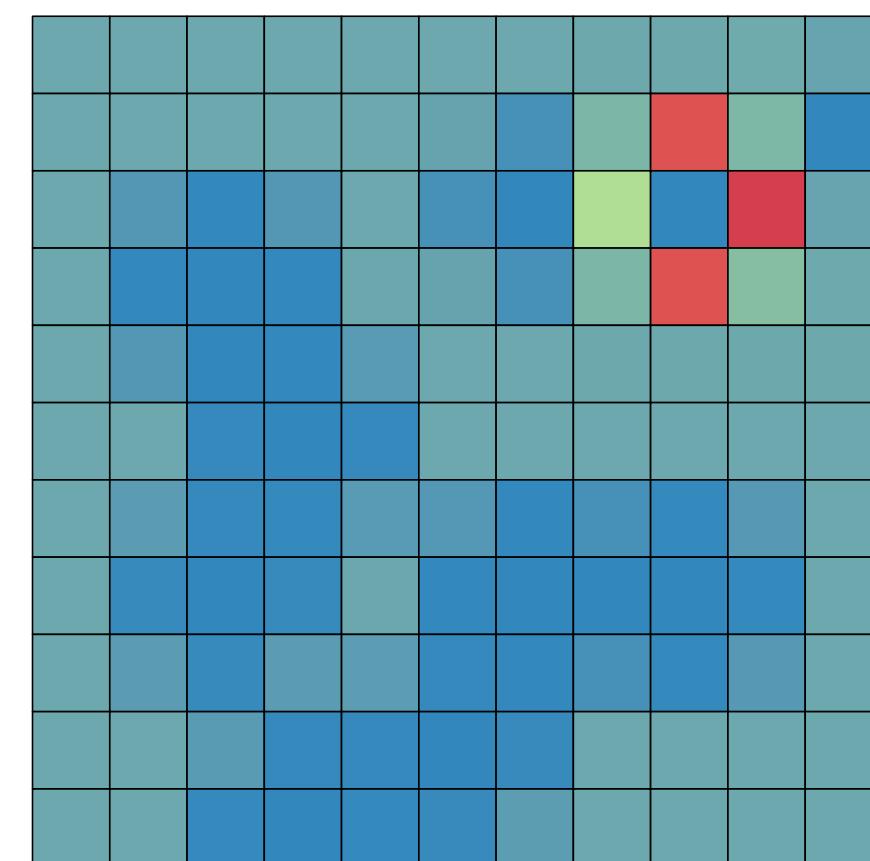
Upper Confidence Bound  
(UCB) Sampling



UCB( $\mathbf{x}$ )

100  
90  
80  
70  
60

Softmax Choice Rule

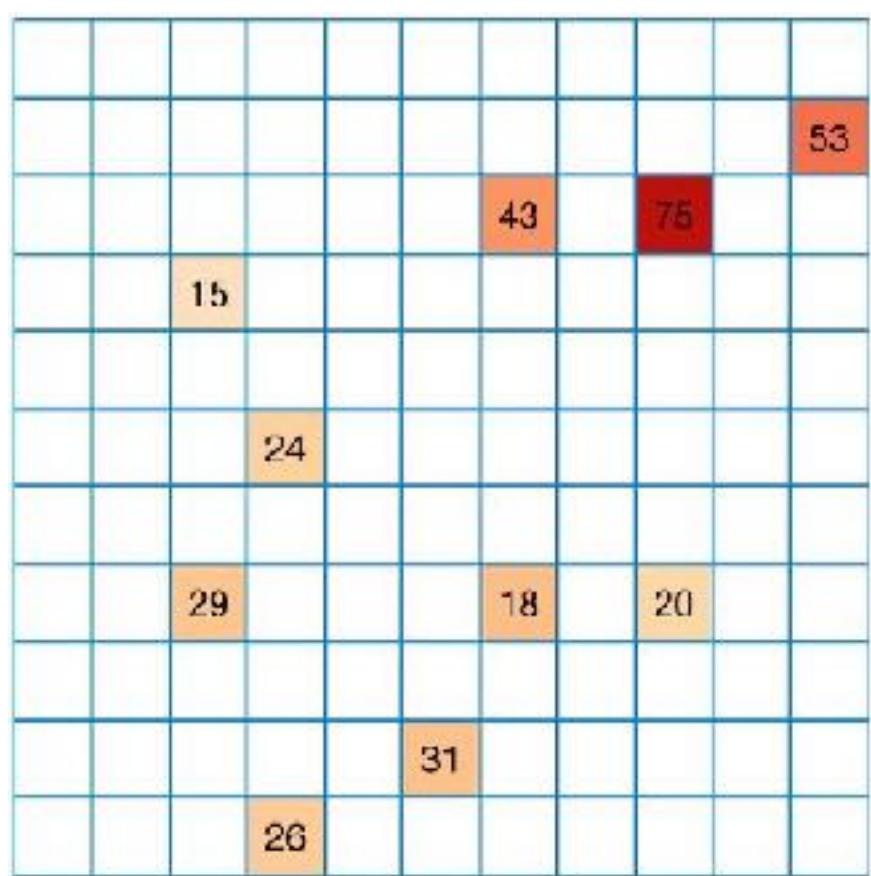


P( $\mathbf{x}$ )

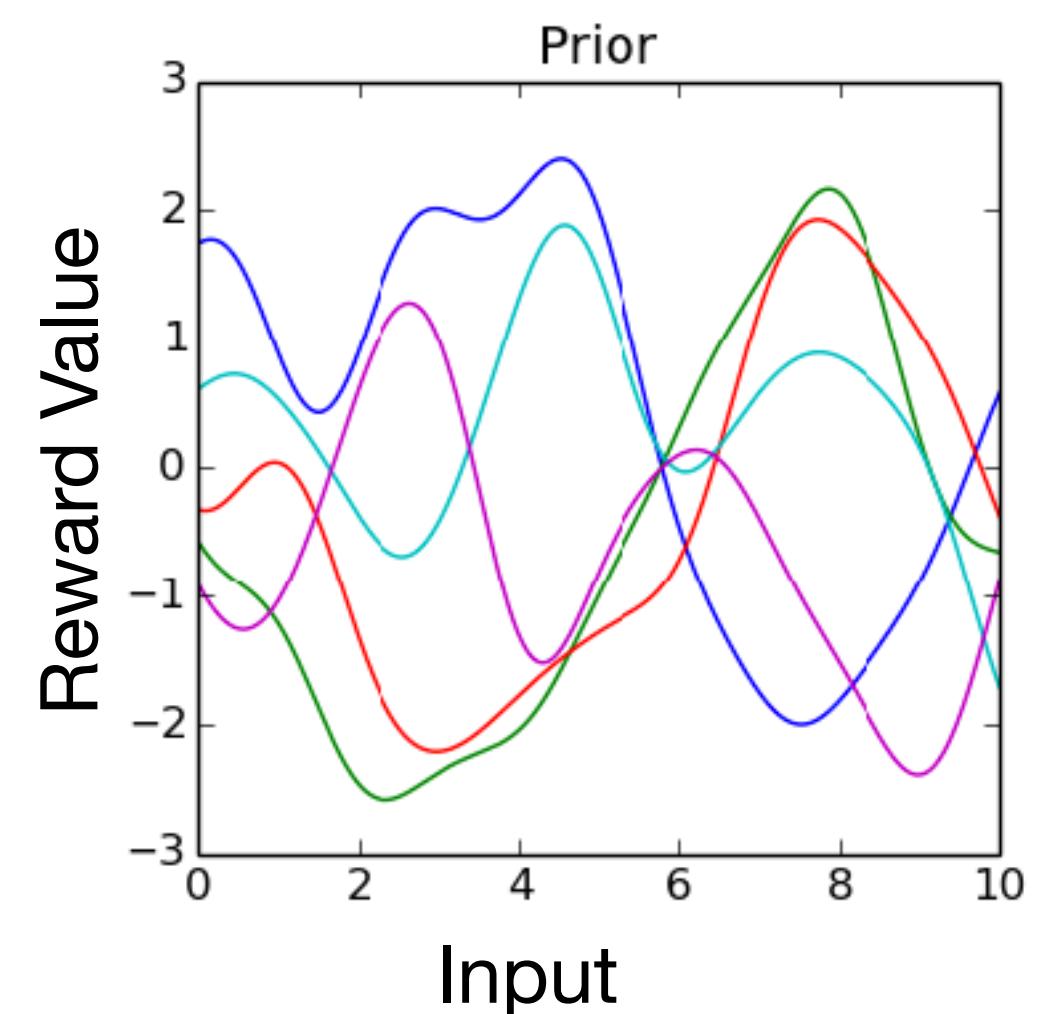
0.100  
0.075  
0.050  
0.025

# GP-UCB Model

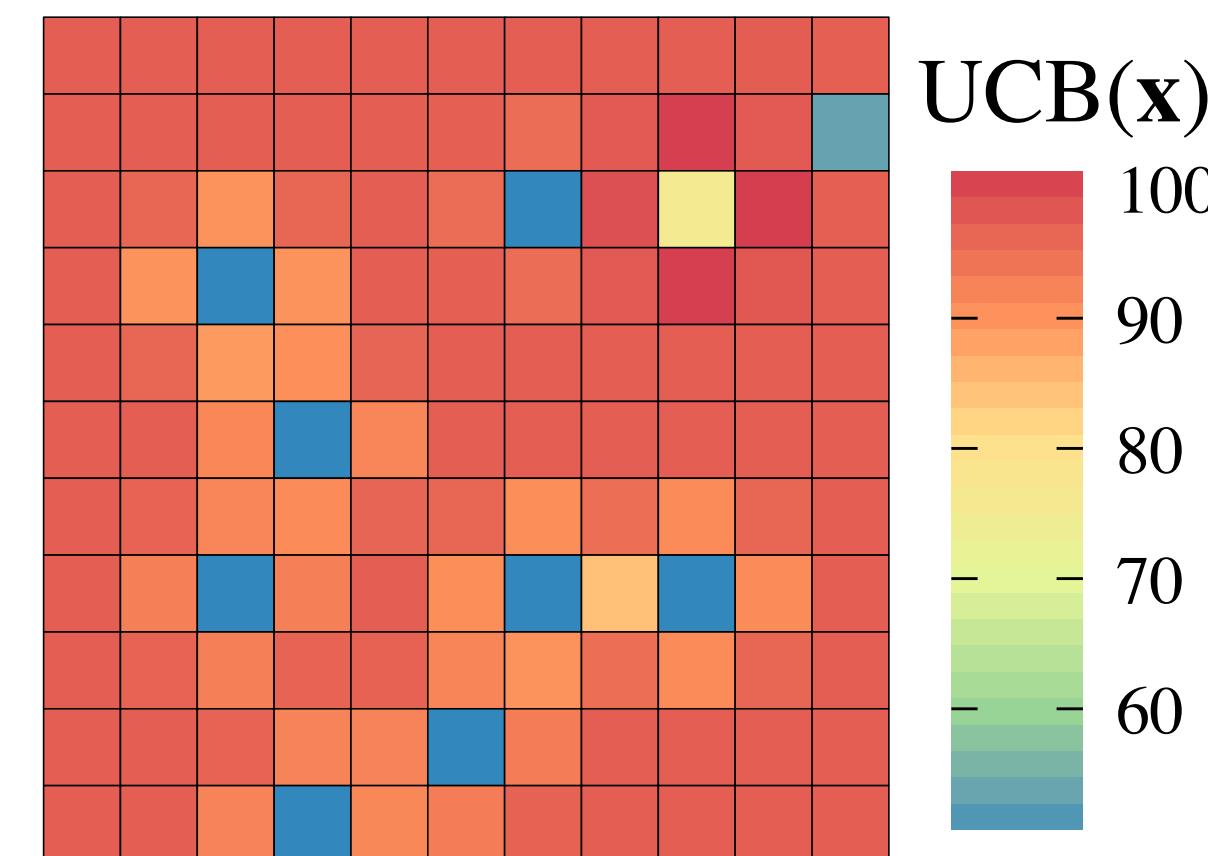
Observations



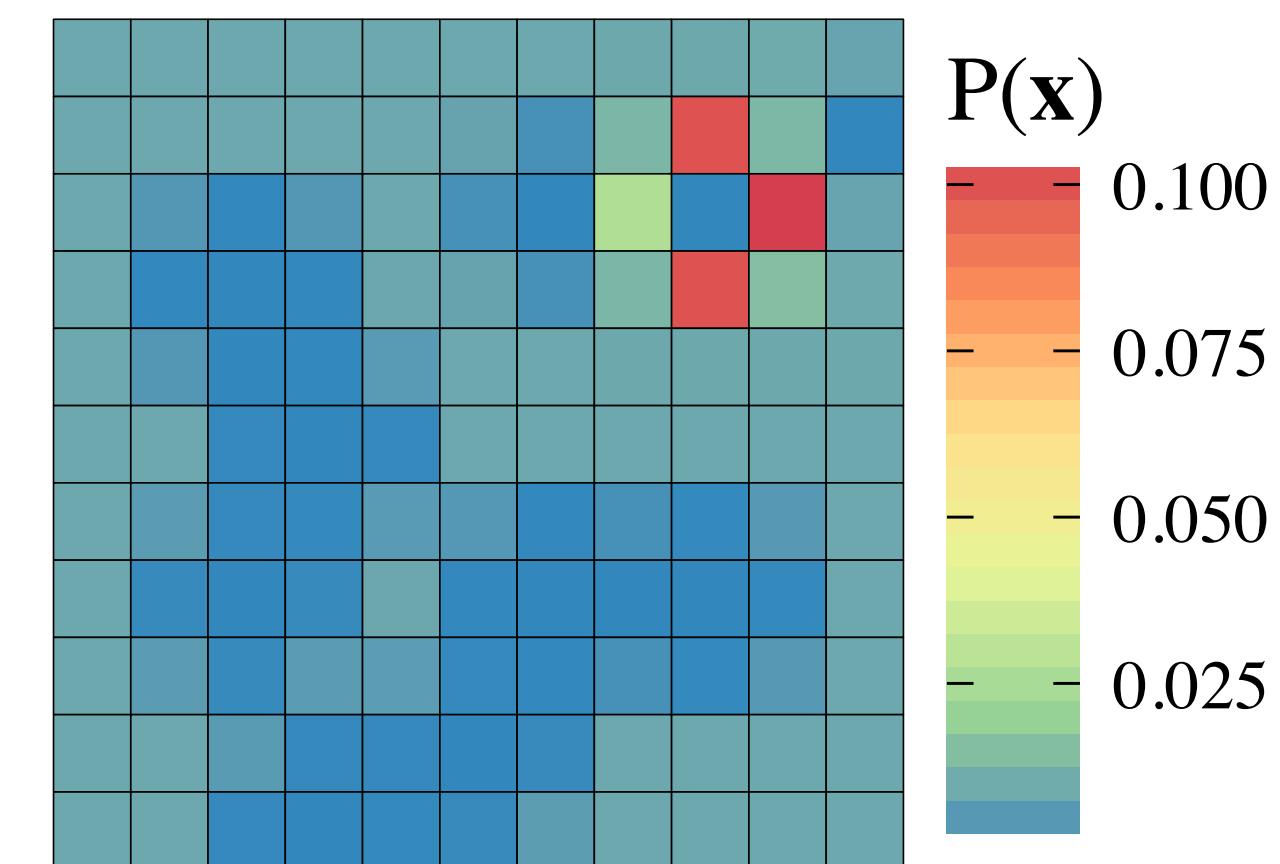
Gaussian Process (GP)



Upper Confidence Bound  
(UCB) Sampling

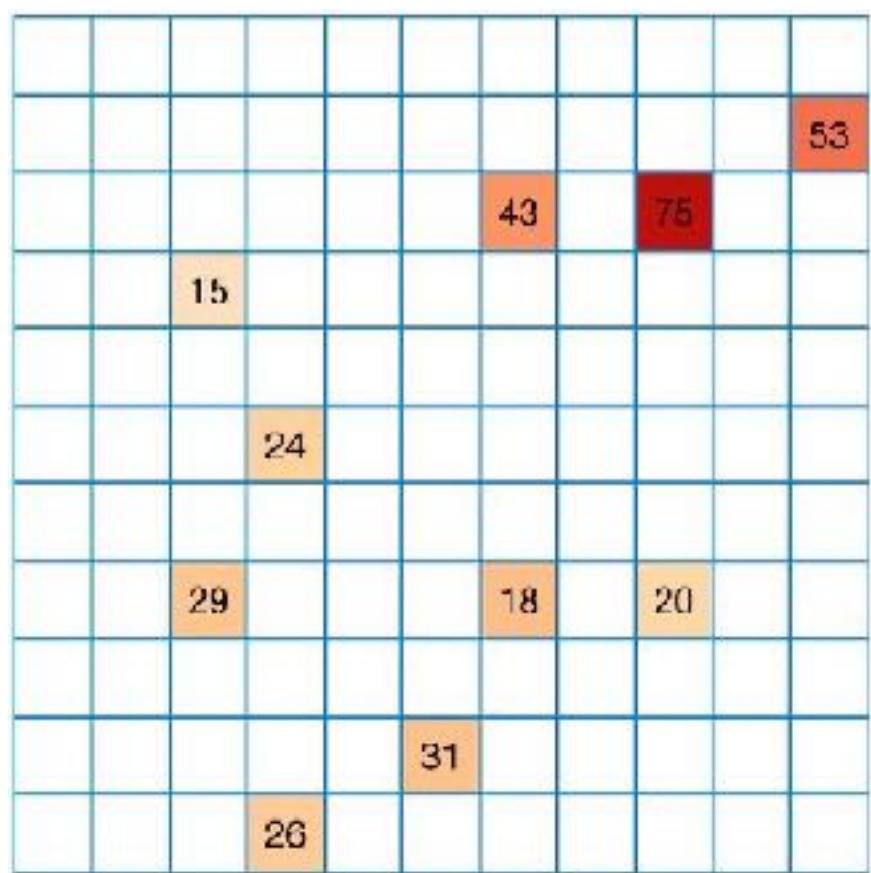


Softmax Choice Rule

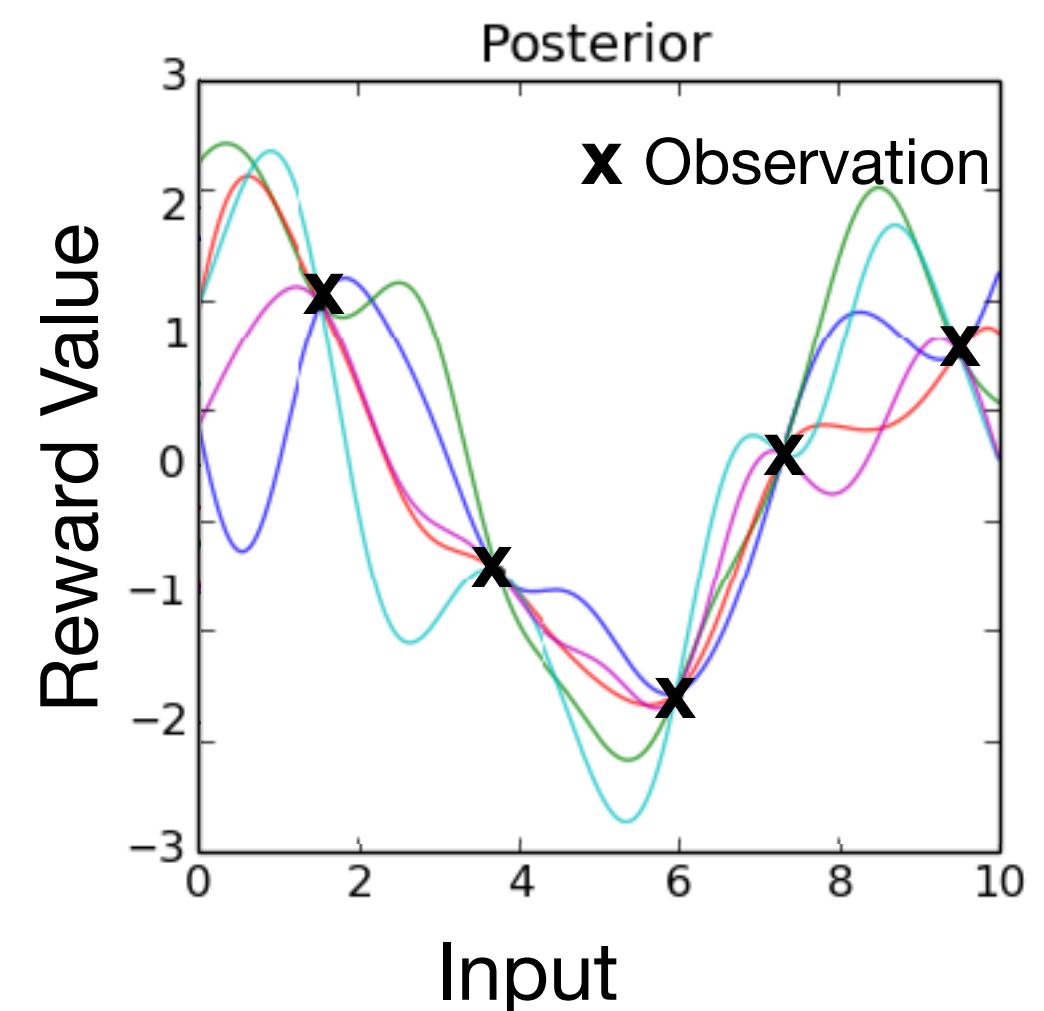


# GP-UCB Model

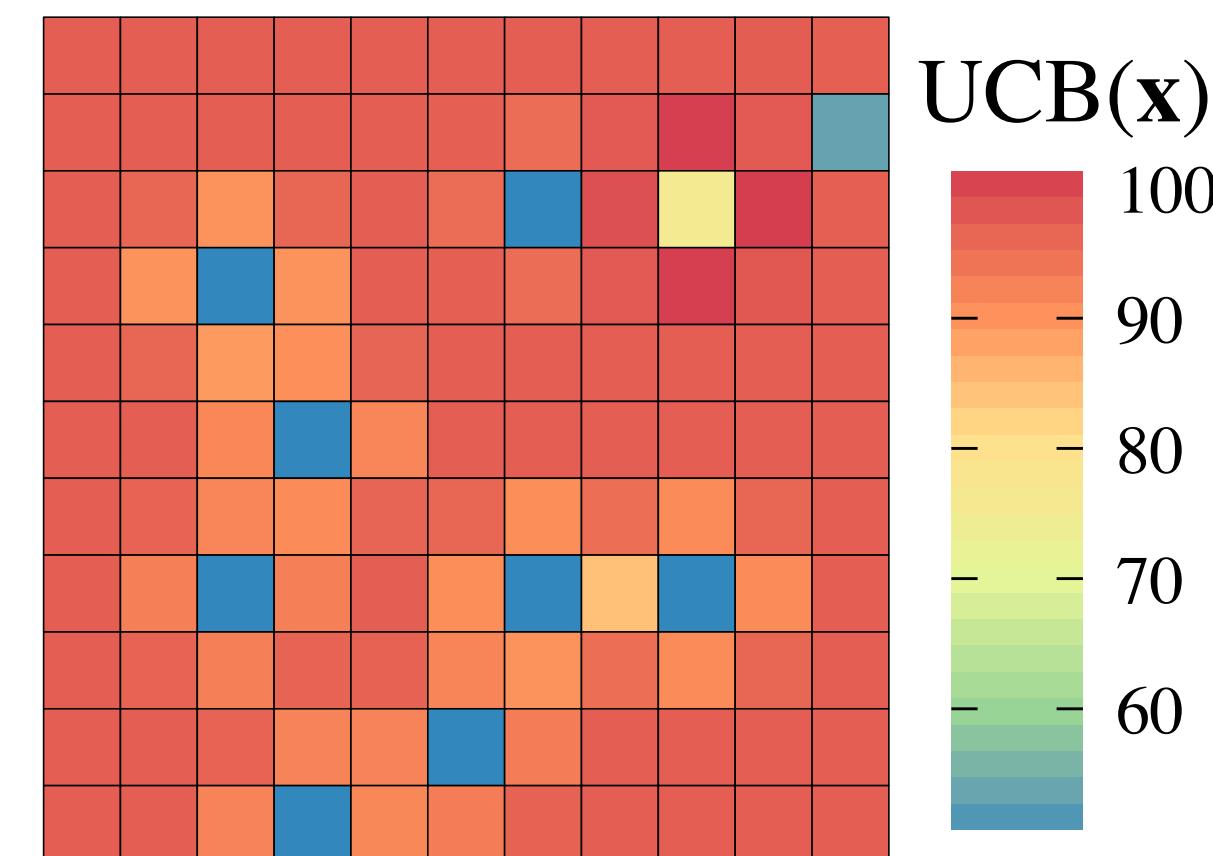
Observations



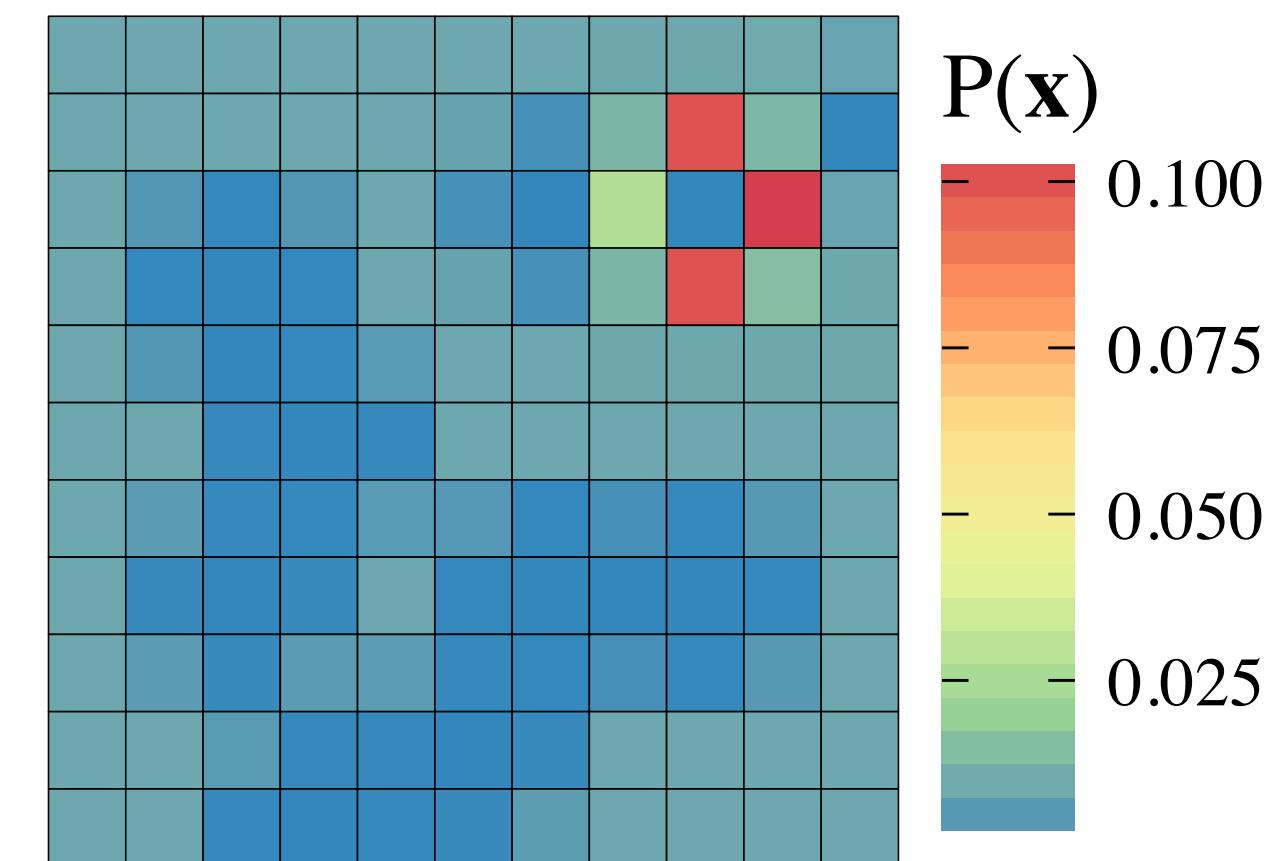
Gaussian Process (GP)



Upper Confidence Bound  
(UCB) Sampling

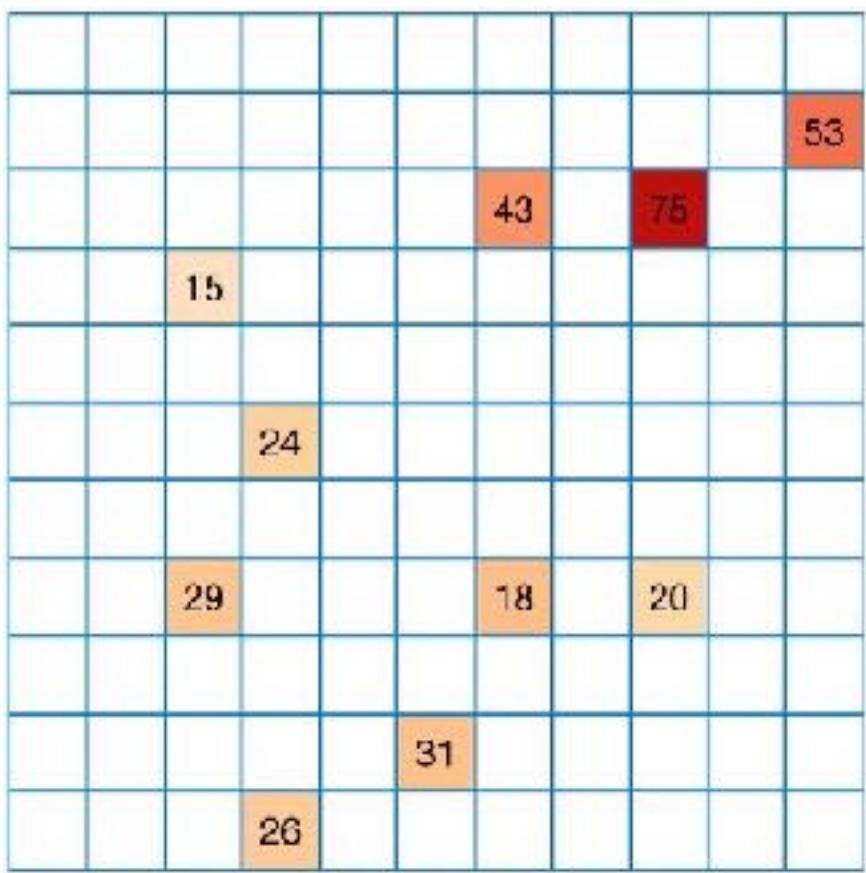


Softmax Choice Rule

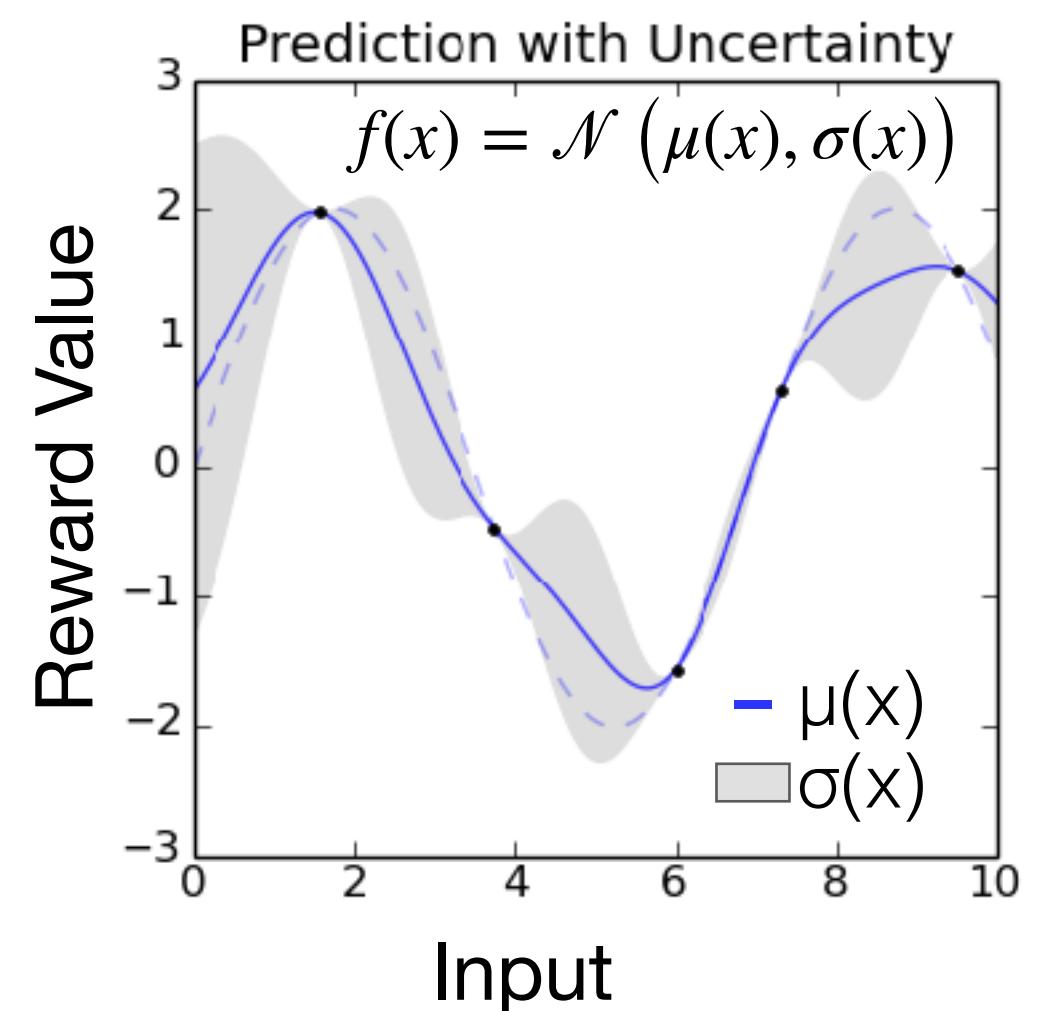


# GP-UCB Model

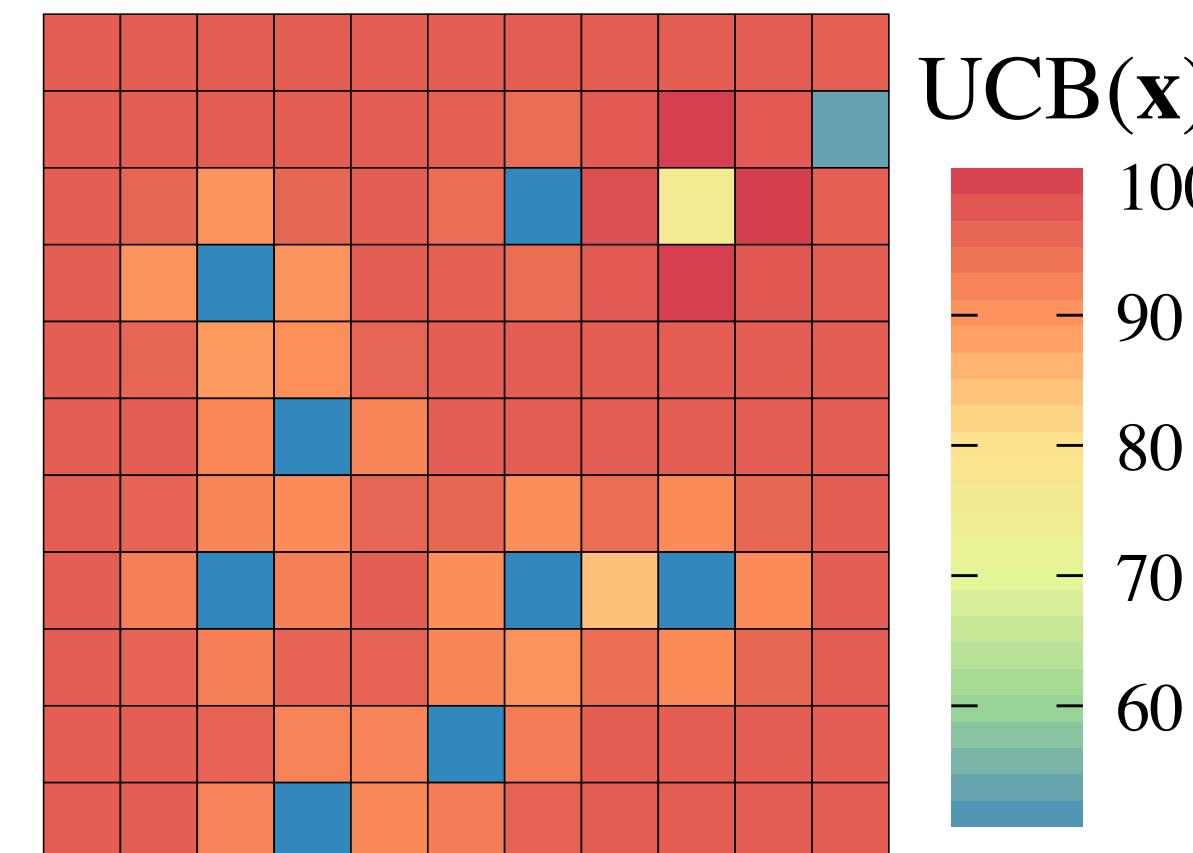
Observations



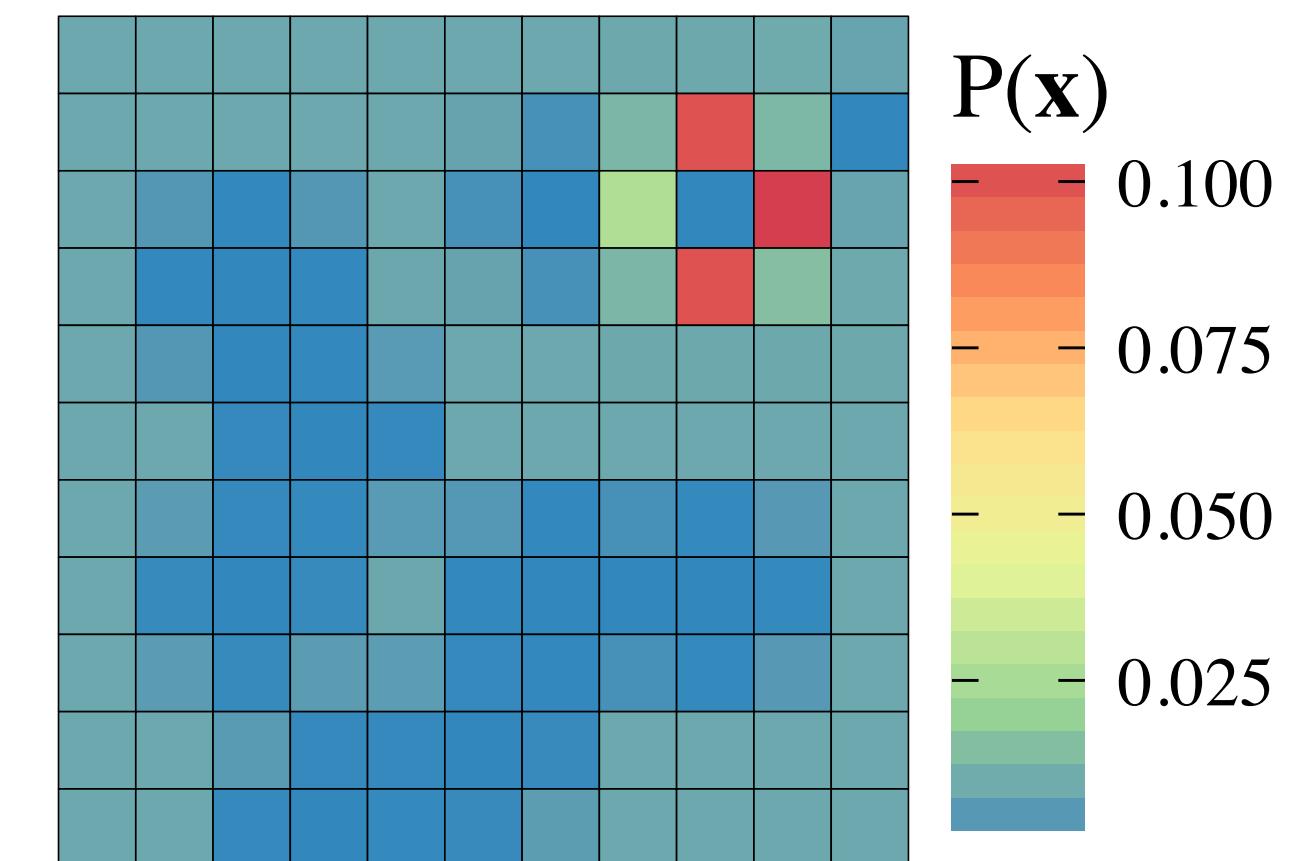
Gaussian Process (GP)



Upper Confidence Bound (UCB) Sampling

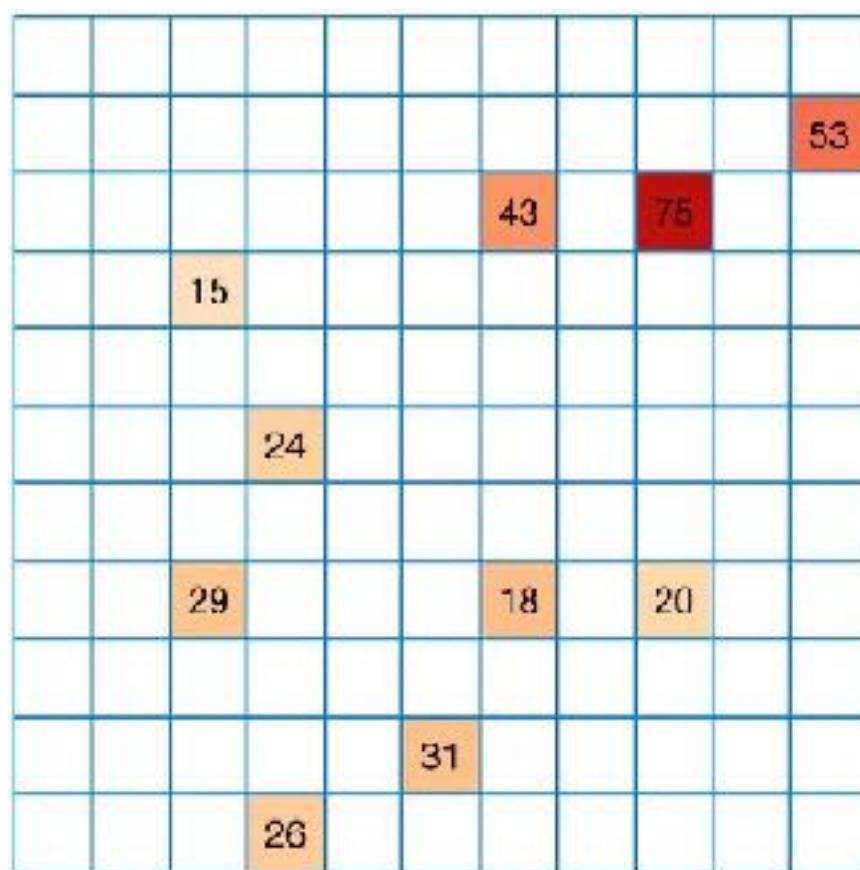


Softmax Choice Rule

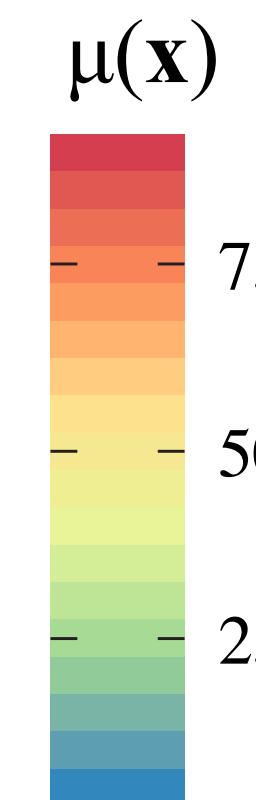
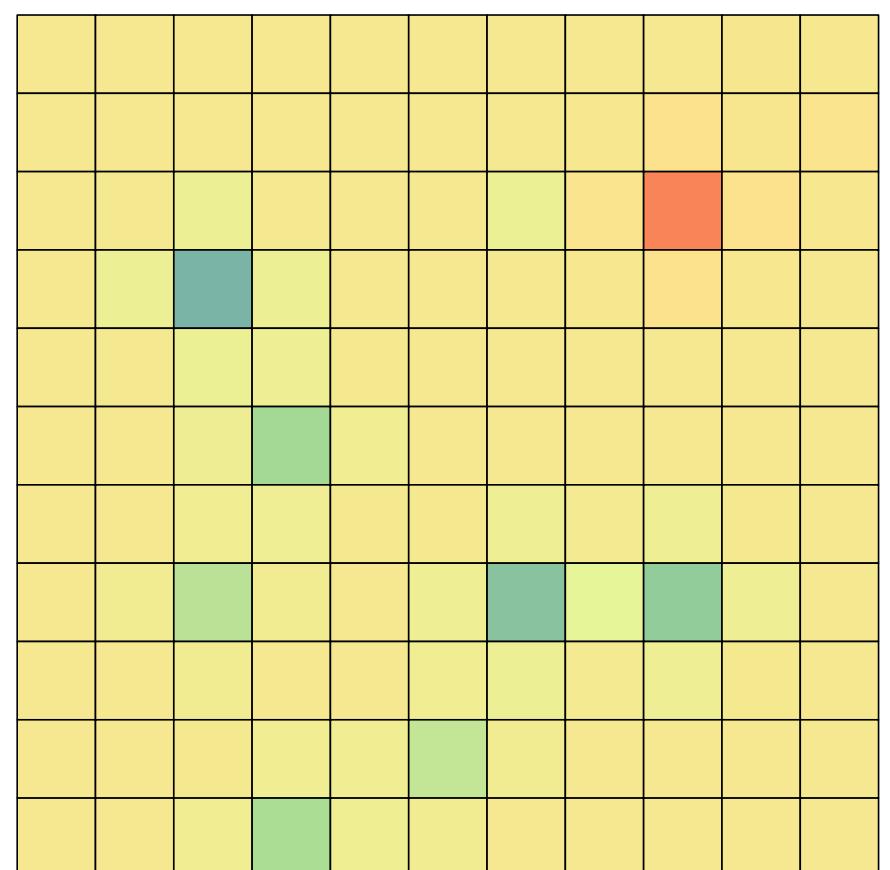


# GP-UCB Model

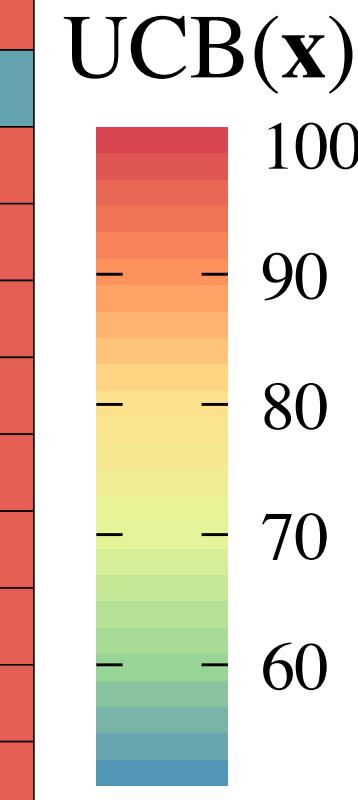
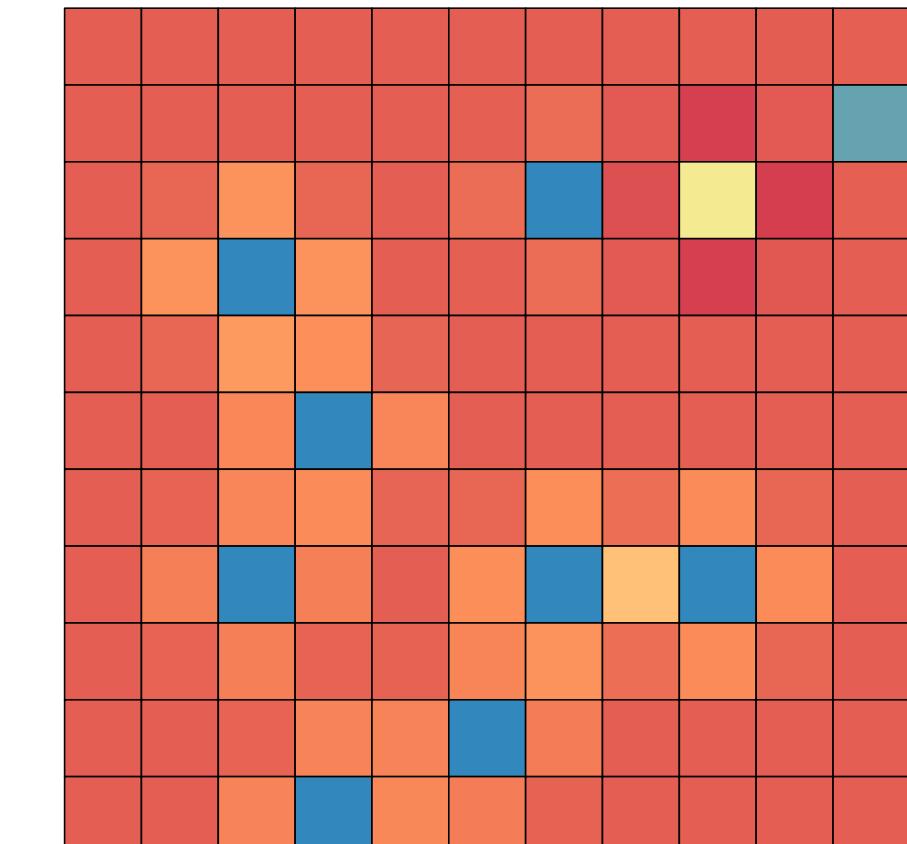
Observations



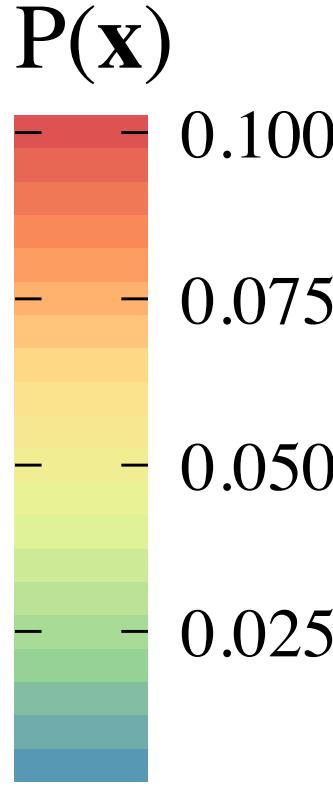
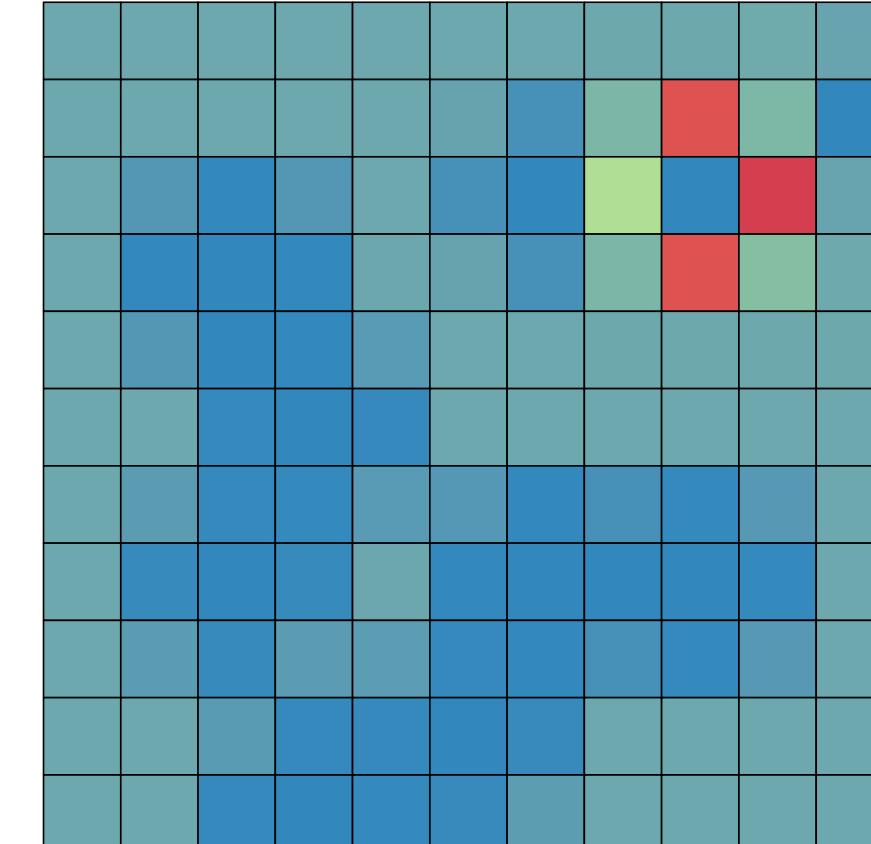
Gaussian Process (GP)



Upper Confidence Bound  
(UCB) Sampling

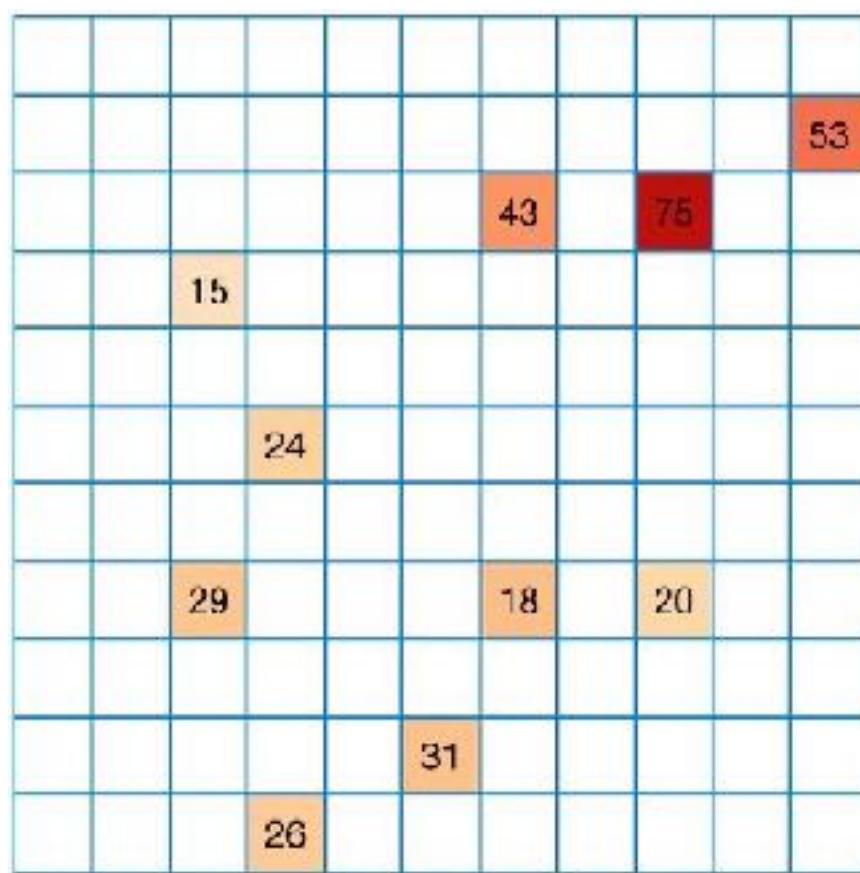


Softmax Choice Rule

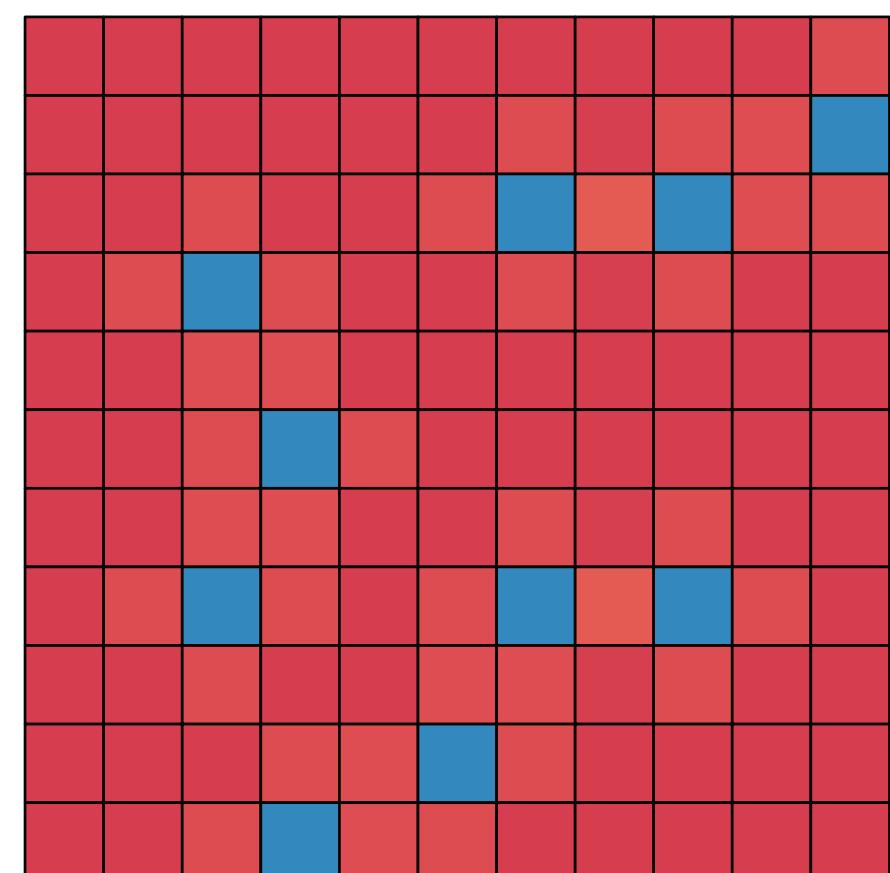


# GP-UCB Model

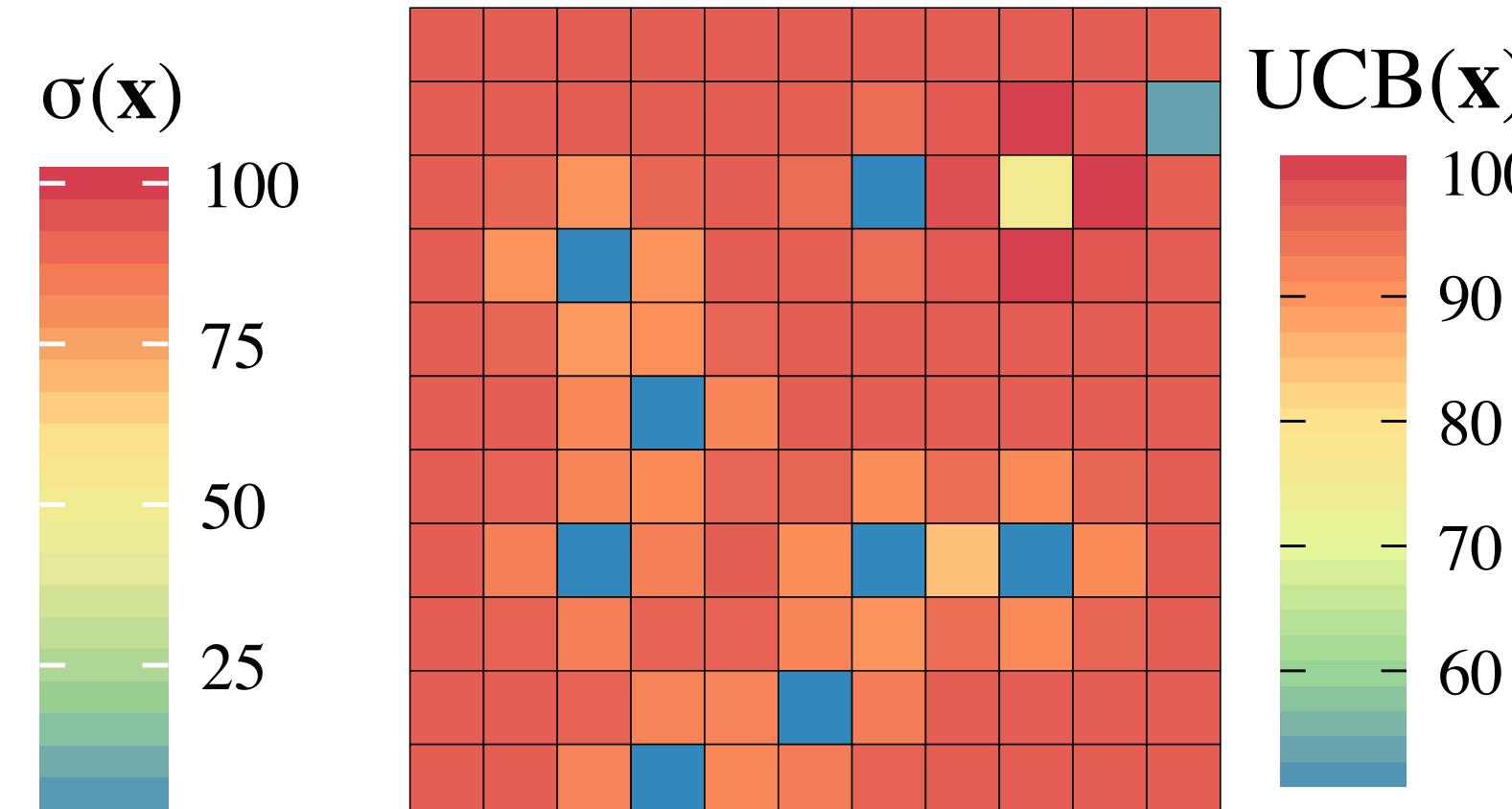
Observations



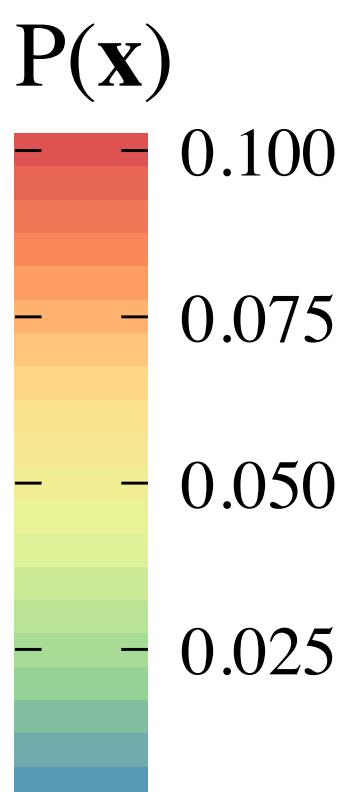
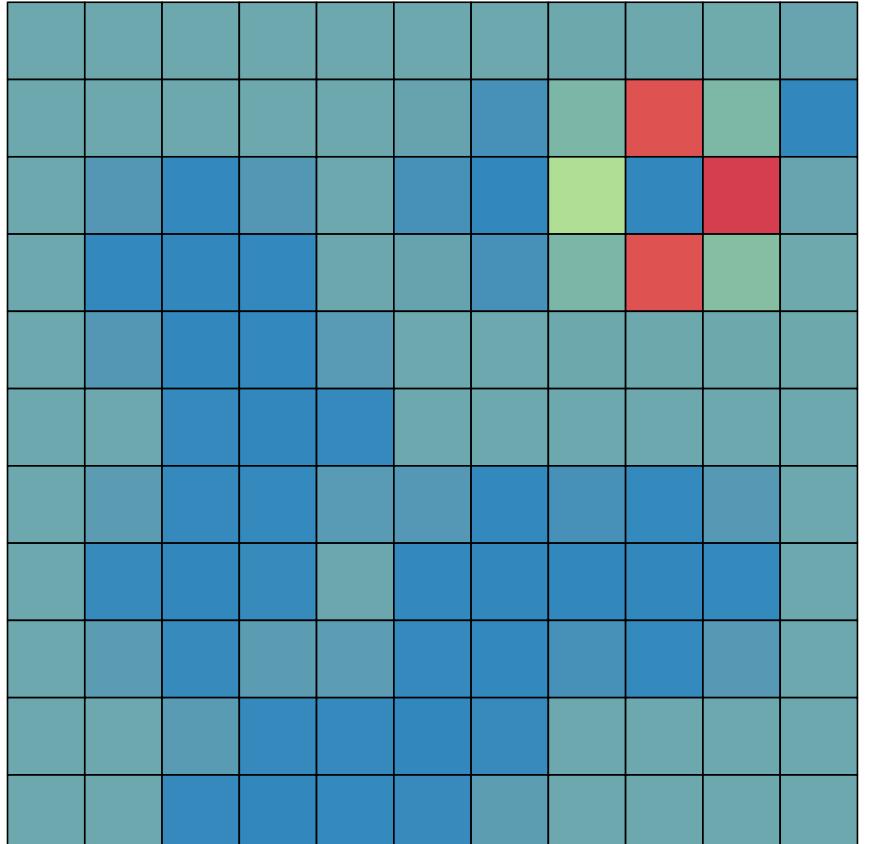
Gaussian Process (GP)



Upper Confidence Bound  
(UCB) Sampling

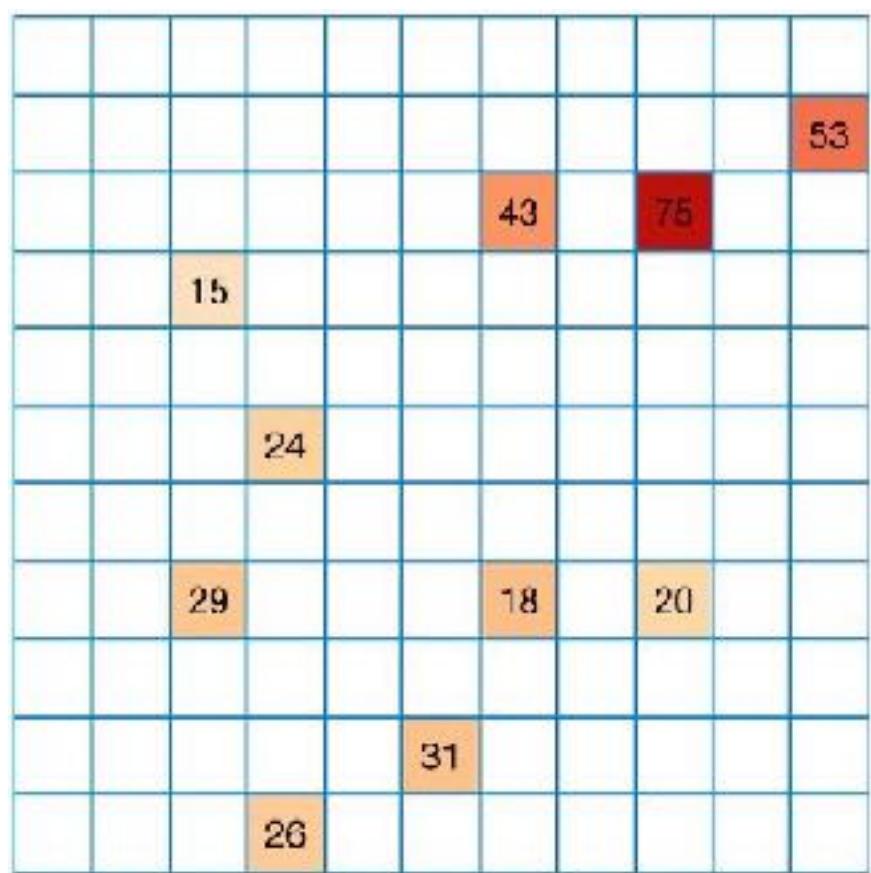


Softmax Choice Rule

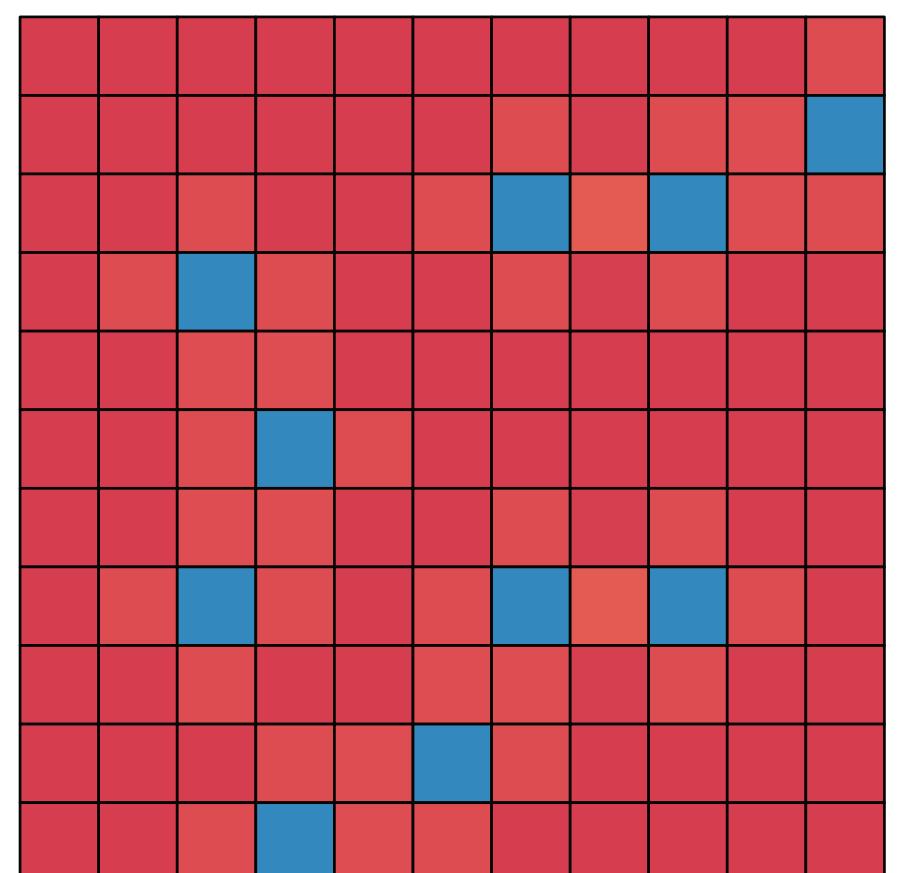


# GP-UCB Model

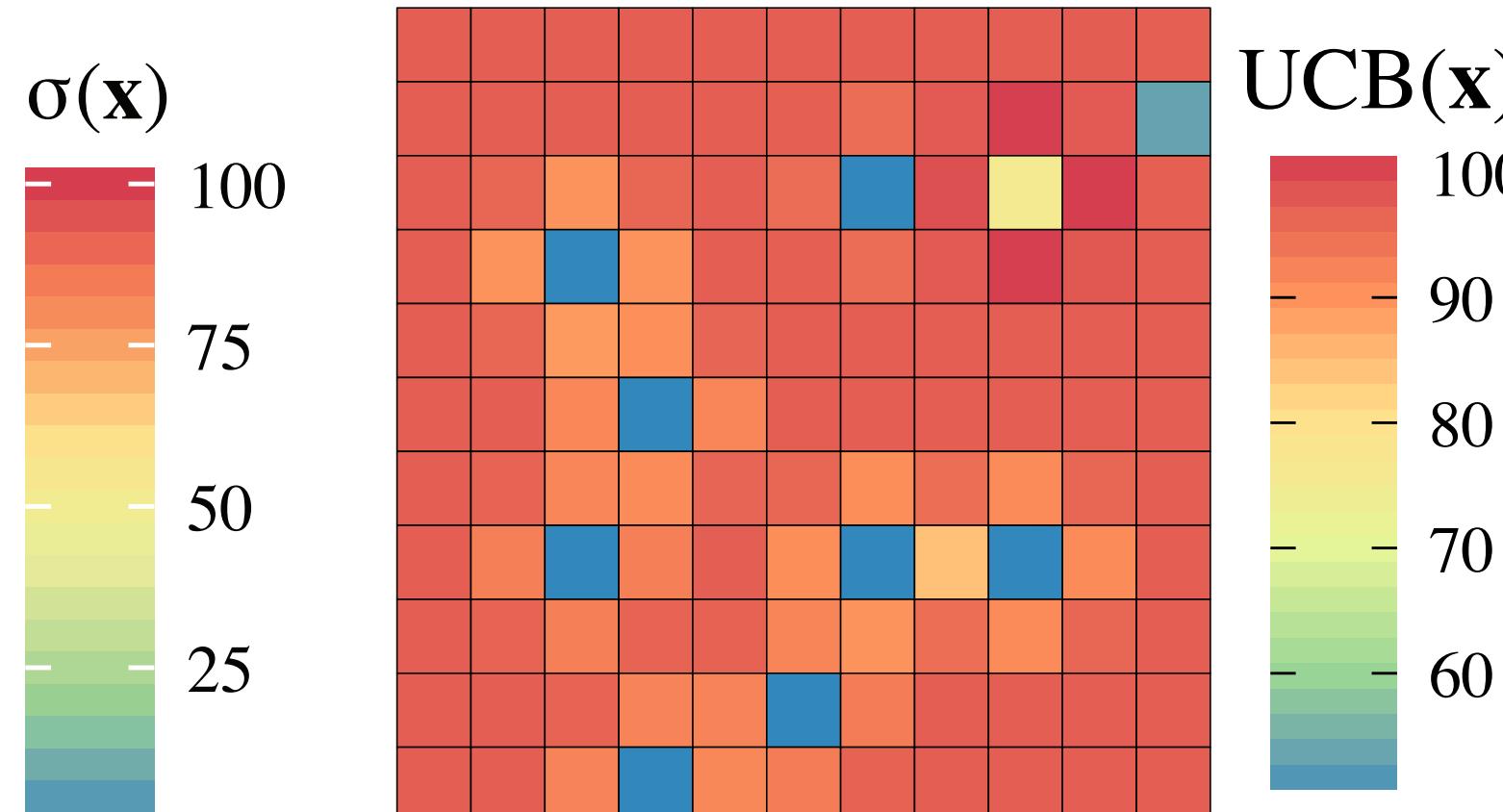
Observations



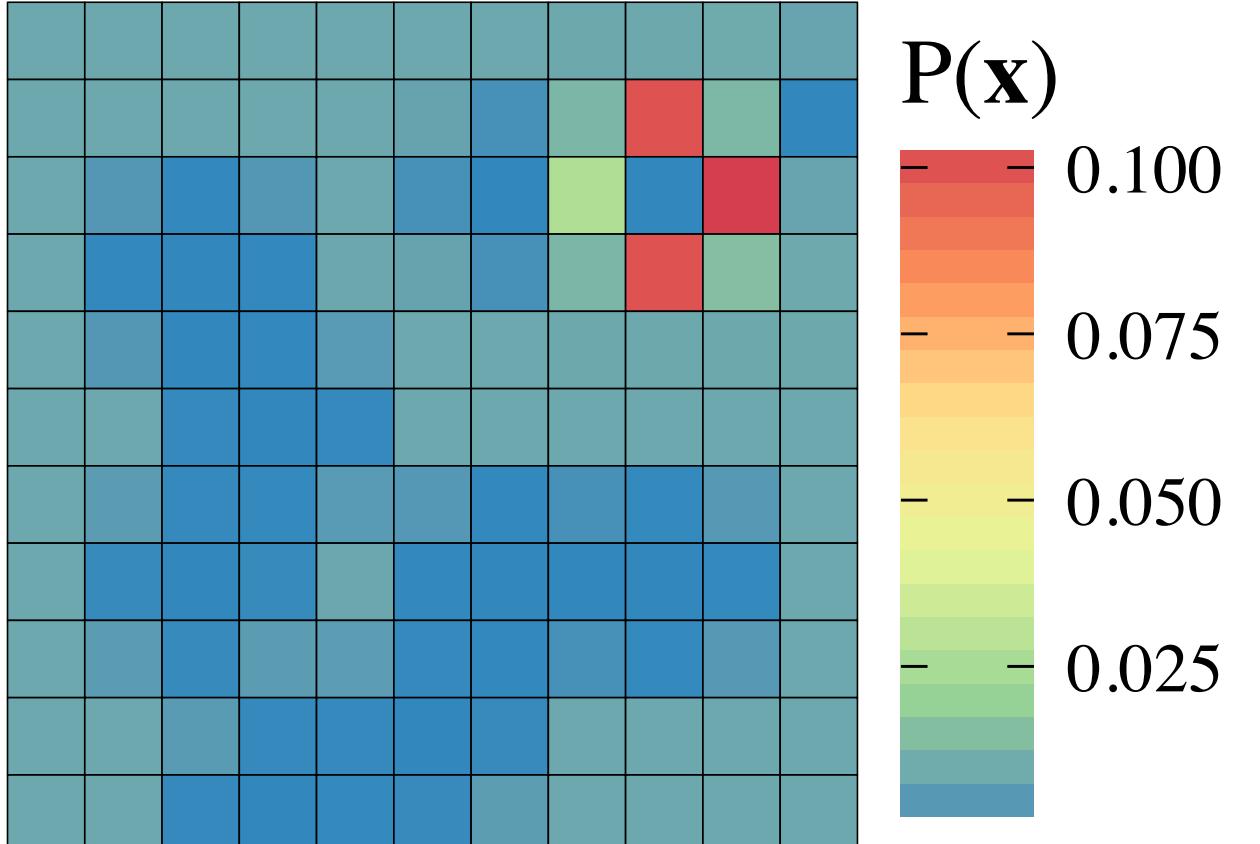
Gaussian Process (GP)



Upper Confidence Bound  
(UCB) Sampling



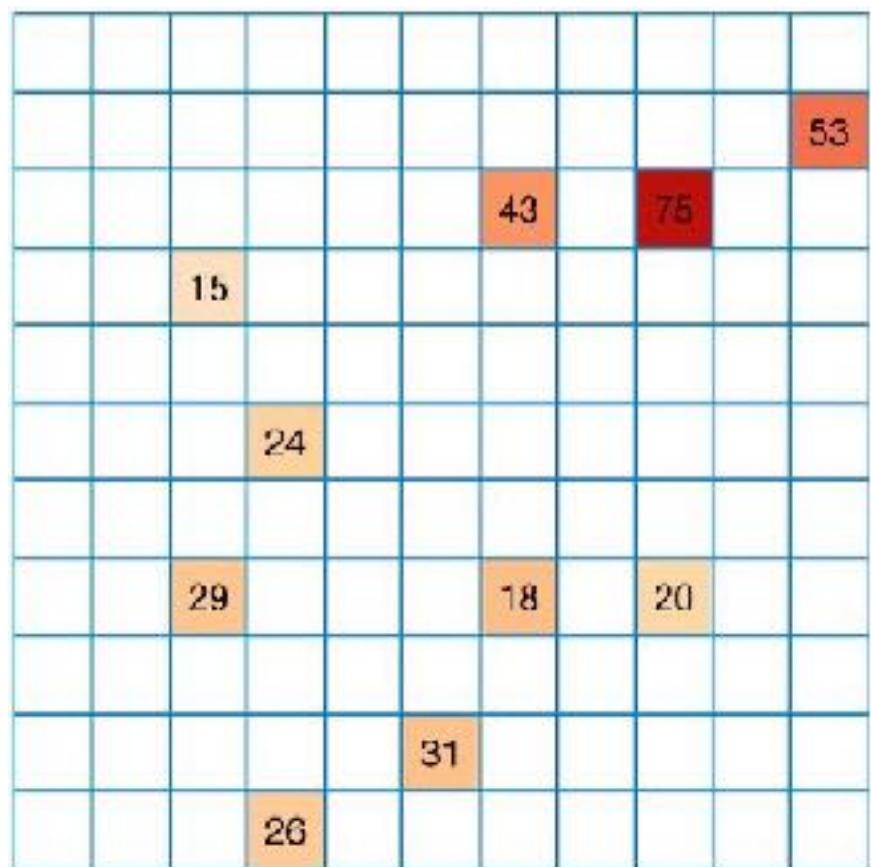
Softmax Choice Rule



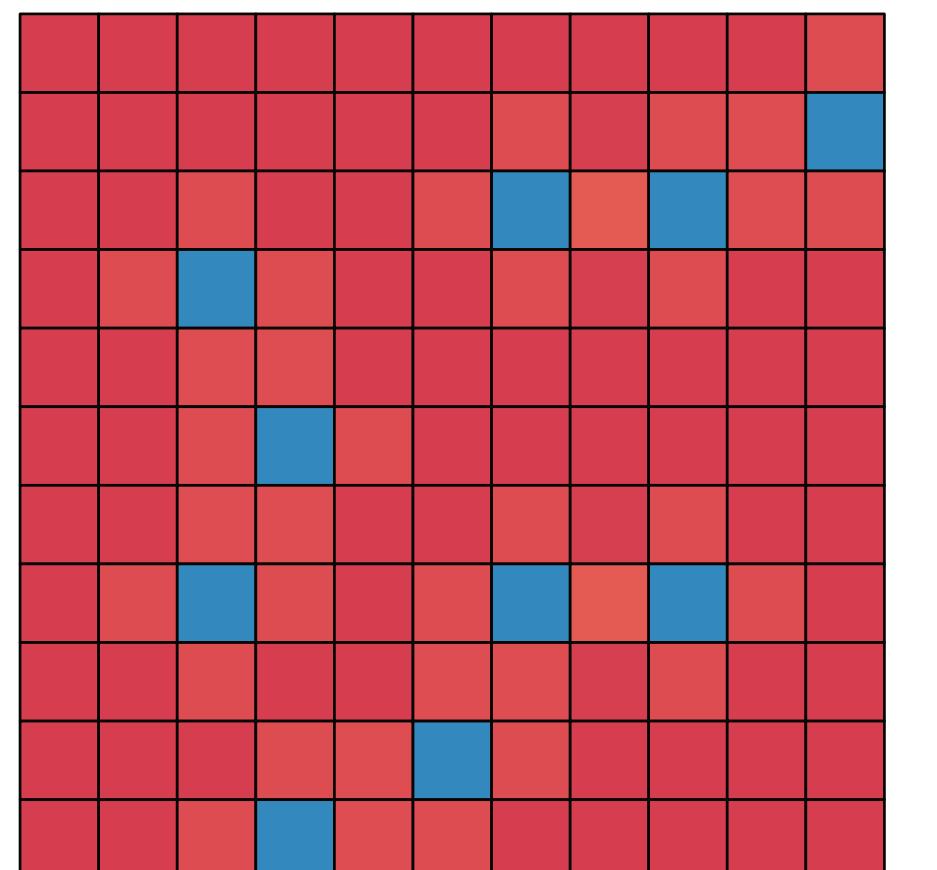
$$k_{RBF}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{||\mathbf{x} - \mathbf{x}'||^2}{2\lambda^2}\right)$$

# GP-UCB Model

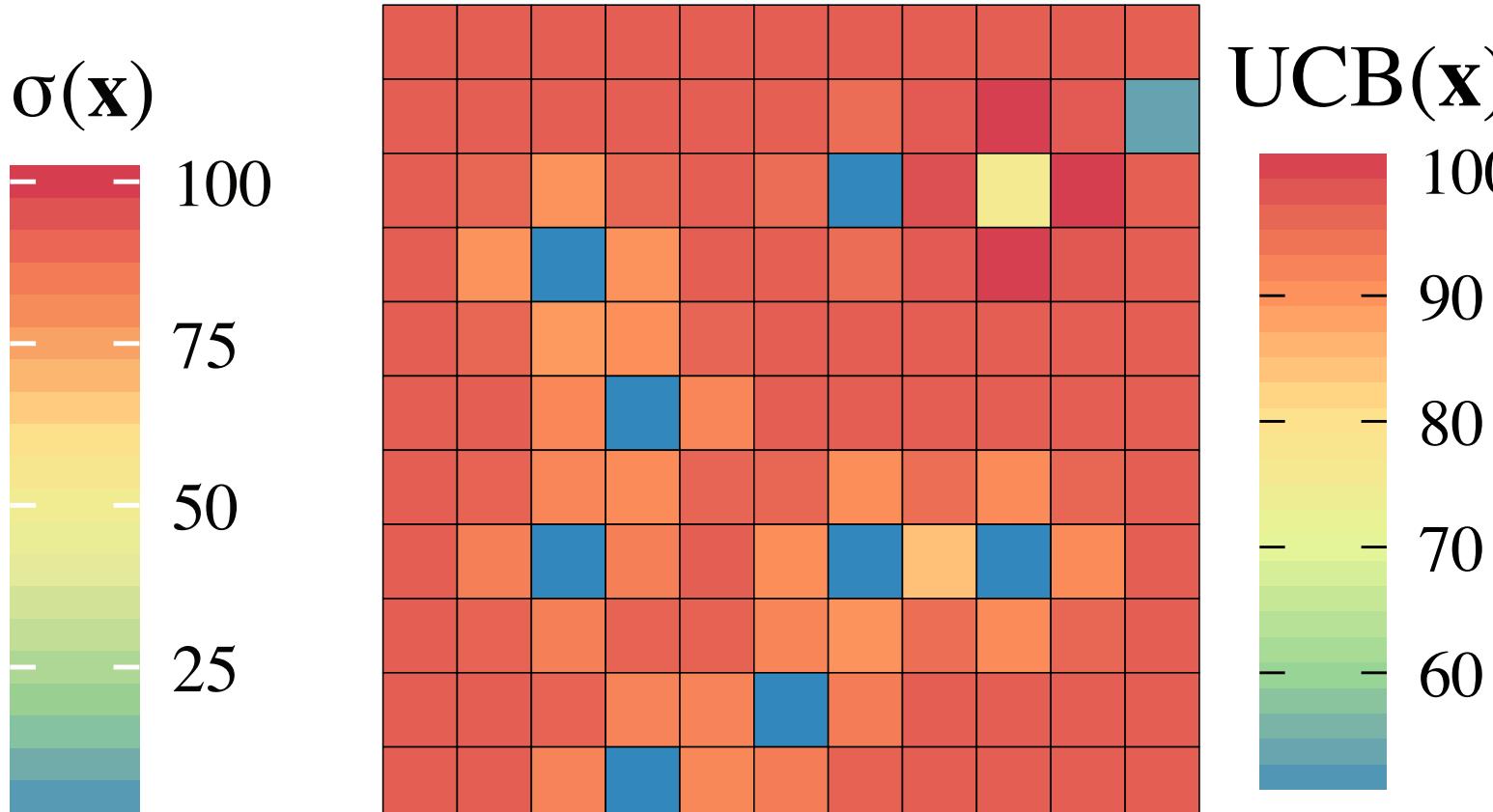
Observations



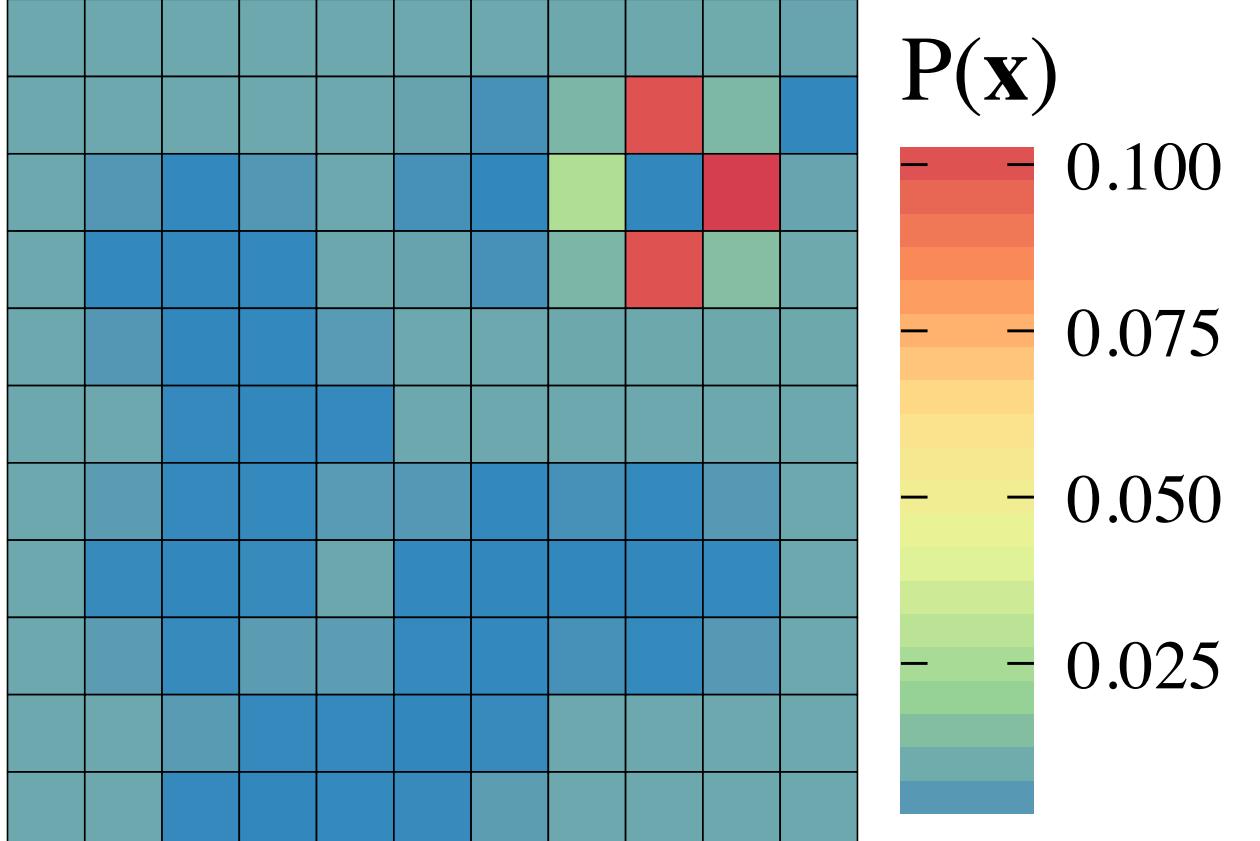
Gaussian Process (GP)



Upper Confidence Bound  
(UCB) Sampling

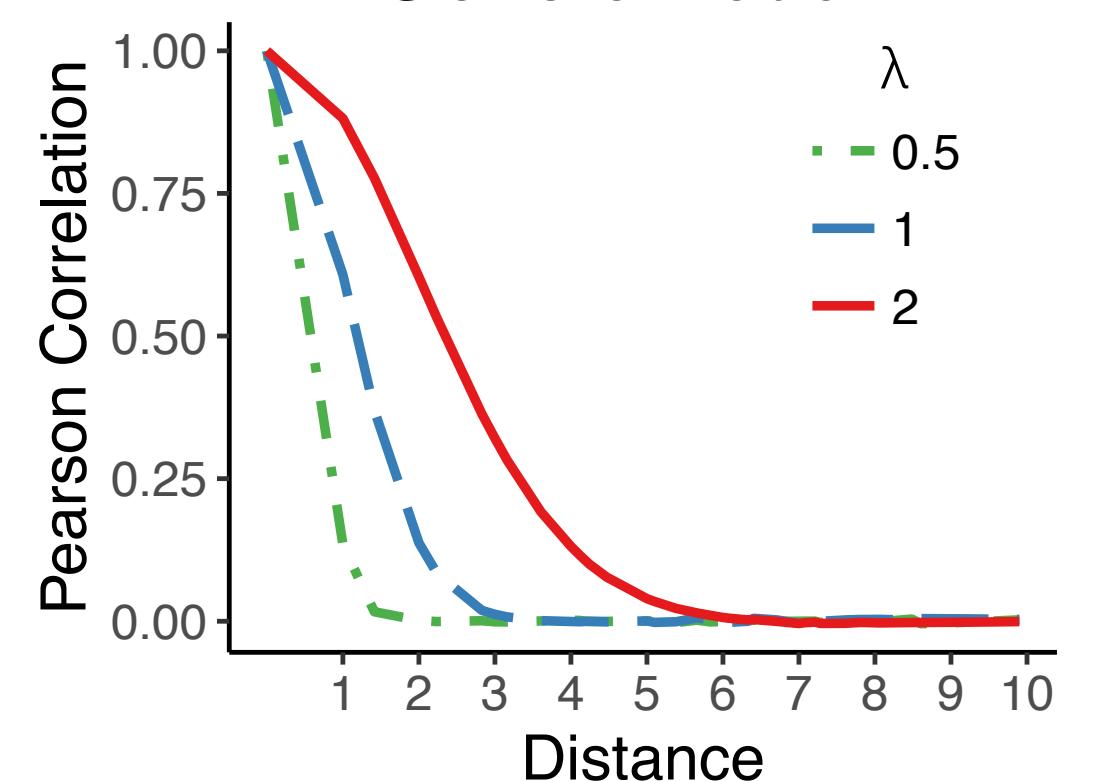


Softmax Choice Rule



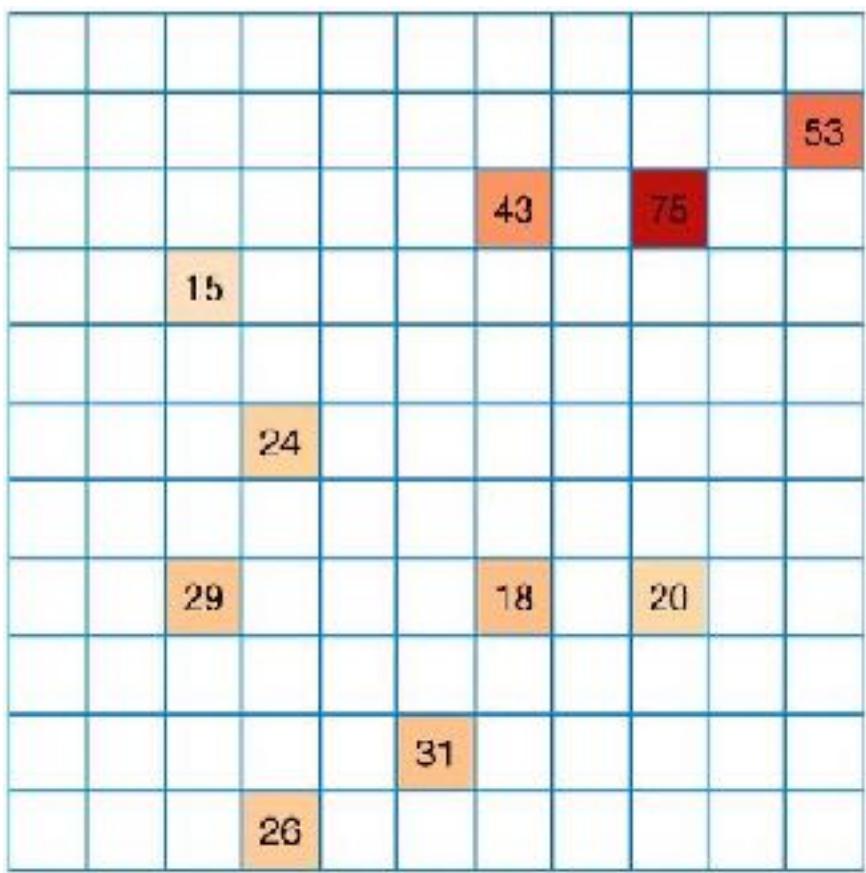
$$k_{RBF}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{||\mathbf{x} - \mathbf{x}'||^2}{2\lambda^2}\right)$$

Generalization  $\lambda$

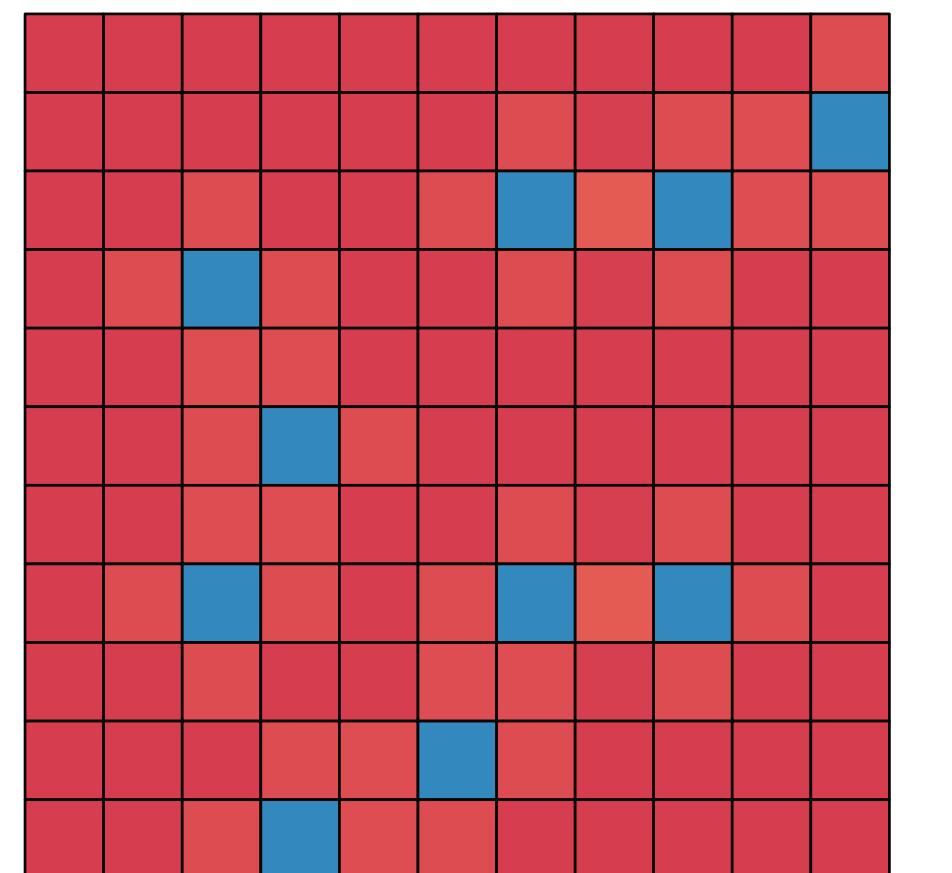


# GP-UCB Model

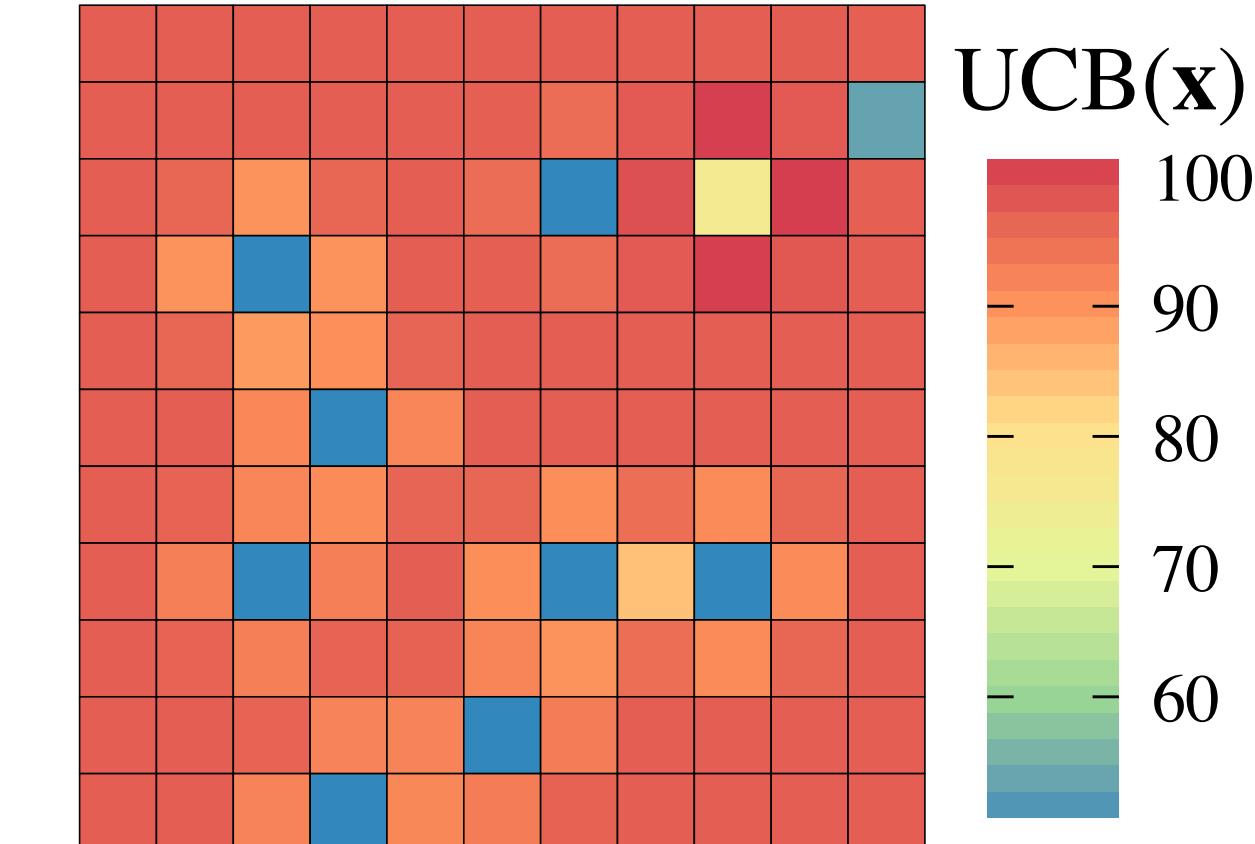
Observations



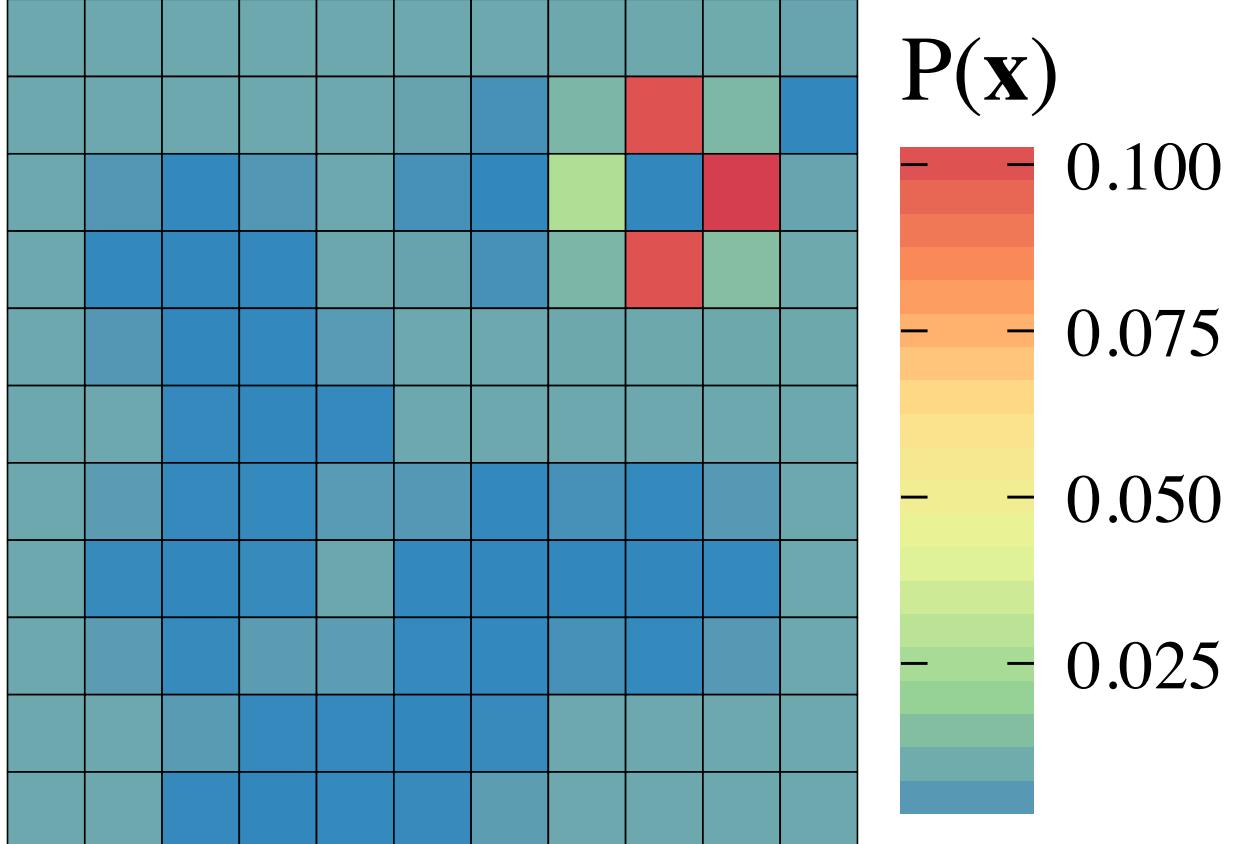
Gaussian Process (GP)



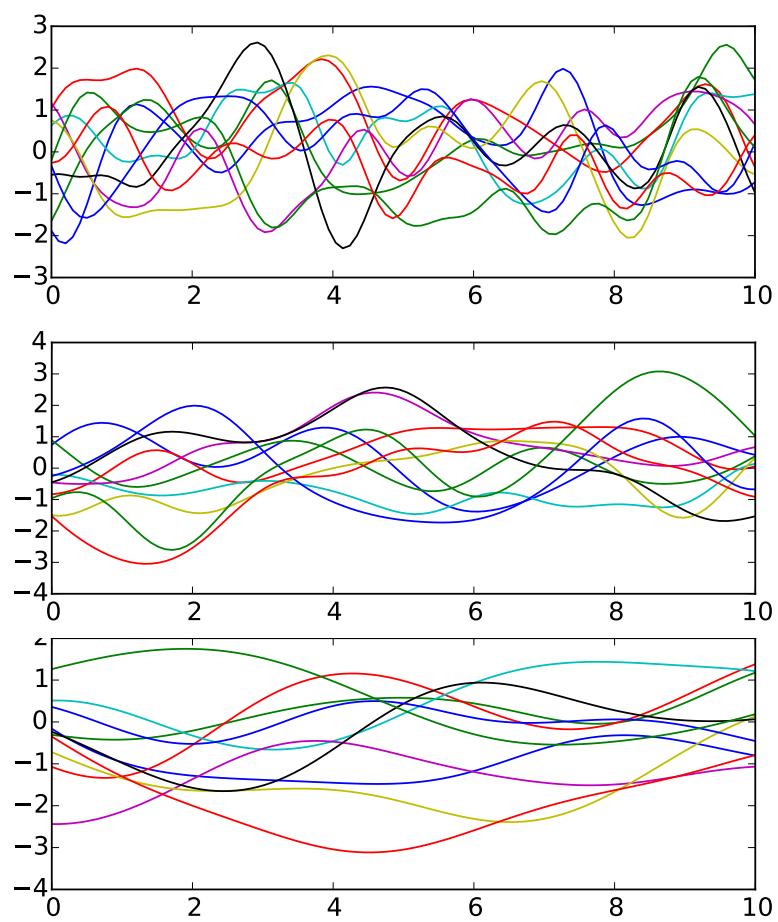
Upper Confidence Bound (UCB) Sampling



Softmax Choice Rule



$\lambda=.5$

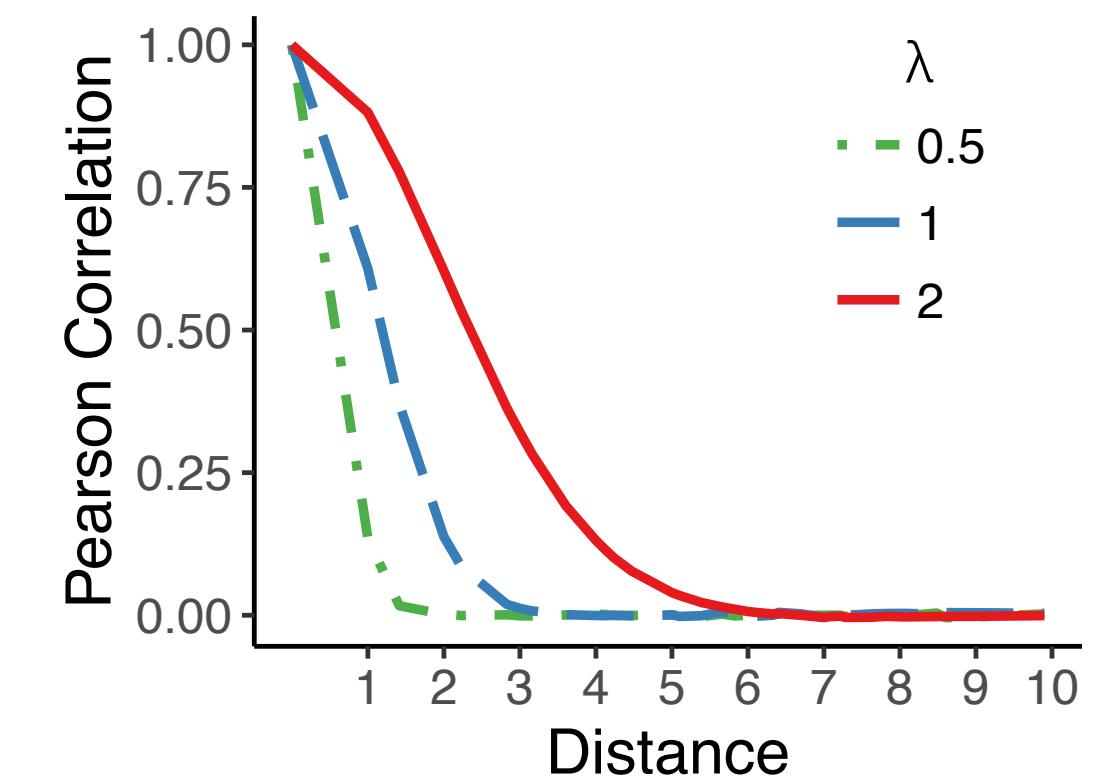


$\lambda=1$

$\lambda=2$

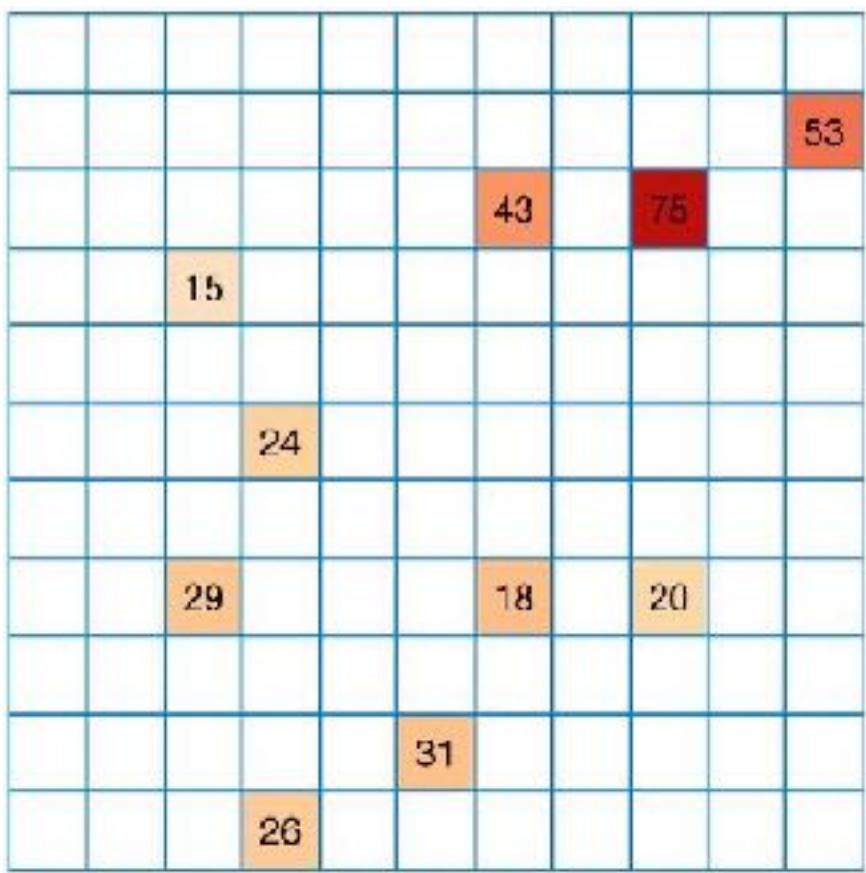
$$k_{RBF}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{||\mathbf{x} - \mathbf{x}'||^2}{2\lambda^2}\right)$$

Generalization  $\lambda$

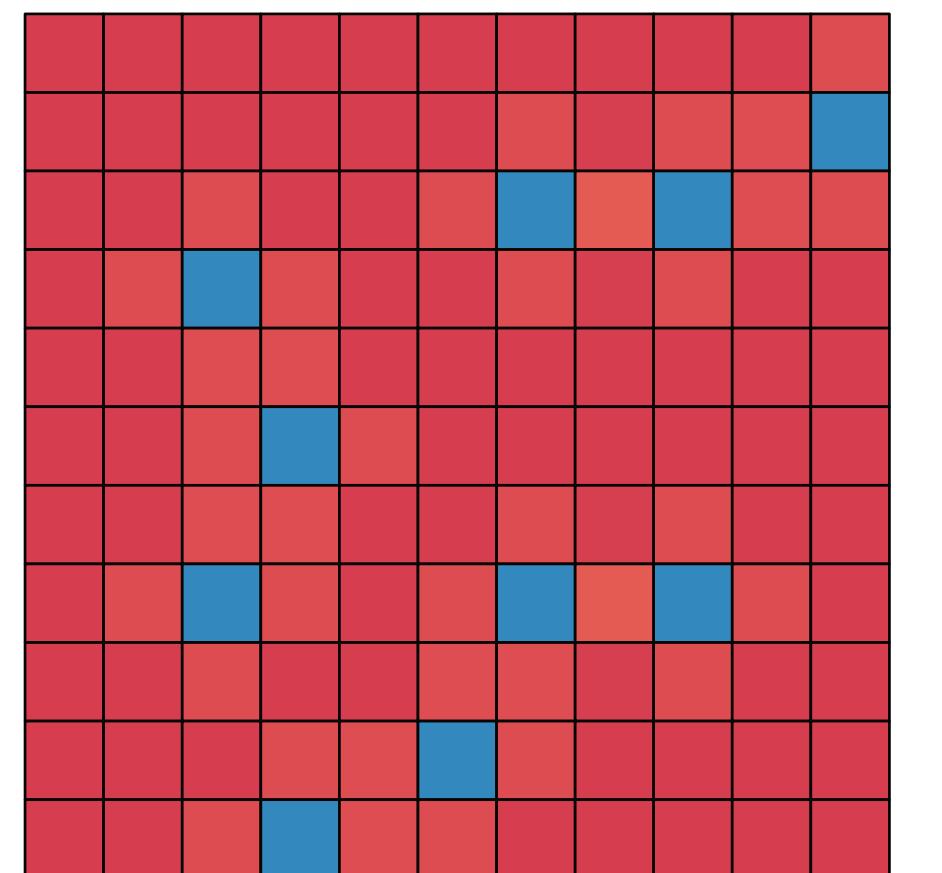


# GP-UCB Model

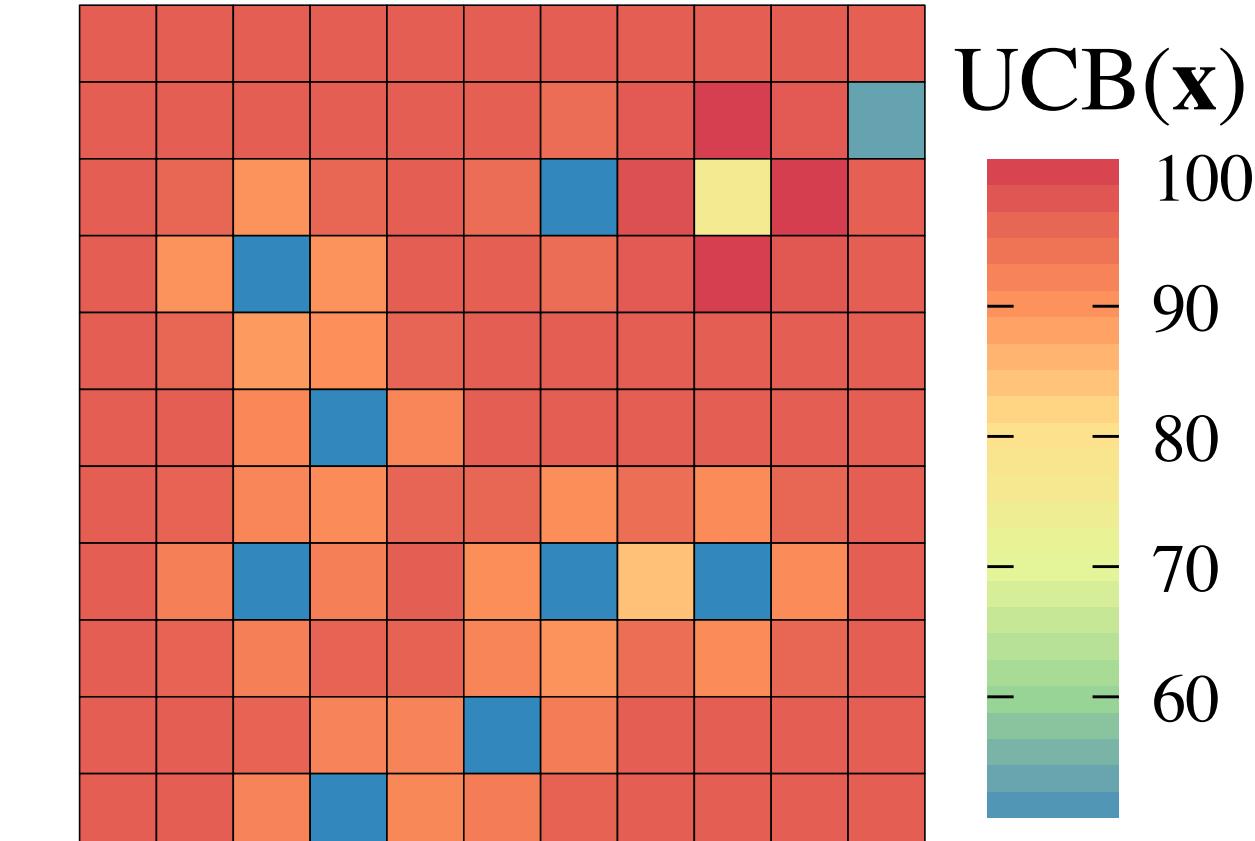
Observations



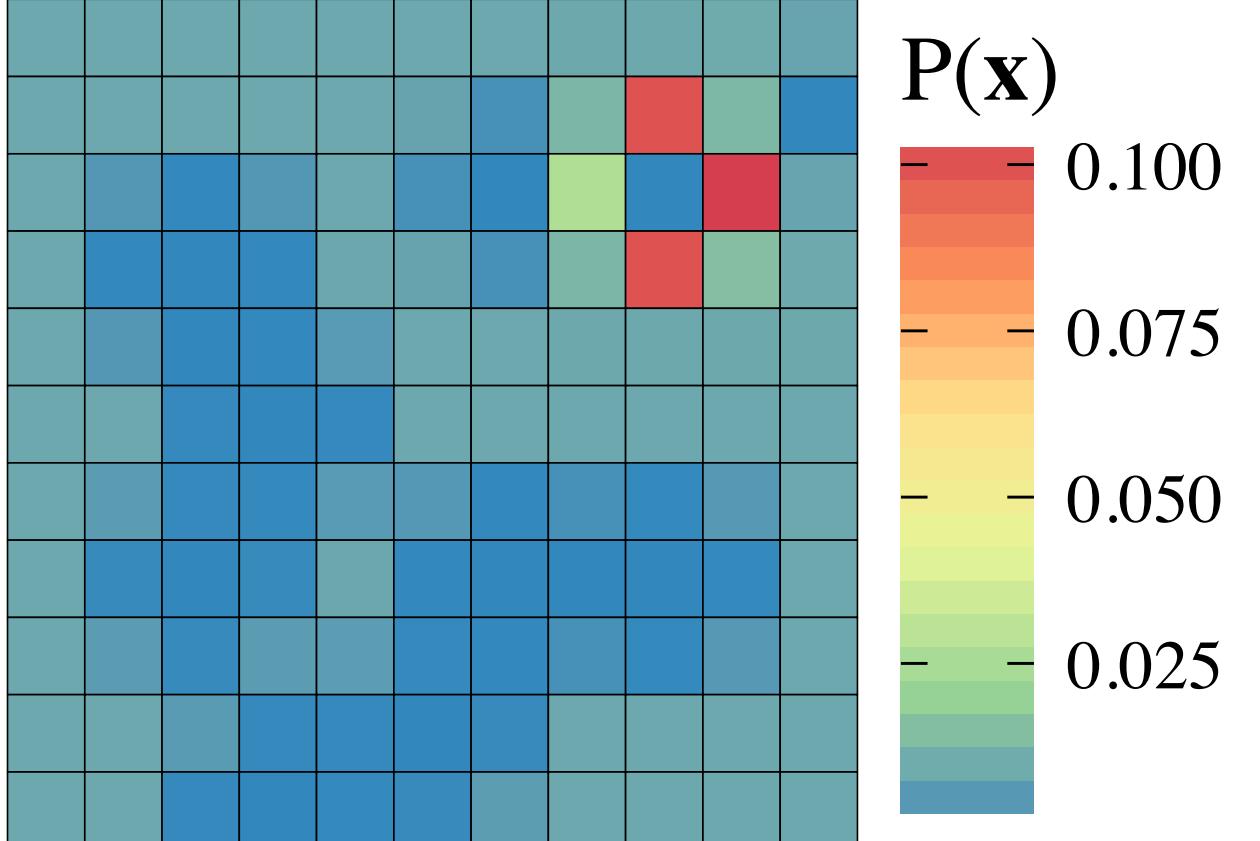
Gaussian Process (GP)



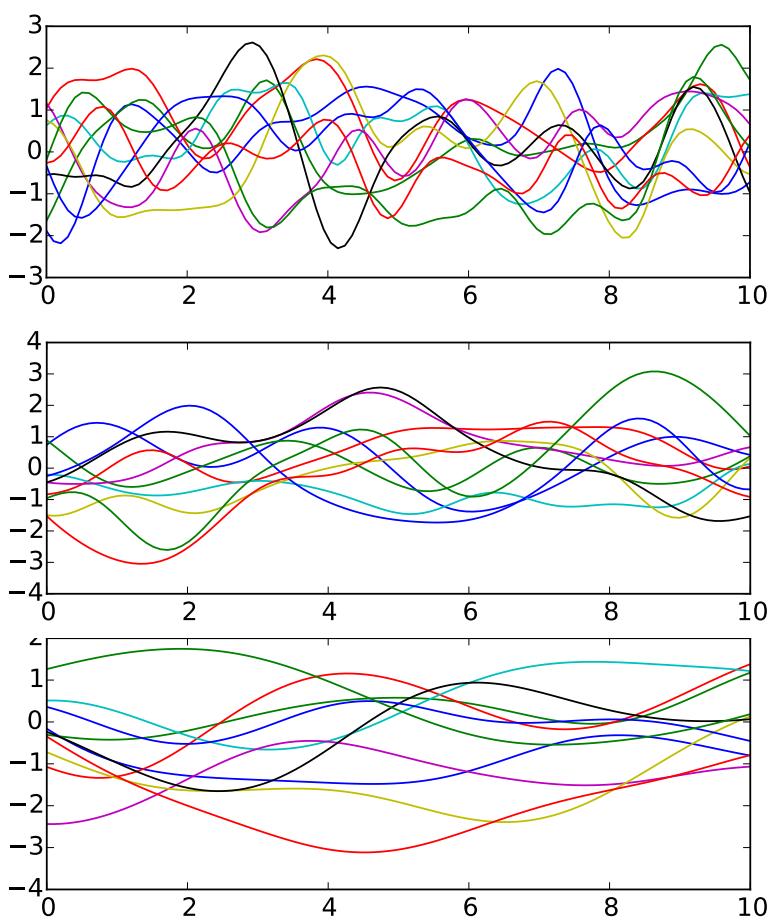
Upper Confidence Bound (UCB) Sampling



Softmax Choice Rule



$\lambda=.5$

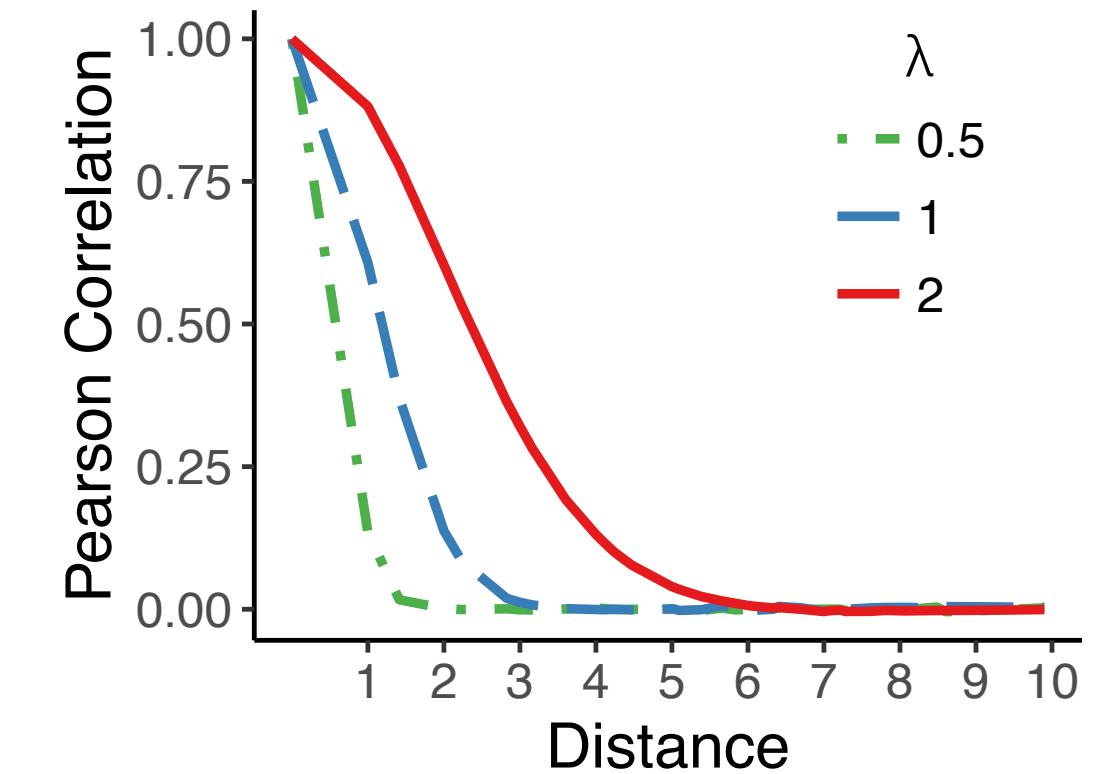


$\lambda=1$

$\lambda=2$

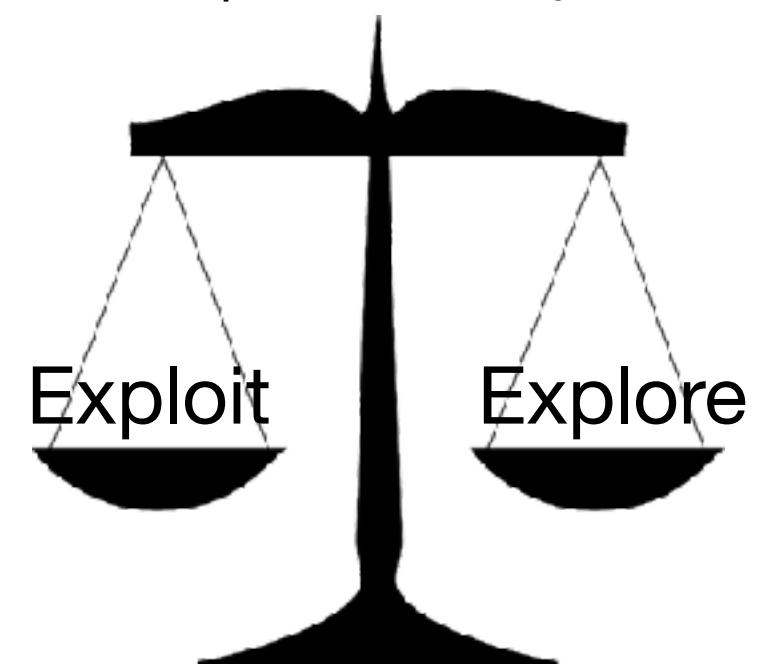
$$k_{RBF}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{||\mathbf{x} - \mathbf{x}'||^2}{2\lambda^2}\right)$$

Generalization  $\lambda$



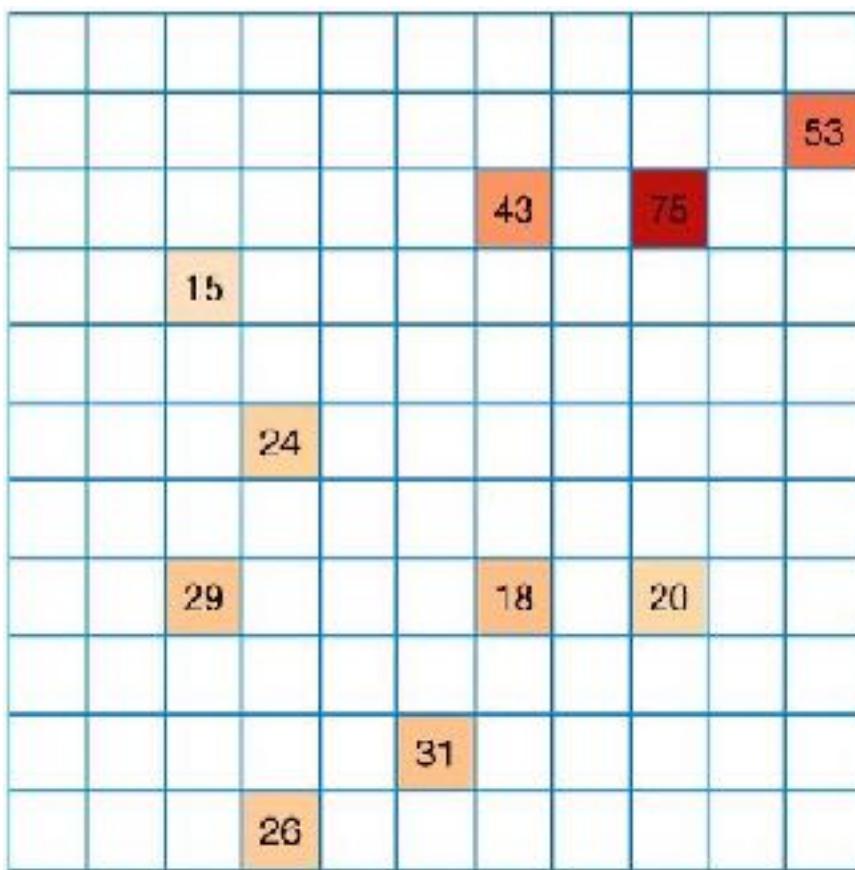
$$UCB(\mathbf{x}_i) = \mu(\mathbf{x}_i) + \beta\sigma(\mathbf{x}_i)$$

Exploration  $\beta$

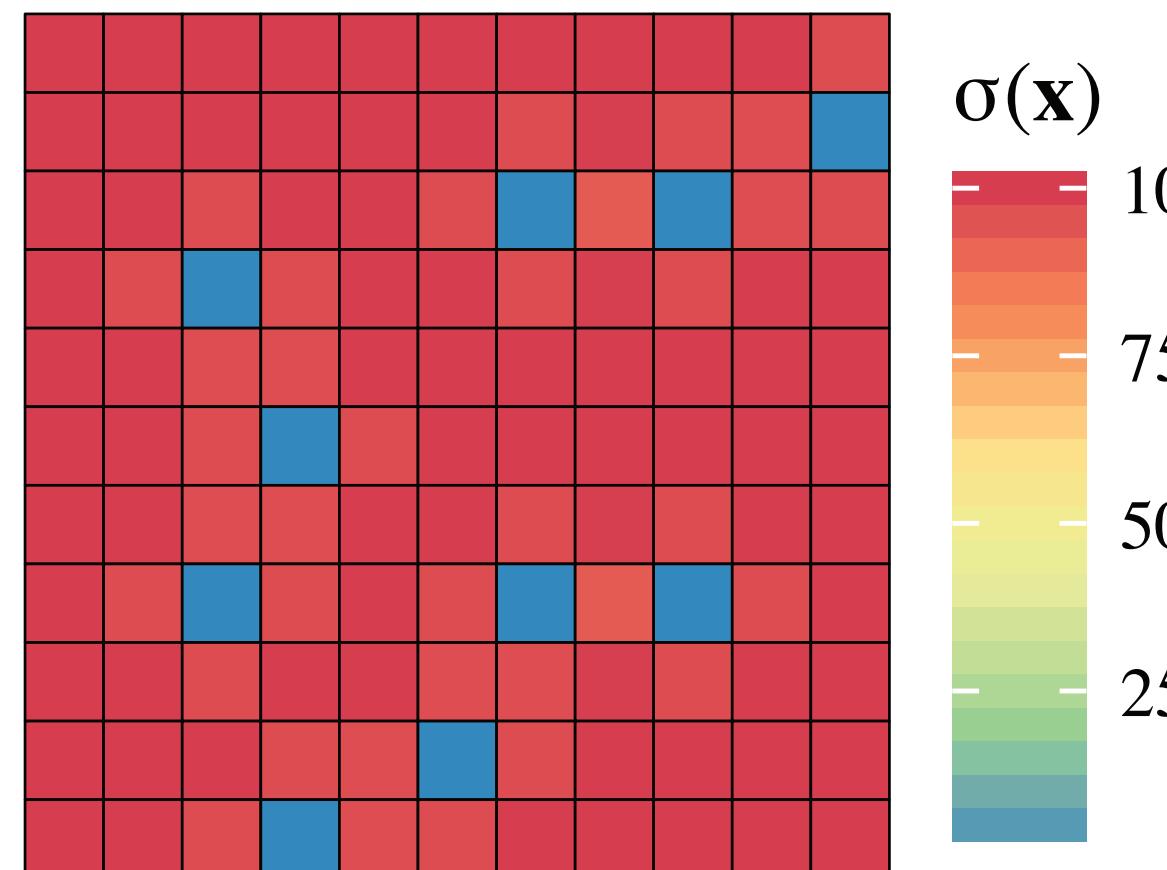


# GP-UCB Model

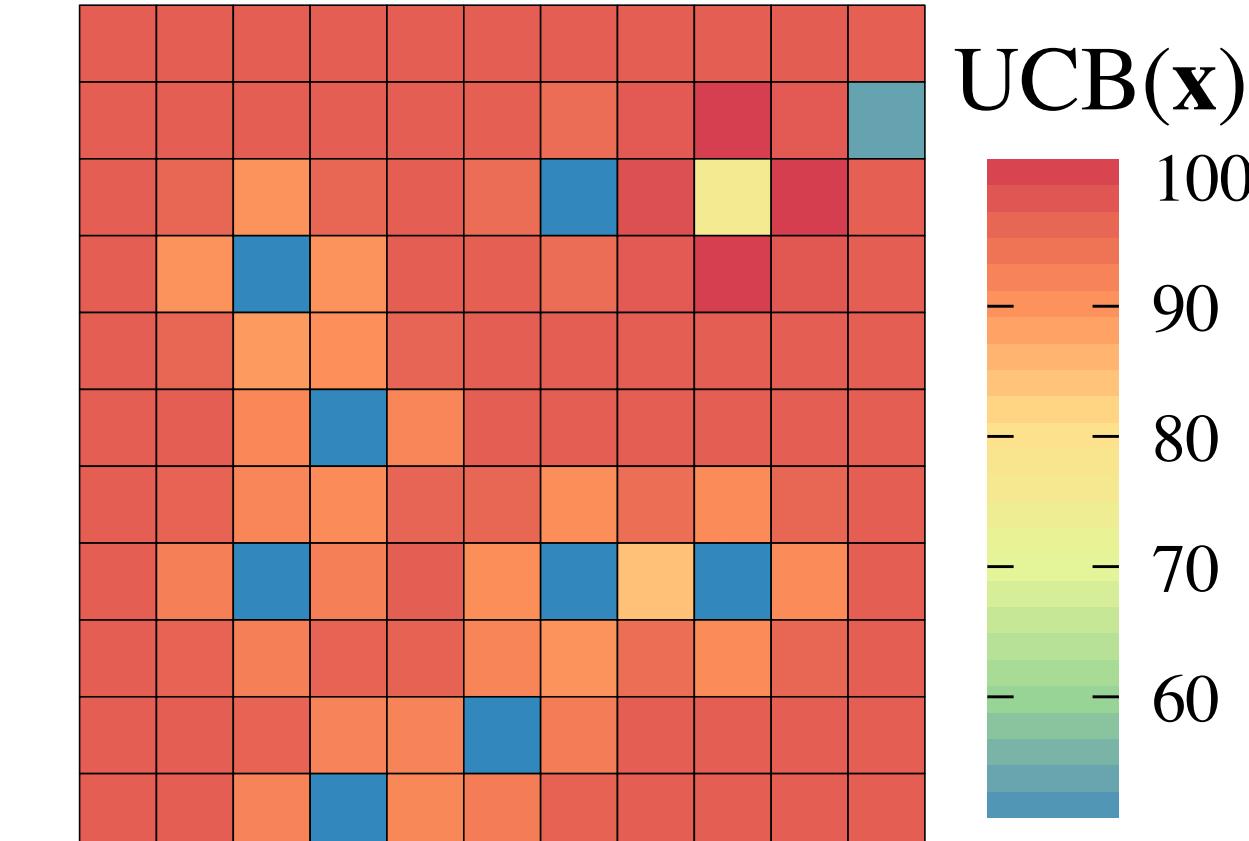
Observations



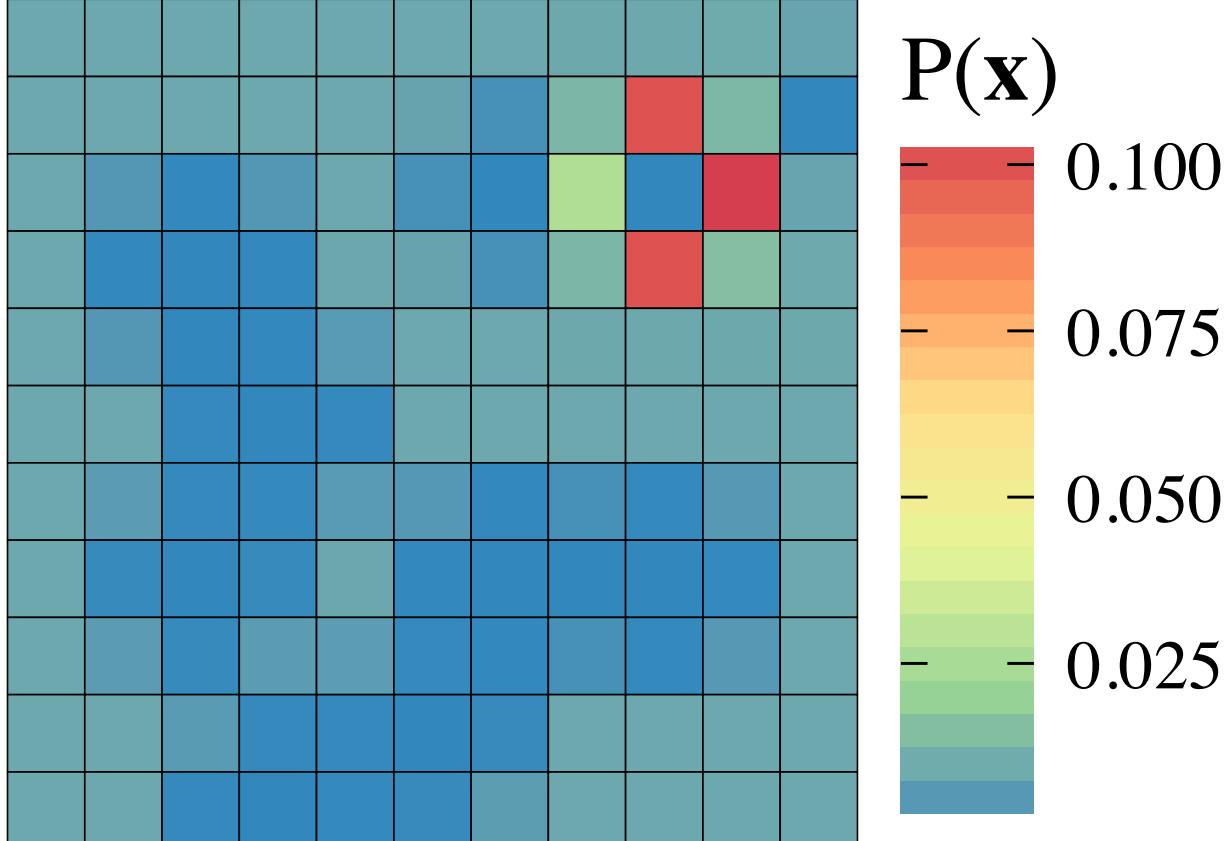
Gaussian Process (GP)



Upper Confidence Bound (UCB) Sampling



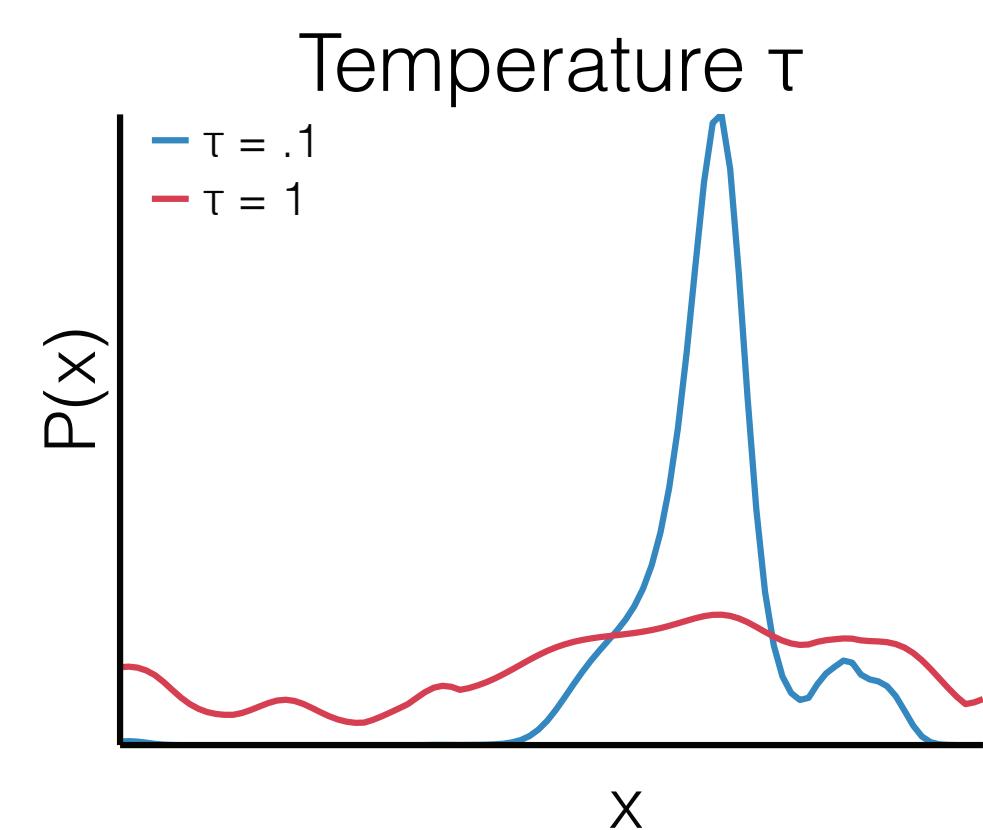
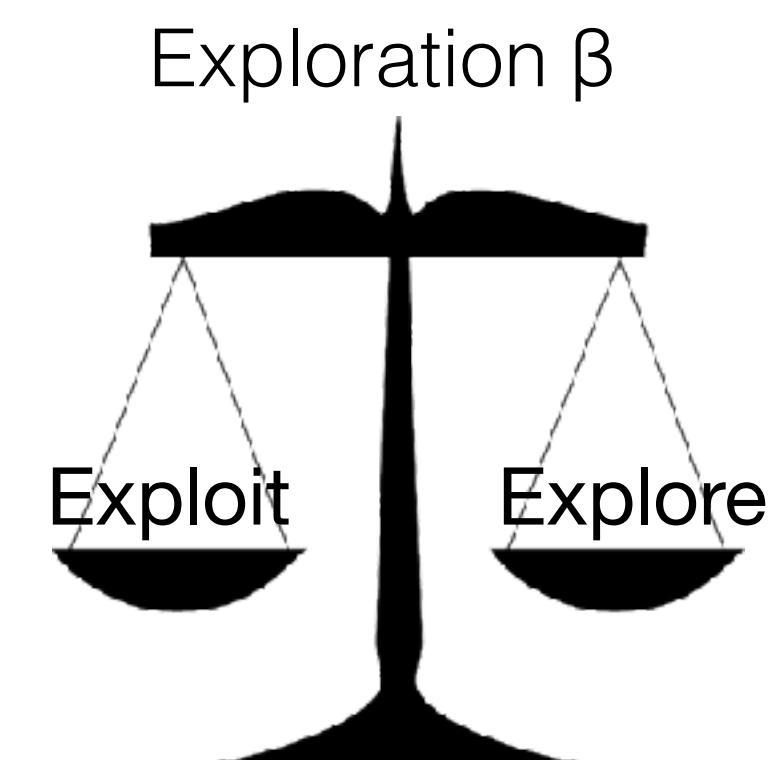
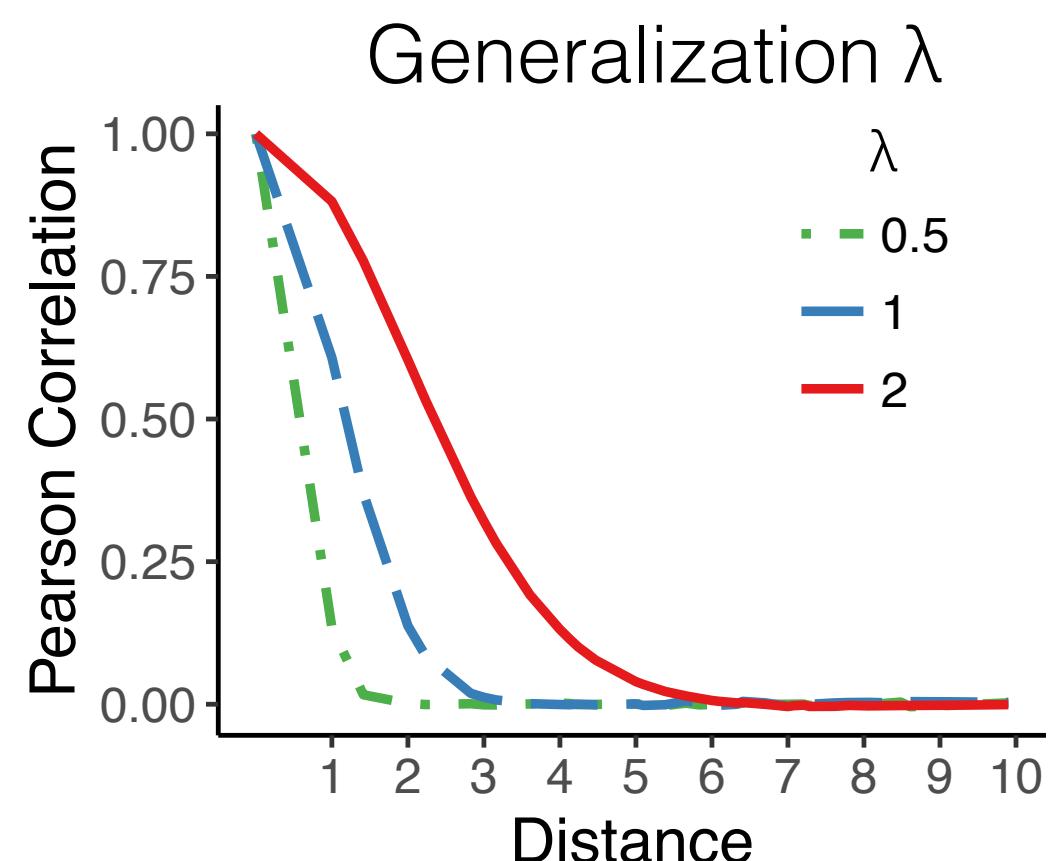
Softmax Choice Rule



$$k_{RBF}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{||\mathbf{x} - \mathbf{x}'||^2}{2\lambda^2}\right)$$

$$UCB(\mathbf{x}_i) = \mu(\mathbf{x}_i) + \beta\sigma(\mathbf{x}_i)$$

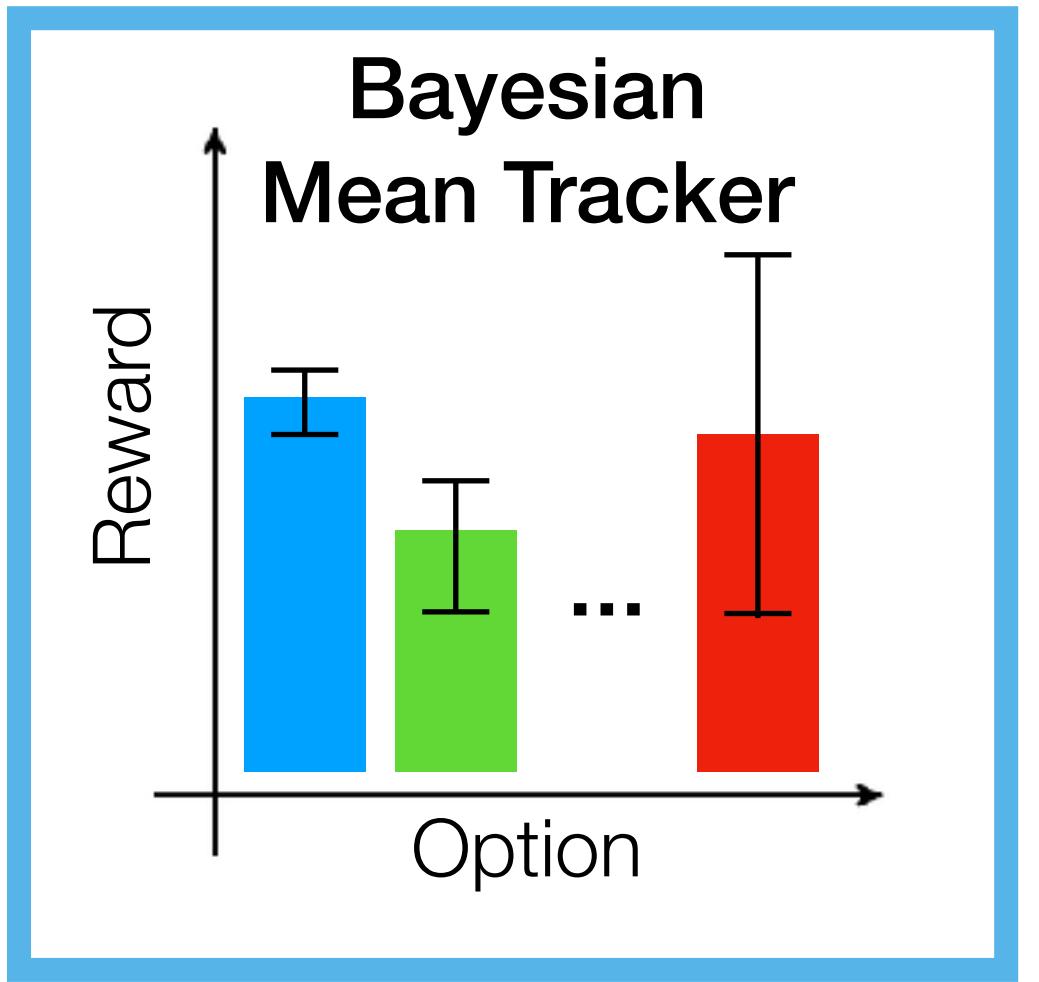
$$P(\mathbf{x}_i) \propto \exp(UCB(\mathbf{x}_i)/\tau)$$



# Model Comparison

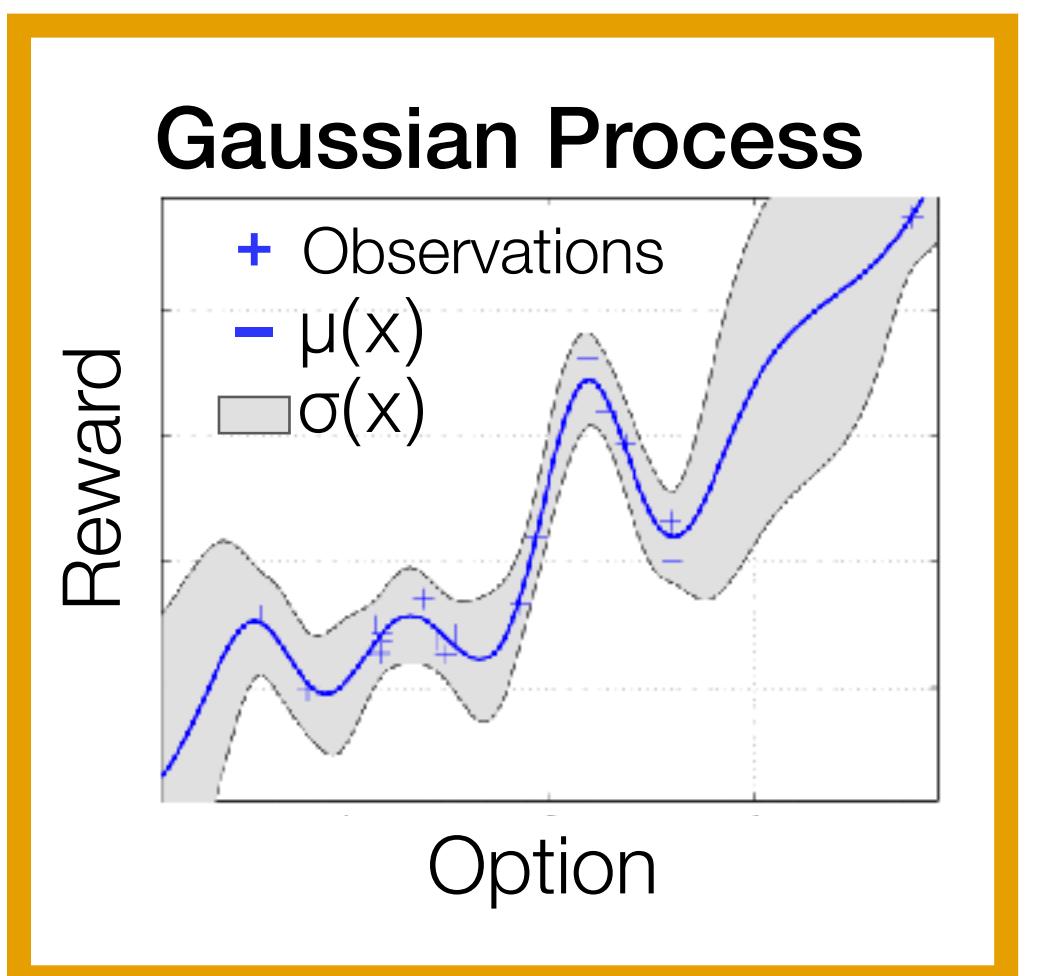
- **Traditional RL model:**

- Learns the value of each option independently
  - e.g., Rescorla-Wagner, Q-learning, Kalman Filter, *Bayesian Mean Tracker* (BMT), etc...
- Can balance explore-exploit dilemma using a variety of sampling strategies, but offers limited guidance about *where to explore*

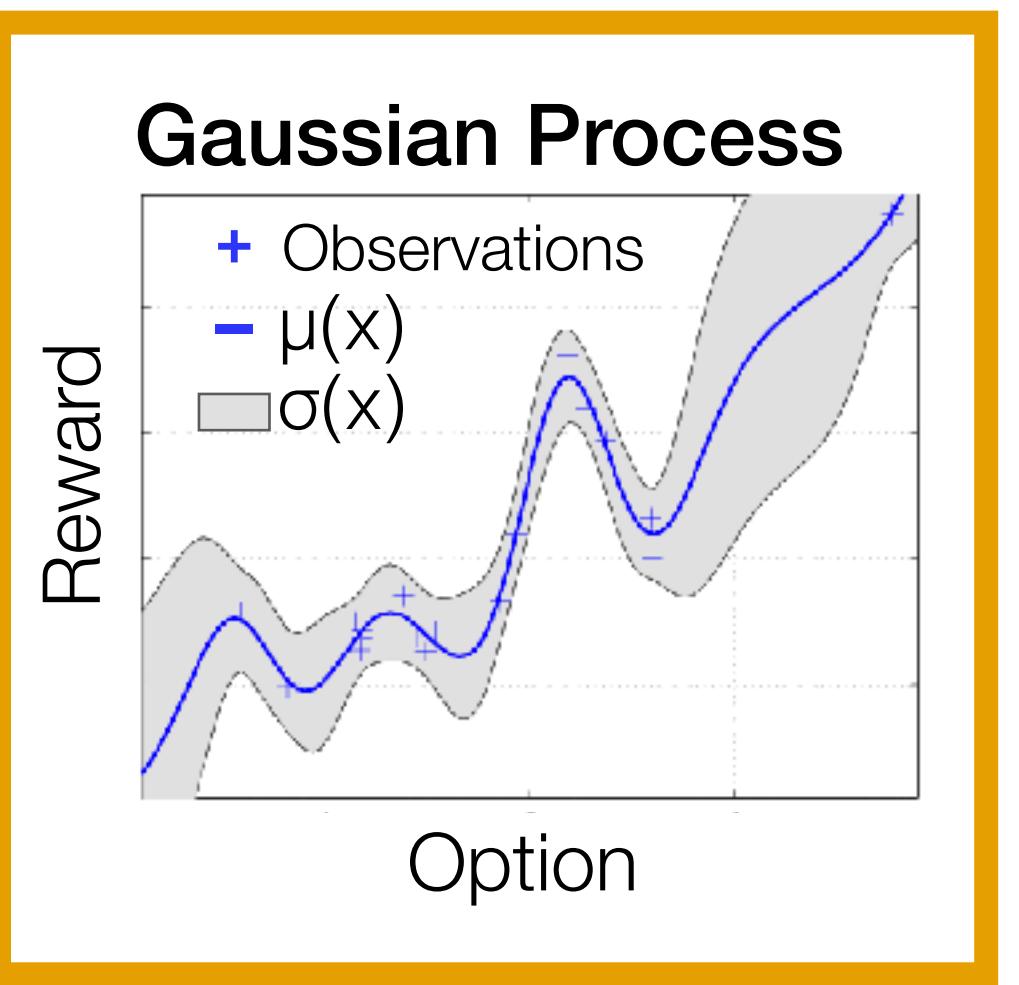
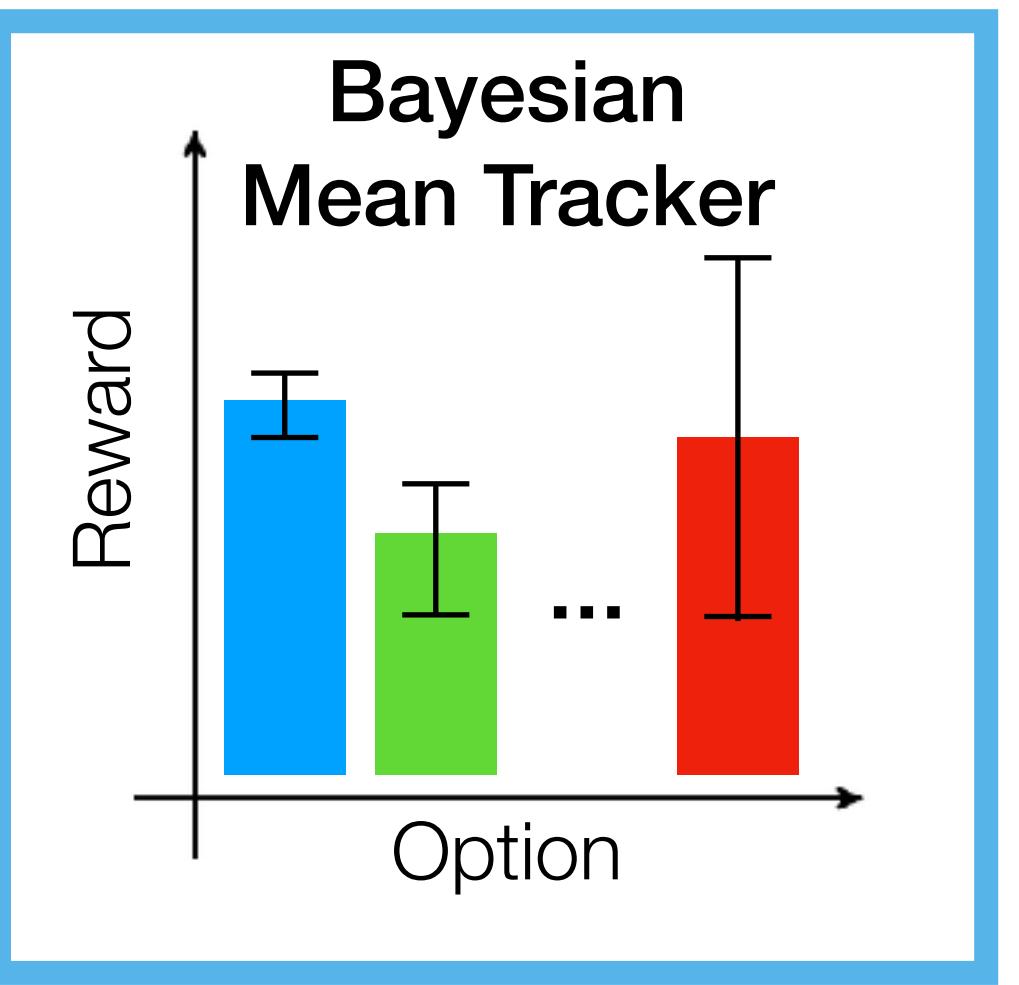


- **Function learning model:**

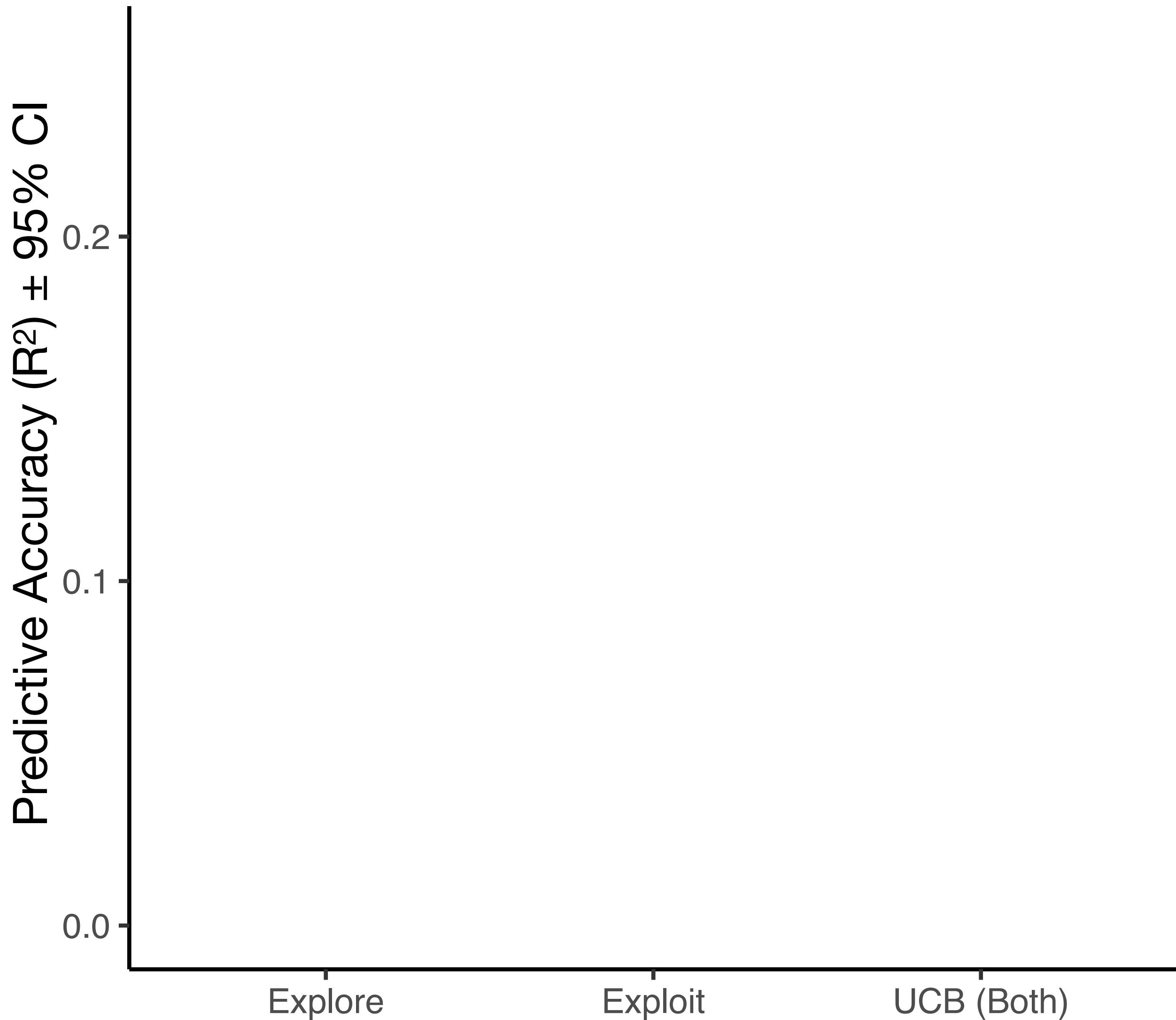
- Uses function approximation to generalize about novel option
  - e.g., Neural Network function approximators, *Gaussian Process* (GP) model, etc...
- Balances explore-exploit using the same sampling strategies as option learning models, but also makes predictions about *where to explore* through generalization



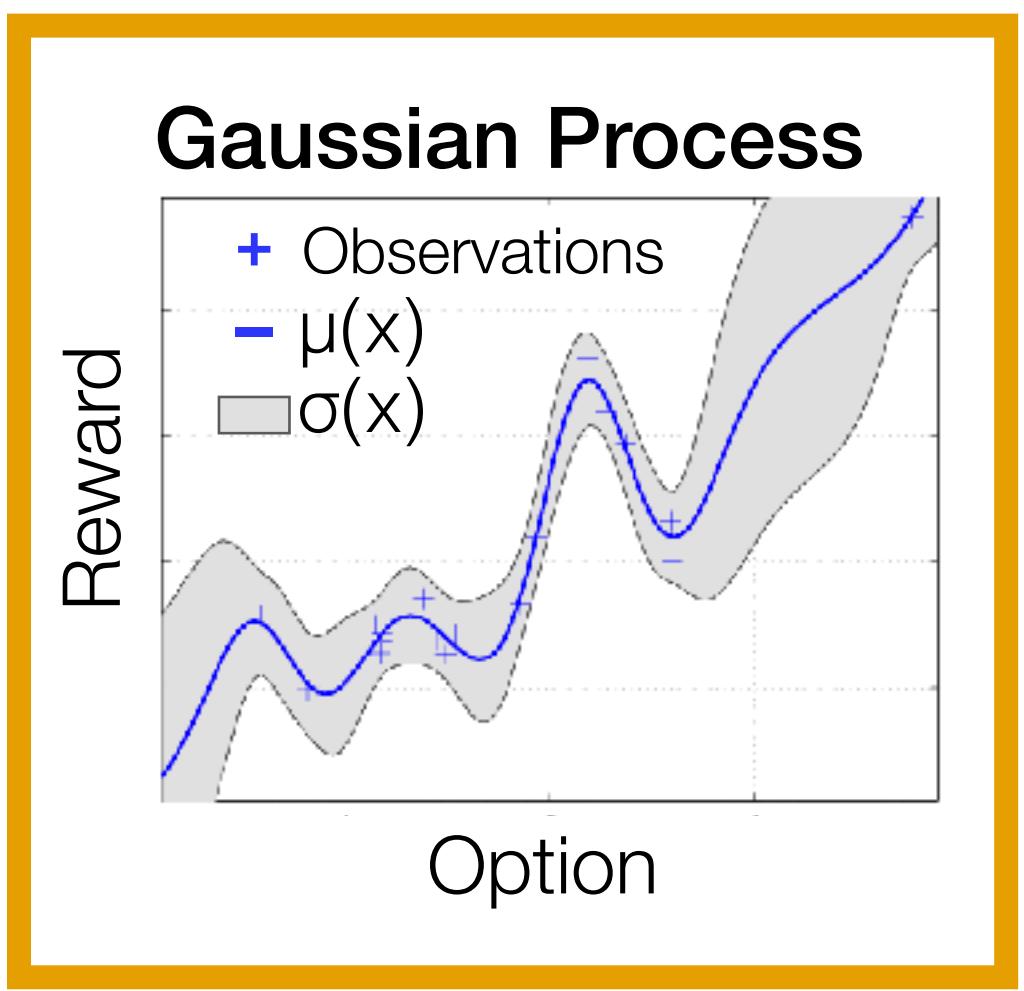
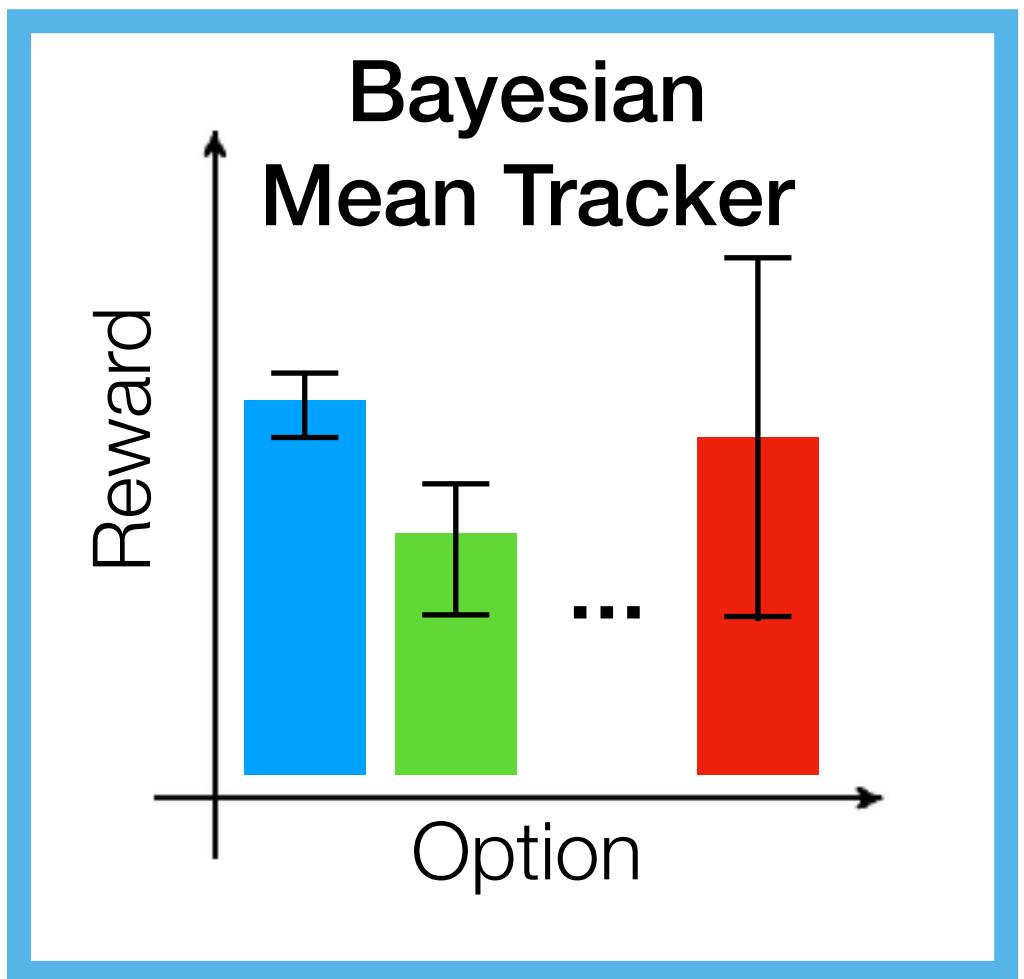
# Model Comparison



# Model Comparison

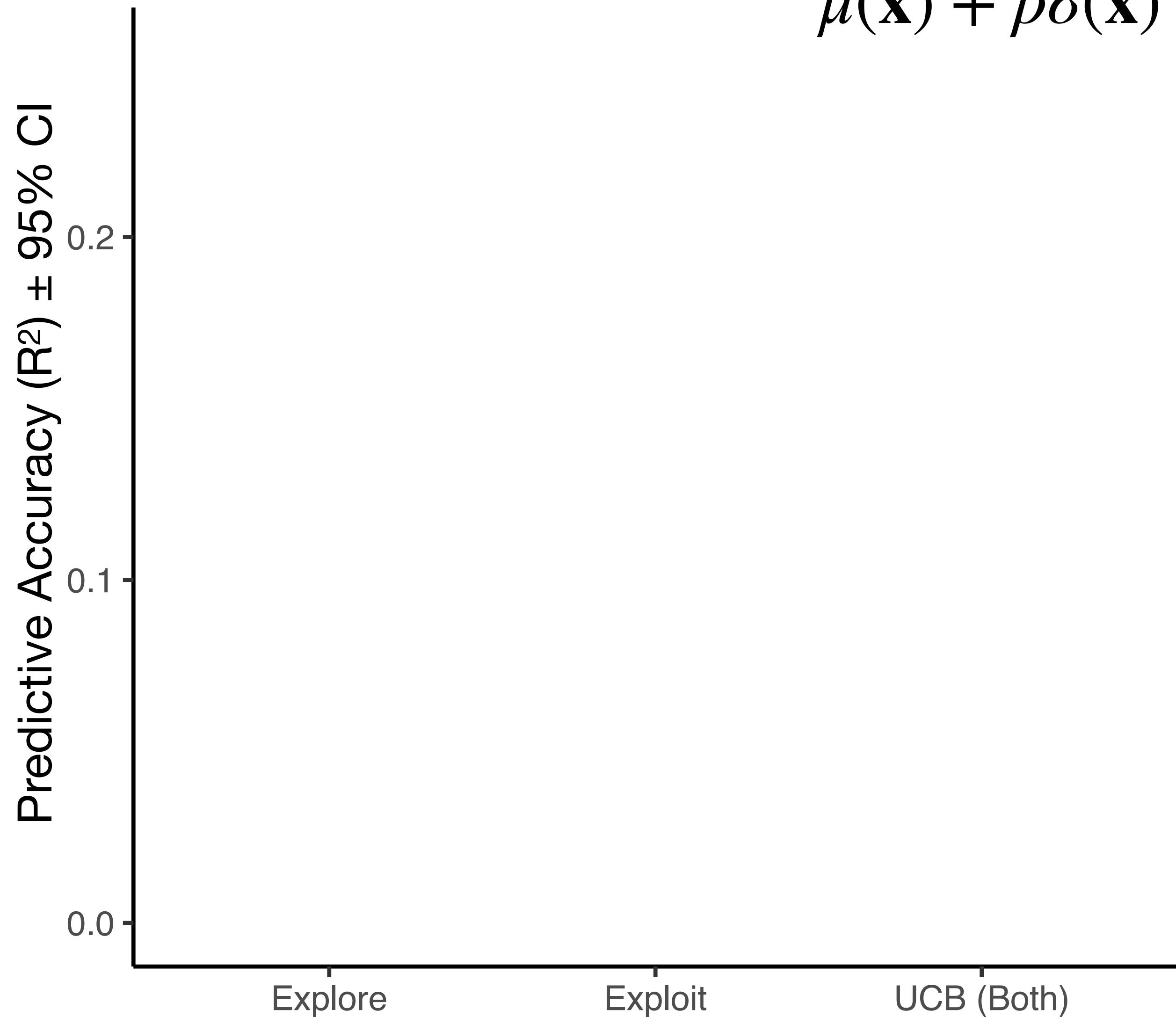


$$R^2 = 1 - \frac{\log \mathcal{L}(\mathcal{M}_k)}{\log \mathcal{L}(\mathcal{M}_{\text{rand}})}$$

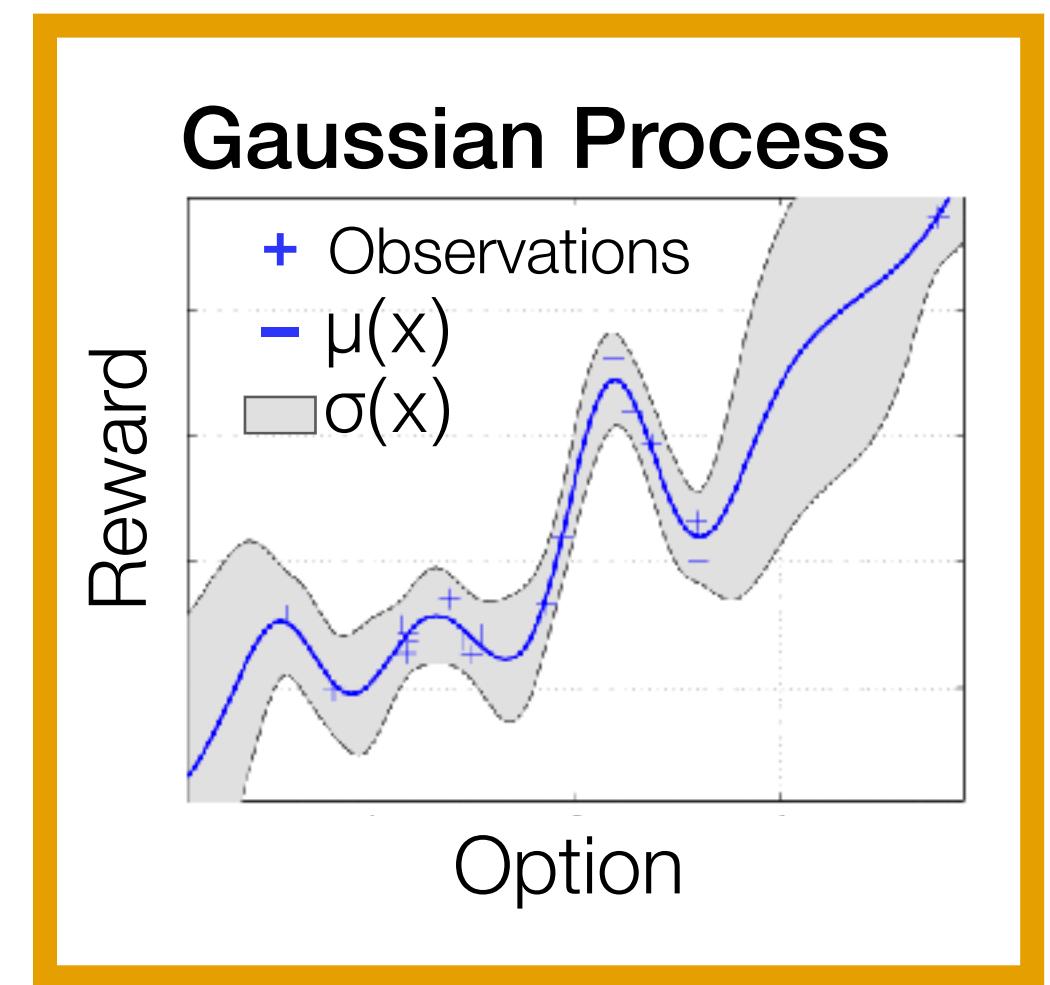
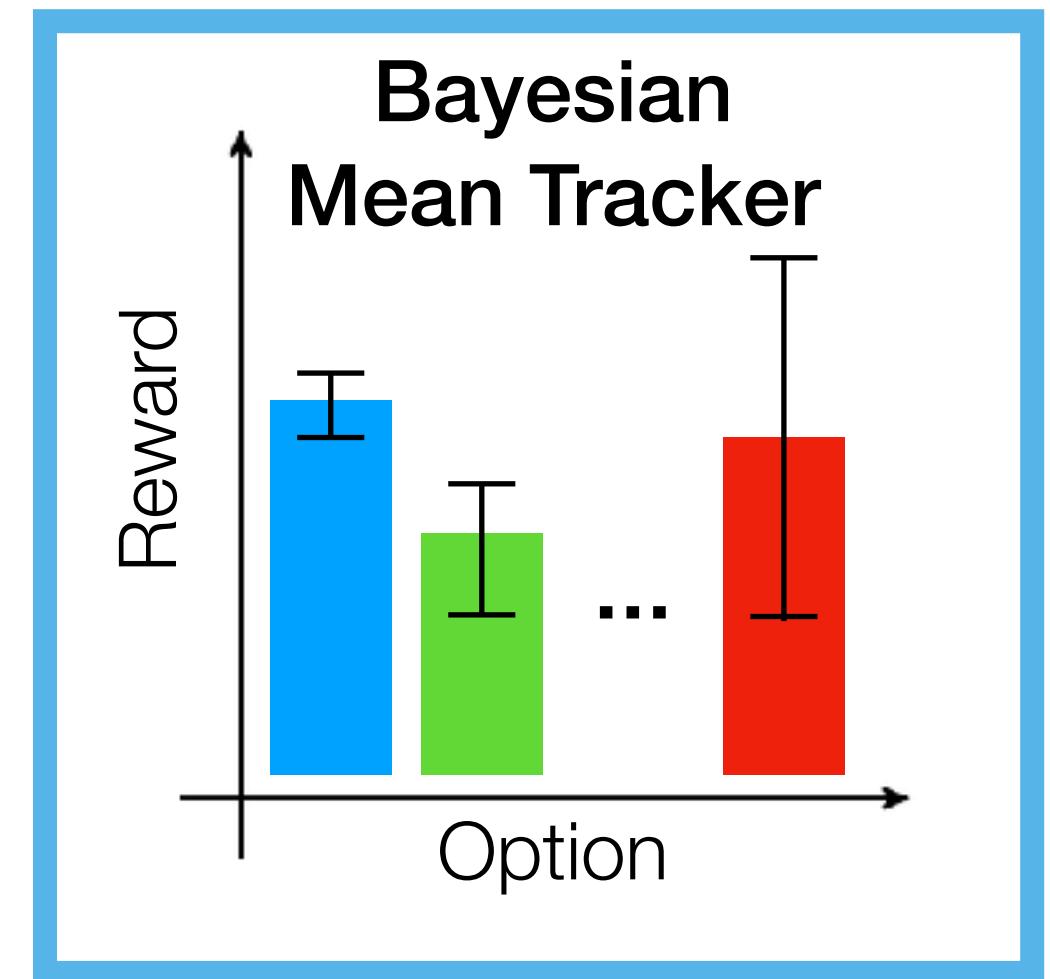


# Model Comparison

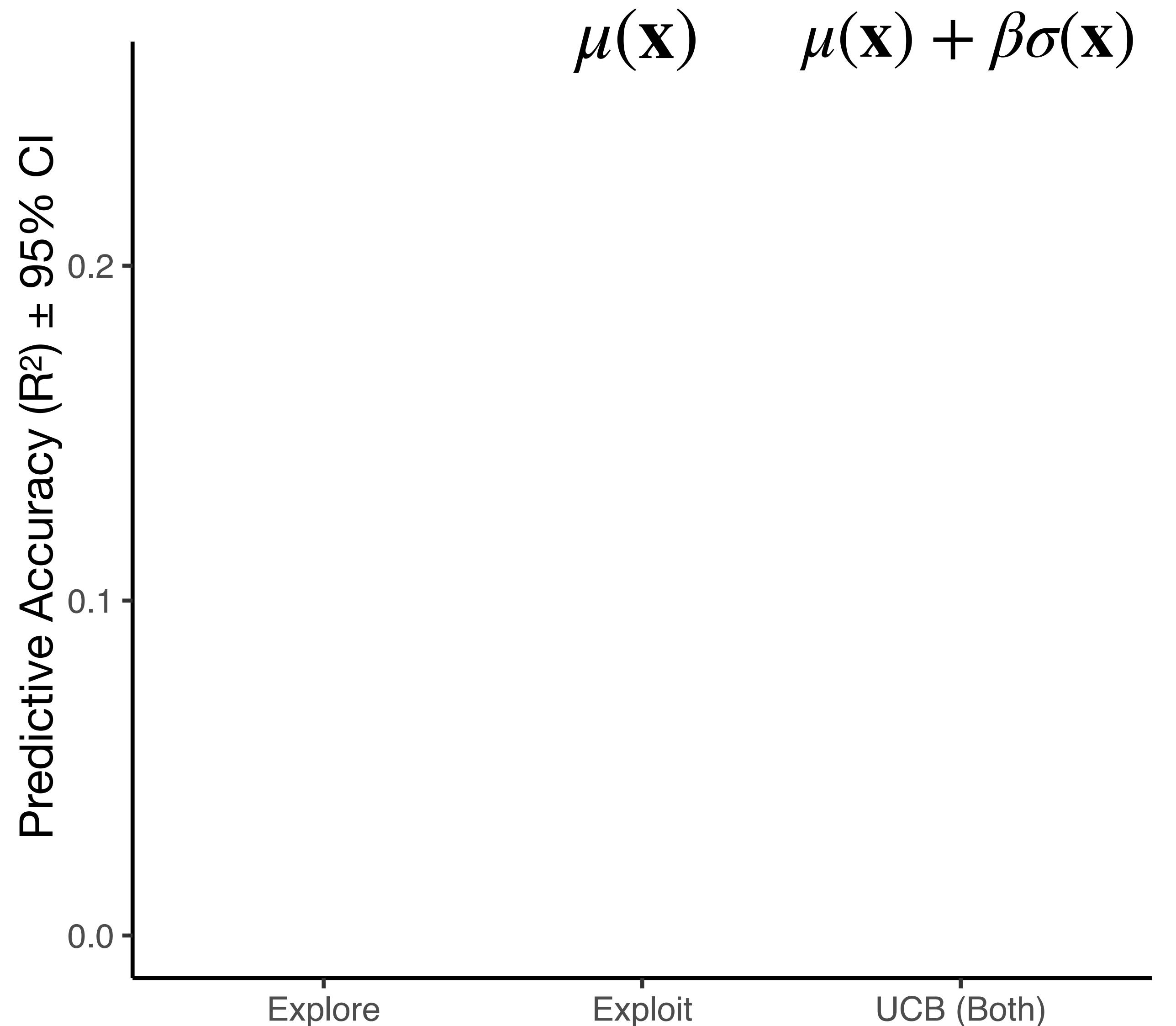
$$\mu(\mathbf{x}) + \beta\sigma(\mathbf{x})$$



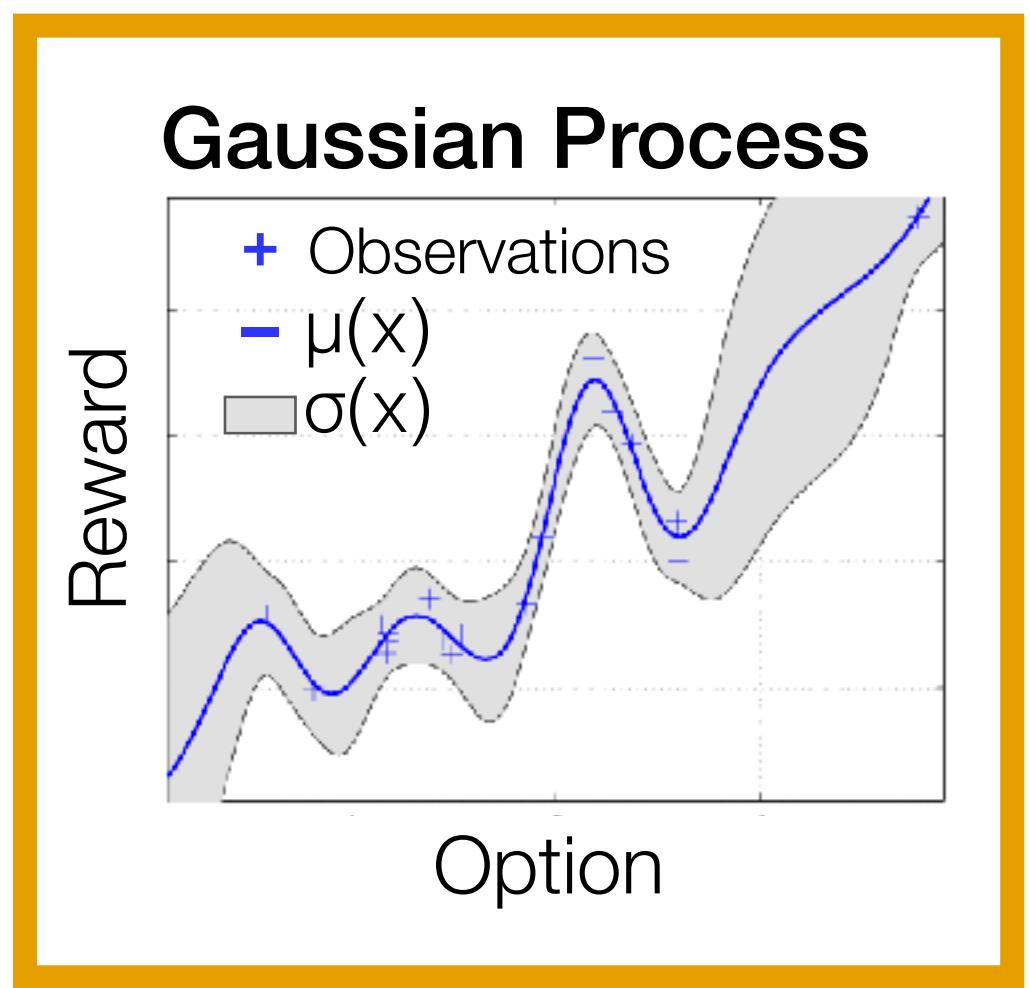
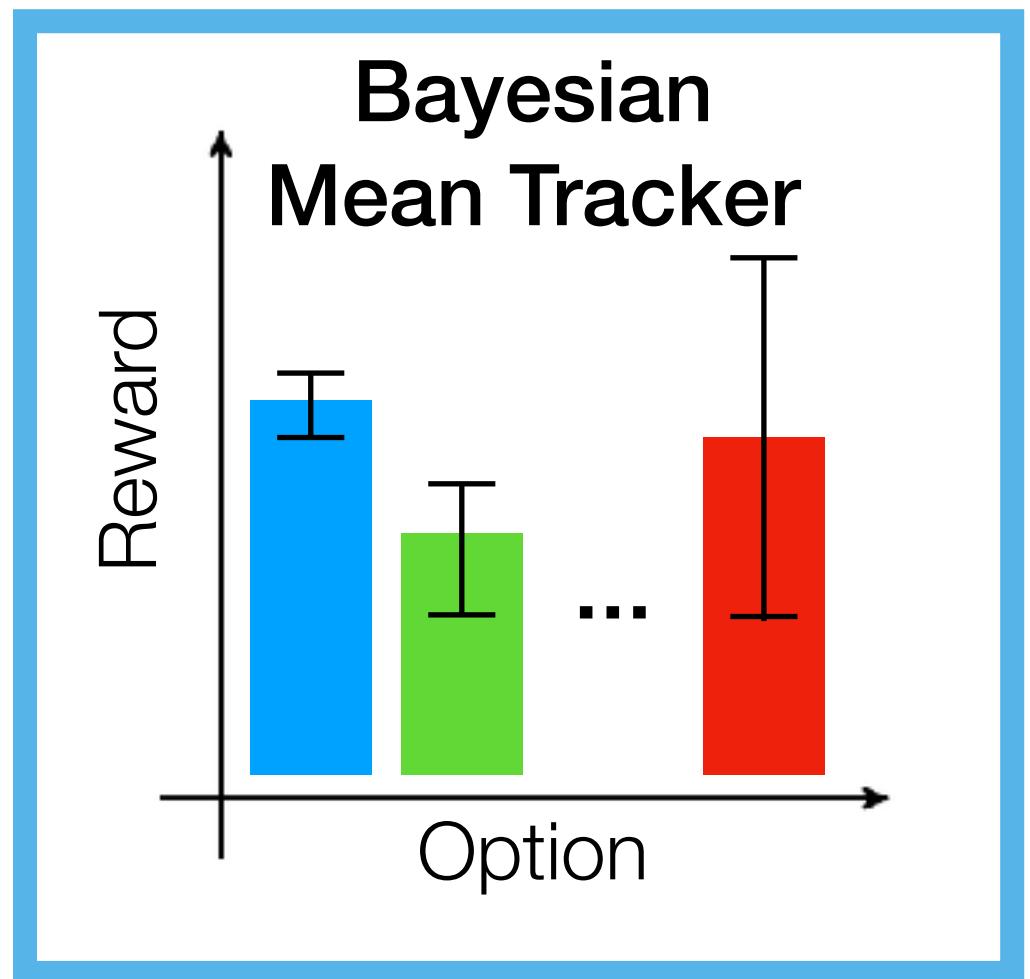
$$R^2 = 1 - \frac{\log \mathcal{L}(\mathcal{M}_k)}{\log \mathcal{L}(\mathcal{M}_{\text{rand}})}$$



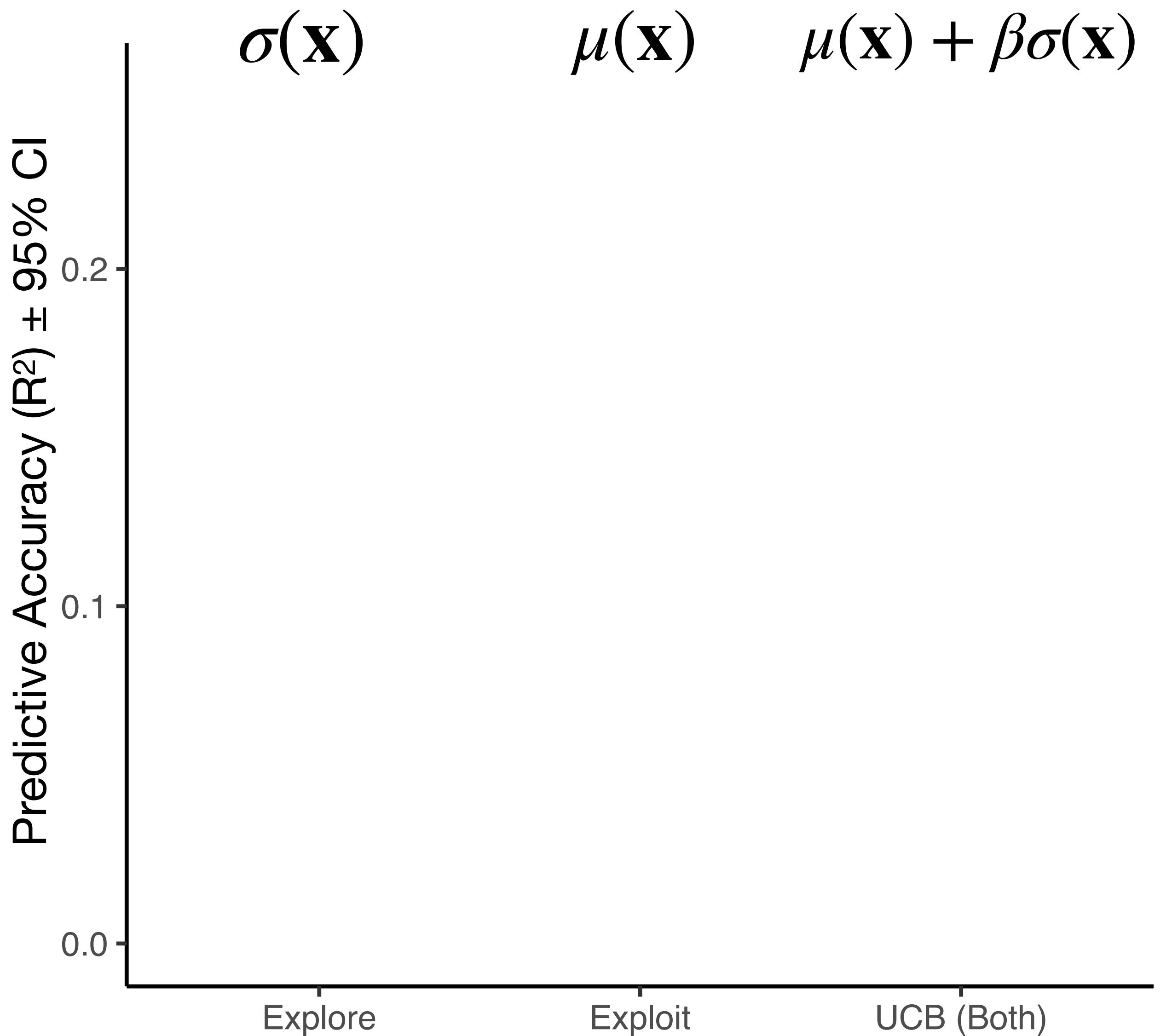
# Model Comparison



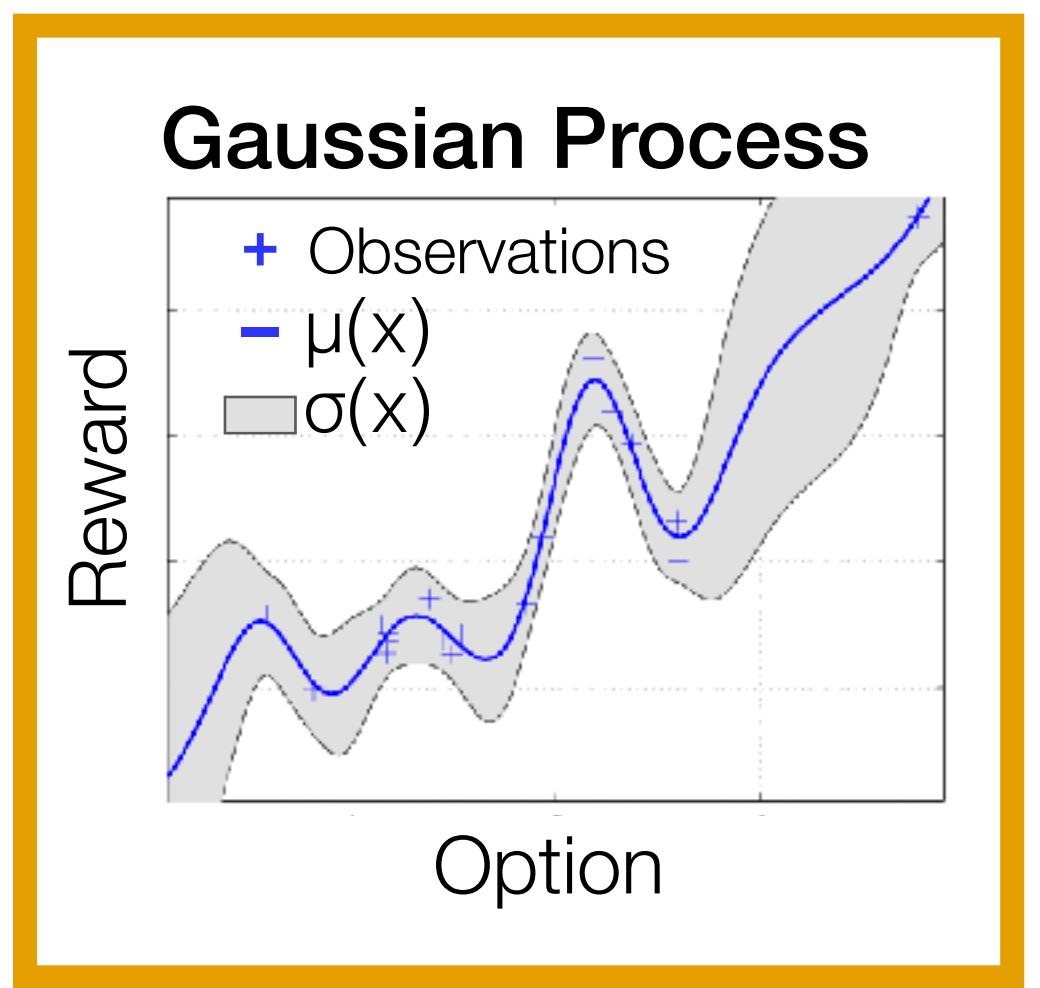
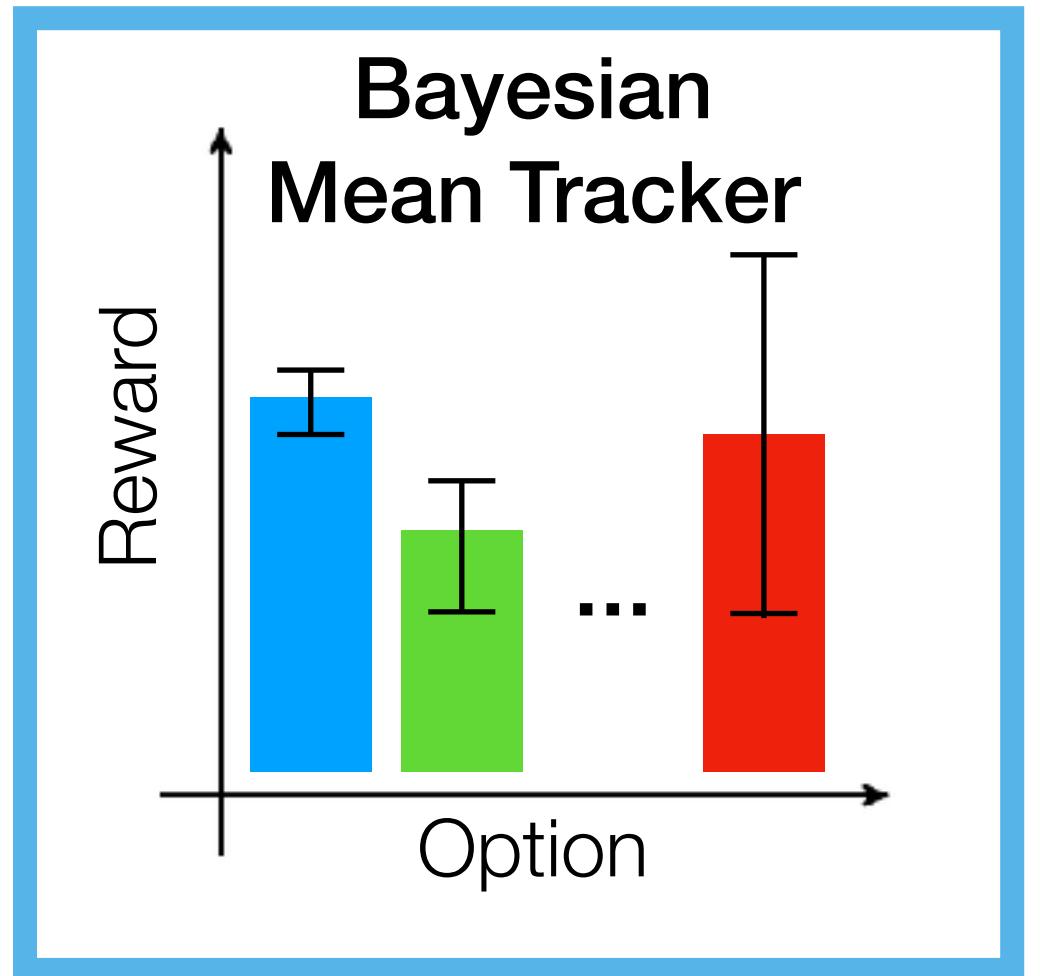
$$R^2 = 1 - \frac{\log \mathcal{L}(\mathcal{M}_k)}{\log \mathcal{L}(\mathcal{M}_{\text{rand}})}$$



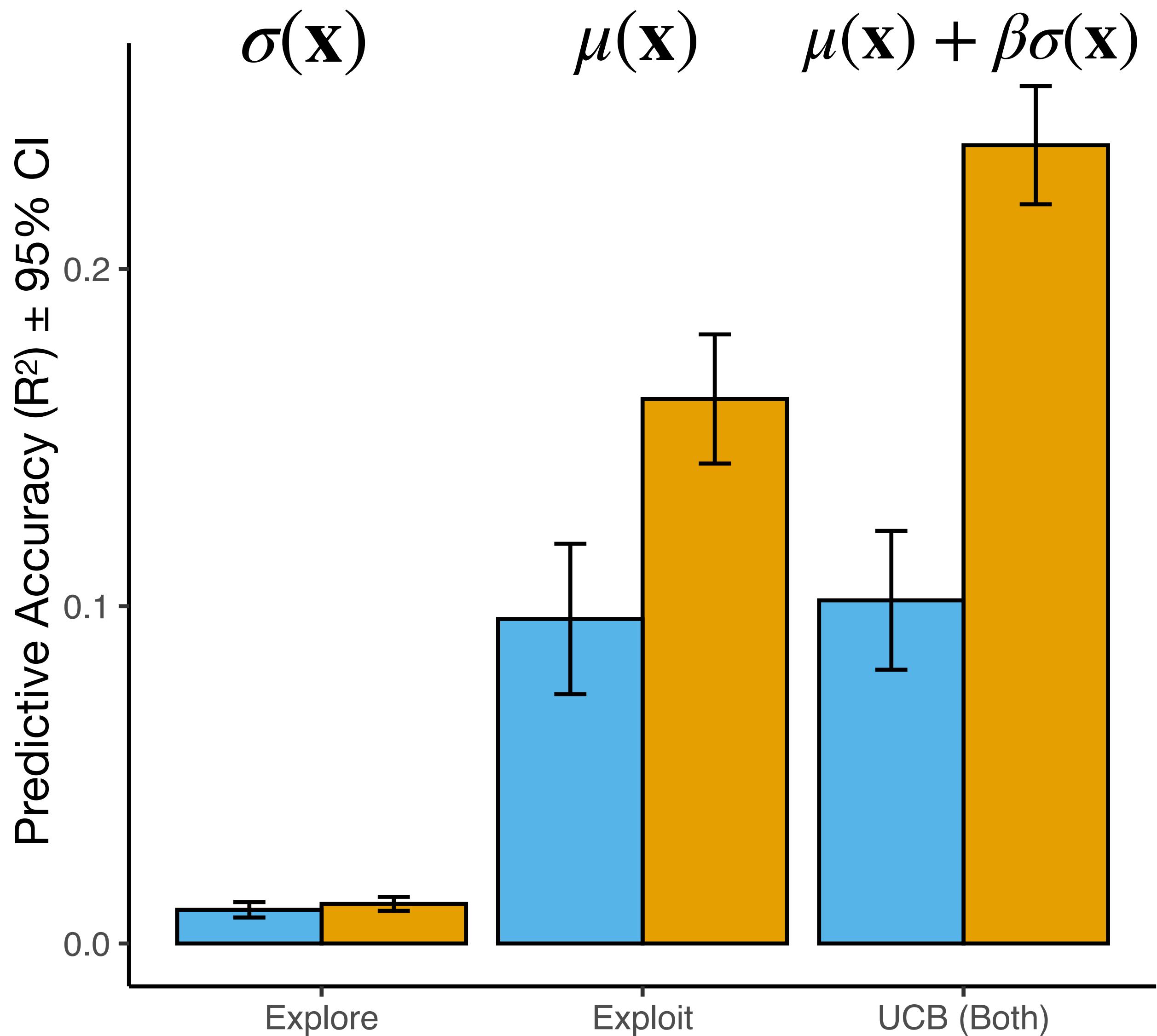
# Model Comparison



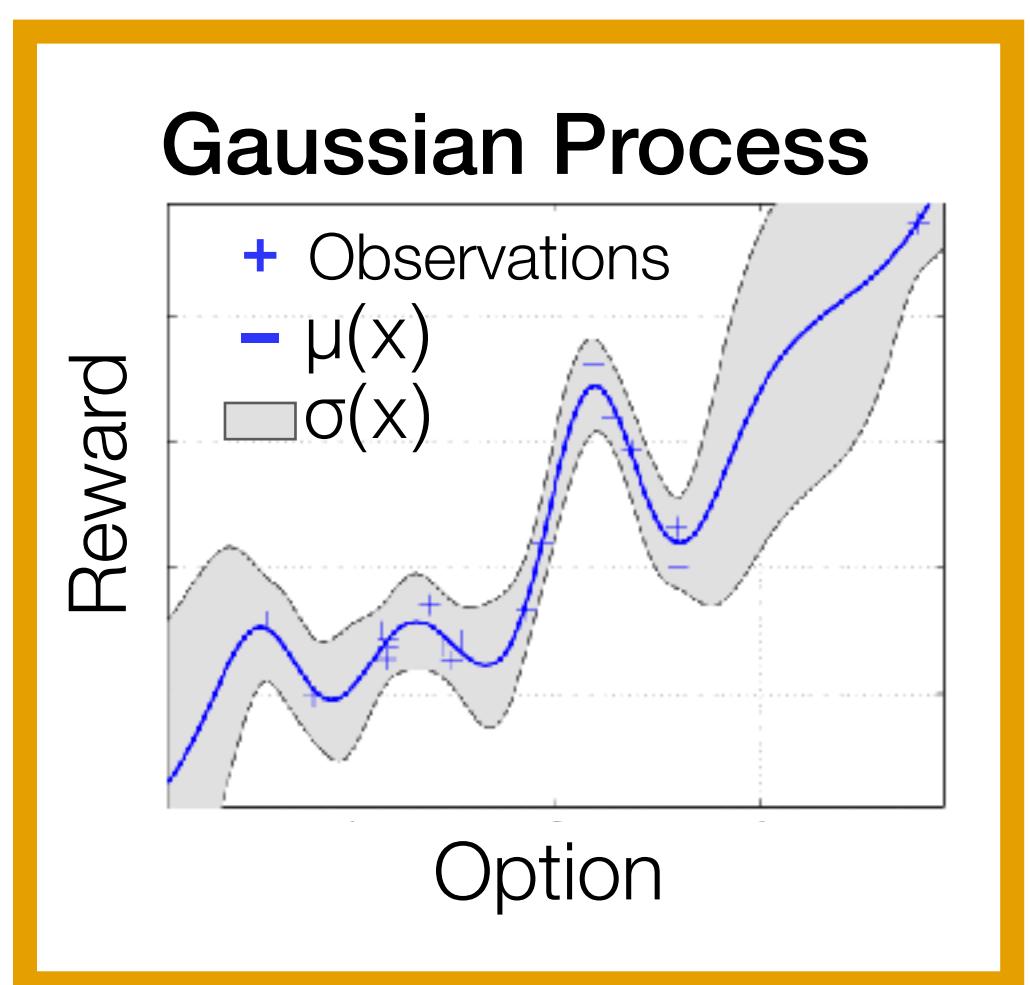
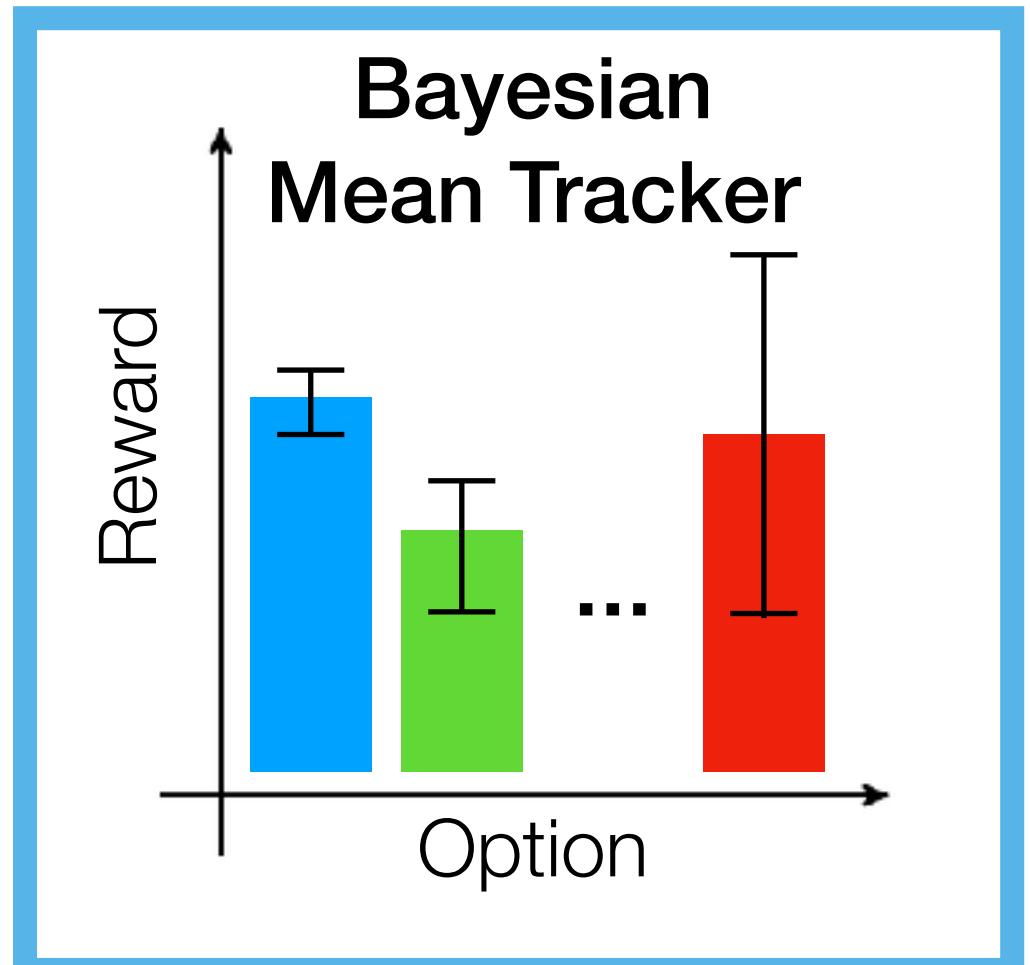
$$R^2 = 1 - \frac{\log \mathcal{L}(\mathcal{M}_k)}{\log \mathcal{L}(\mathcal{M}_{\text{rand}})}$$



# Model Comparison



$$R^2 = 1 - \frac{\log \mathcal{L}(\mathcal{M}_k)}{\log \mathcal{L}(\mathcal{M}_{\text{rand}})}$$



# Establishing a new paradigm

## 1. Generalization guides exploration

Wu, Schulz, Nelson, Speekenbrink & Meder (*Cogsci* 2017)

Wu, Schulz, Nelson, Speekenbrink & Meder (*Nature Human Behaviour* 2018)

## 2. Learning like a child

Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

Meder, Wu, Schulz & Ruggeri (*Dev Sci* in press)

## 3. Graph-structured Generalization

Wu, Schulz & Gershman (*Cogsci* 2019)

Wu, Schulz & Gershman (*CCN* 2019)

Wu, Schulz & Gershman (*Comput Brain Behav* 2020)

## 4. Search in abstract conceptual spaces

Wu, Schulz, Garvert, Meder & Schuck (*Cogsci* 2018)

Wu, Schulz, Garvert, Meder & Schuck (*PLOS Comp Bio* 2020)

## 5. Safe exploration

Schulz, Wu, Huys, Krause & Speekenbrink (*Cognitive Science* 2018)

## 6. Clinically depressed populations

Schefft, Wu, Meder, Köhler & Schulz (*in prep*)

## 7. Social search in VR

Wu, Ho, Kahl, Leuker, Meder & Kurvers (*bioRxiv* 2021)

# Establishing a new paradigm

## 1. Generalization guides exploration

Wu, Schulz, Nelson, Speekenbrink & Meder (*Cogsci* 2017)

Wu, Schulz, Nelson, Speekenbrink & Meder (*Nature Human Behaviour* 2018)

## 2. Learning like a child

Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

Meder, Wu, Schulz & Ruggeri (*Dev Sci* in press)

## 3. Graph-structured Generalization

Wu, Schulz & Gershman (*Cogsci* 2019)

Wu, Schulz & Gershman (*CCN* 2019)

Wu, Schulz & Gershman (*Comput Brain Behav* 2020)

## 4. Search in abstract conceptual spaces

Wu, Schulz, Garvert, Meder & Schuck (*Cogsci* 2018)

Wu, Schulz, Garvert, Meder & Schuck (*PLOS Comp Bio* 2020)

## 5. Safe exploration

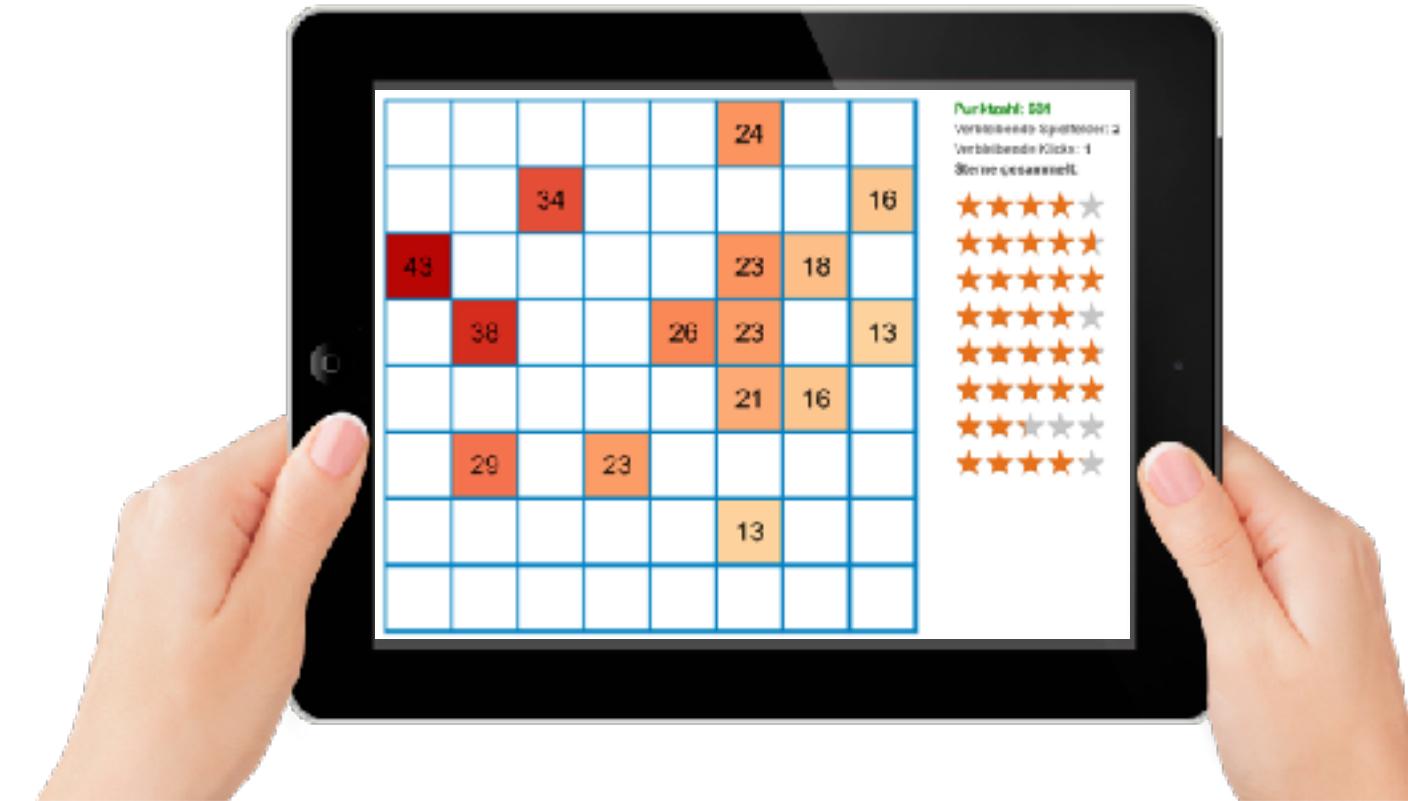
Schulz, Wu, Huys, Krause & Speekenbrink (*Cognitive Science* 2018)

## 6. Clinically depressed populations

Schefft, Wu, Meder, Köhler & Schulz (*in prep*)

## 7. Social search in VR

Wu, Ho, Kahl, Leuker, Meder & Kurvers (*bioRxiv* 2021)



Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

# Establishing a new paradigm

## 1. Generalization guides exploration

Wu, Schulz, Nelson, Speekenbrink & Meder (*Cogsci* 2017)

Wu, Schulz, Nelson, Speekenbrink & Meder (*Nature Human Behaviour* 2018)

## 2. Learning like a child

Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

Meder, Wu, Schulz & Ruggeri (*Dev Sci* in press)

## 3. Graph-structured Generalization

Wu, Schulz & Gershman (*Cogsci* 2019)

Wu, Schulz & Gershman (*CCN* 2019)

Wu, Schulz & Gershman (*Comput Brain Behav* 2020)

## 4. Search in abstract conceptual spaces

Wu, Schulz, Garvert, Meder & Schuck (*Cogsci* 2018)

Wu, Schulz, Garvert, Meder & Schuck (*PLOS Comp Bio* 2020)

## 5. Safe exploration

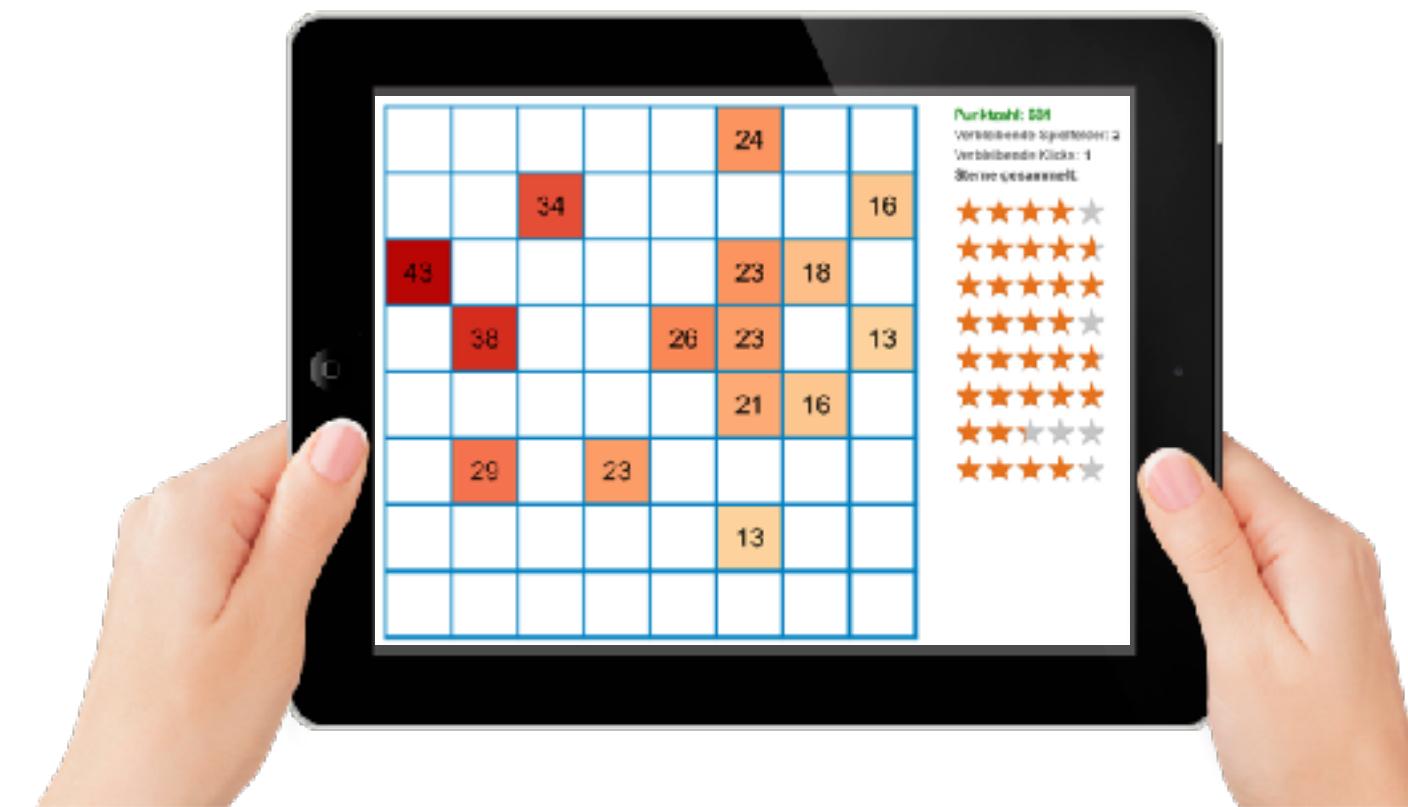
Schulz, Wu, Huys, Krause & Speekenbrink (*Cognitive Science* 2018)

## 6. Clinically depressed populations

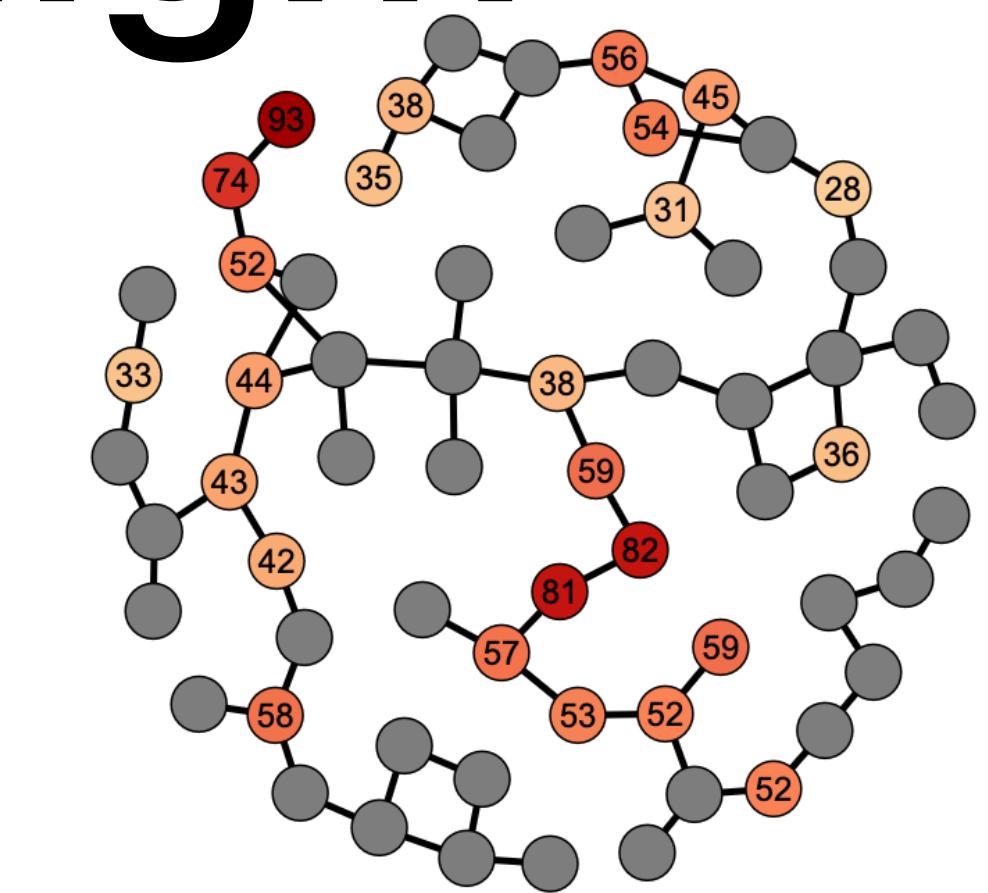
Schefft, Wu, Meder, Köhler & Schulz (*in prep*)

## 7. Social search in VR

Wu, Ho, Kahl, Leuker, Meder & Kurvers (*bioRxiv* 2021)



Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)



Wu, Schulz & Gershman (*CBB* 2020)

# Establishing a new paradigm.

# 1. Generalization guides exploration

Wu, Schulz, Nelson, Speekenbrink & Meder (*Cogsci* 2017)

Wu, Schulz, Nelson, Speekenbrink & Meder (*Nature Human Behaviour* 2018)

# 2. Learning like a child

Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

Meder, Wu, Schulz & Ruggeri (*Dev Sci* in press)

## 3. Graph-structured Generalization

Wu, Schulz & Gershman (*Cogsci* 2019)

Wu, Schulz & Gershman (CCN 2019)

Wu, Schulz & Gershman (Comput Brain Behav 2020)

## 4. Search in abstract conceptual spaces

Wu, Schulz, Garvert, Meder & Schuck (*Coascj* 2018)

Wu, Schulz, Garvert, Meder & Schuck (bioRxiv 2019);  
 Wu, Schulz, Garvert, Meder & Schuck (PLOS Comp Bio 2020)

# 5. Safe exploration

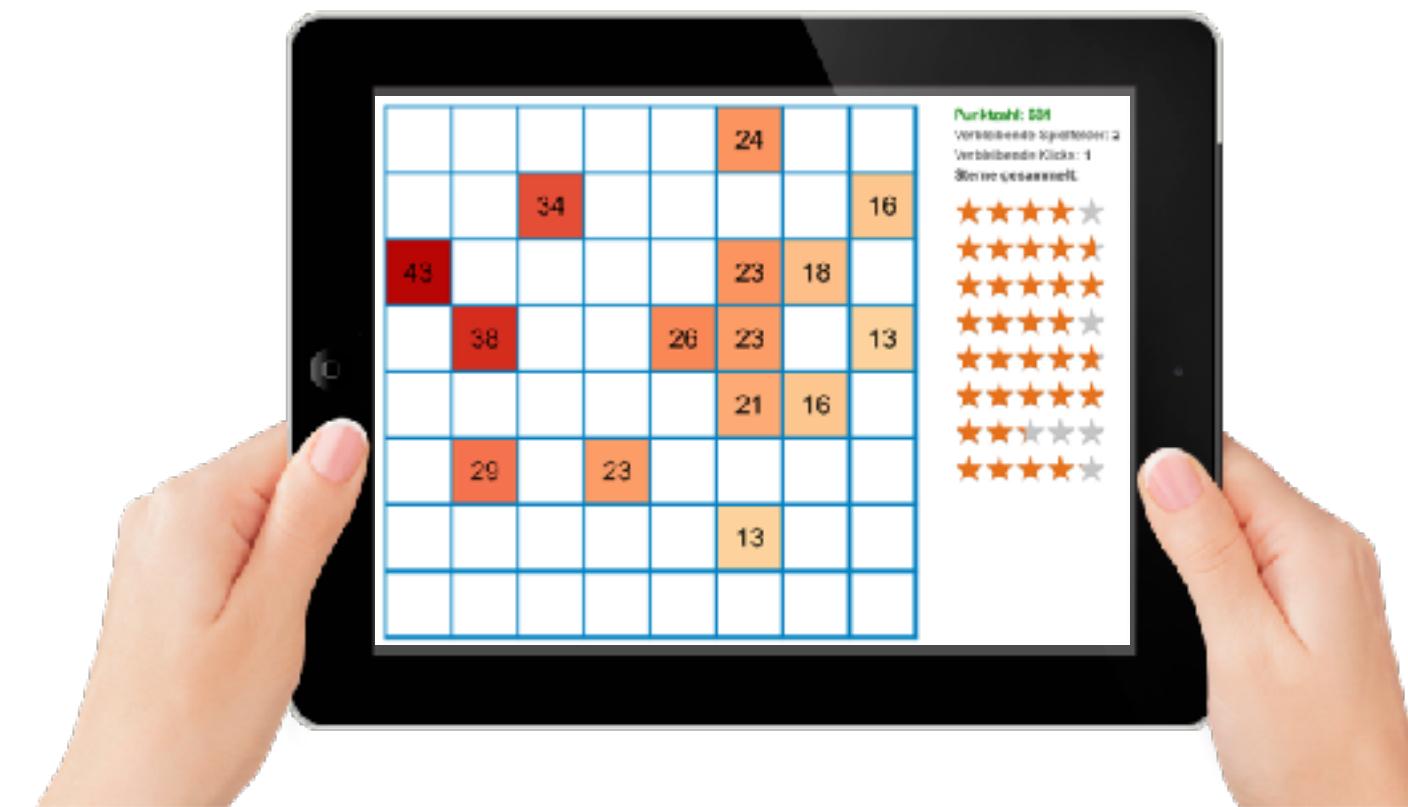
Schulz, Wu, Huys, Krause & Speekenbrink (*Cognitive Science* 2018)

## 6. Clinically depressed populations

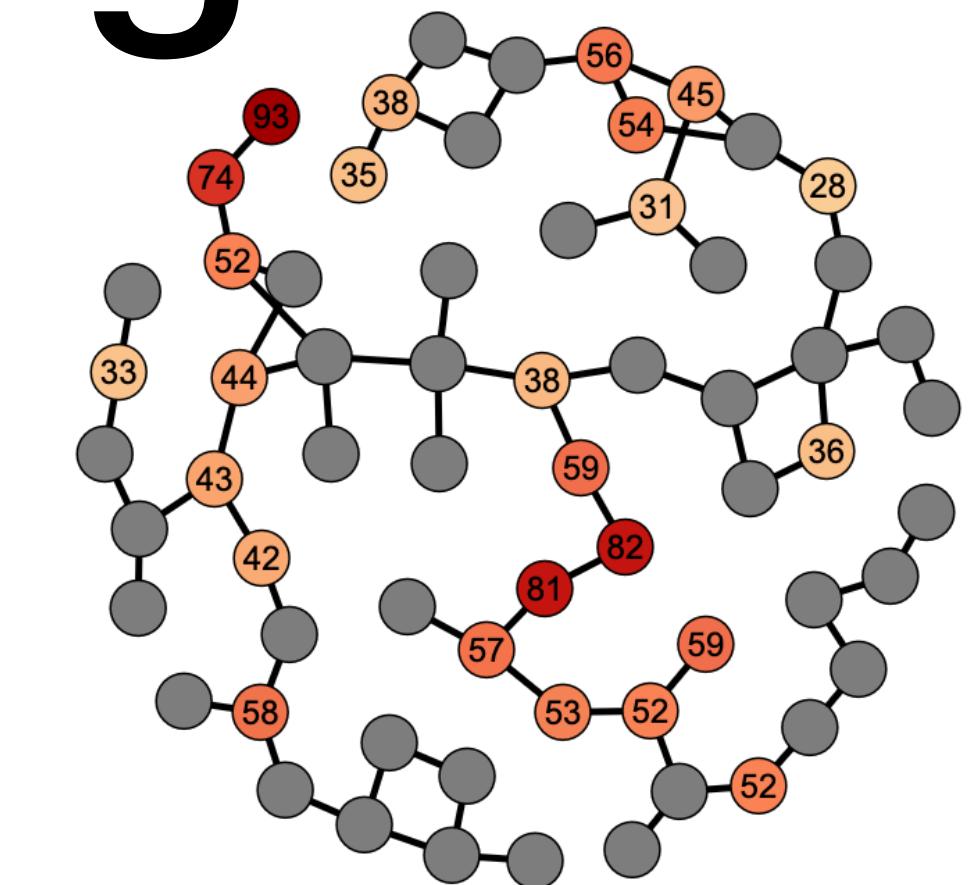
Schefft, Wu, Meder, Köhler & Schulz (in prep)

# 7. Social search in VR

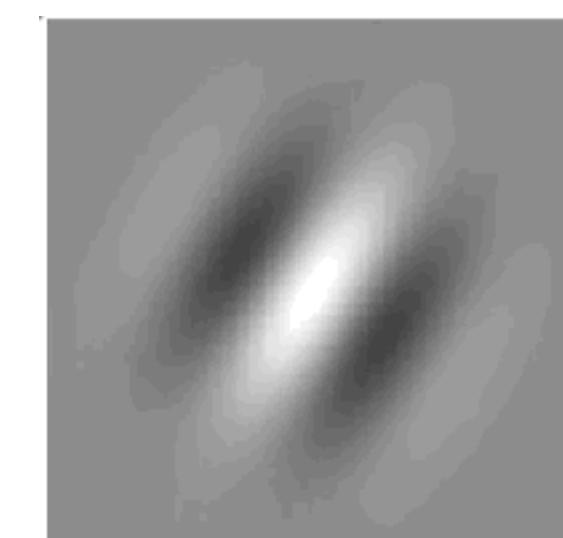
WU Ho, Kahl, Leuker, Meder & Kurvers (*bioRxiv*, 2021)



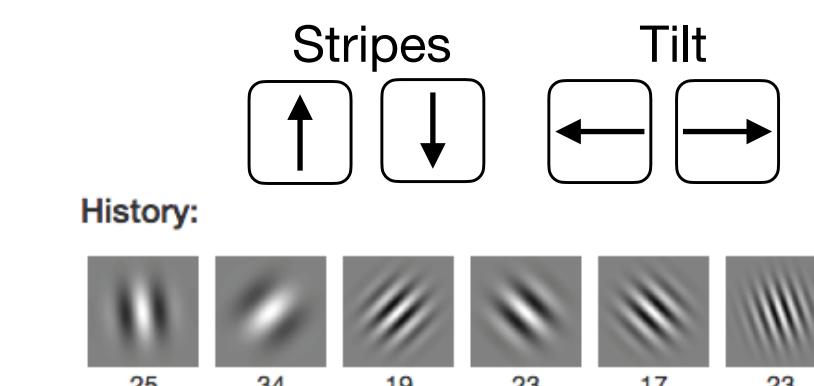
Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)



Wu, Schulz & Gershman (*CBB* 2020)



Current Score: 141  
Trials Remaining: 14  
Rounds Remaining: 10



Wu, Schulz, Garvert, Meder & Schuck  
*(PLOS Comp Bio 2020)*

# Establishing a new paradigm

## 1. Generalization guides exploration

Wu, Schulz, Nelson, Speekenbrink & Meder (*Cogsci* 2017)

Wu, Schulz, Nelson, Speekenbrink & Meder (*Nature Human Behaviour* 2018)

## 2. Learning like a child

Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

Meder, Wu, Schulz & Ruggeri (*Dev Sci* in press)

## 3. Graph-structured Generalization

Wu, Schulz & Gershman (*Cogsci* 2019)

Wu, Schulz & Gershman (*CCN* 2019)

Wu, Schulz & Gershman (*Comput Brain Behav* 2020)

## 4. Search in abstract conceptual spaces

Wu, Schulz, Garvert, Meder & Schuck (*Cogsci* 2018)

Wu, Schulz, Garvert, Meder & Schuck (*PLOS Comp Bio* 2020)

## 5. Safe exploration

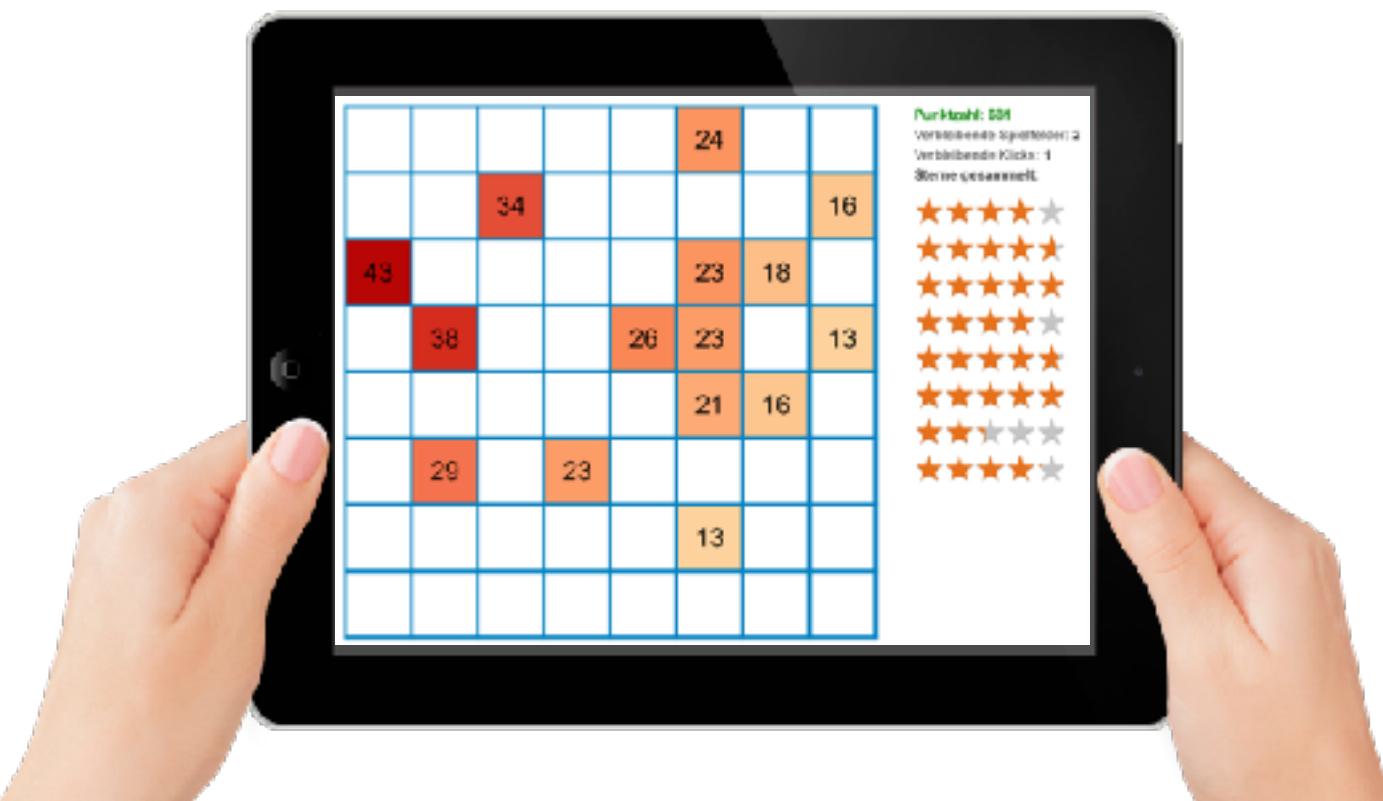
Schulz, Wu, Huys, Krause & Speekenbrink (*Cognitive Science* 2018)

## 6. Clinically depressed populations

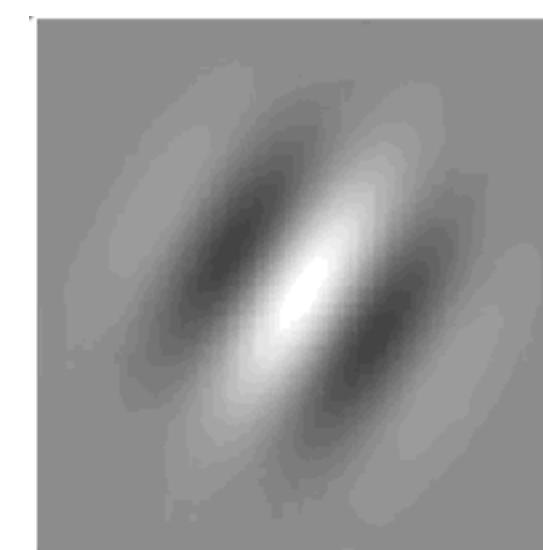
Schefft, Wu, Meder, Köhler & Schulz (*in prep*)

## 7. Social search in VR

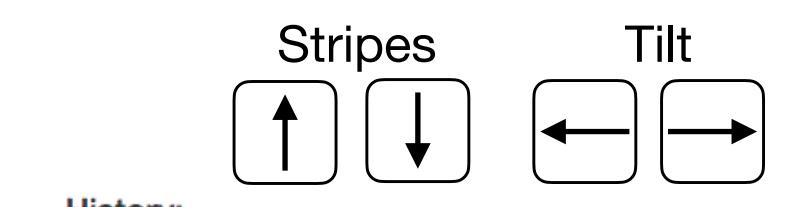
Wu, Ho, Kahl, Leuker, Meder & Kurvers (*bioRxiv* 2021)



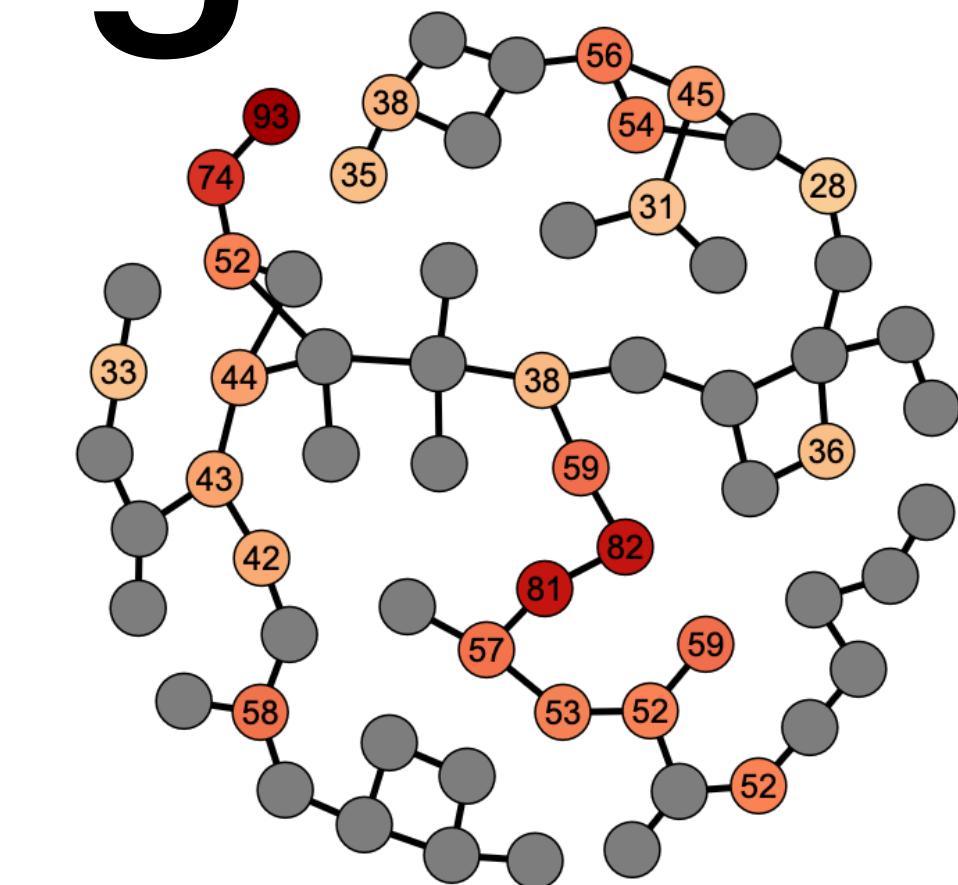
Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)



Current Score: 141  
Trials Remaining: 14  
Rounds Remaining: 10



Wu, Schulz, Garvert, Meder & Schuck  
(*PLOS Comp Bio* 2020)



Wu, Schulz & Gershman (*CBB* 2020)



Schulz, Wu, Huys, Krause & Speekenbrink  
(*Cognitive Science* 2018)



# Establishing a new paradigm

## 1. Generalization guides exploration

Wu, Schulz, Nelson, Speekenbrink & Meder (*Cogsci* 2017)

Wu, Schulz, Nelson, Speekenbrink & Meder (*Nature Human Behaviour* 2018)

## 2. Learning like a child

Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

Meder, Wu, Schulz & Ruggeri (*Dev Sci* in press)

## 3. Graph-structured Generalization

Wu, Schulz & Gershman (*Cogsci* 2019)

Wu, Schulz & Gershman (*CCN* 2019)

Wu, Schulz & Gershman (*Comput Brain Behav* 2020)

## 4. Search in abstract conceptual spaces

Wu, Schulz, Garvert, Meder & Schuck (*Cogsci* 2018)

Wu, Schulz, Garvert, Meder & Schuck (*PLOS Comp Bio* 2020)

## 5. Safe exploration

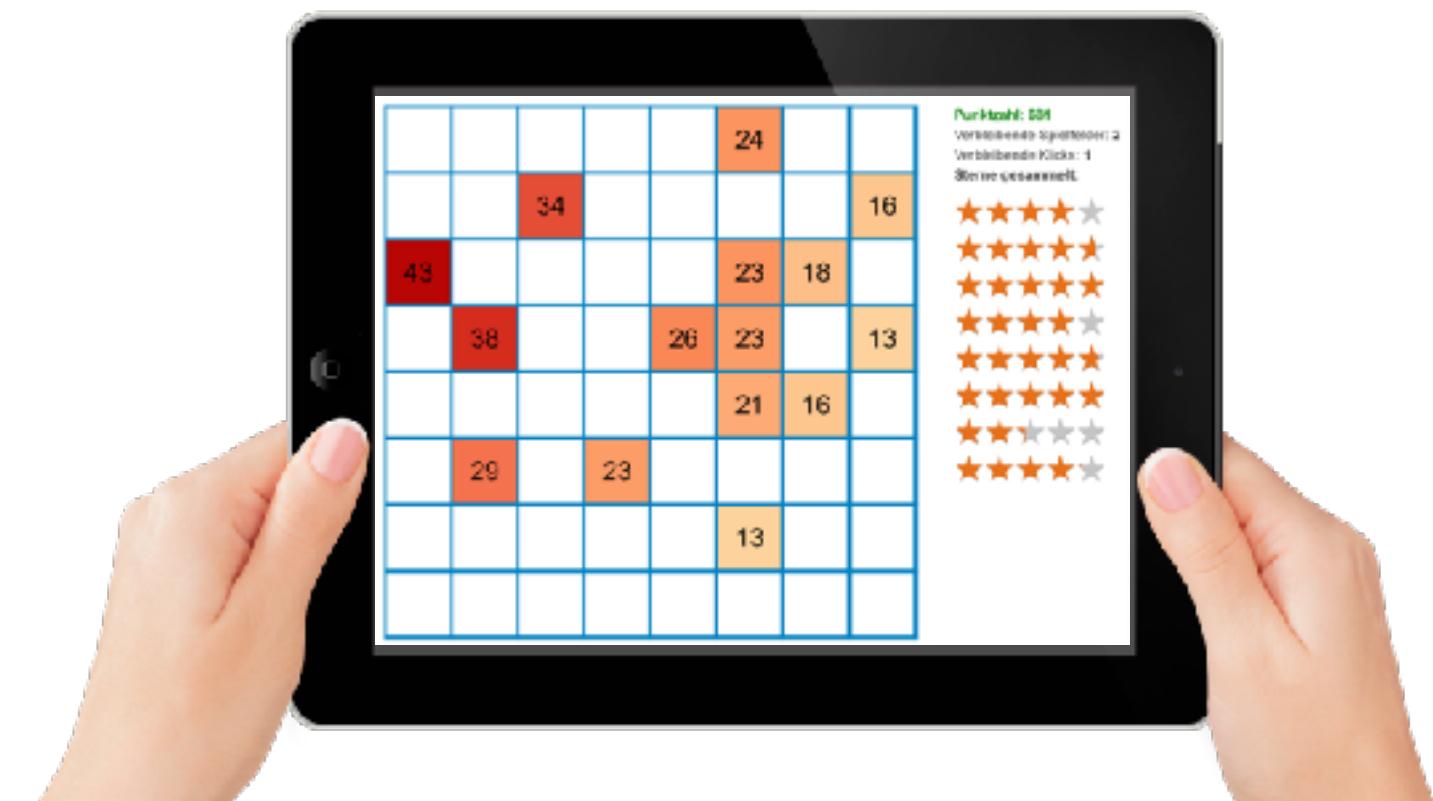
Schulz, Wu, Huys, Krause & Speekenbrink (*Cognitive Science* 2018)

## 6. Clinically depressed populations

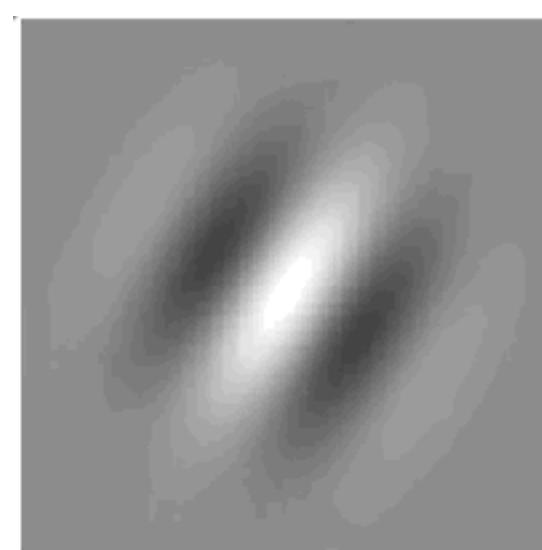
Schefft, Wu, Meder, Köhler & Schulz (*in prep*)

## 7. Social search in VR

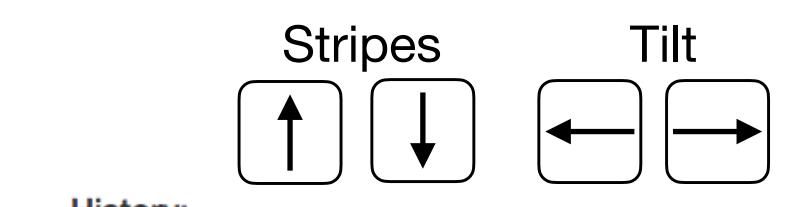
Wu, Ho, Kahl, Leuker, Meder & Kurvers (*bioRxiv* 2021)



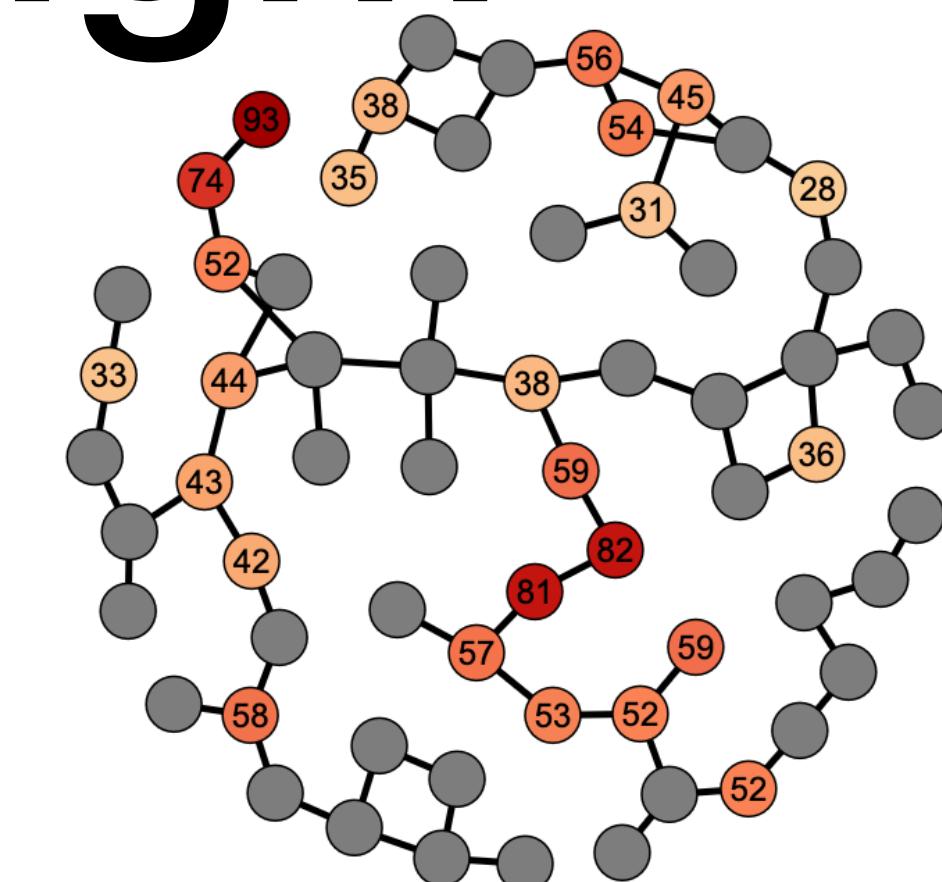
Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)



Current Score: 141  
Trials Remaining: 14  
Rounds Remaining: 10



Wu, Schulz, Garvert, Meder & Schuck  
(*PLOS Comp Bio* 2020)



Wu, Schulz & Gershman (*CBB* 2020)



Schulz, Wu, Huys, Krause & Speekenbrink  
(*Cognitive Science* 2018)



# Establishing a new paradigm

## 1. Generalization guides exploration

Wu, Schulz, Nelson, Speekenbrink & Meder (*Cogsci* 2017)

Wu, Schulz, Nelson, Speekenbrink & Meder (*Nature Human Behaviour* 2018)

## 2. Learning like a child

Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

Meder, Wu, Schulz & Ruggeri (*Dev Sci* in press)

## 3. Graph-structured Generalization

Wu, Schulz & Gershman (*Cogsci* 2019)

Wu, Schulz & Gershman (*CCN* 2019)

Wu, Schulz & Gershman (*Comput Brain Behav* 2020)

## 4. Search in abstract conceptual spaces

Wu, Schulz, Garvert, Meder & Schuck (*Cogsci* 2018)

Wu, Schulz, Garvert, Meder & Schuck (*PLOS Comp Bio* 2020)

## 5. Safe exploration

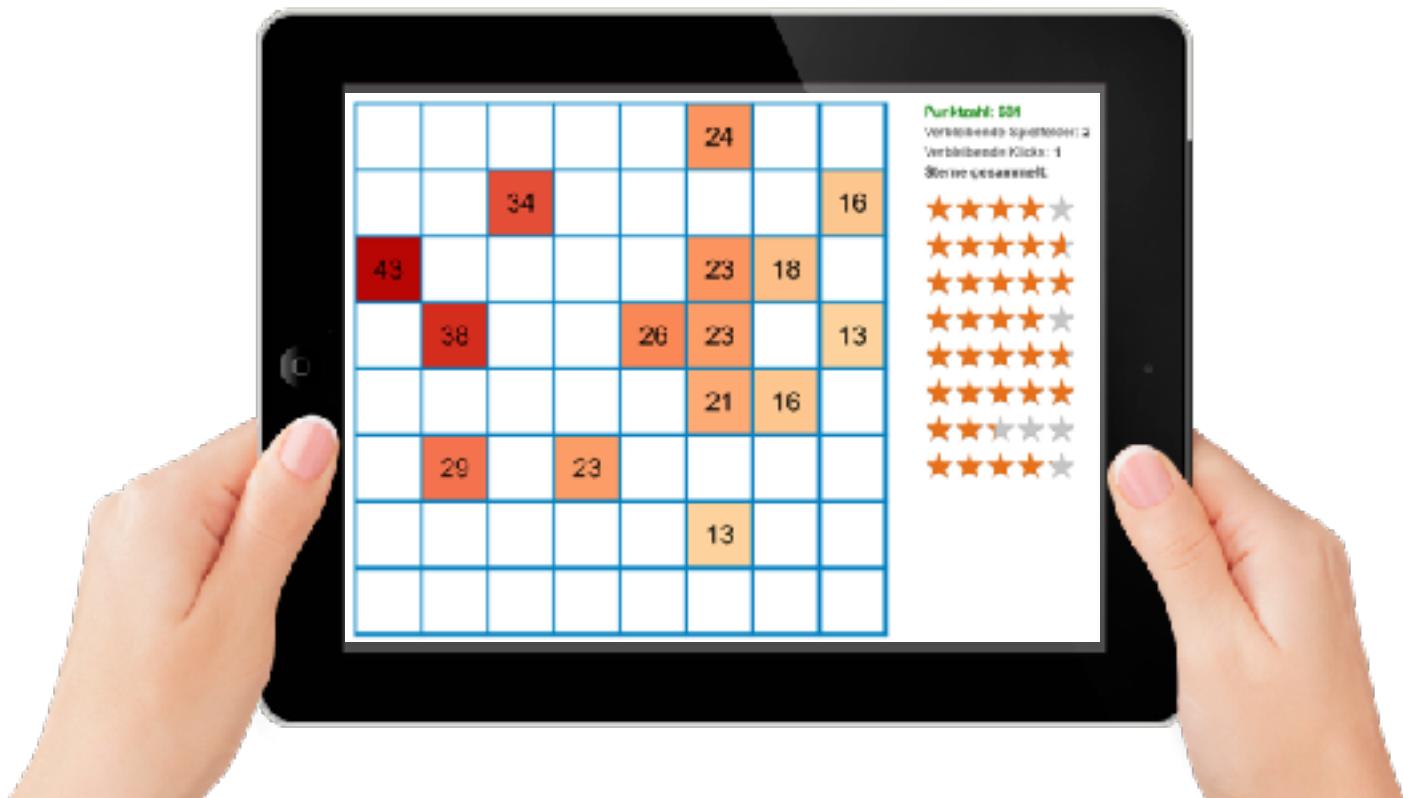
Schulz, Wu, Huys, Krause & Speekenbrink (*Cognitive Science* 2018)

## 6. Clinically depressed populations

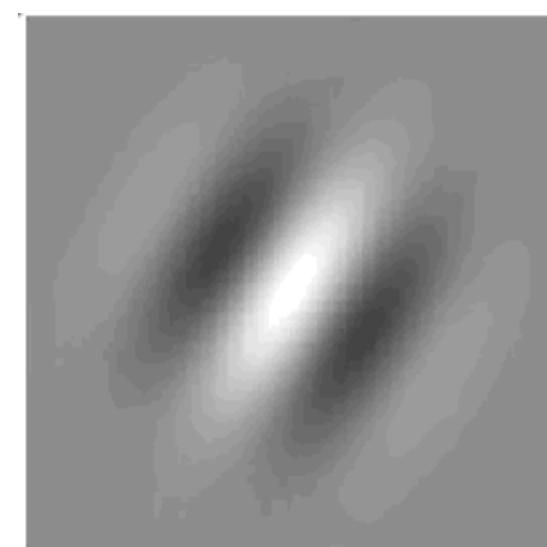
Schefft, Wu, Meder, Köhler & Schulz (*in prep*)

## 7. Social search in VR

Wu, Ho, Kahl, Leuker, Meder & Kurvers (*bioRxiv* 2021)

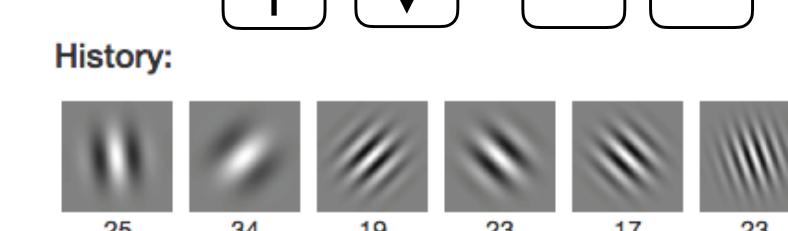


Schulz, Wu, Ruggeri & Meder (*PsychSci* 2019)

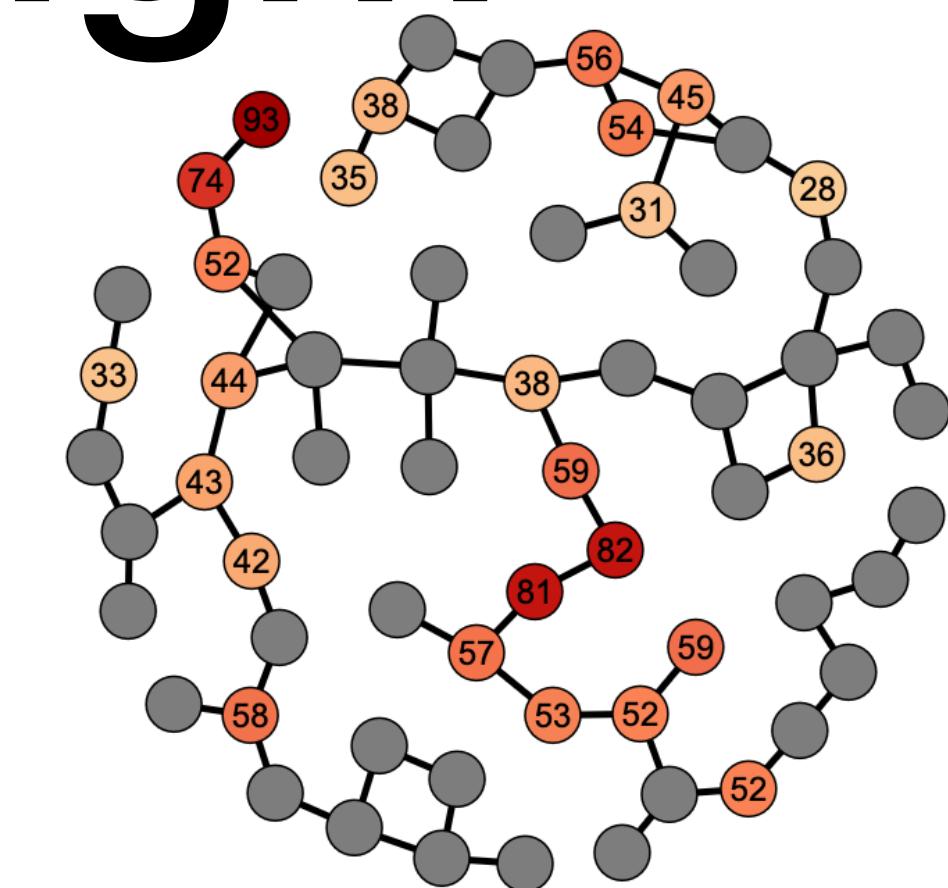


Current Score: 141  
Trials Remaining: 14  
Rounds Remaining: 10

Stripes      Tilt  
↑ ↓ ← →



Wu, Schulz, Garvert, Meder & Schuck  
(*PLOS Comp Bio* 2020)



Wu, Schulz & Gershman (*CBB* 2020)



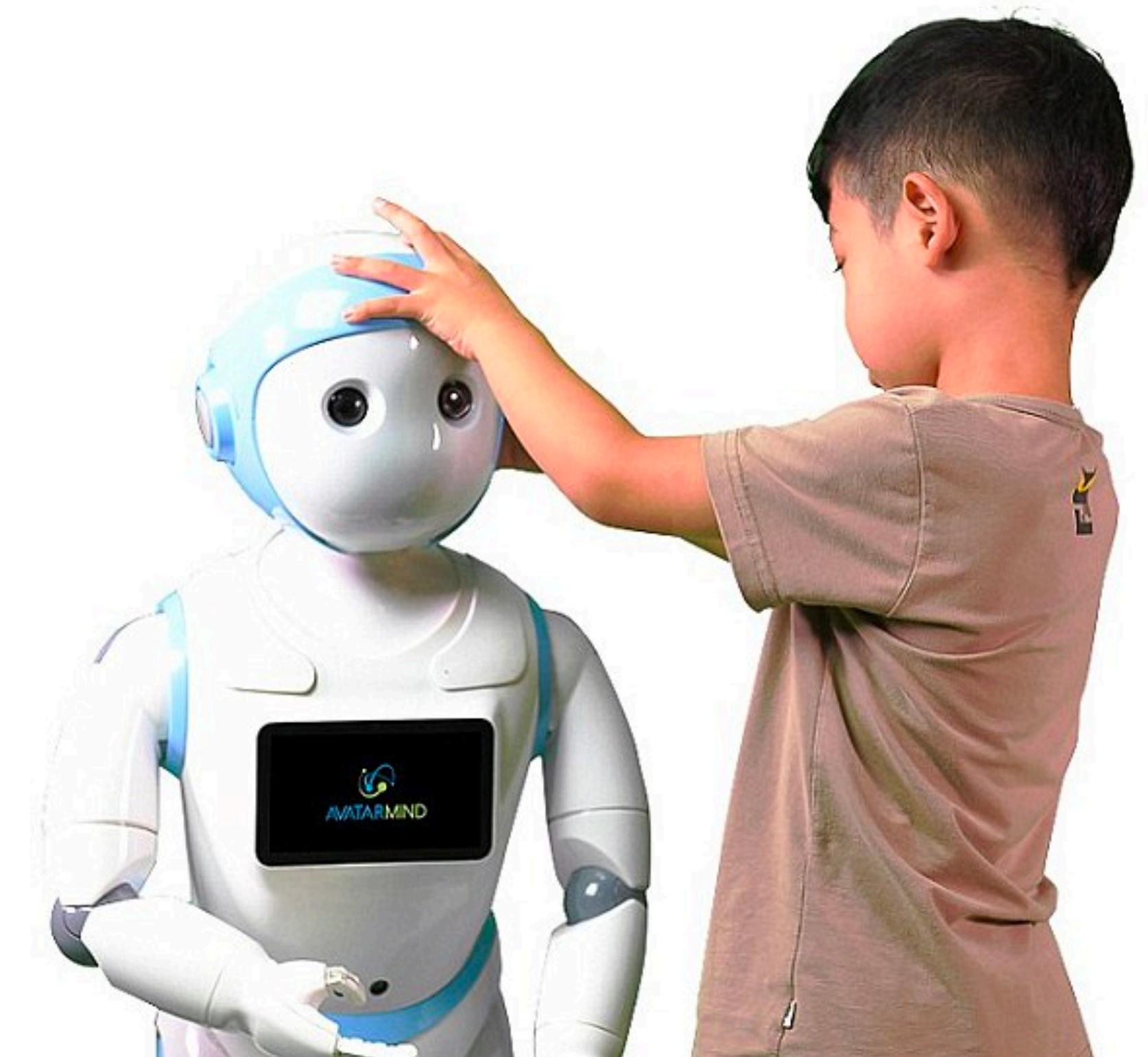
Schulz, Wu, Huys, Krause & Speekenbrink  
(*Cognitive Science* 2018)

# Part 2

How to learn like a child

# What AI can learn from Children

- Josh Tenenbaum (MIT): Children are the only known information processing system that demonstrably and reproducibly develop into intelligent systems
- Turing (1950) suggested we should build AI that learns like a child
- How do children learn differently from adults?



# What AI can learn from Children

- Josh Tenenbaum (MIT): Children are the only known information processing system that demonstrably and reproducibly develop into intelligent systems
- Turing (1950) suggested we should build AI that learns like a child
- How do children learn differently from adults?
- One robust finding is they are highly variable!



video by Francis Vachon

# What explains the extensive variability found in children's search behavior?

- High temperature sampling hypothesis:
  - Children initially perform high-temperature search, which gradually “cool offs” as they grow older (Gopnik et al., *Curr Dir Psych Sci* 2017)



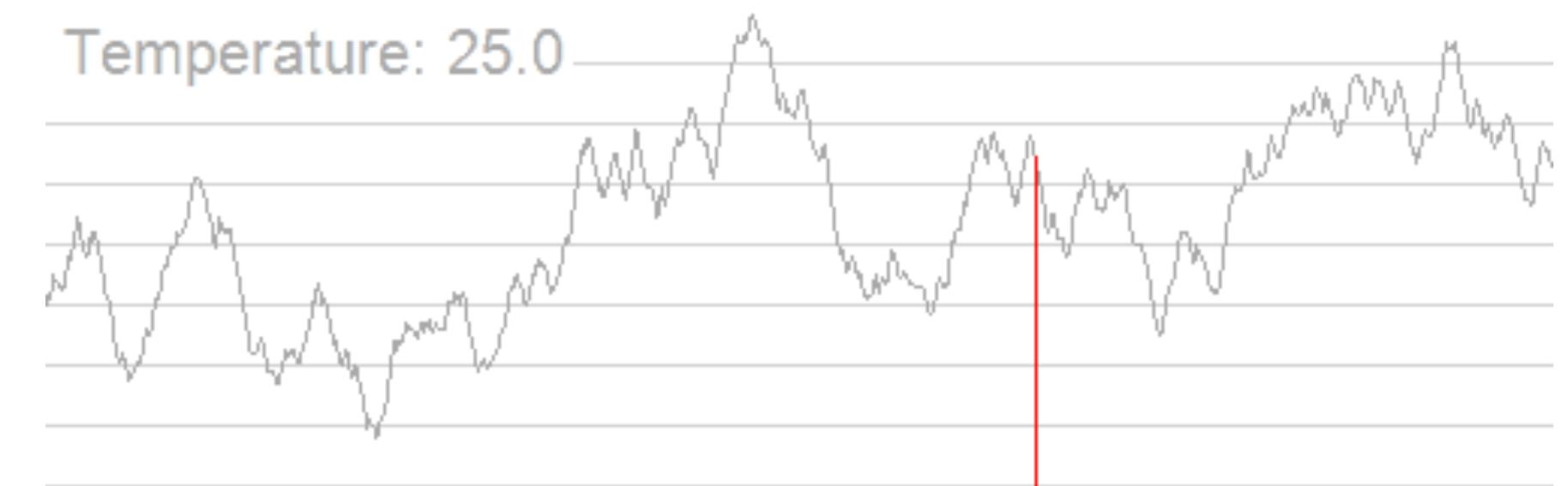
# What explains the extensive variability found in children's search behavior?

- High temperature sampling hypothesis:
  - Children initially perform high-temperature search, which gradually “cool offs” as they grow older (Gopnik et al., *Curr Dir Psych Sci* 2017)



# Sources of Developmental Differences

- Higher **temperature** sampling that cools off over age? (Gopnik et al., 2017)



# Sources of Developmental Differences

- Higher **temperature** sampling that cools off over age? (Gopnik et al., 2017)



# Sources of Developmental Differences

- Higher **temperature** sampling that cools off over age? (Gopnik et al., 2017)

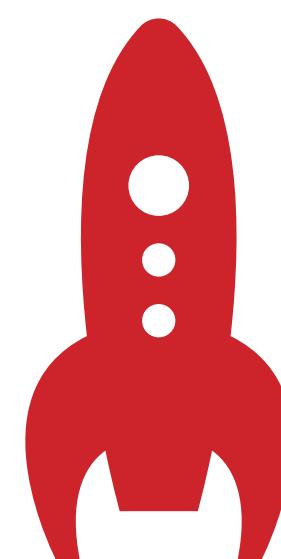
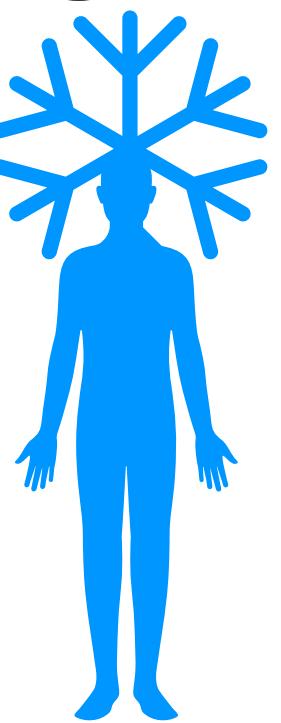


# Sources of Developmental Differences

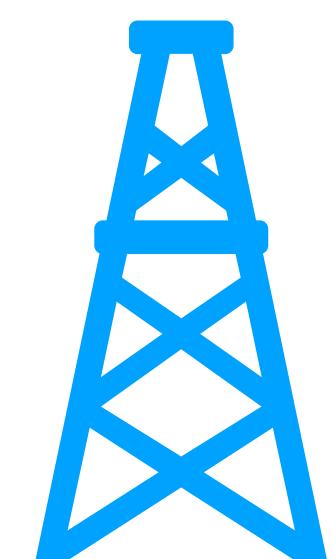
- Higher **temperature** sampling that cools off over age? (Gopnik et al., 2017)
- Changes in **directed exploration** rather than random exploration? (Wilson et al., 2014)



less random  
sampling



less directed  
exploration

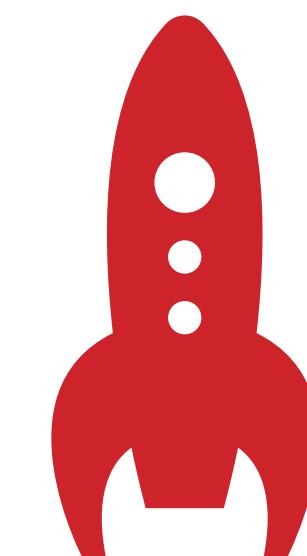
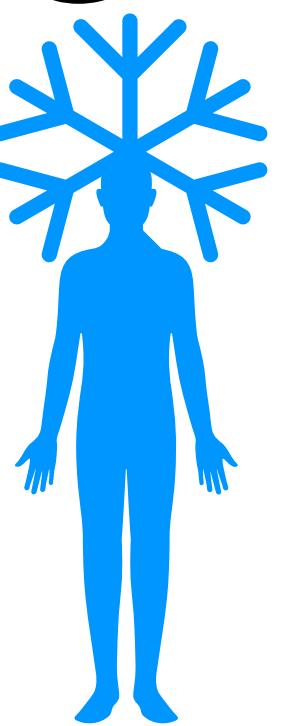


# Sources of Developmental Differences

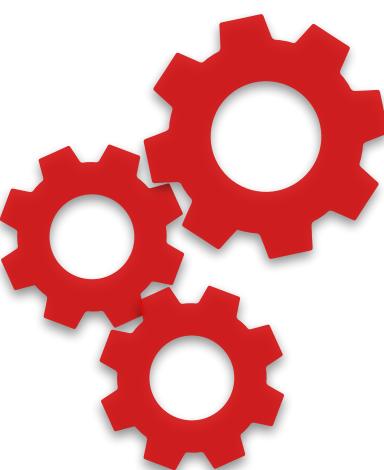
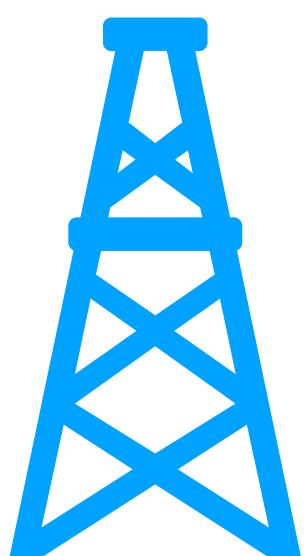
- Higher **temperature** sampling that cools off over age? (Gopnik et al., 2017)
- Changes in **directed exploration** rather than random exploration? (Wilson et al., 2014)
- Refinement of cognitive representations and processes supporting **generalization**? (Blanco et al., 2016)



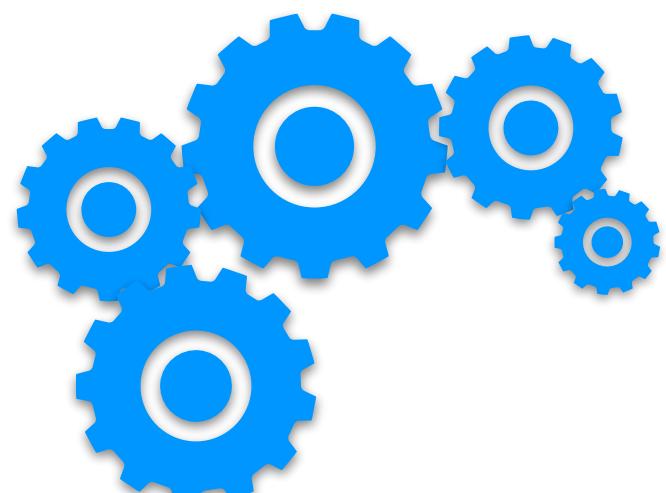
less random  
sampling



less directed  
exploration



broader  
generalization

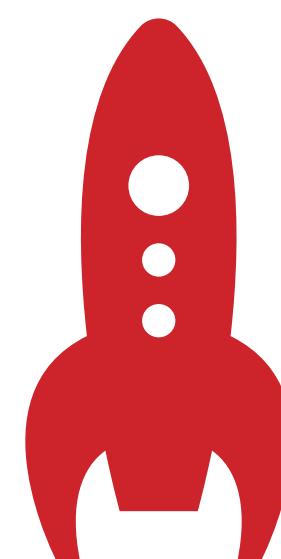
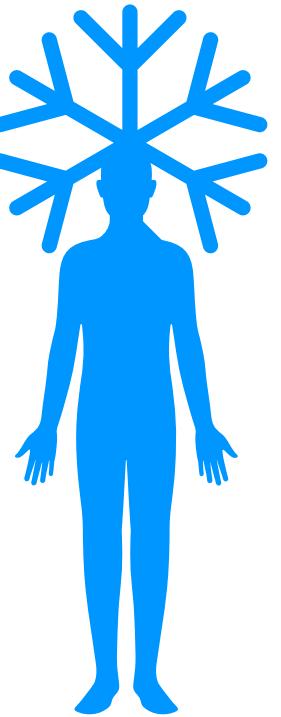


# Sources of Developmental Differences

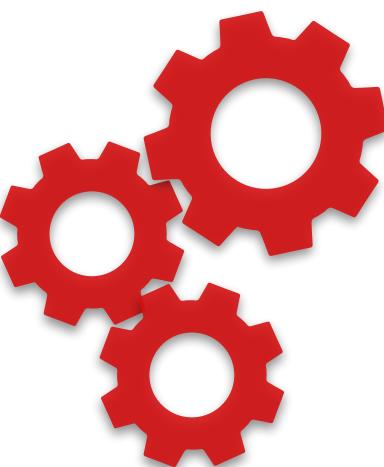
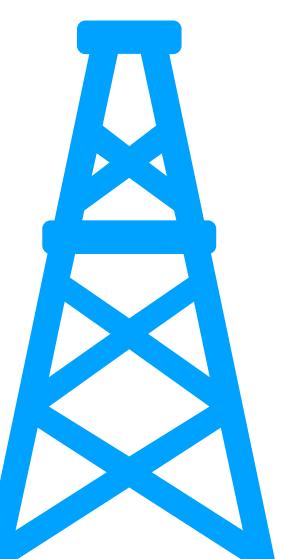
- Higher **temperature** sampling that cools off over age? (Gopnik et al., 2017)
- Changes in **directed exploration** rather than random exploration? (Wilson et al., 2014)
- Refinement of cognitive representations and processes supporting **generalization**? (Blanco et al., 2016)

 $\tau$ 

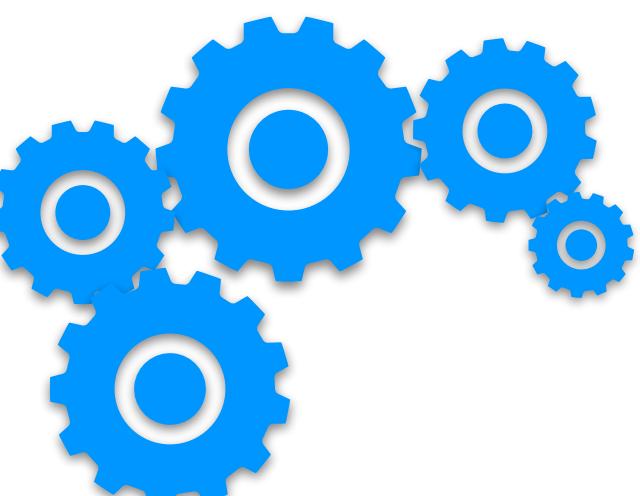
less random  
sampling

 $\beta$ 

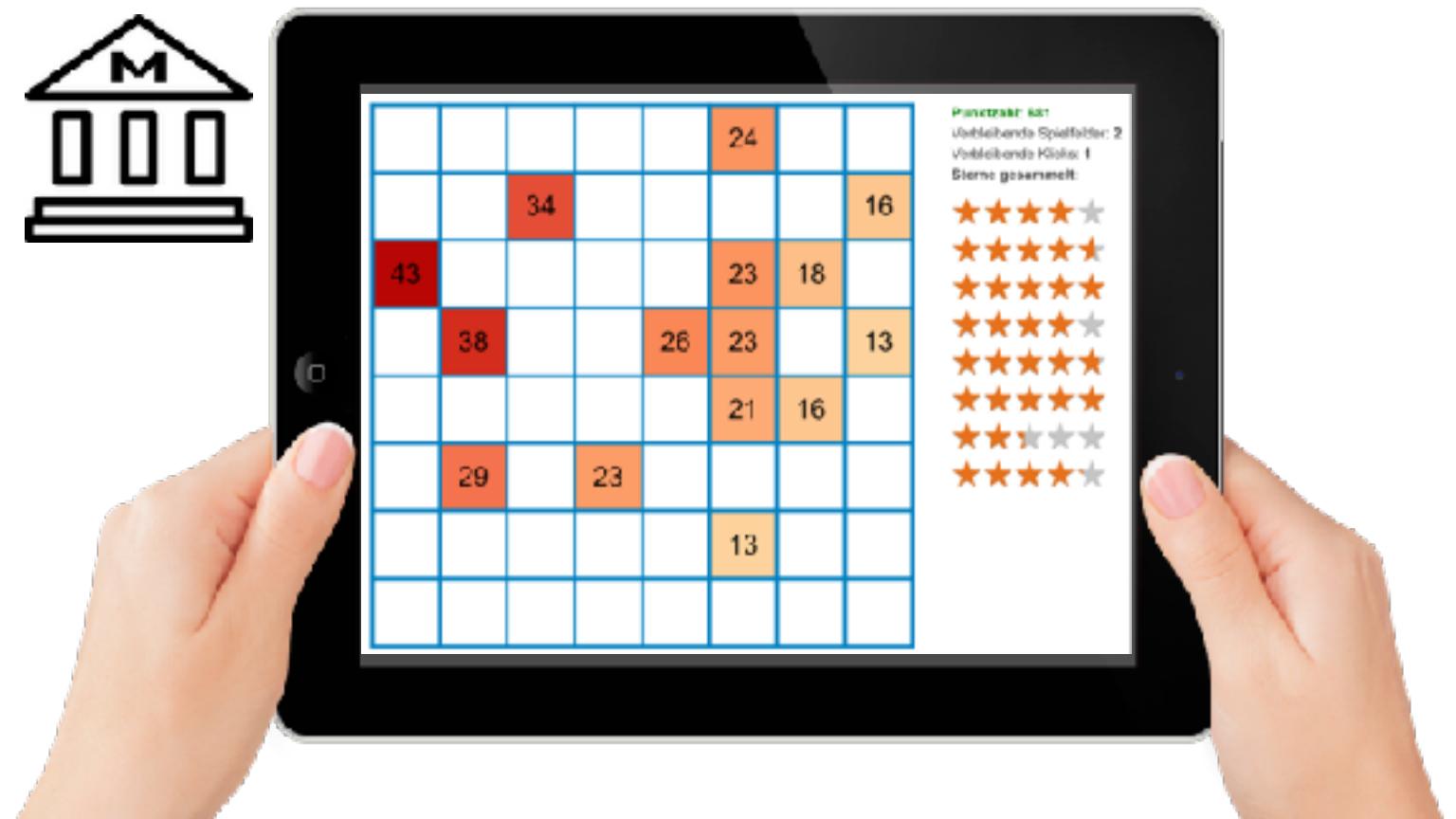
less directed  
exploration

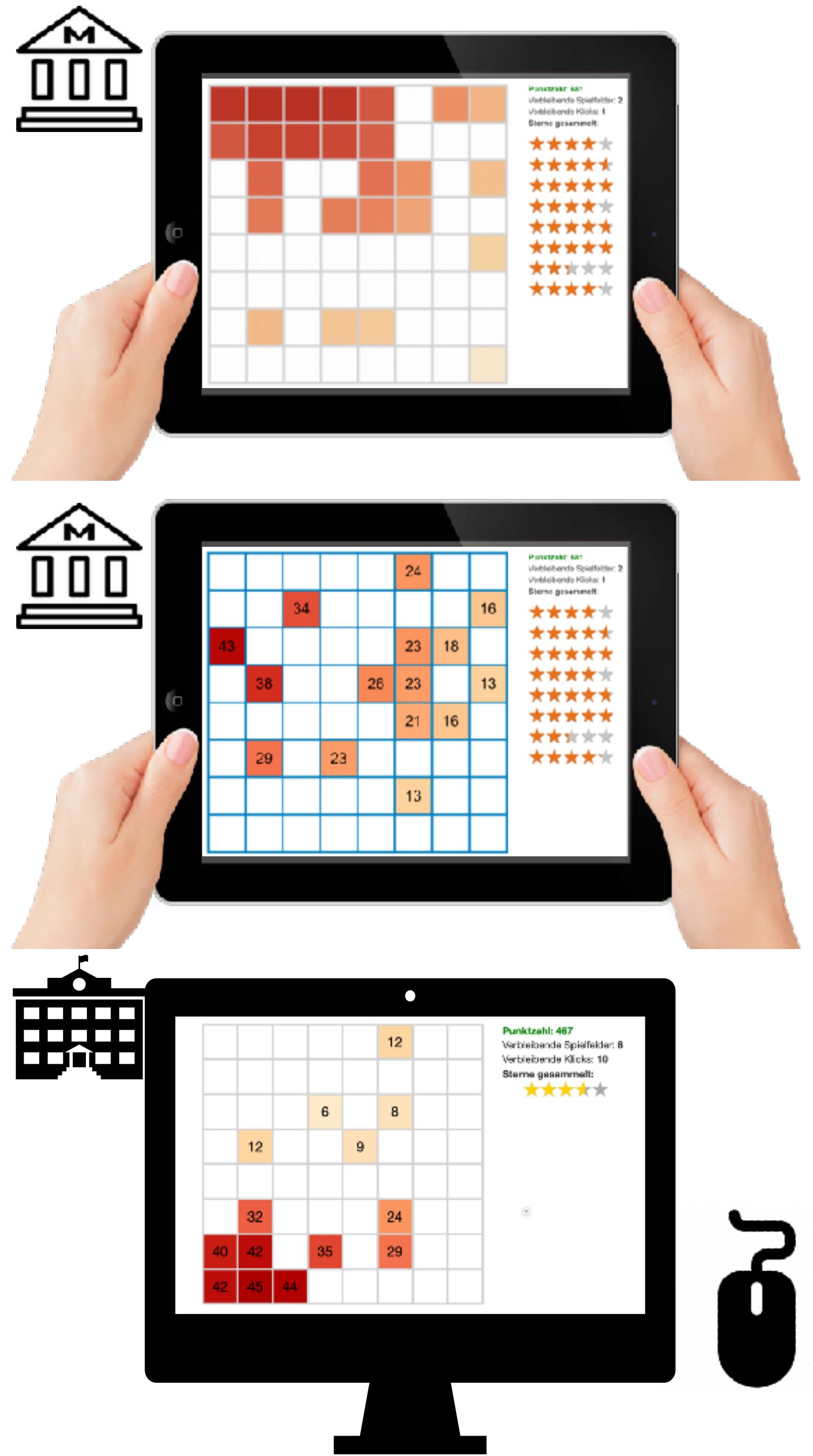
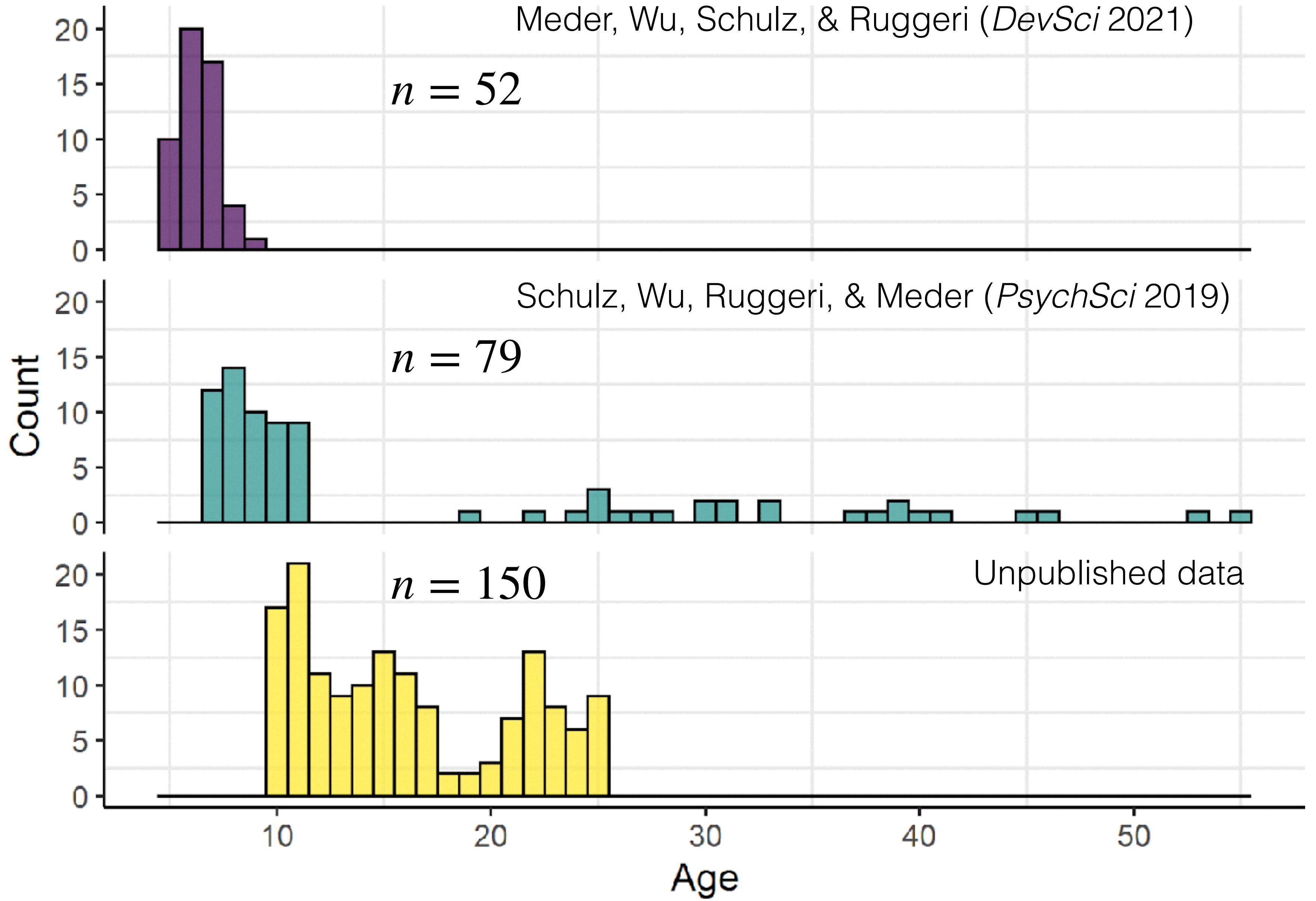
 $\lambda$ 

broader  
generalization

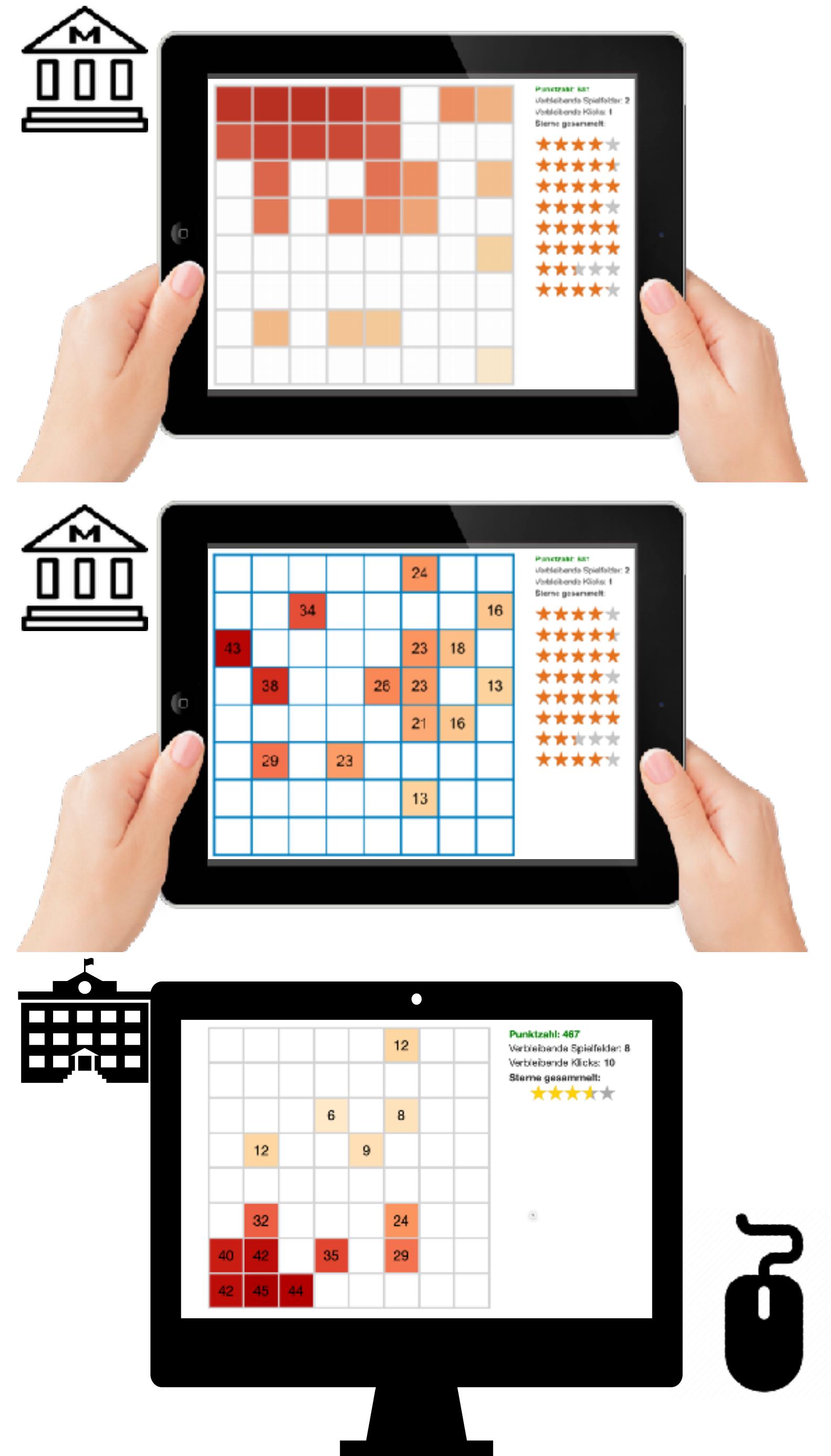
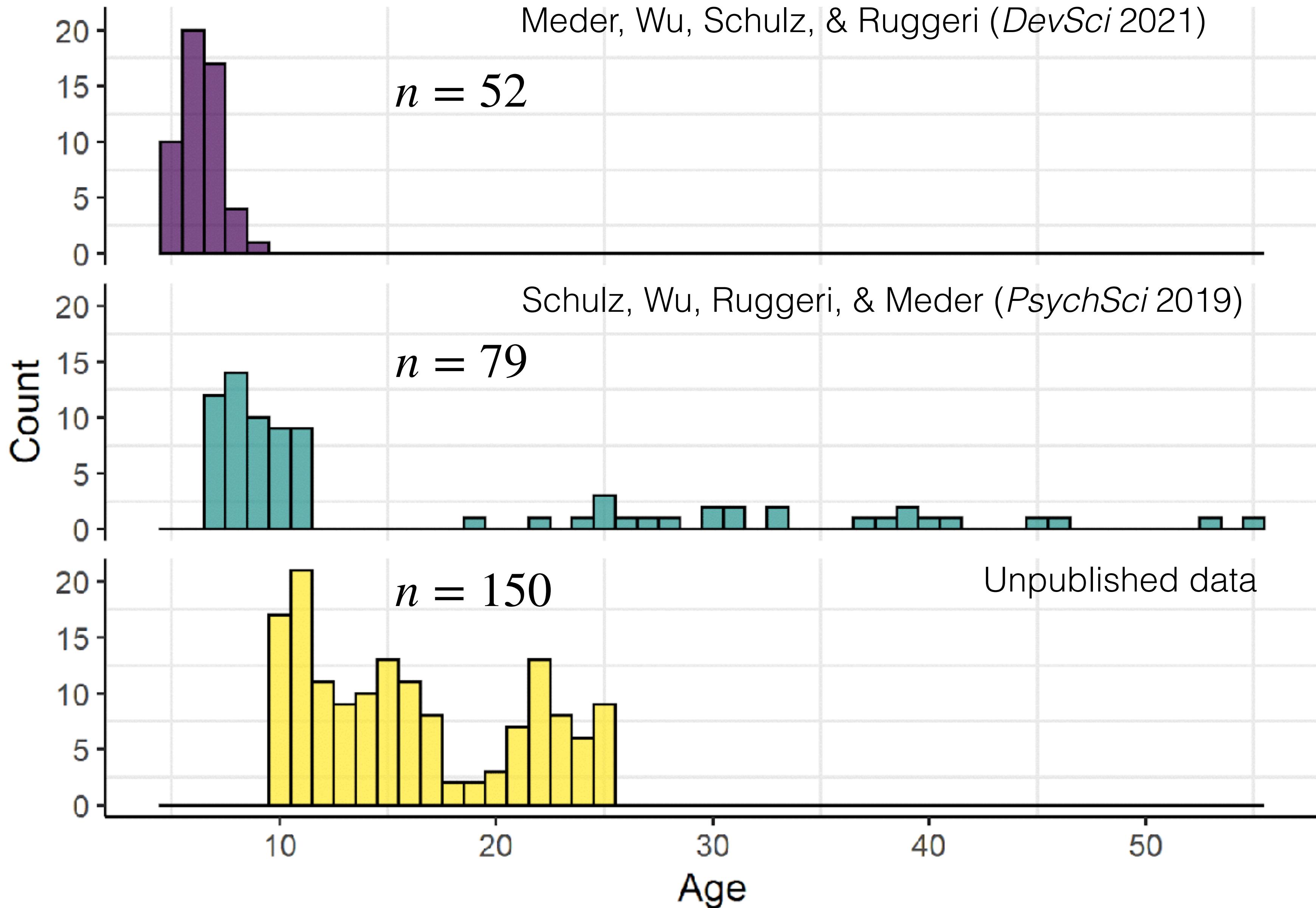


Schulz, Wu, Ruggeri, & Meder (*PsychSci* 2019)



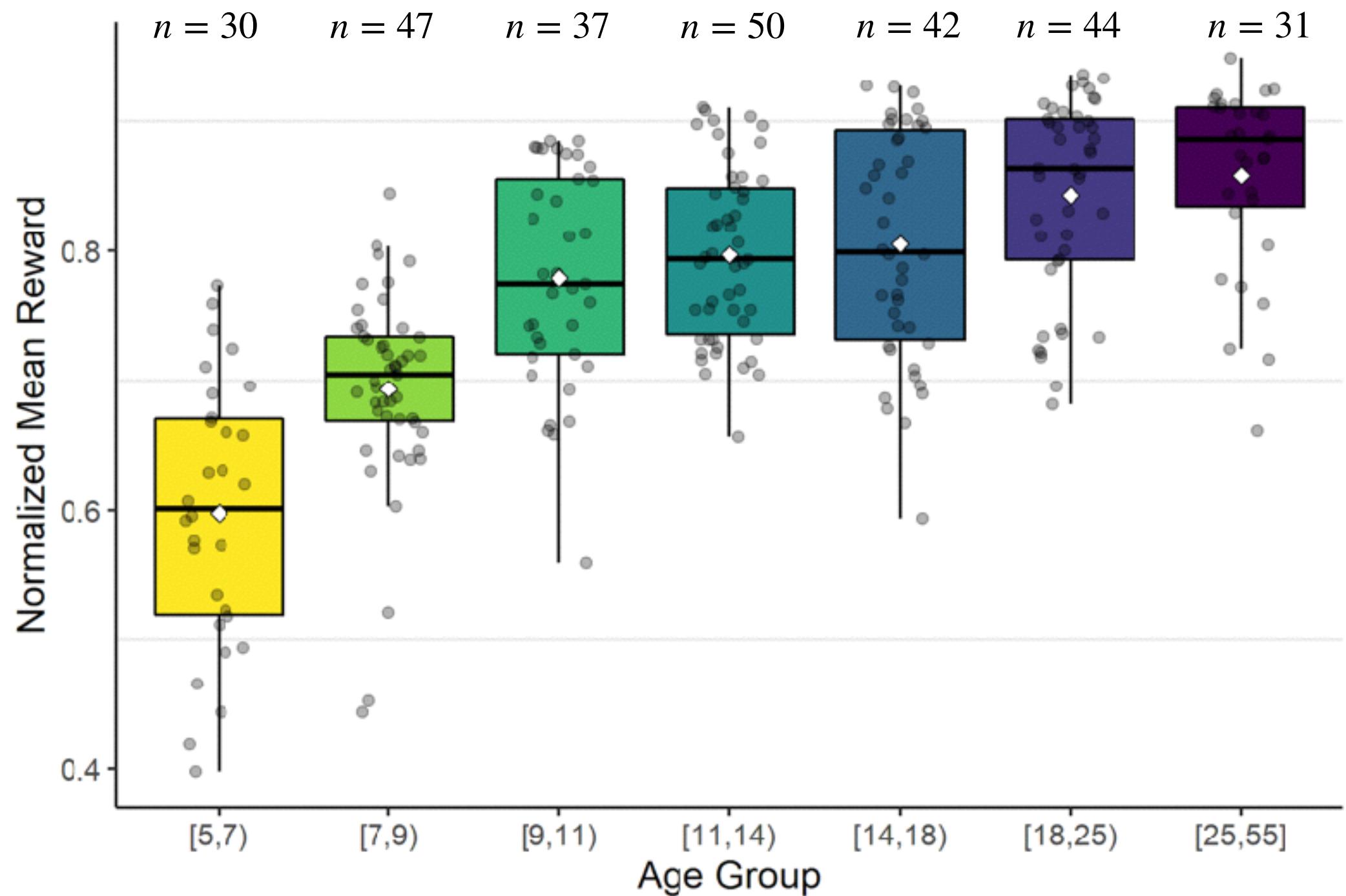


Filtered to use the same sets of environments, same grid size, and same number of trials per round



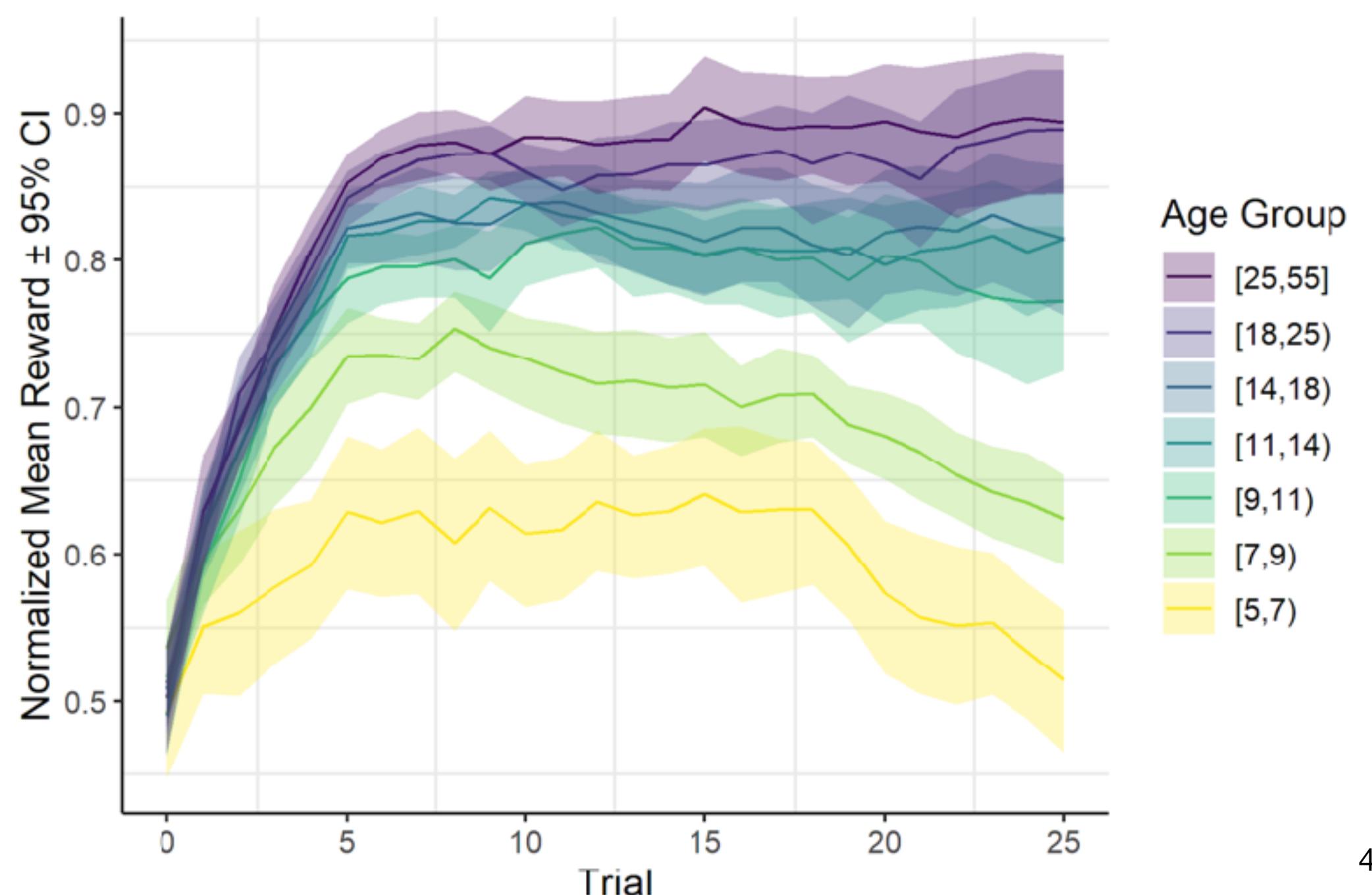
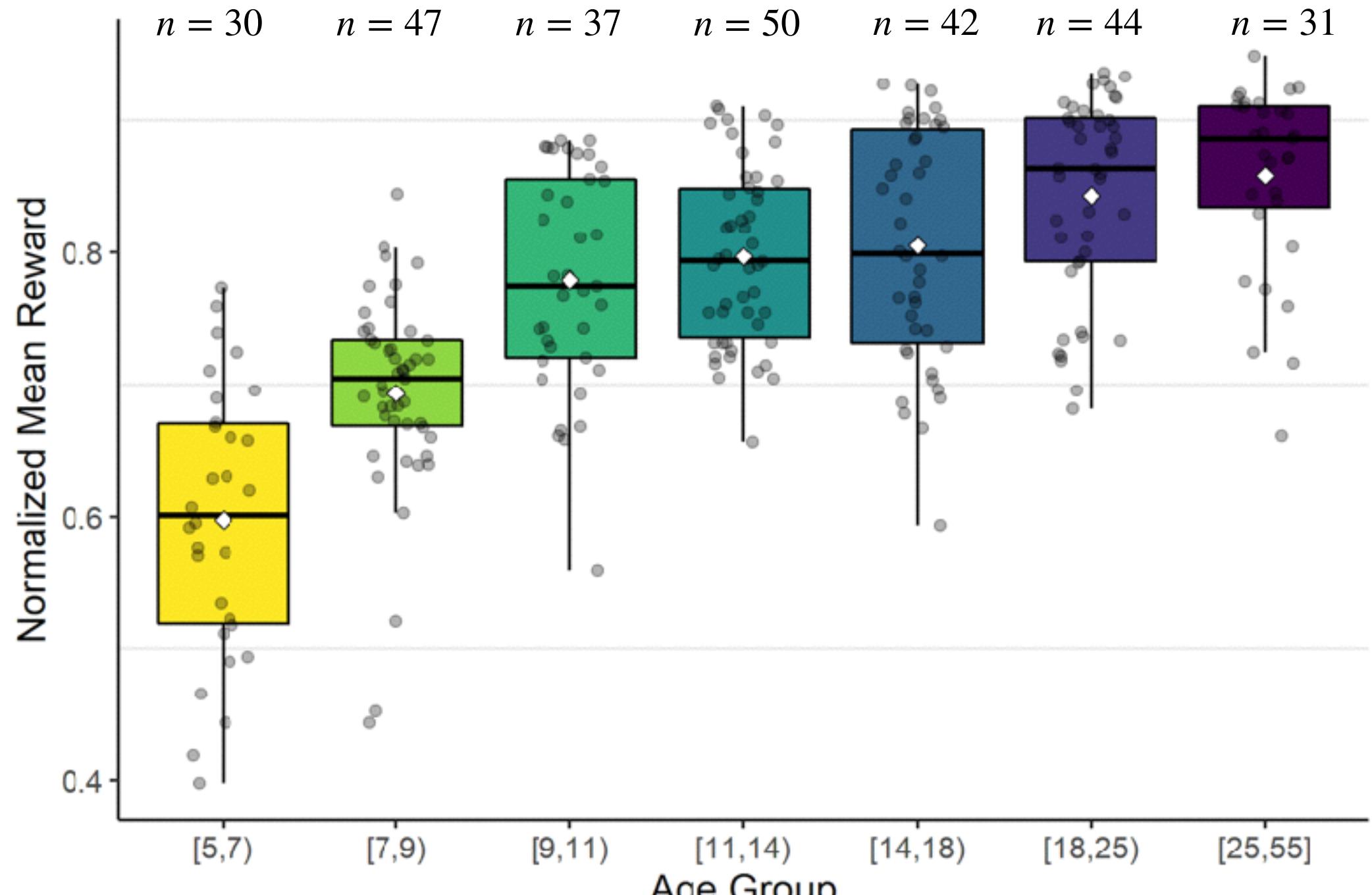
# Behavioral results

- Performance increases over age



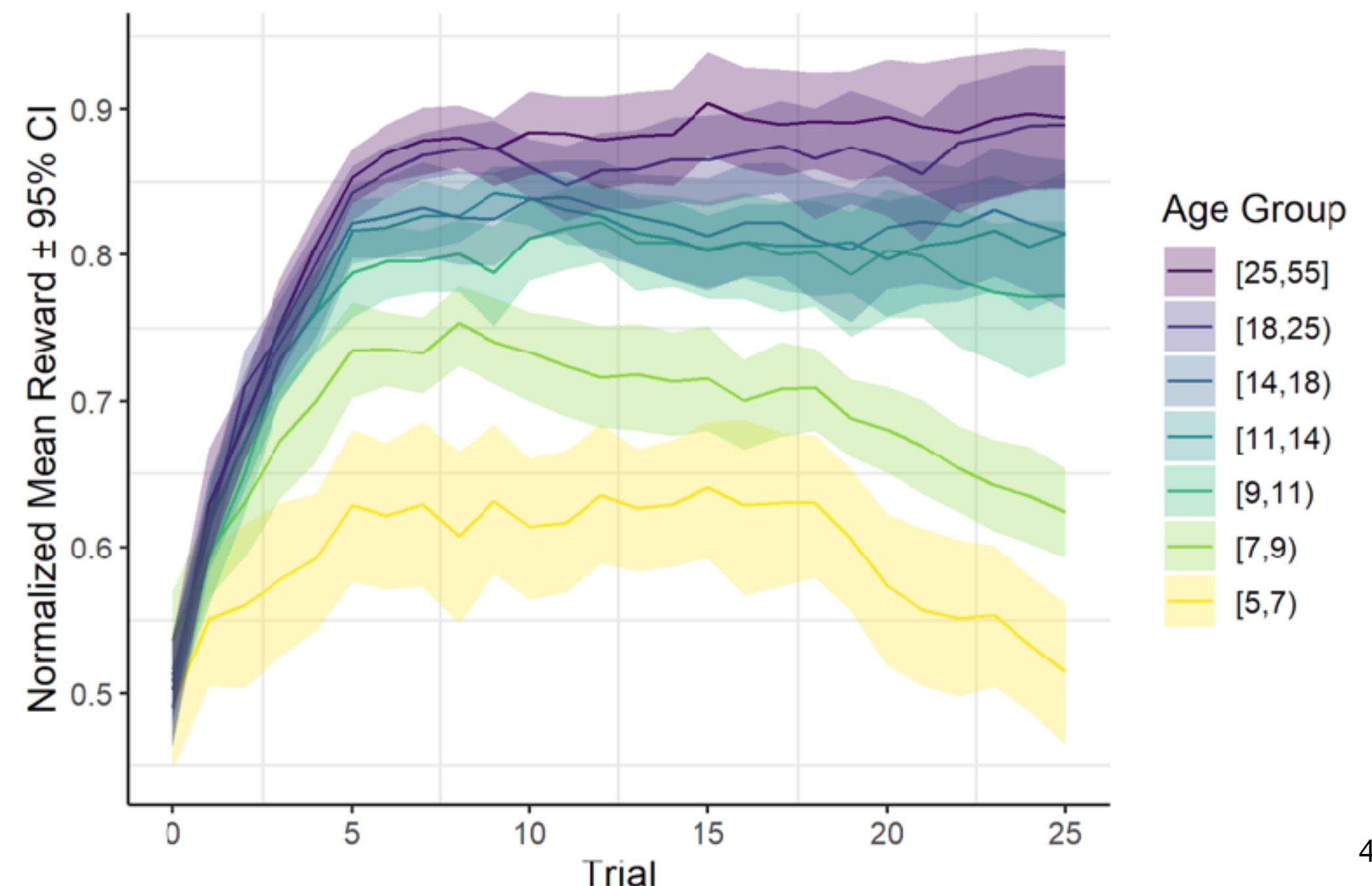
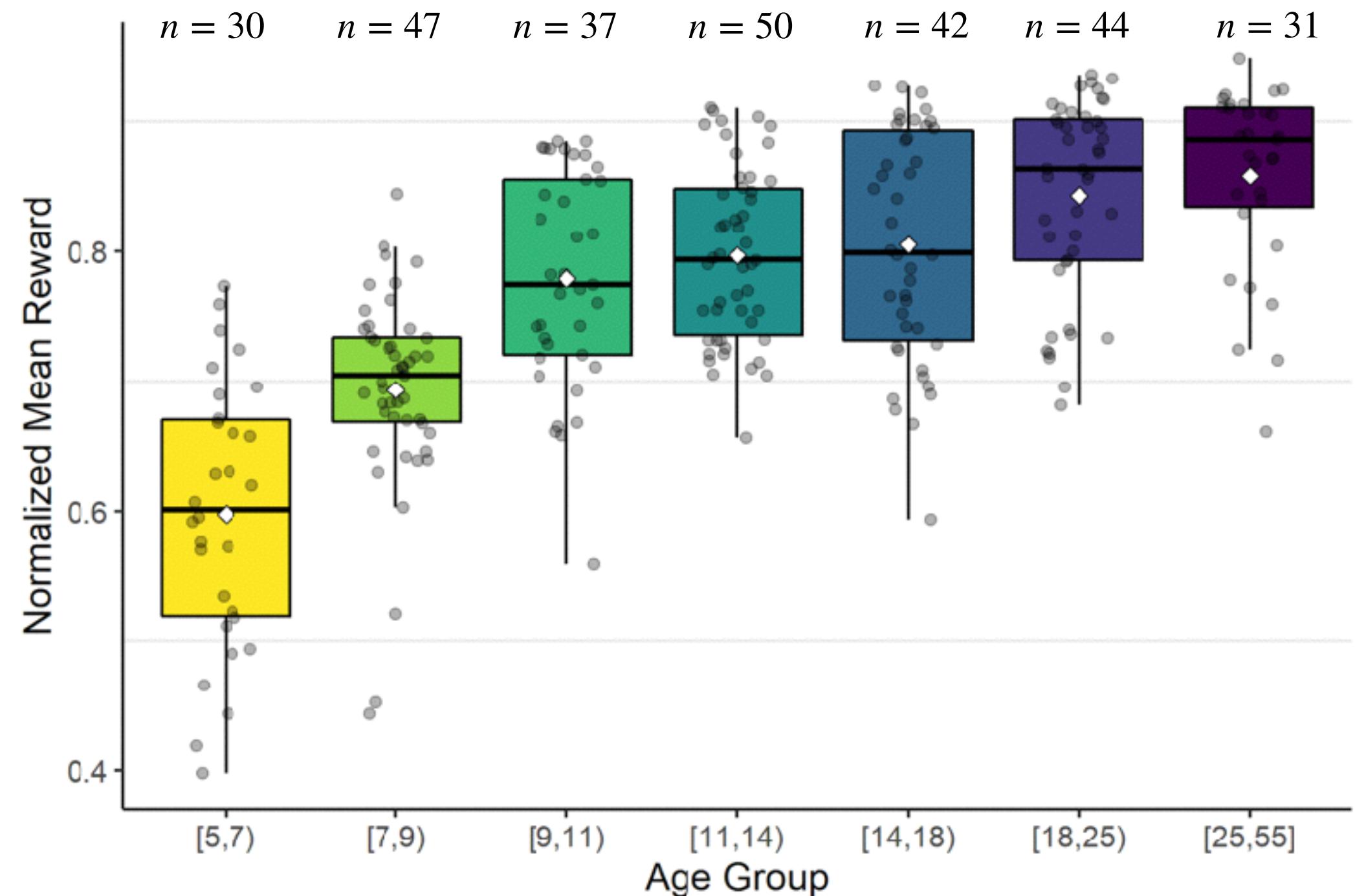
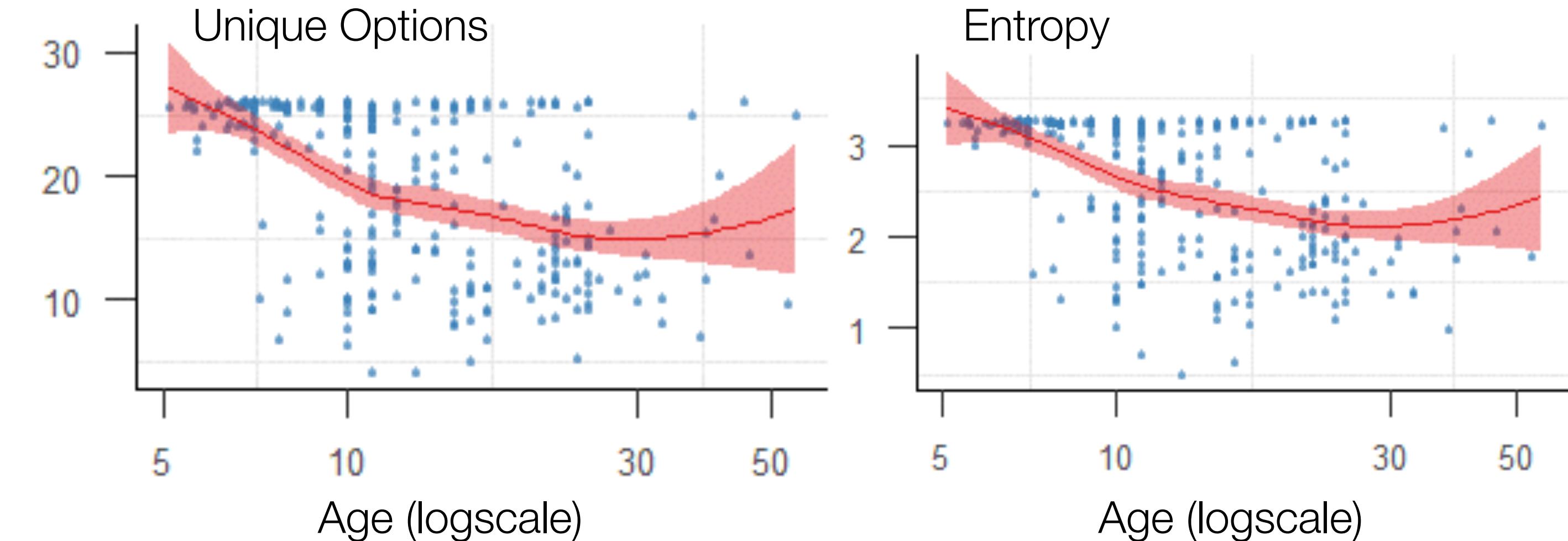
# Behavioral results

- Performance increases over age
- Age-related differences are already evident in the first few trials
  - Older subjects have steeper learning curves
  - Younger children have decaying learning curves, consistent with over-exploration, i.e., more unique options and higher entropy of choices



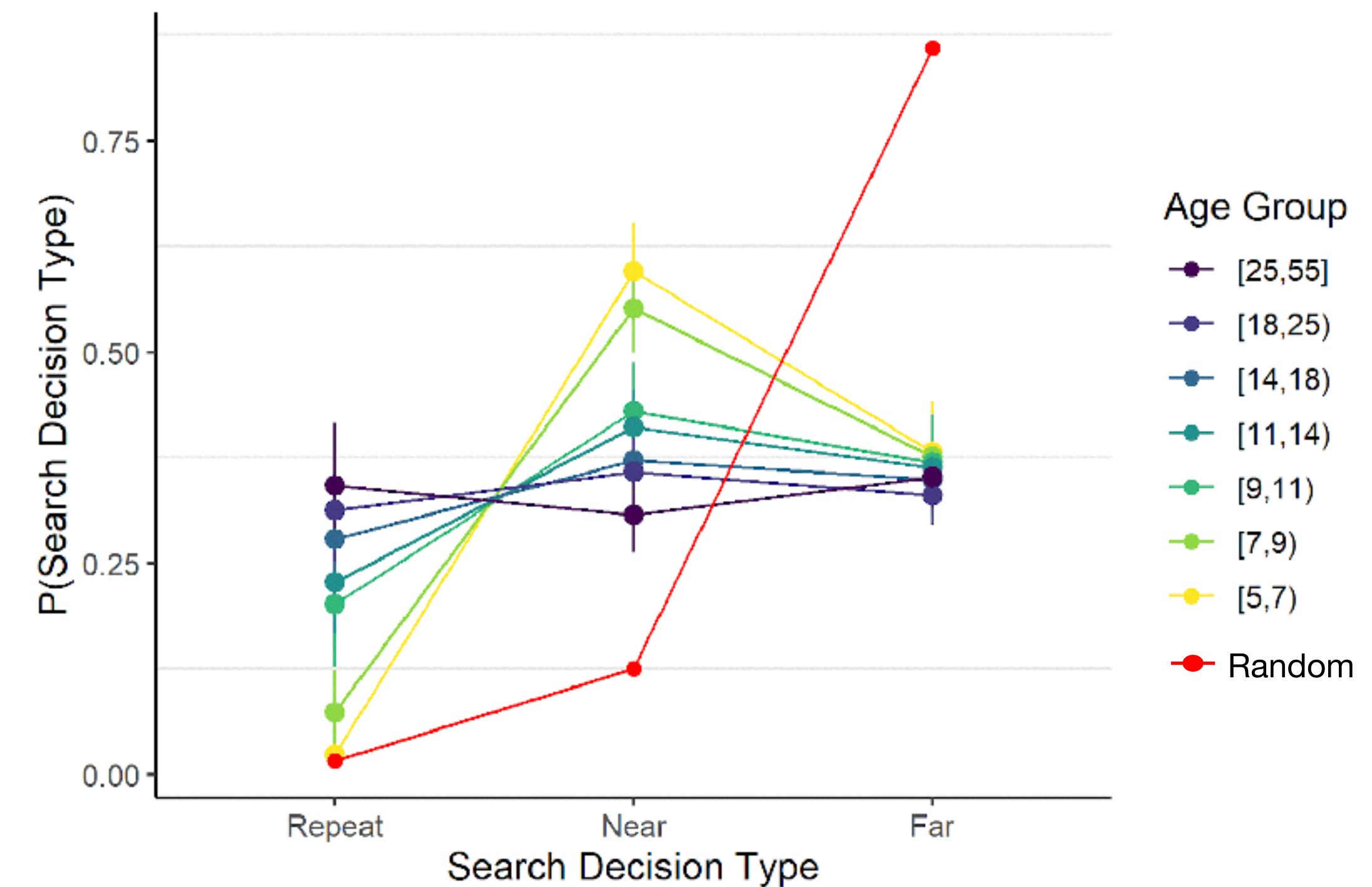
# Behavioral results

- Performance increases over age
- Age-related differences are already evident in the first few trials
  - Older subjects have steeper learning curves
  - Younger children have decaying learning curves, consistent with over-exploration, i.e., more unique options and higher entropy of choices



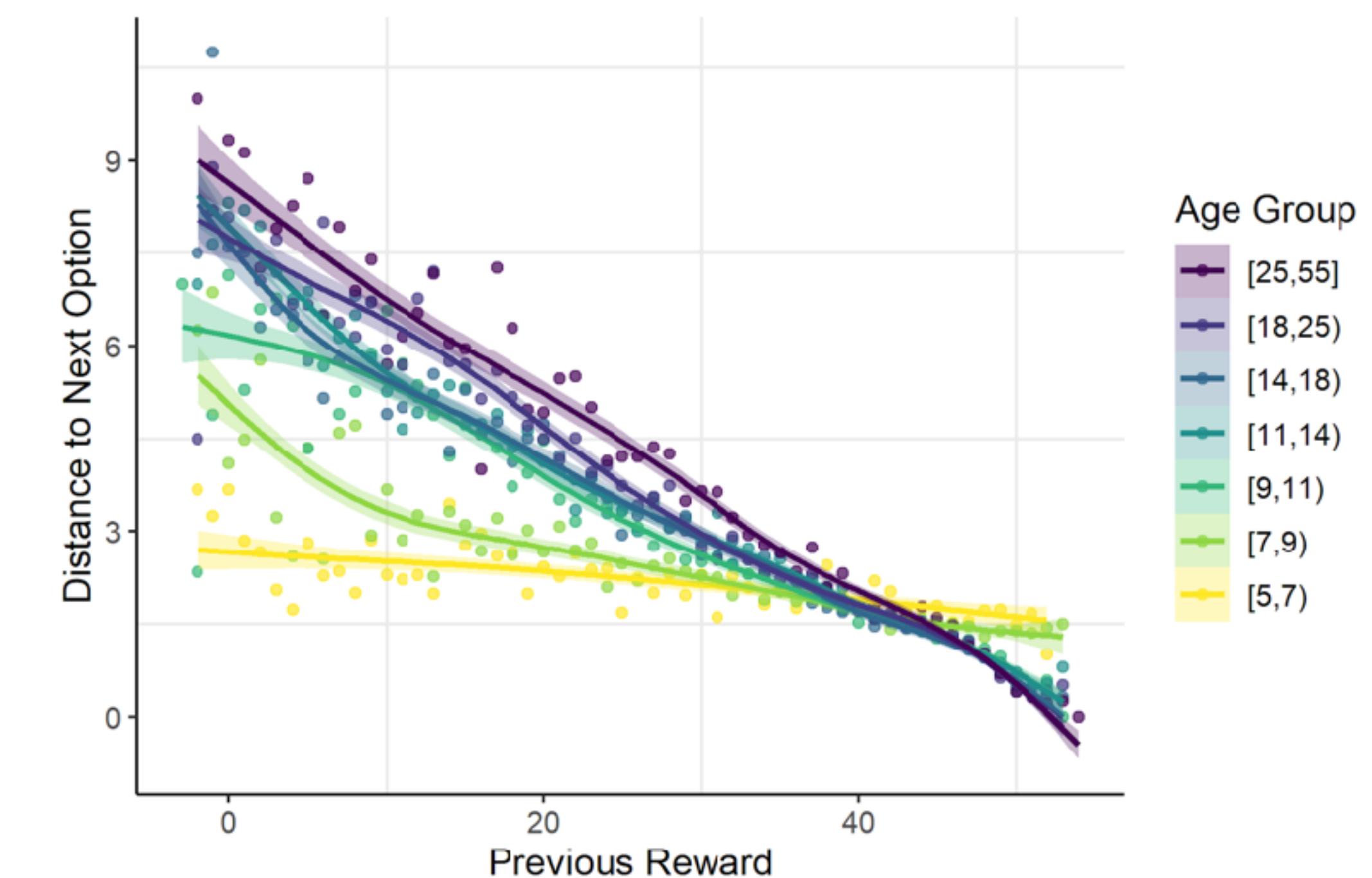
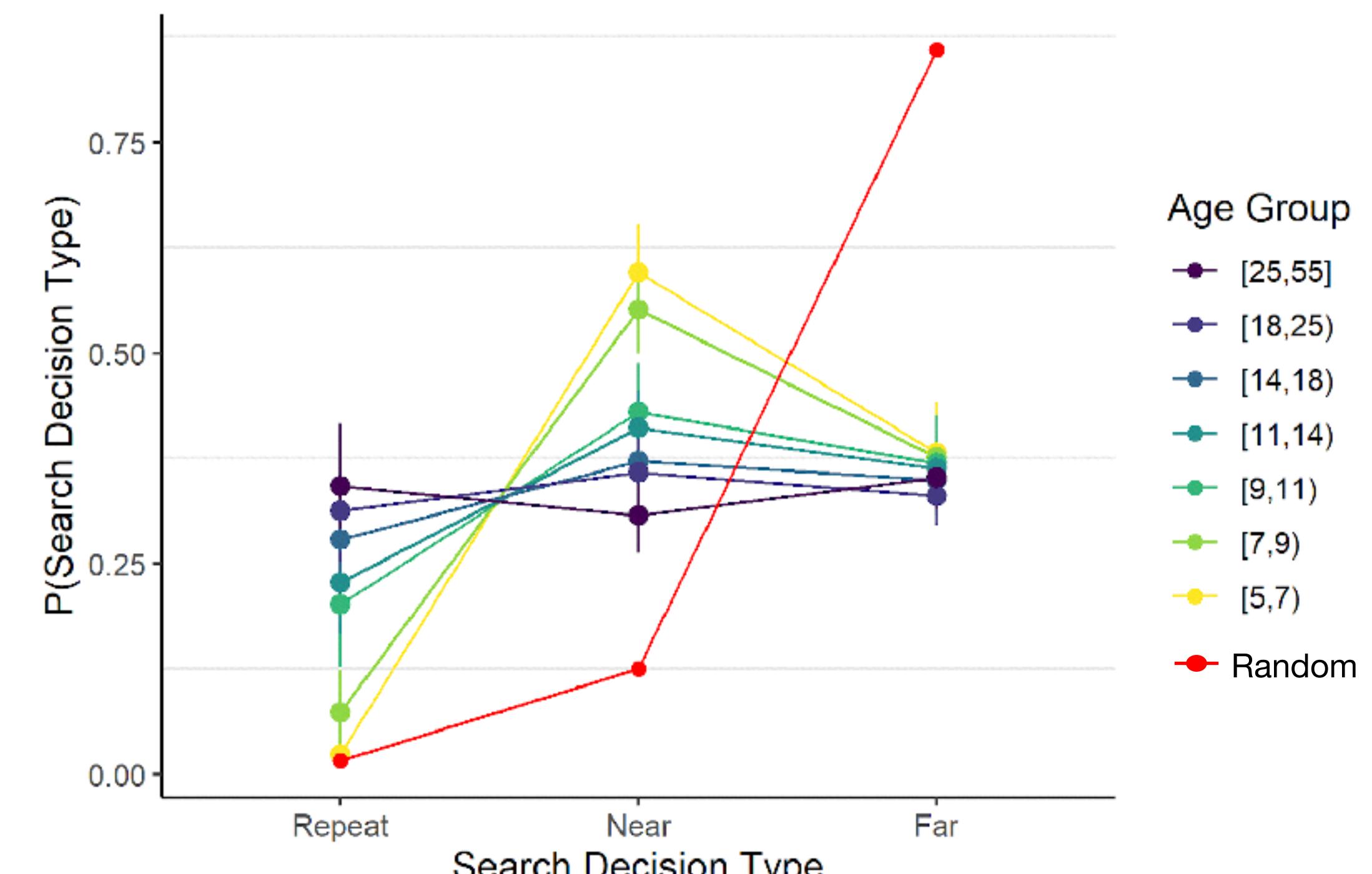
# Behavioral results

- Categorized decisions as either:
  - Repeat (same as last choice)
  - Near (neighboring option)
  - Far (any other choice)
- $P(\text{near}) > P(\text{repeat})$  for younger children, but reaches parity in adults



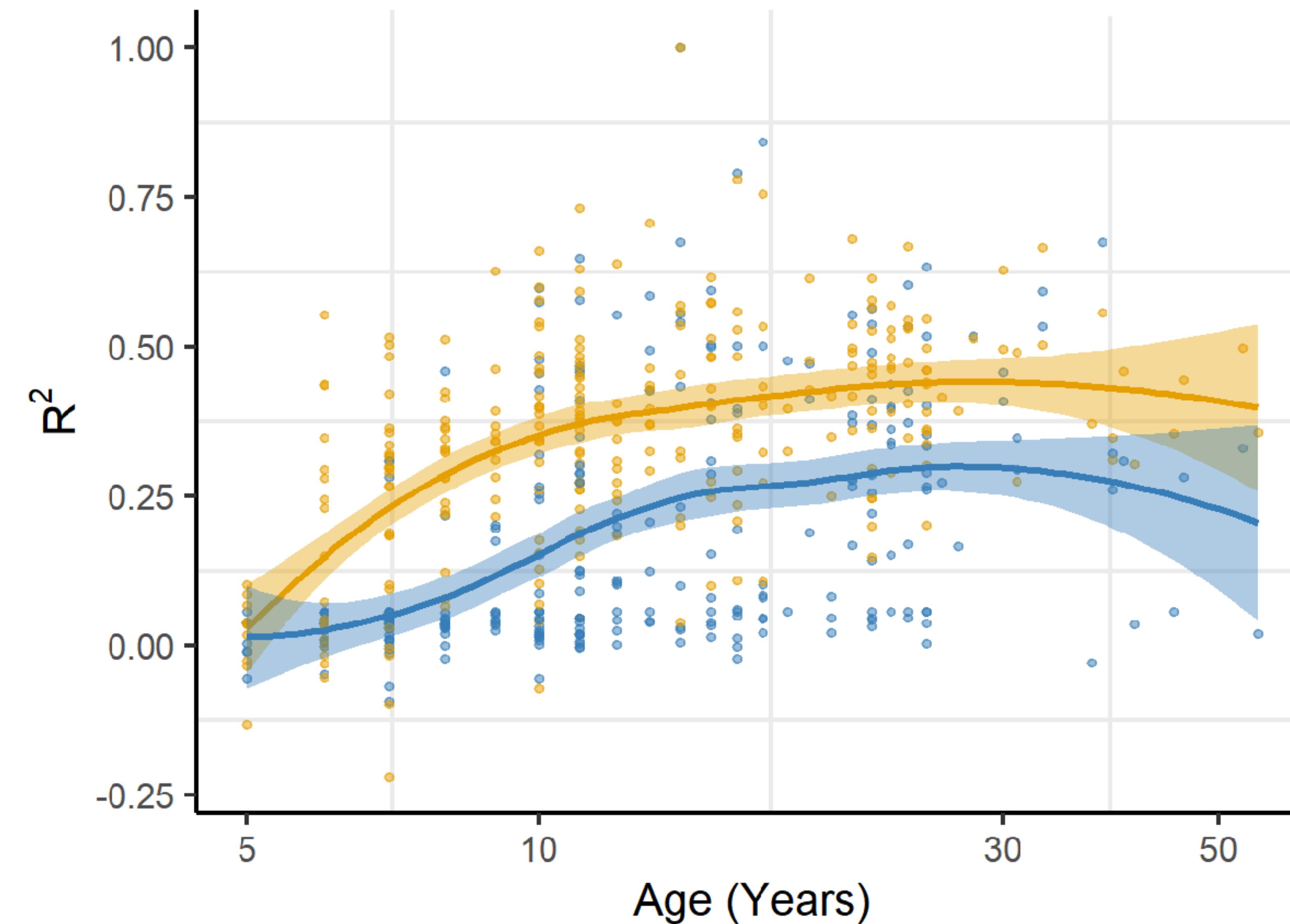
# Behavioral results

- Categorized decisions as either:
  - Repeat (same as last choice)
  - Near (neighboring option)
  - Far (any other choice)
- $P(\text{near}) > P(\text{repeat})$  for younger children, but reaches parity in adults
- Younger children are also less responsive in adapting search distance to reward outcomes
  - Over the lifespan, this develops into a linear relationship, resembling a gradual form of win-stay lose-shift

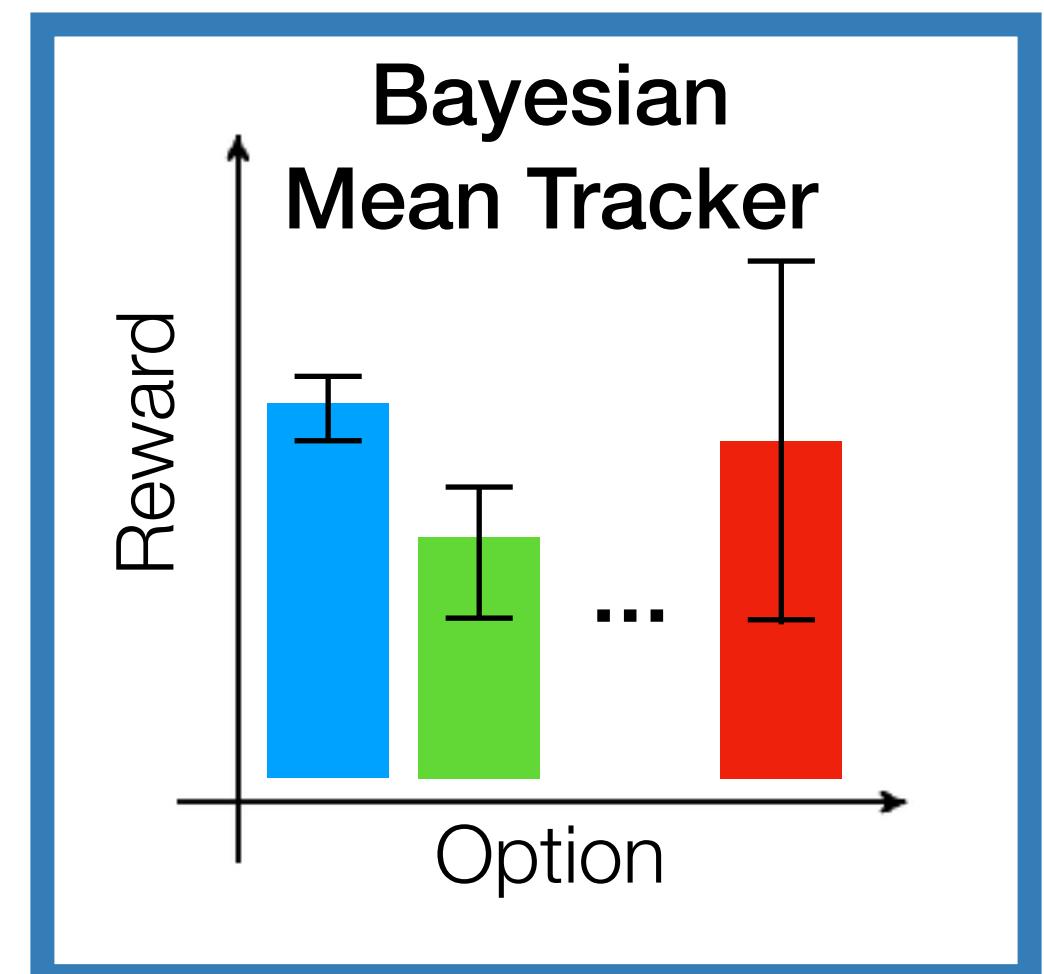
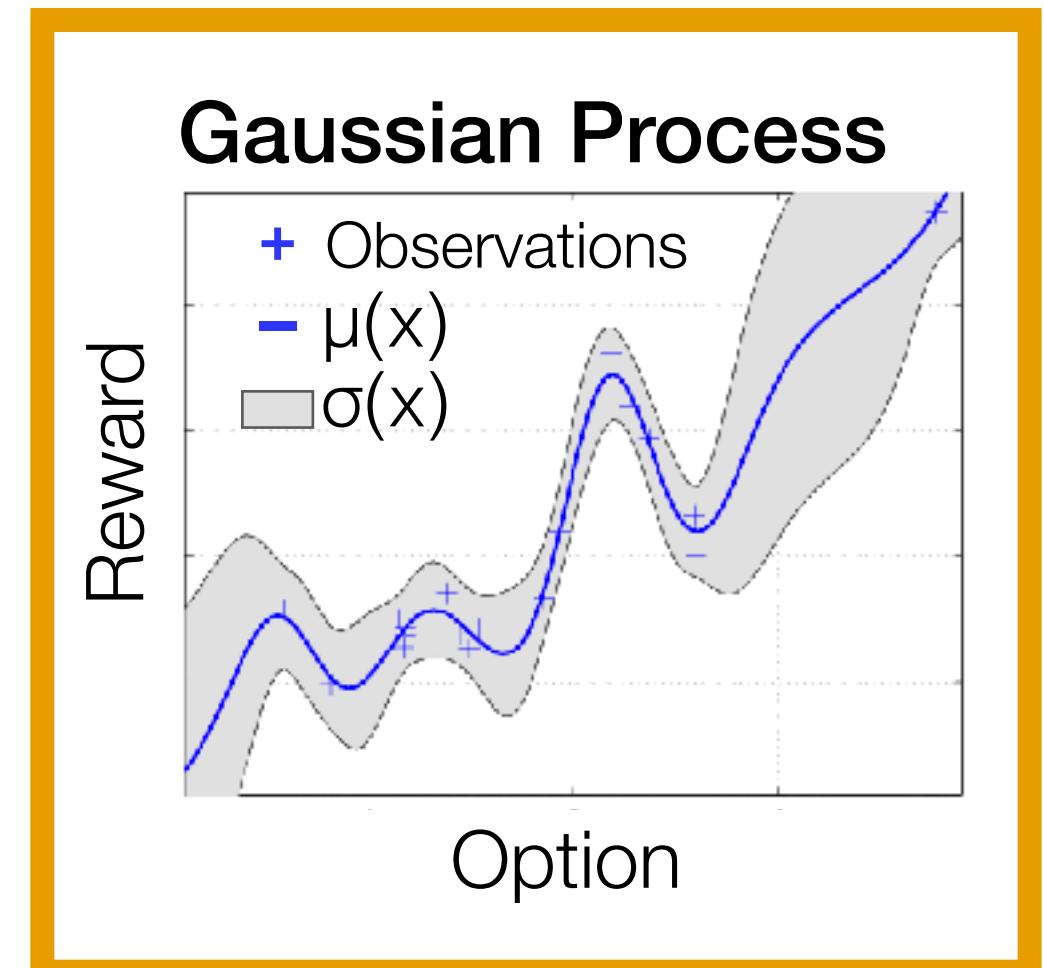


# Model results

$$R^2 = 1 - \frac{\log \mathcal{L}(\mathcal{M}_k)}{\log \mathcal{L}(\mathcal{M}_{\text{rand}})}$$

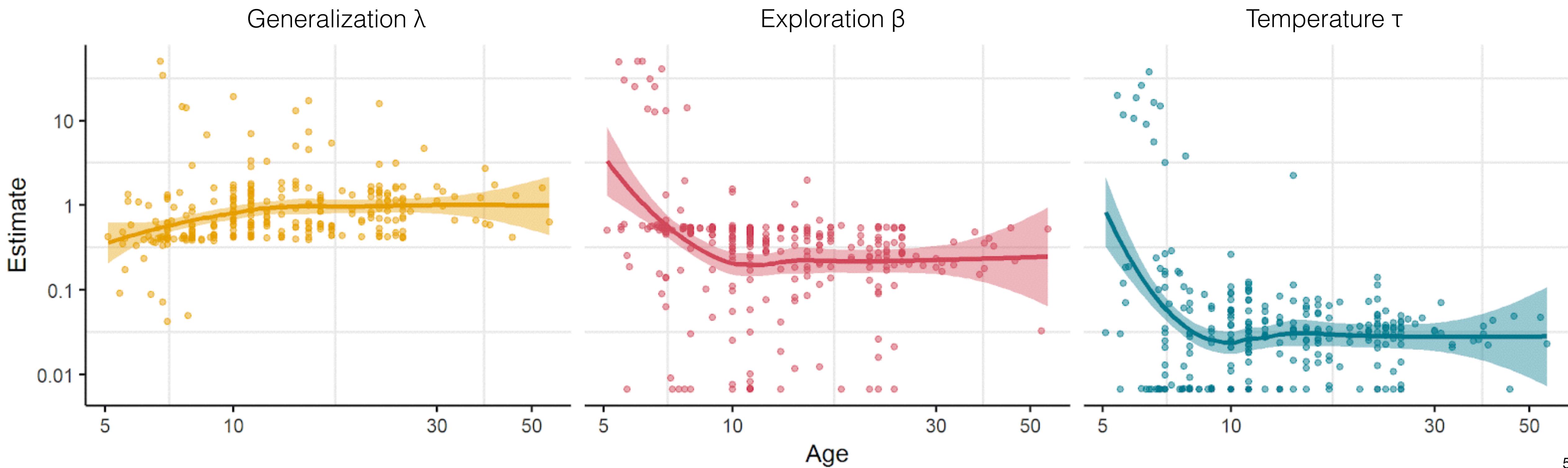


Model  
GP  
BMT

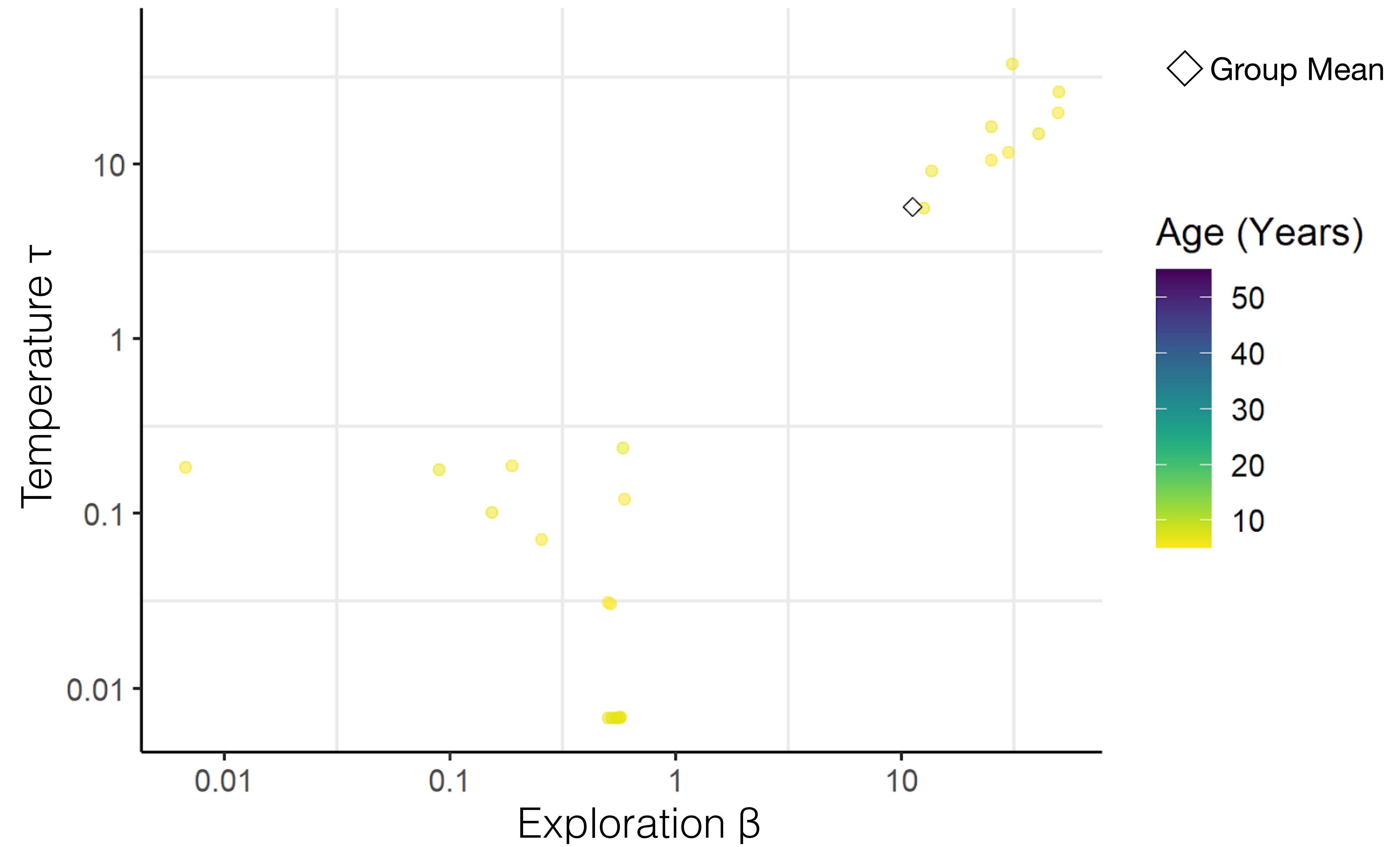


# Parameter estimates

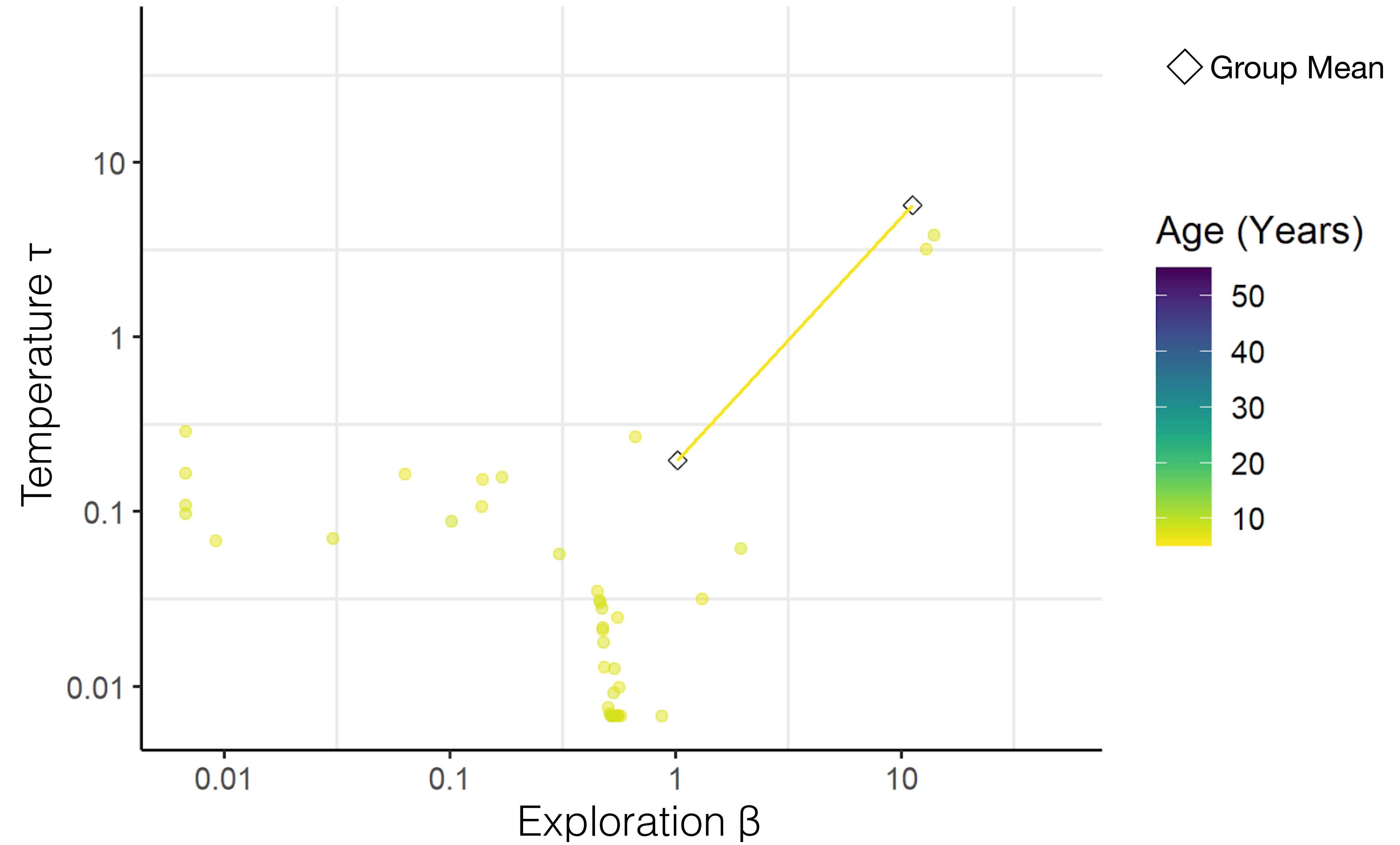
- Small uptick in generalization  $\lambda$
- Large decrease in both uncertainty-directed exploration  $\beta$  and temperature  $\tau$



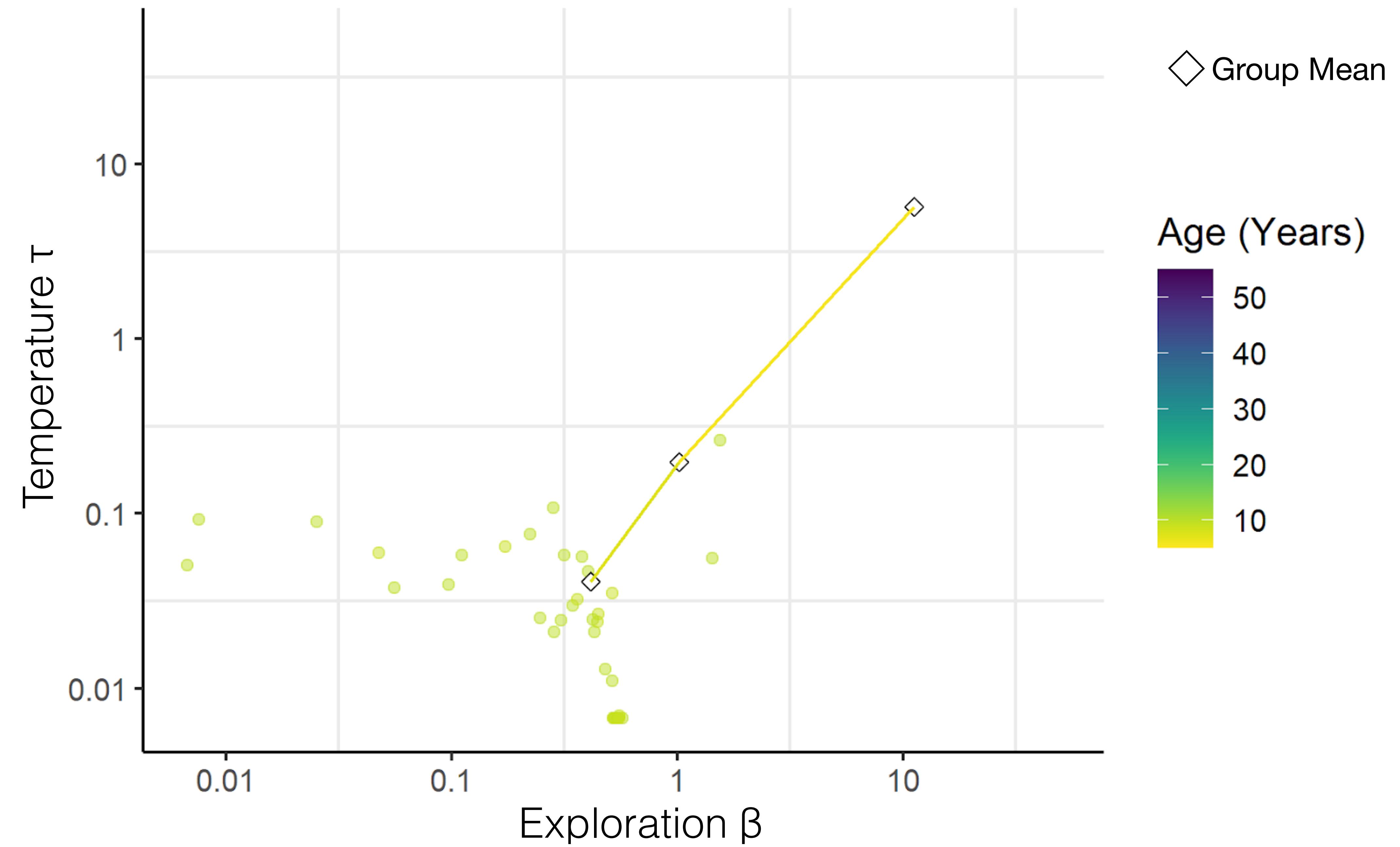
Age Group [5,7)



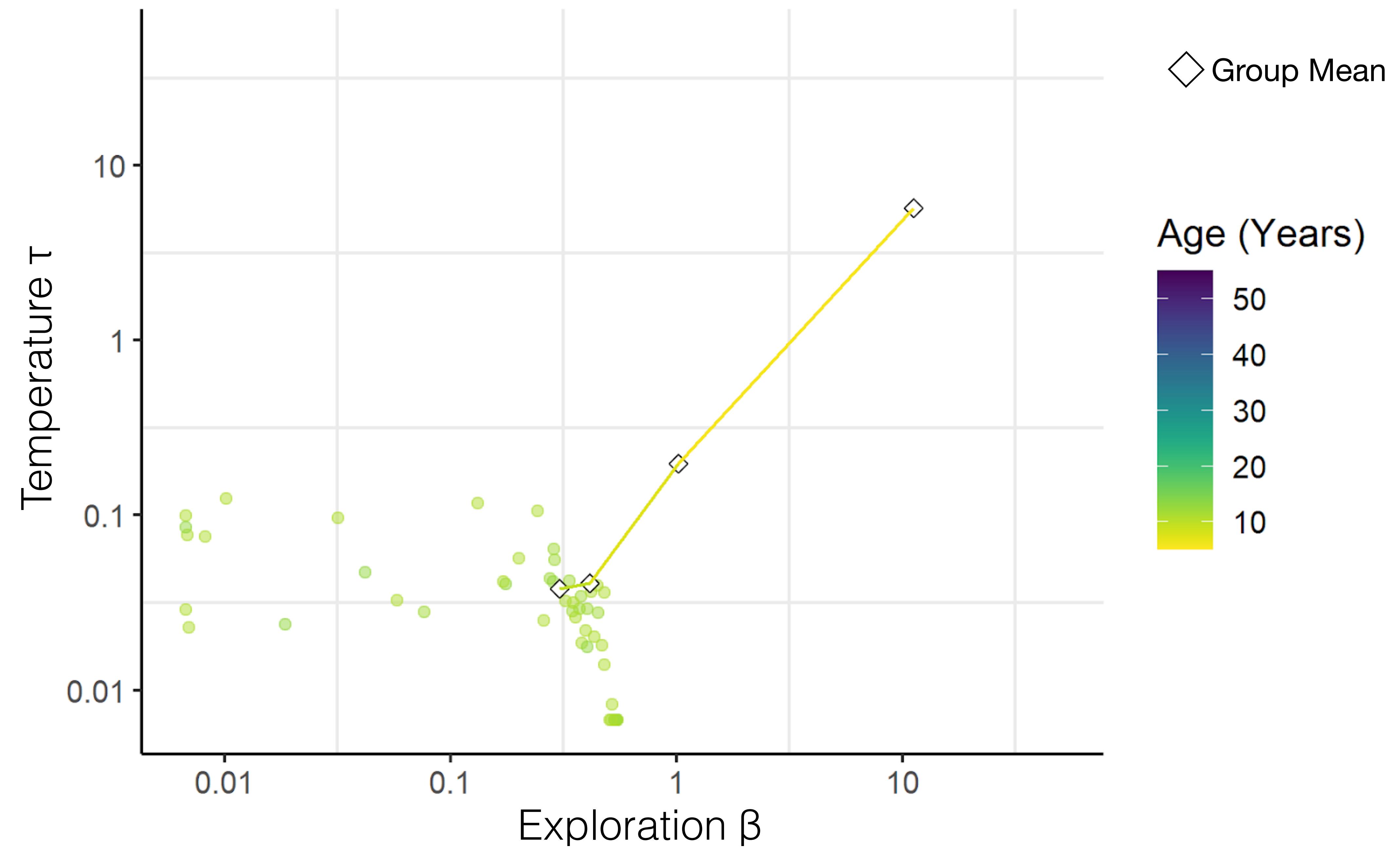
Age Group [7,9)



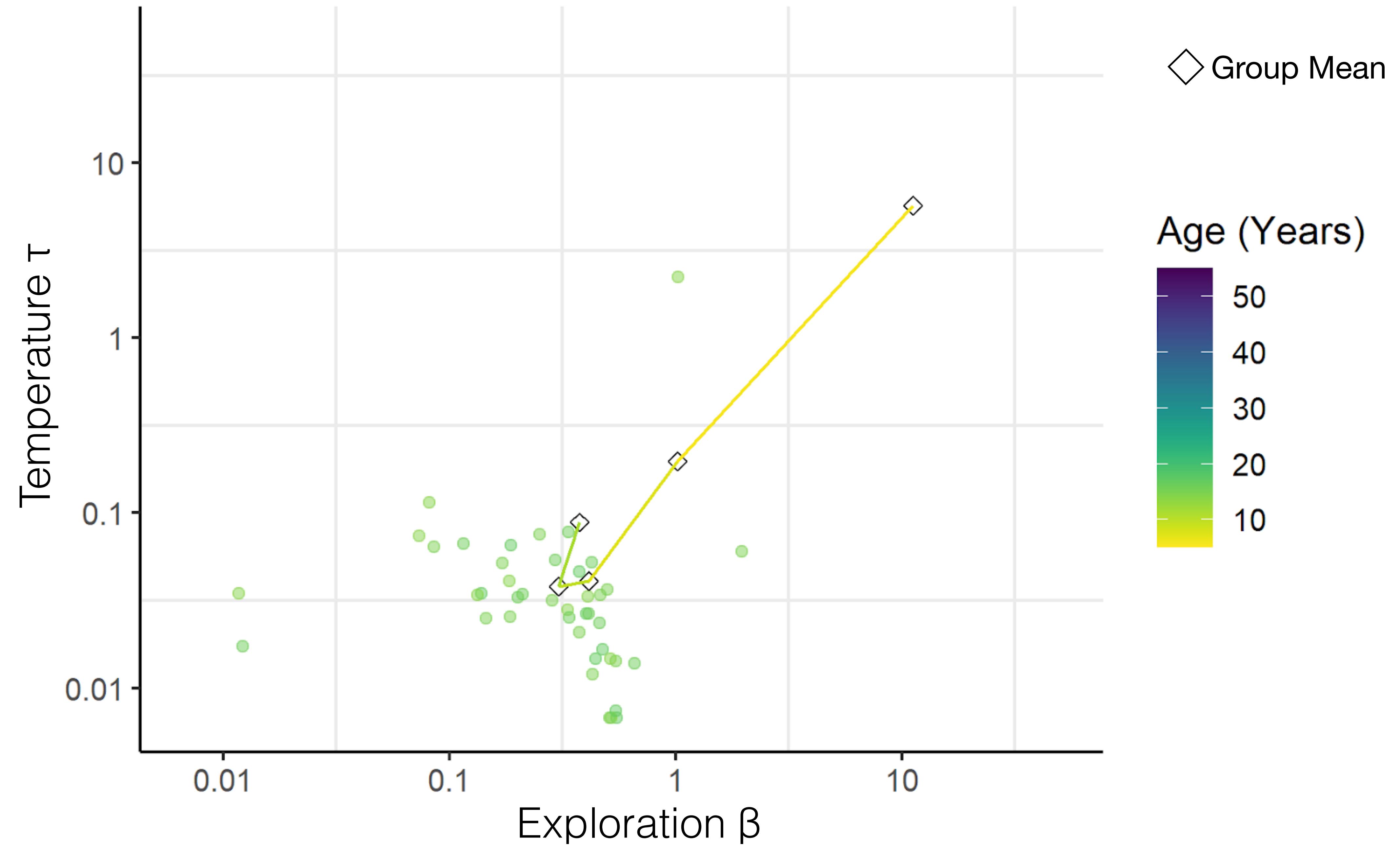
Age Group [9,11)



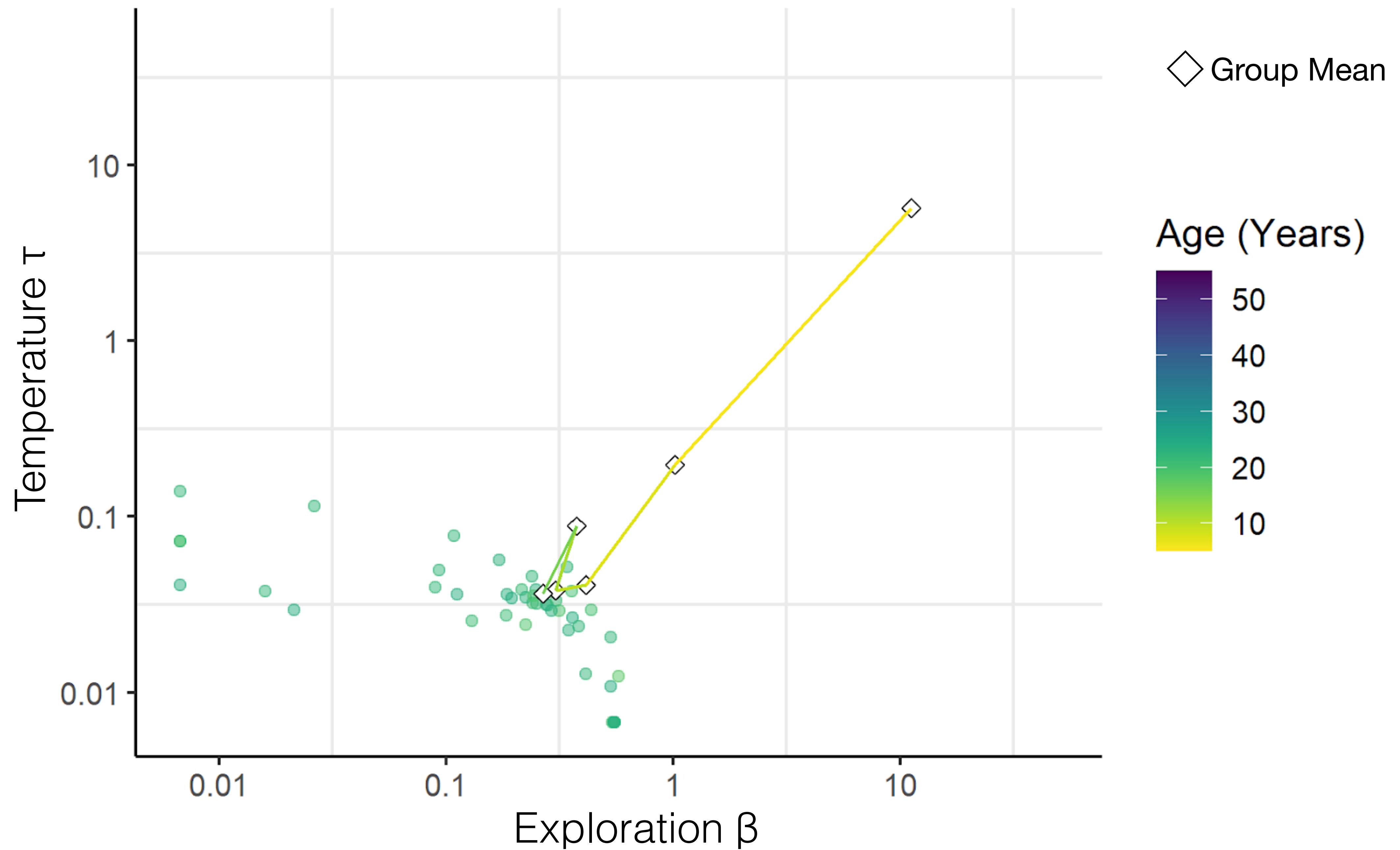
Age Group [11, 14)



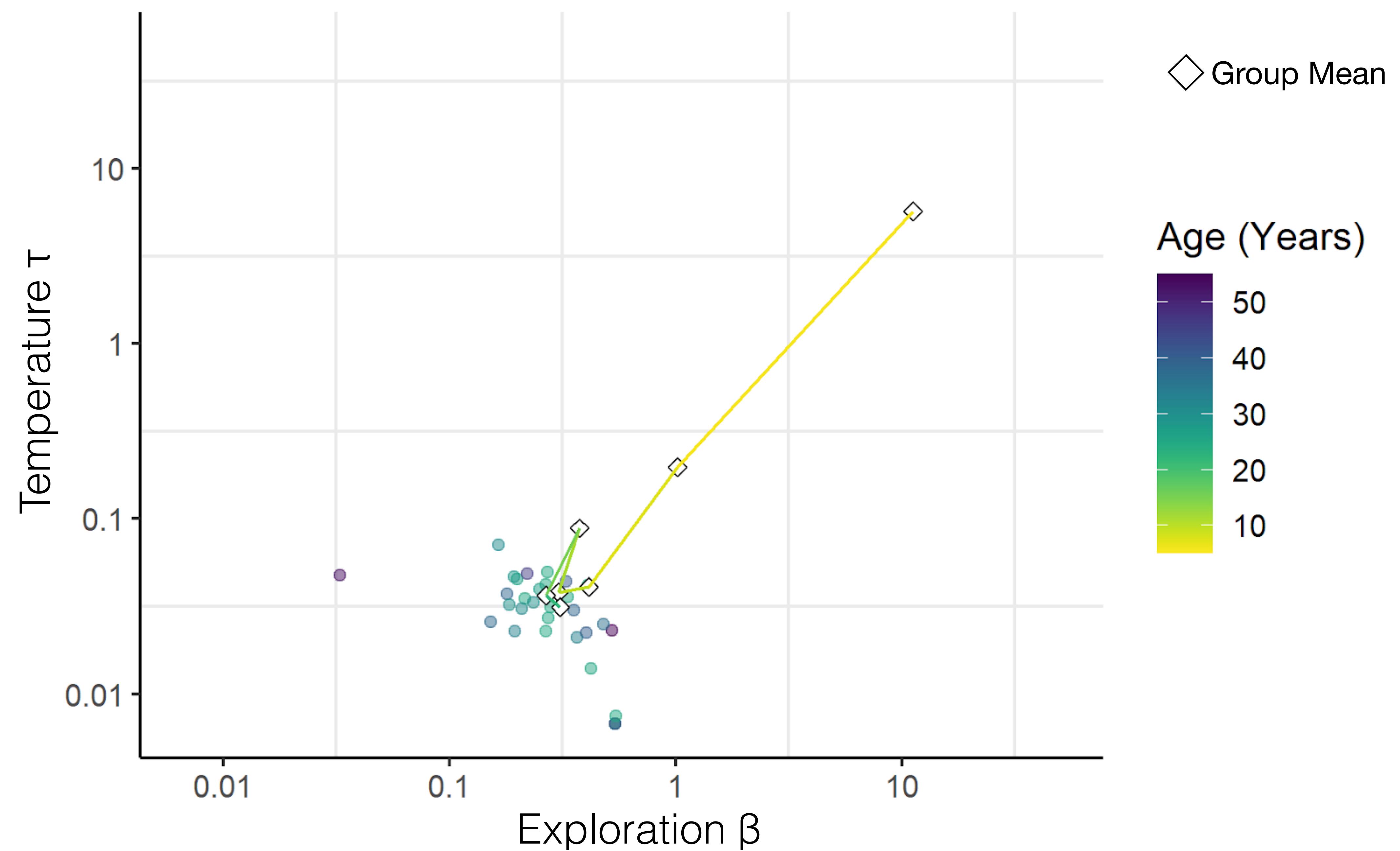
Age Group [14, 18)



# Age Group [18,25)



Age Group [25,55]

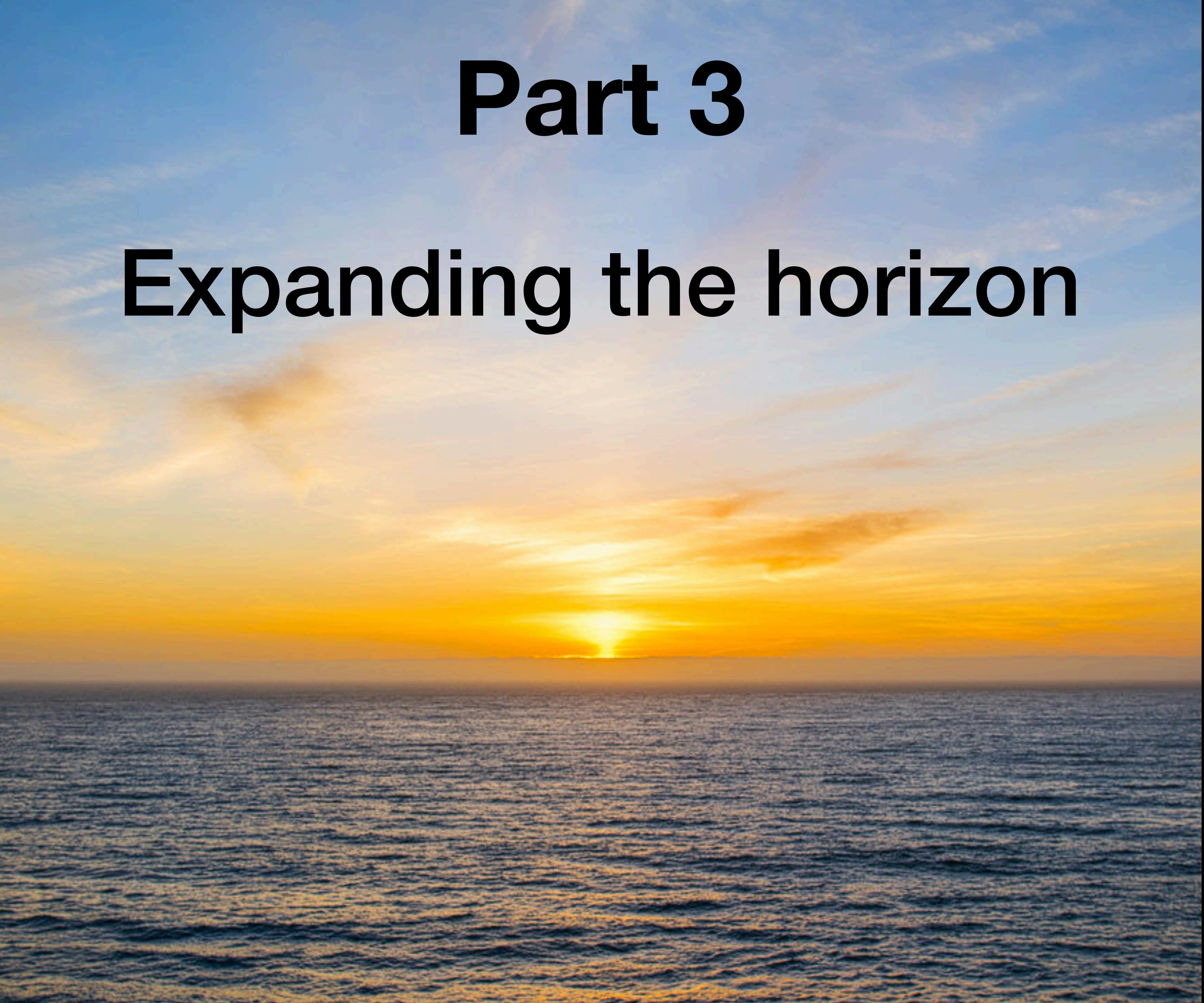


# Summary and Future Directions

- The strategic use of uncertainty-directed exploration and predictive generalization comes online already at a very young age (5-7 year olds)
  - Simultaneous reduction in both directed and random exploration in early childhood
  - Children are not just more random, but also hungrier for information
- While there is an uptick in random exploration during adolescence, this is relatively minor compared to changes in childhood
  - Consistent with theories that increased exploration in adolescence is largely driven by social rather than cognitive factors
- Future work can use model simulations to examine which is the best normative developmental trajectory through model space

# **Part 3**

## **Expanding the horizon**



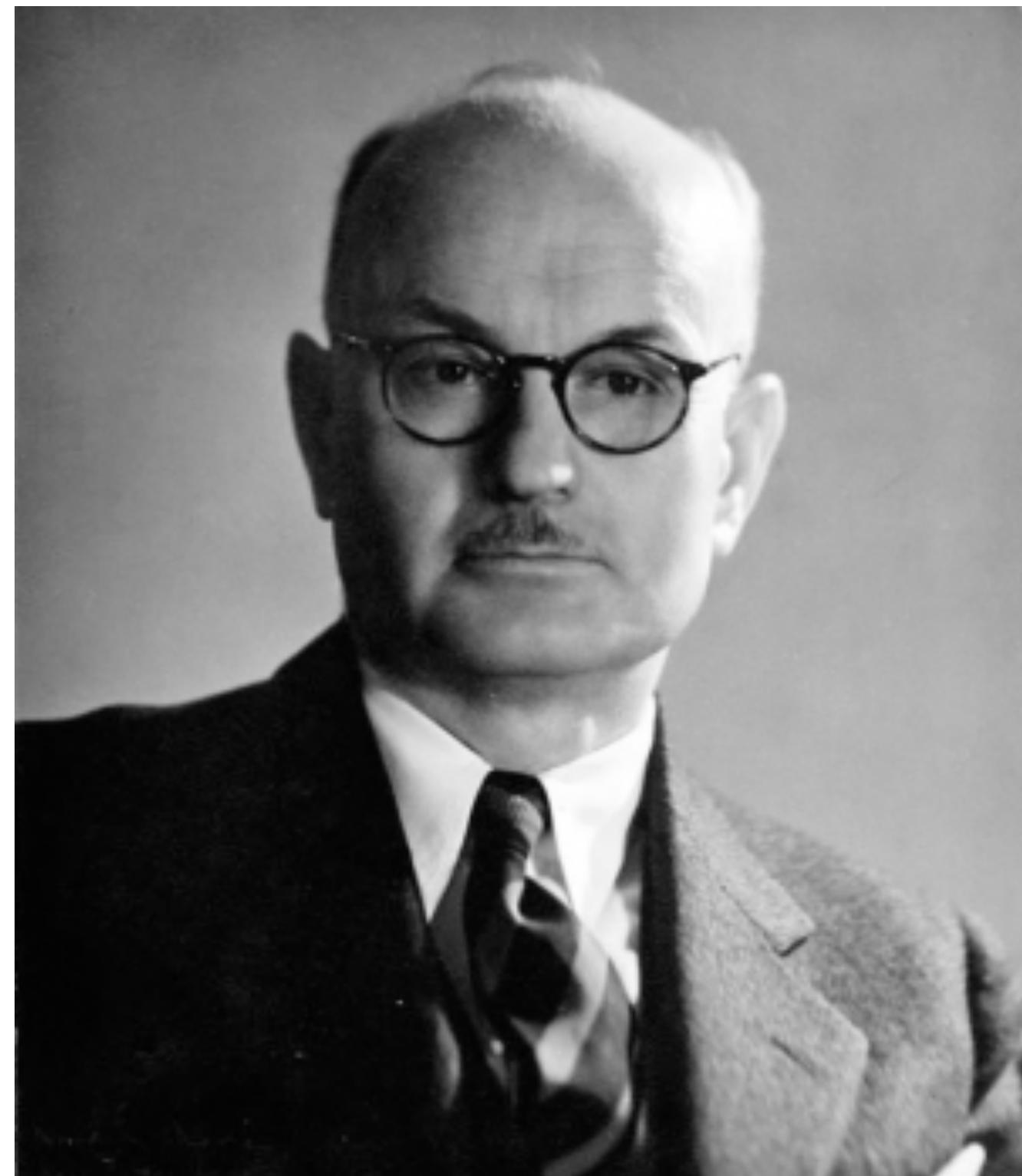
# Part 3

## Expanding the horizon

**Generalization in Conceptual and  
Structured domains**

# Cognitive Maps for Navigation

## Cognitive Maps



Tolman (*Psych Rev*, 1948)

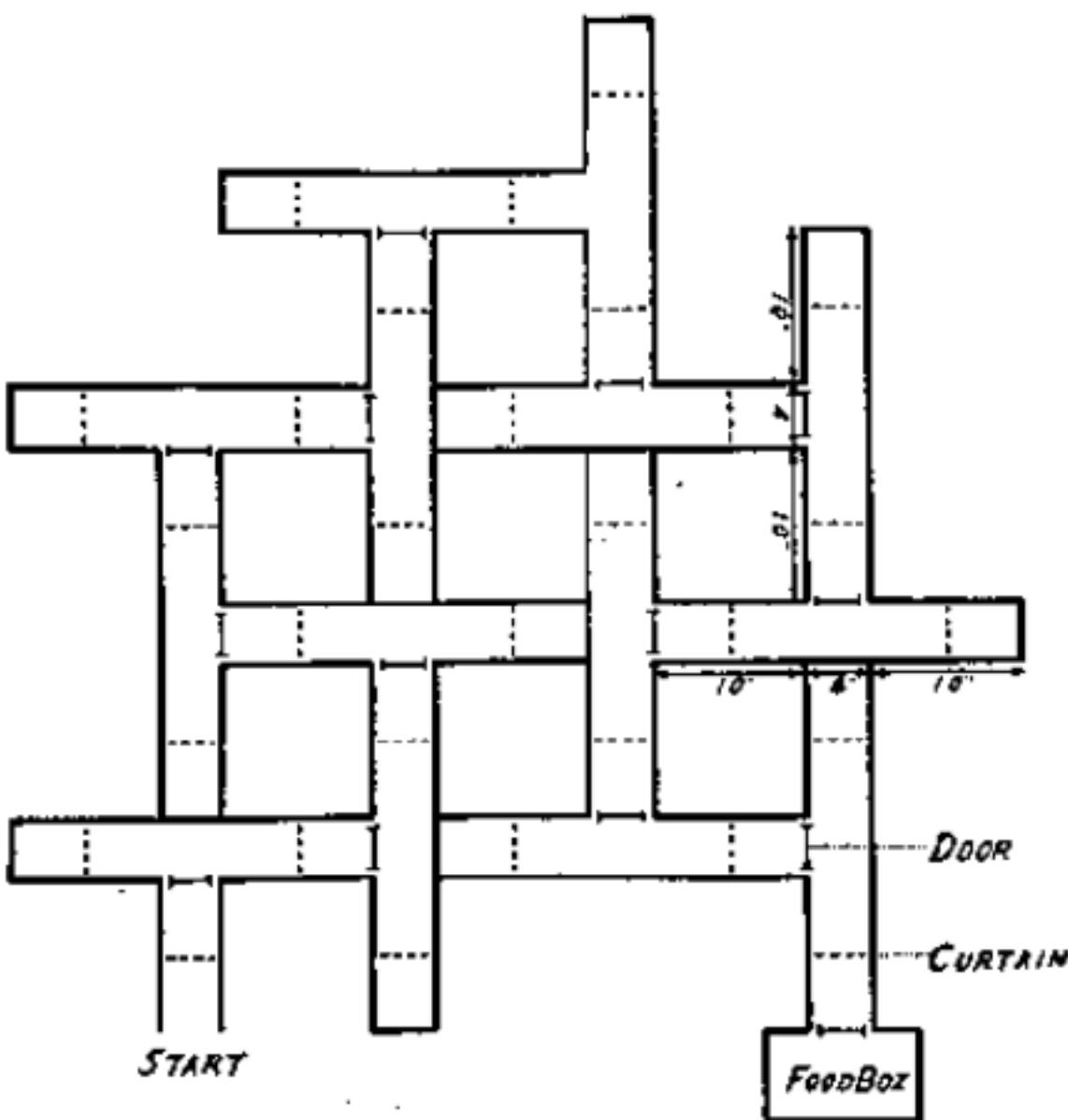


FIG. 1  
(From M. H. Elliott, The effect of change of reward on the maze performance of rats. *Univ. Calif. Publ. Psychol.*, 1928, 4, p. 20.)

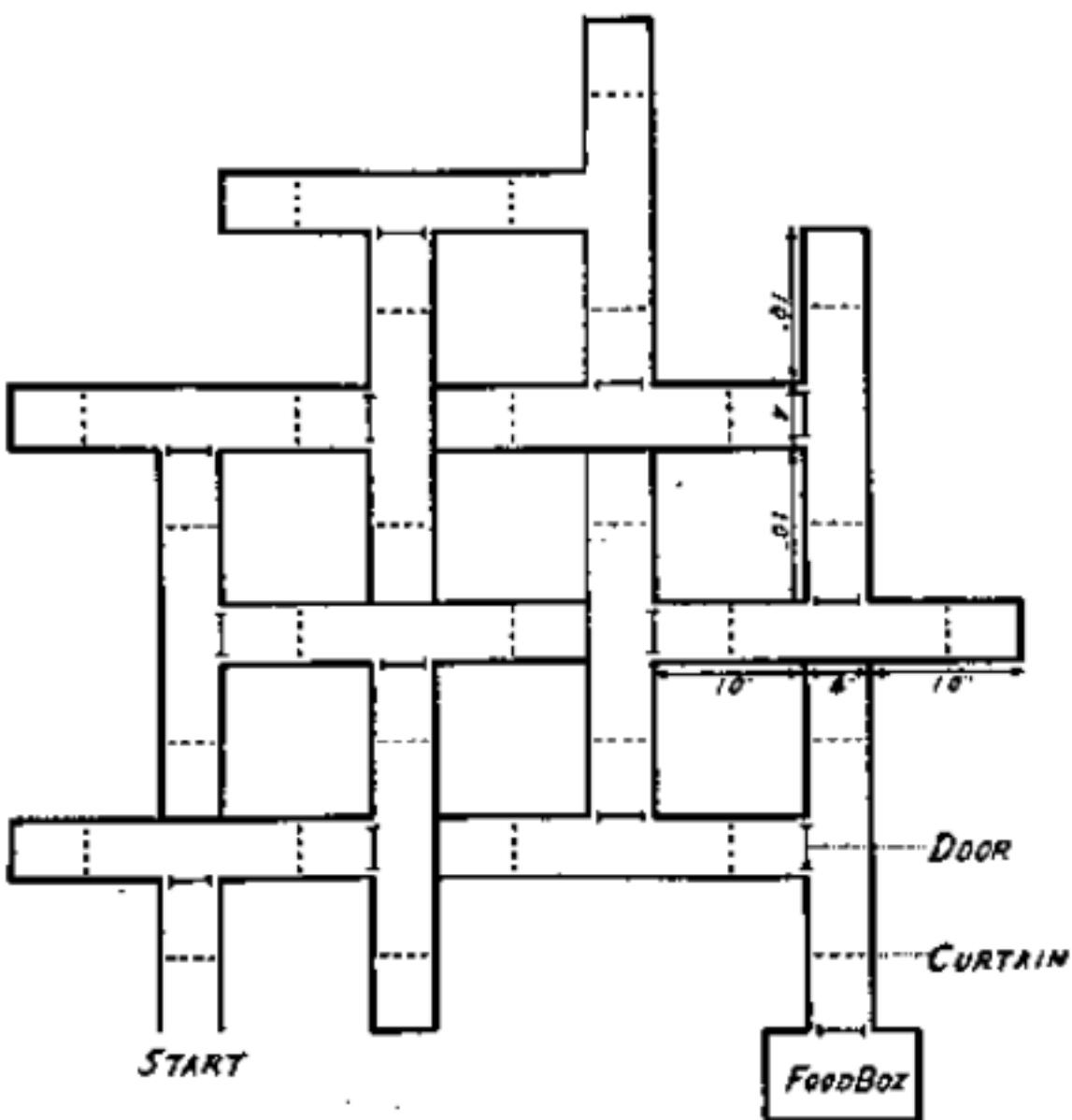
... “in the course of learning something like a field map of the environment gets established in the rat's brain”

# Cognitive Maps for Navigation

## Cognitive Maps



Tolman (*Psych Rev*, 1948)



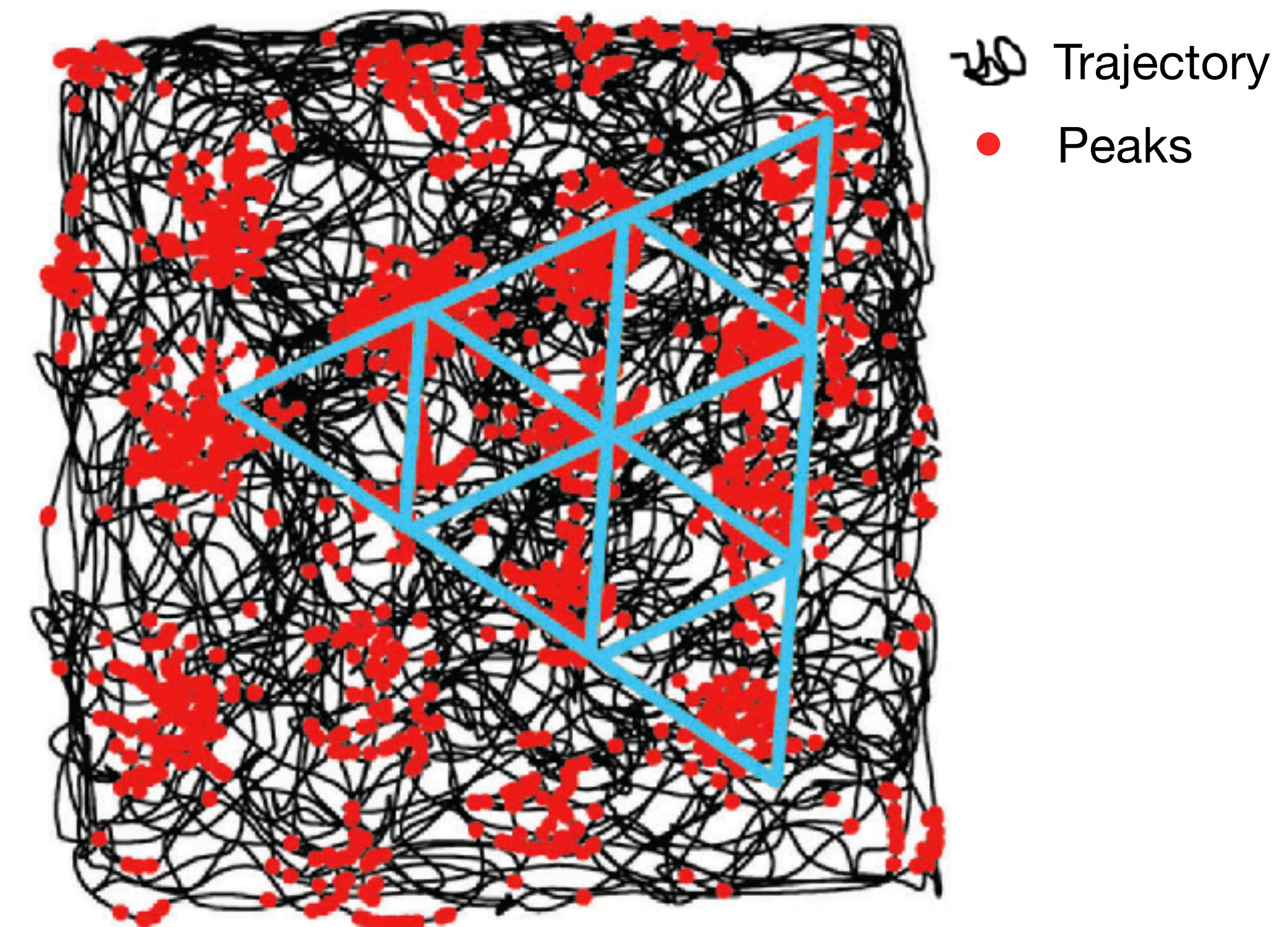
Plan of maze  
14-Unit T-Alley Maze

FIG. 1

(From M. H. Elliott, *The effect of change of reward on the maze performance of rats*. *Univ. Calif. Publ. Psychol.*, 1928, 4, p. 20.)

... “in the course of learning something like a field map of the environment gets established in the rat's brain”

## Grid Cells

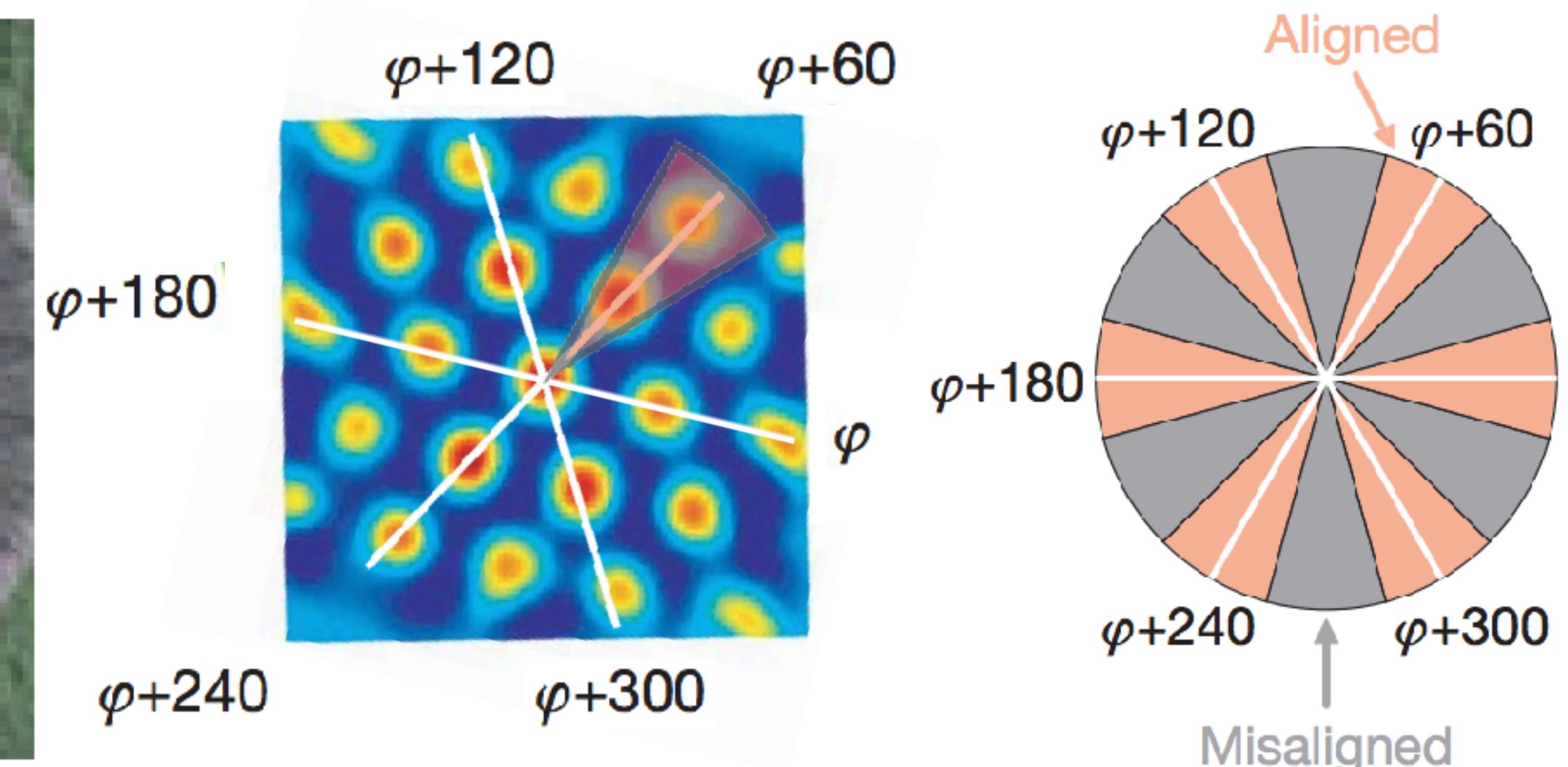


Hafting et al., (2005)

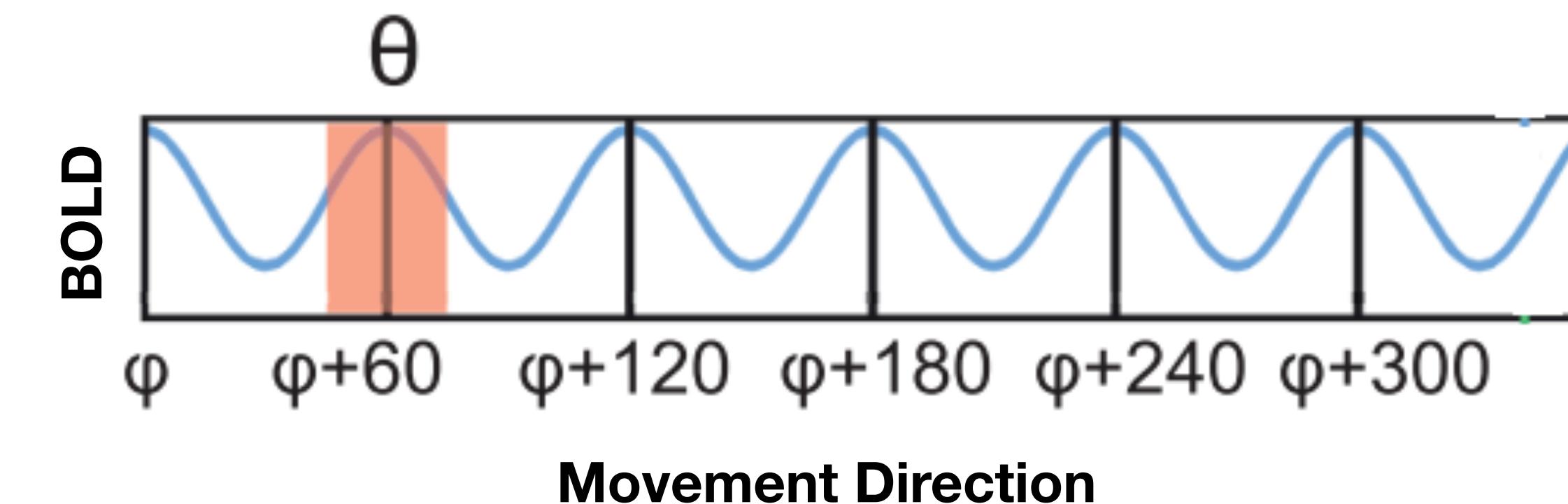
Moser, Rowland, & Moser (2015)

# We can measure trajectories in human navigation

Spatial trajectory (Birds eye)



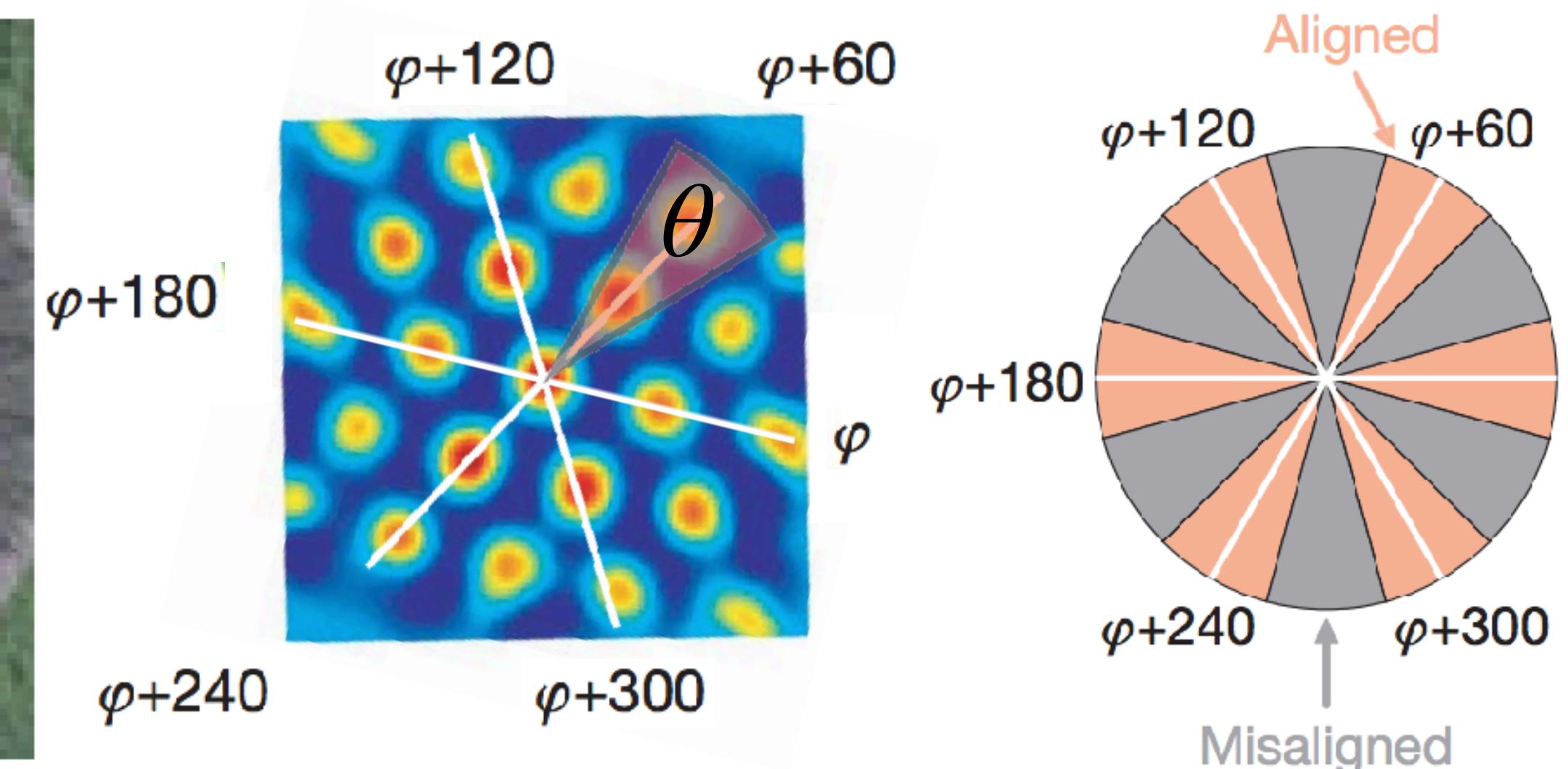
Participant perspective



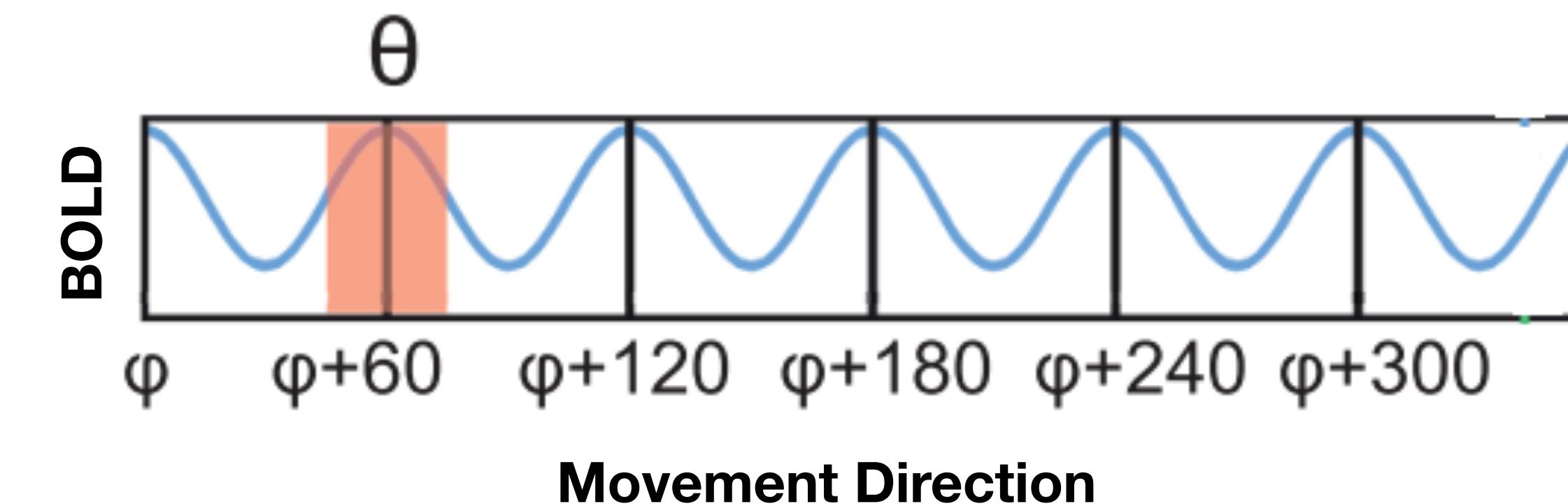
Doeller, Barry, & Burgess (*Nature*, 2010)

# We can measure trajectories in human navigation

Spatial trajectory (Birds eye)



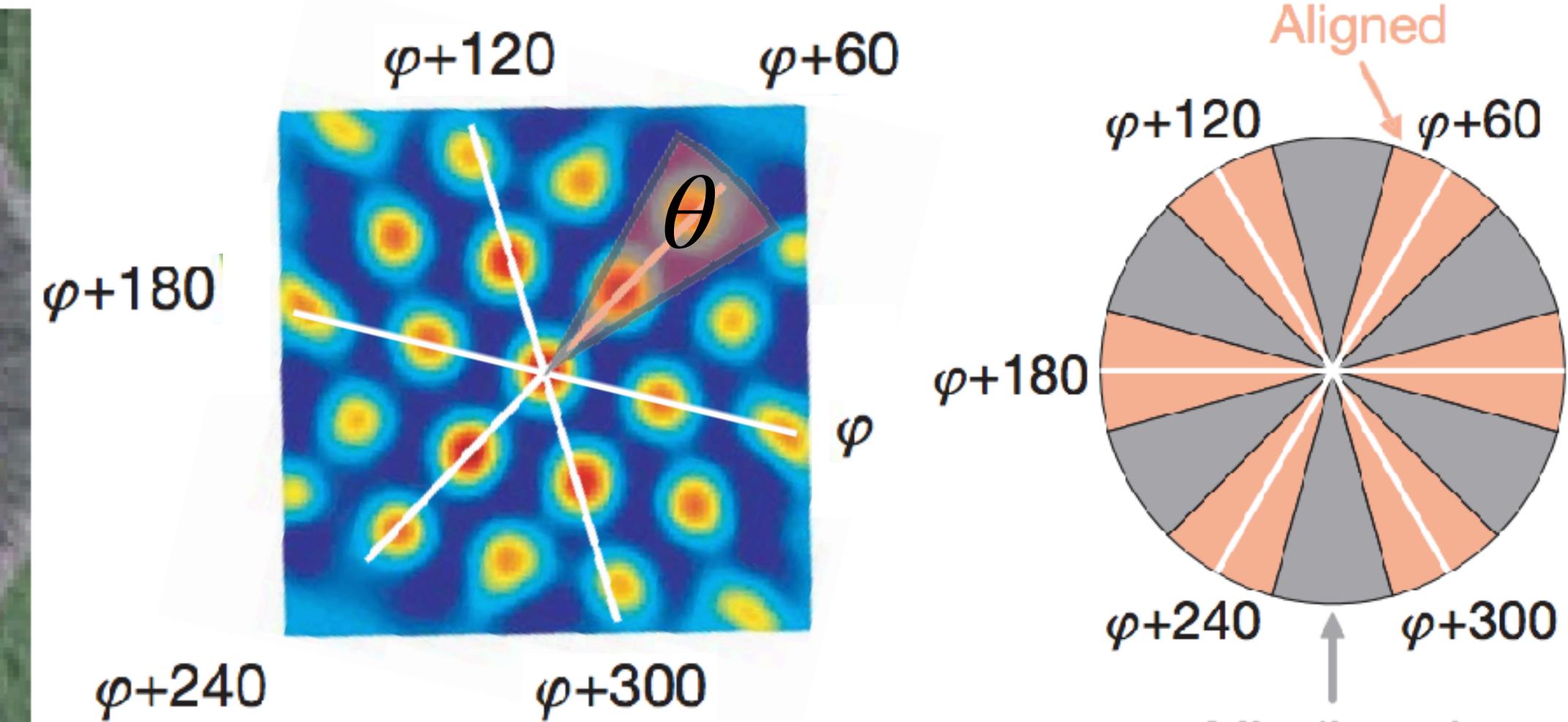
Participant perspective



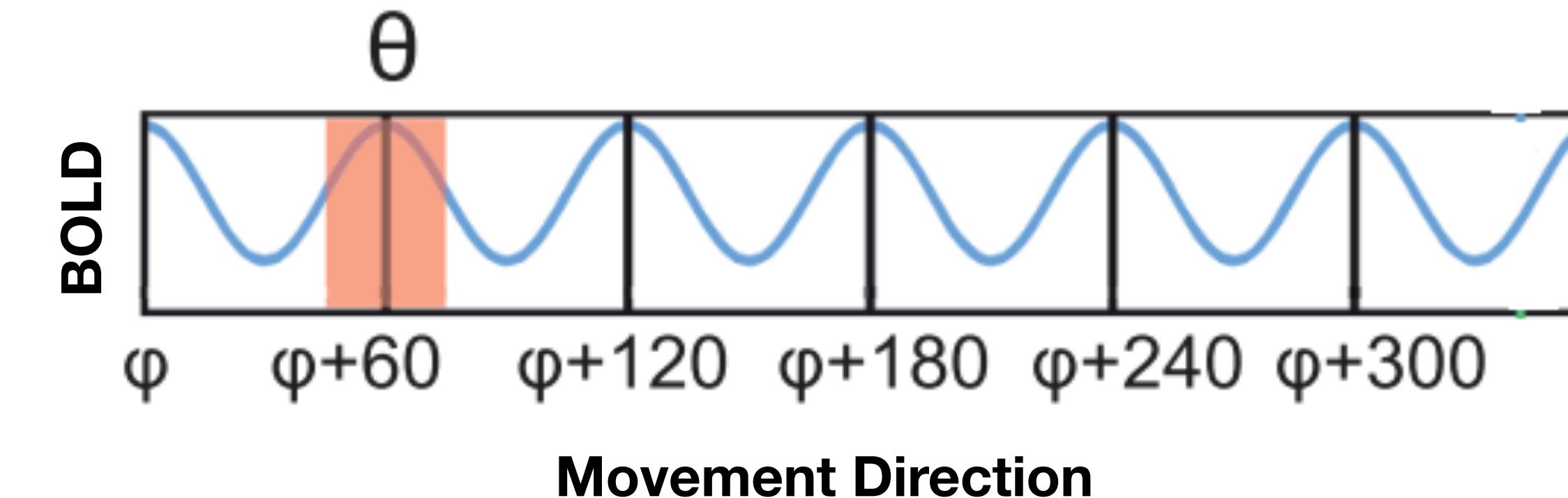
Doeller, Barry, & Burgess (*Nature*, 2010)

# We can measure trajectories in human navigation

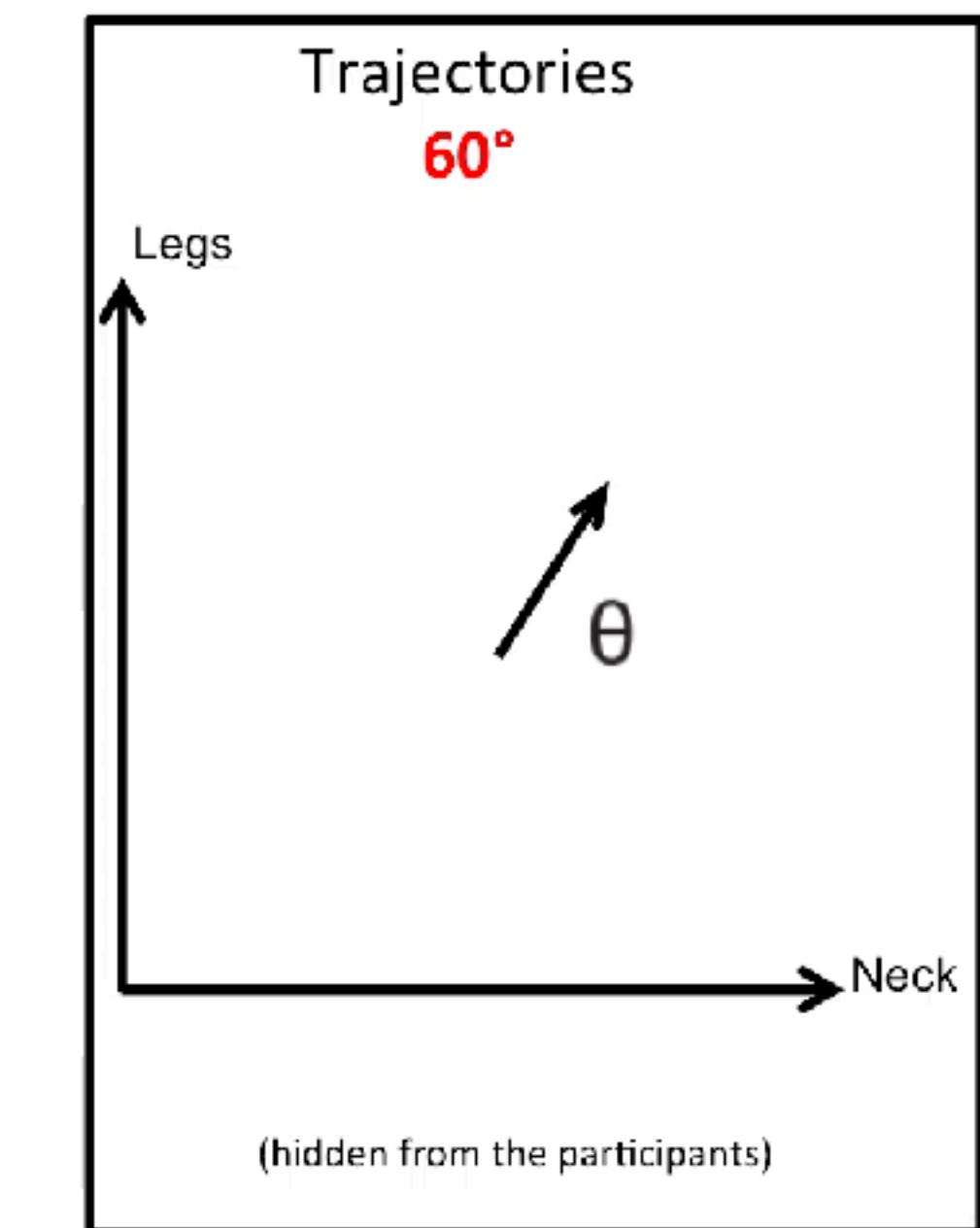
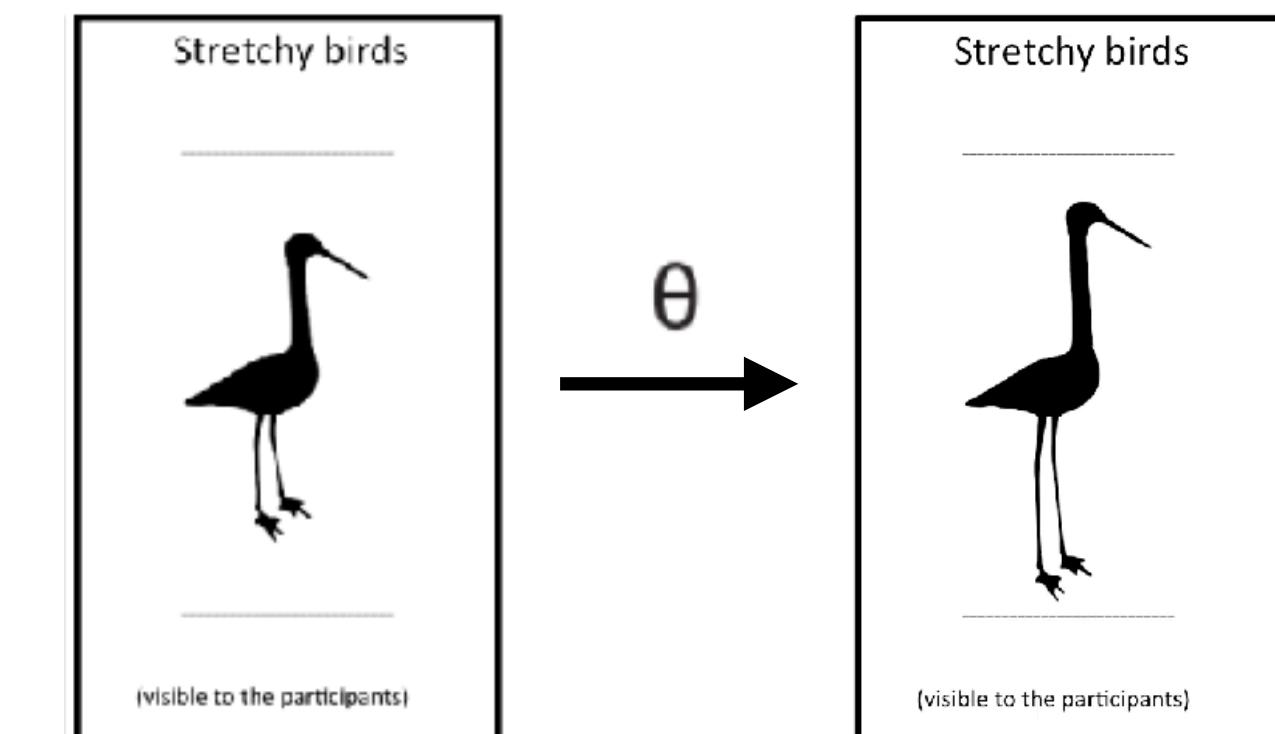
Spatial trajectory (Birds eye)



Participant perspective

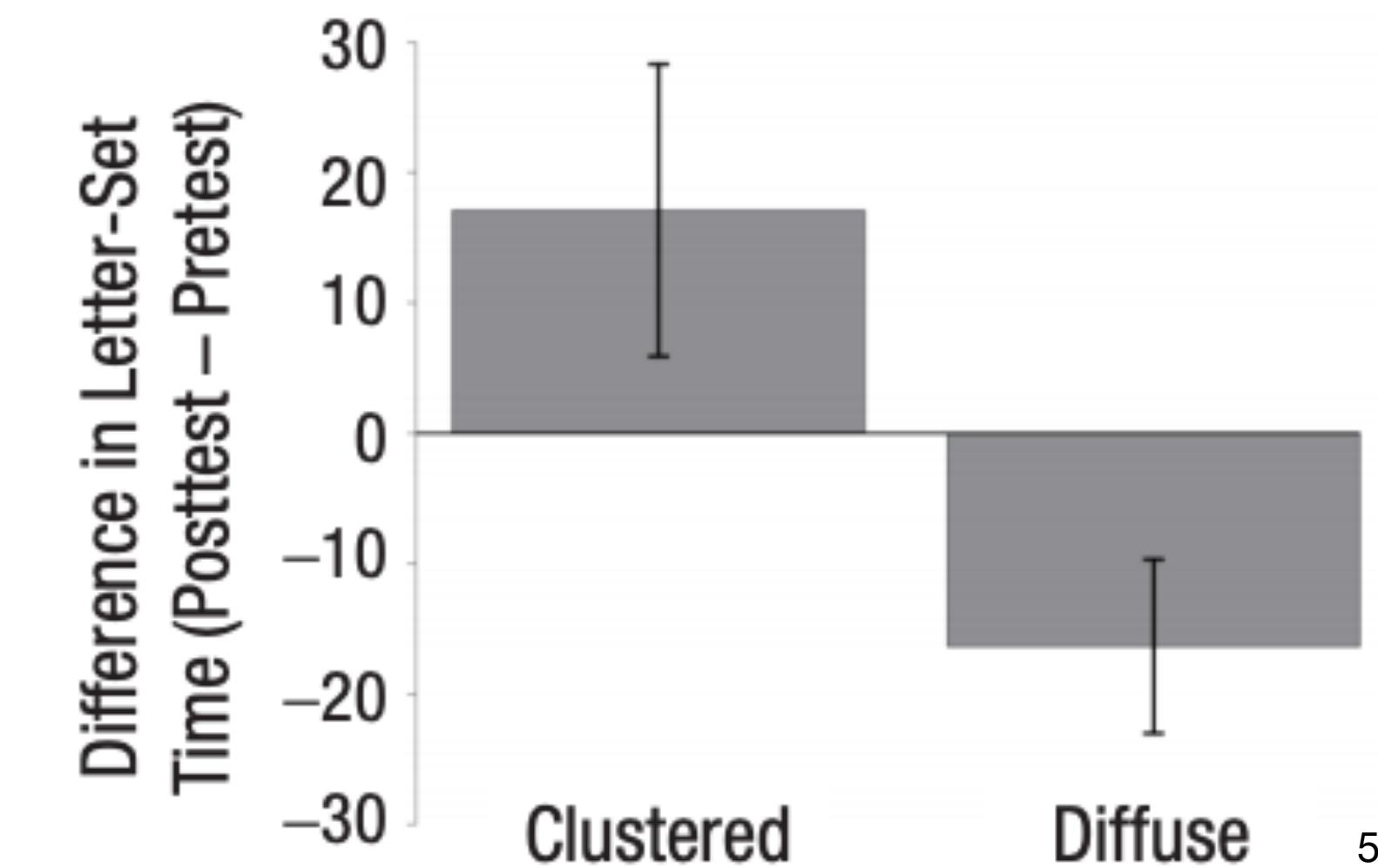
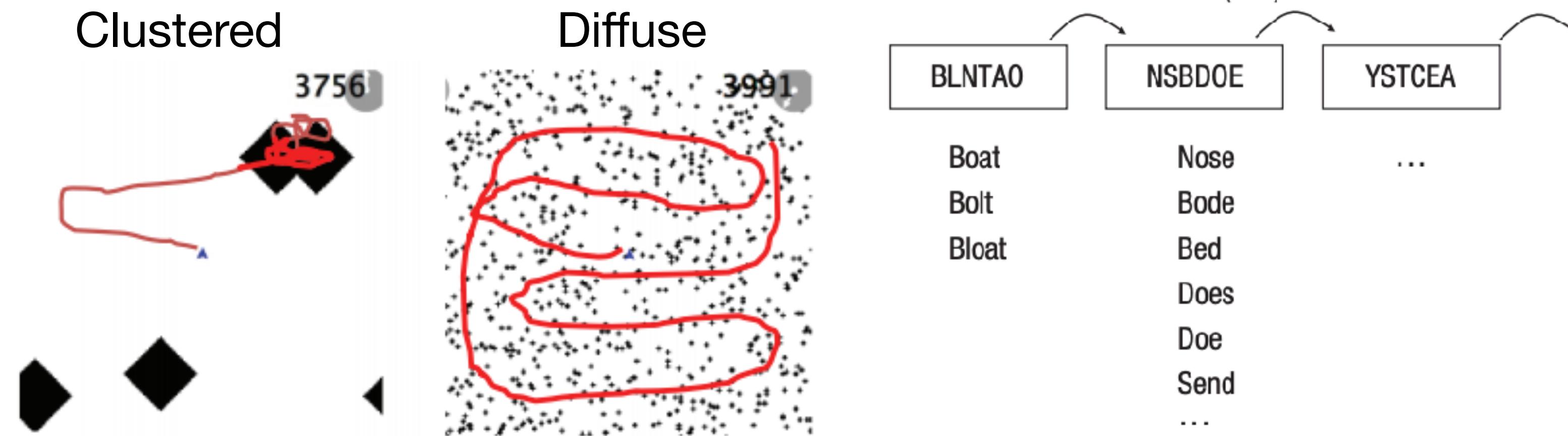
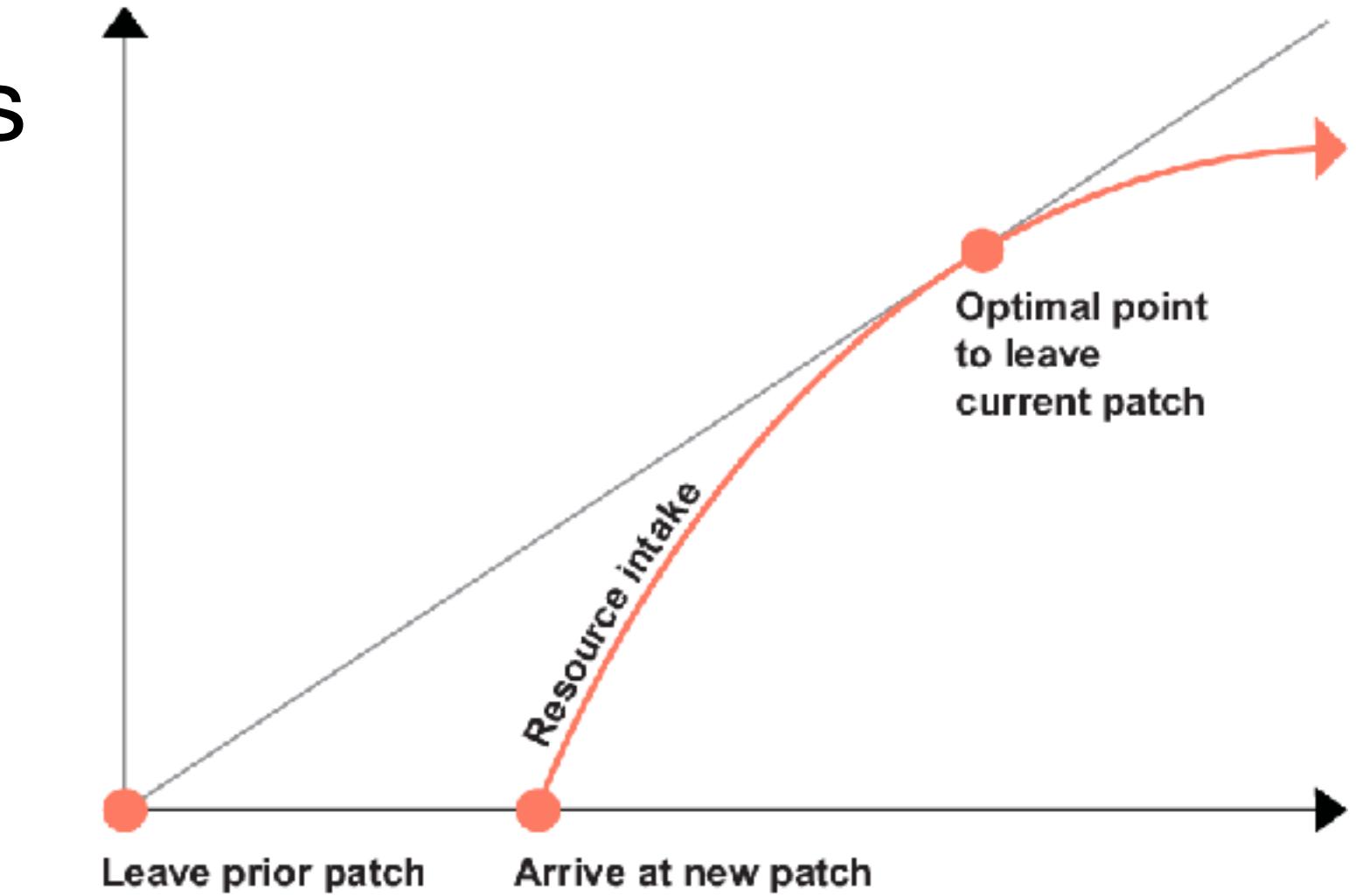


Also in non-spatial domains!



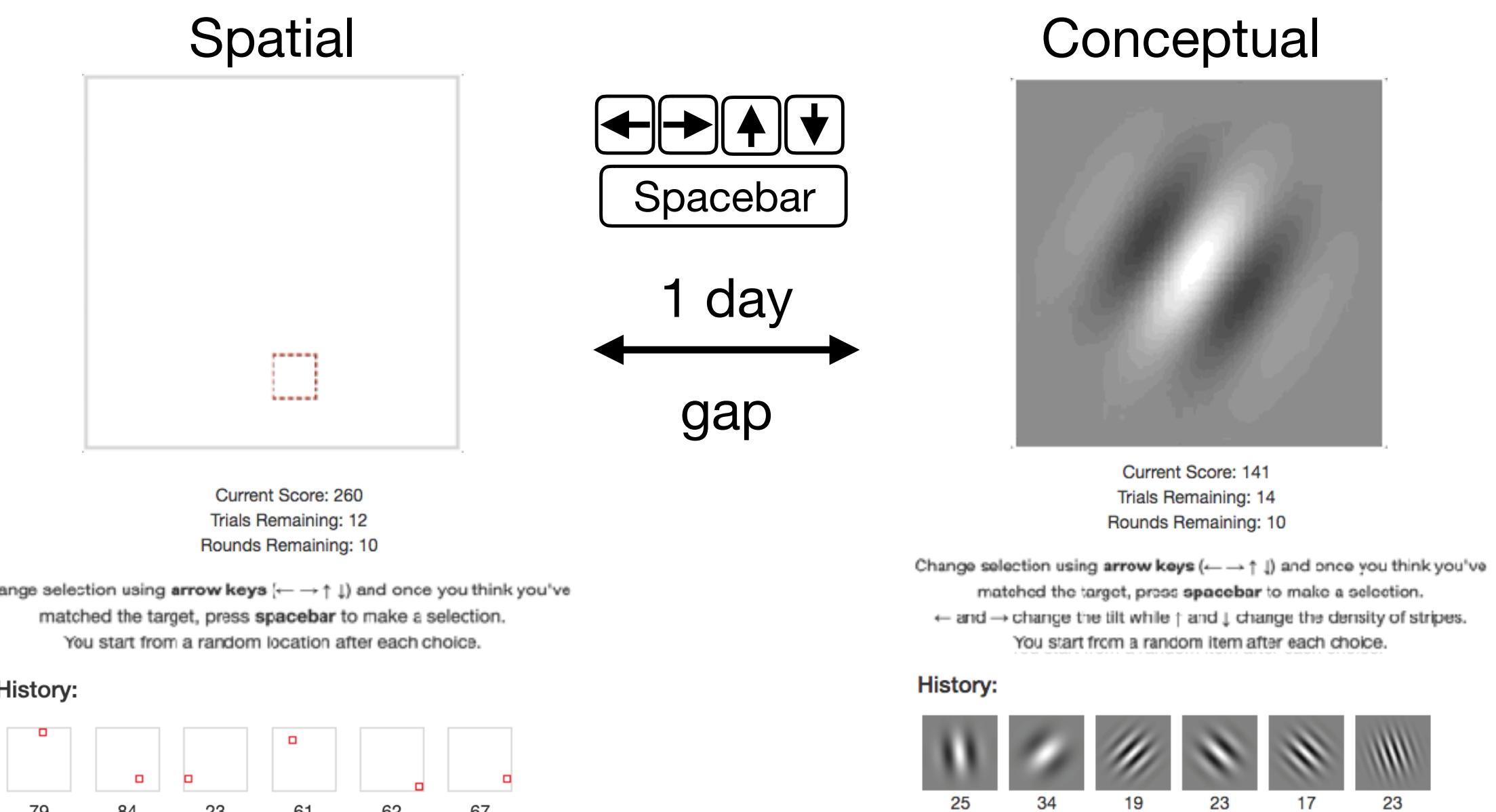
# Spatial Rewards Influence Semantic Foraging

- Search in external and internal spaces follow similar principles of optimal foraging  
Charnov (1976); Pirolli & Card (1999)
- The distribution of resources in a spatial foraging task can influence semantic search patterns in a word generation task  
Hills, Todd, & Goldstone (2008)
- “Exaptation” of spatial cognition to other domains  
Hills (2006); Hills, Todd, & Goldstone (2008)



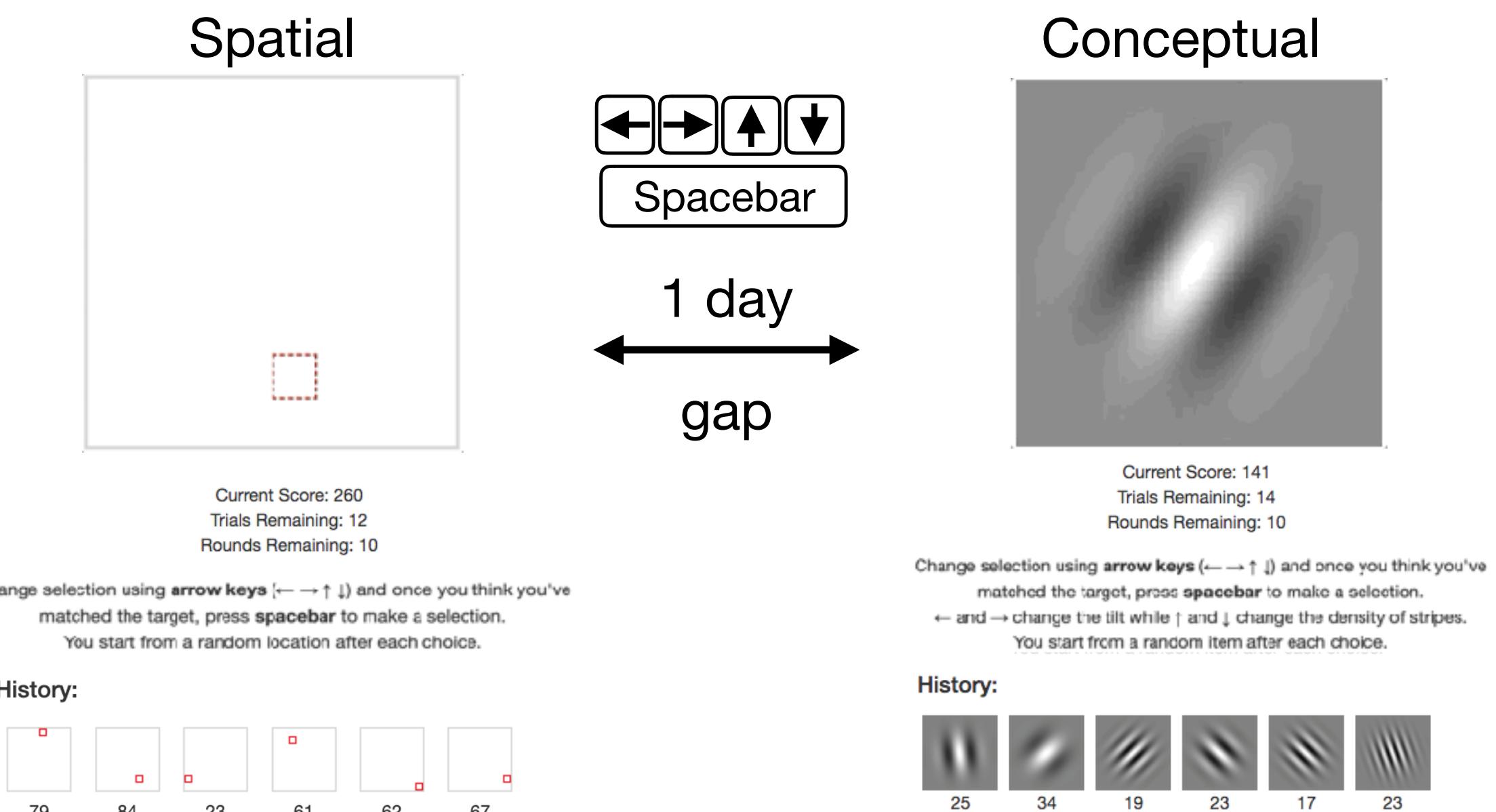
# Connecting Spatial and Conceptual Search

- Since there is evidence for a common neural representation for both spatial and conceptual navigation, what are the downstream implications for behavior?
- Are there domain general principles for generalization (about novel stimuli) and exploration (in new environments)?
- Within-subject experiment, where participants used either spatial or conceptual features to guide the search for rewards



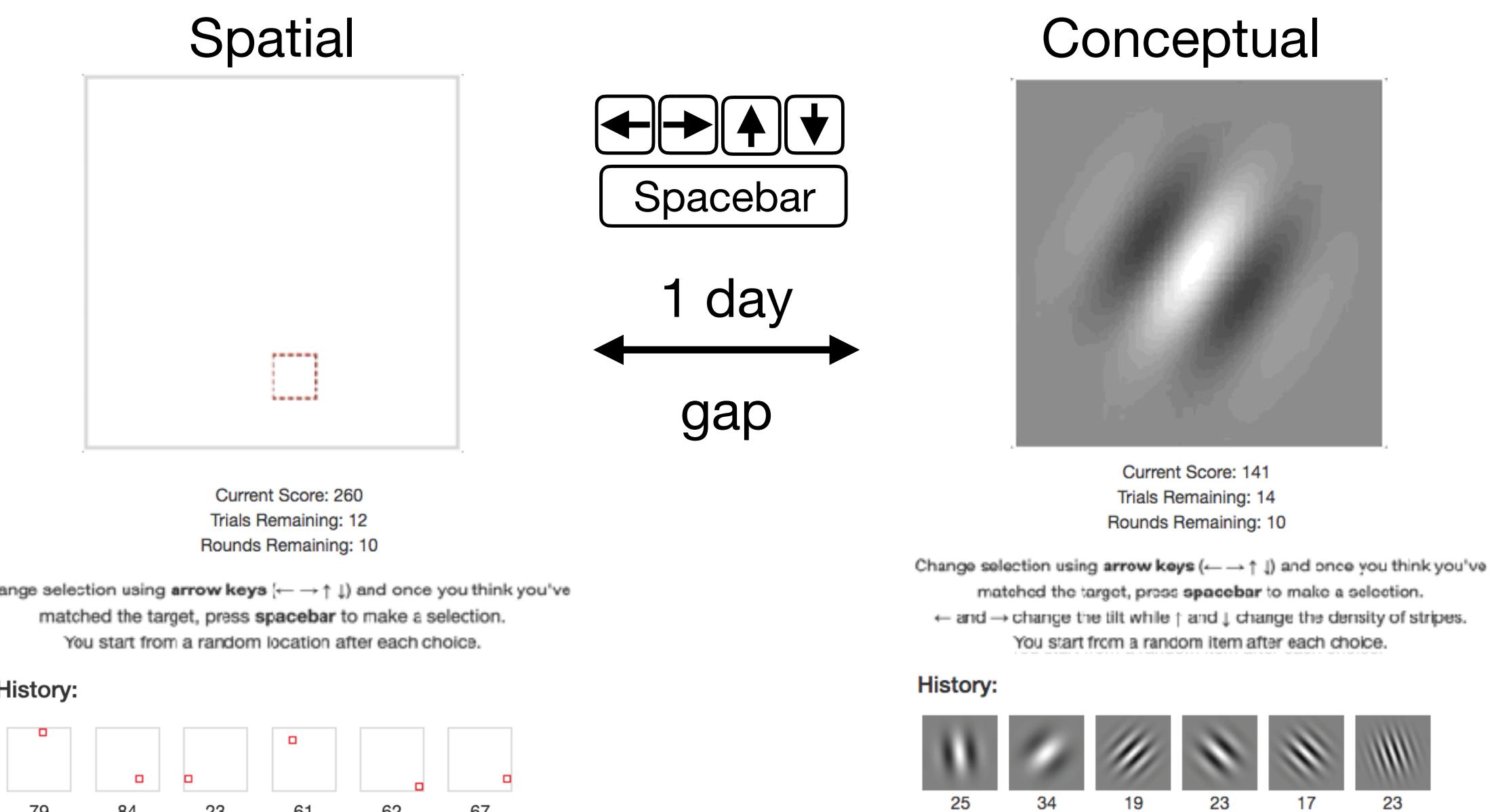
# Connecting Spatial and Conceptual Search

- Since there is evidence for a common neural representation for both spatial and conceptual navigation, what are the downstream implications for behavior?
- Are there domain general principles for generalization (about novel stimuli) and exploration (in new environments)?
- Within-subject experiment, where participants used either spatial or conceptual features to guide the search for rewards



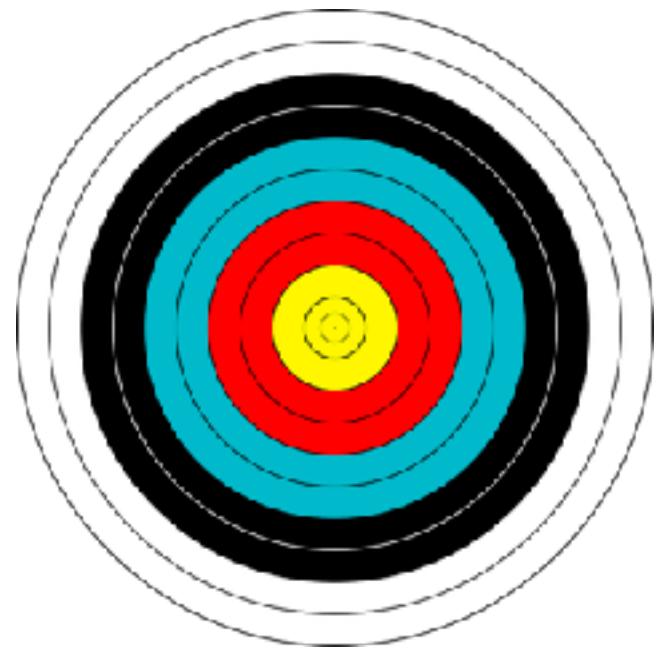
# Connecting Spatial and Conceptual Search

- Since there is evidence for a common neural representation for both spatial and conceptual navigation, what are the downstream implications for behavior?
- Are there domain general principles for generalization (about novel stimuli) and exploration (in new environments)?
- Within-subject experiment, where participants used either spatial or conceptual features to guide the search for rewards



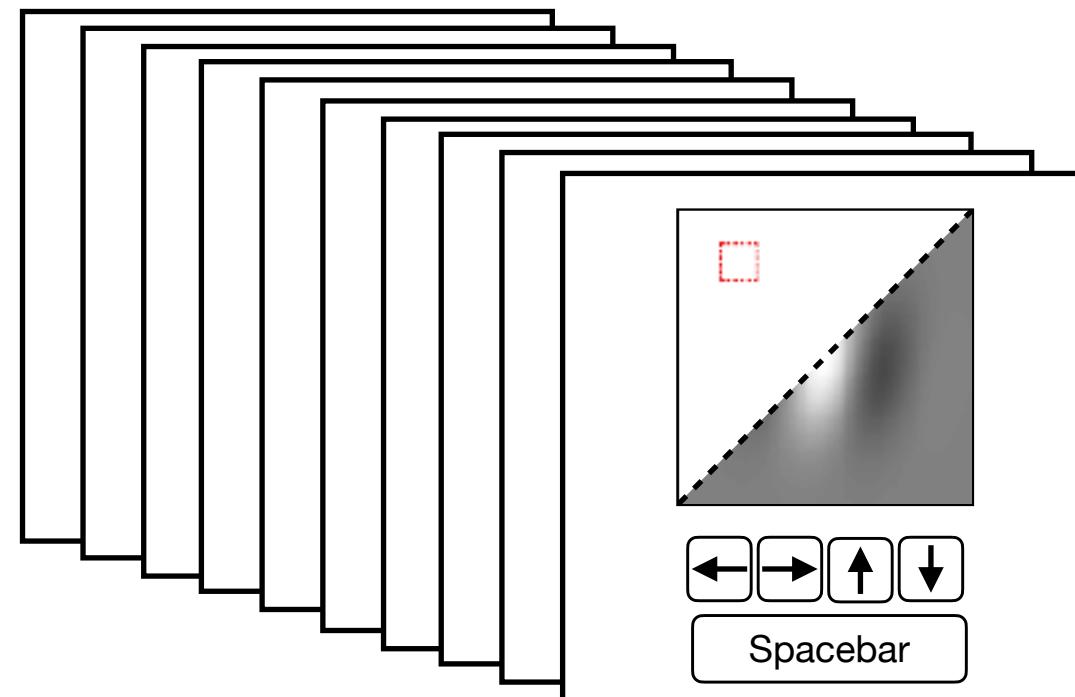
# Task Design

## Training Phase



Comprehension  
Questions

## Main Search Task



10 rounds  
20 trials in each

## Judgment and Confidence

Match target stimuli until learning criterion reached

### Current Selection



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

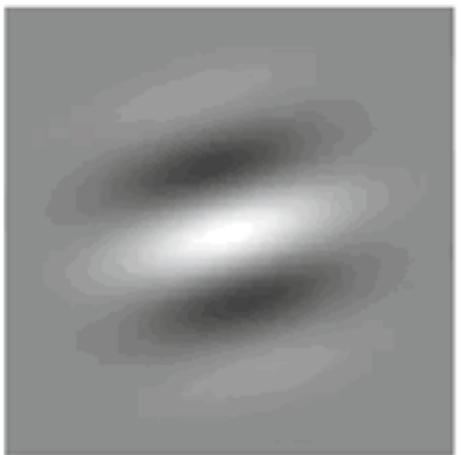
Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.  
You start from a random location after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target



### Current Selection



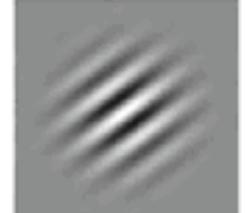
Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.  
 $\leftarrow$  and  $\rightarrow$  change the tilt while  $\uparrow$  and  $\downarrow$  change the density of stripes.

You start from a random item after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target

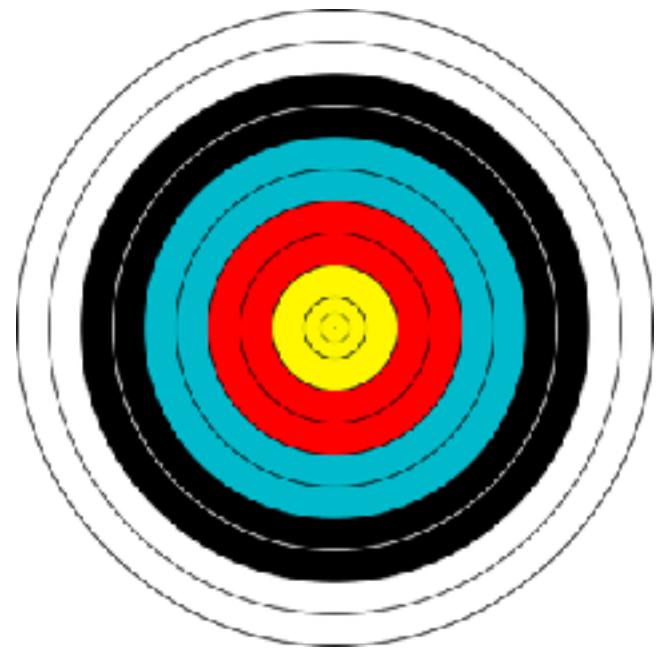


At least 32 trials AND a run of 9 out of 10 correct

Bonus  
Round

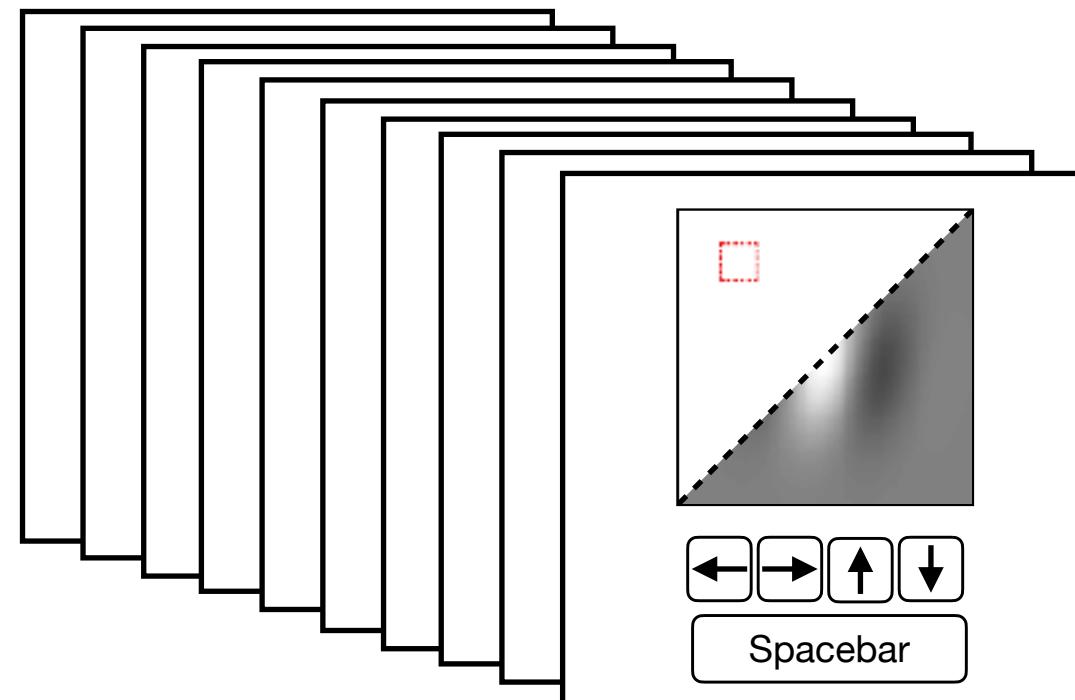
# Task Design

## Training Phase



Comprehension Questions

## Main Search Task



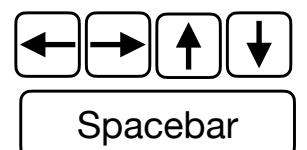
10 rounds  
20 trials in each

## Judgment and Confidence

Bonus Round

Match target stimuli until learning criterion reached

### Current Selection



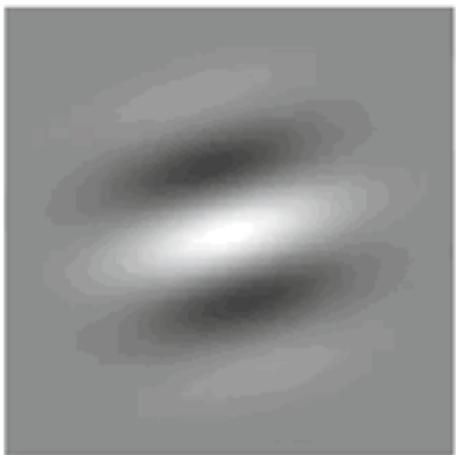
Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

Change selection using **arrow keys** ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press **spacebar** to make a selection.  
You start from a random location after each choice.  
The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target



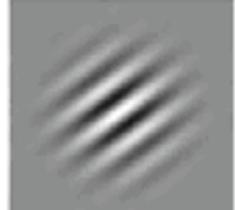
### Current Selection



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

Change selection using **arrow keys** ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press **spacebar** to make a selection.  
 $\leftarrow$  and  $\rightarrow$  change the tilt while  $\uparrow$  and  $\downarrow$  change the density of stripes.  
You start from a random item after each choice.  
The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

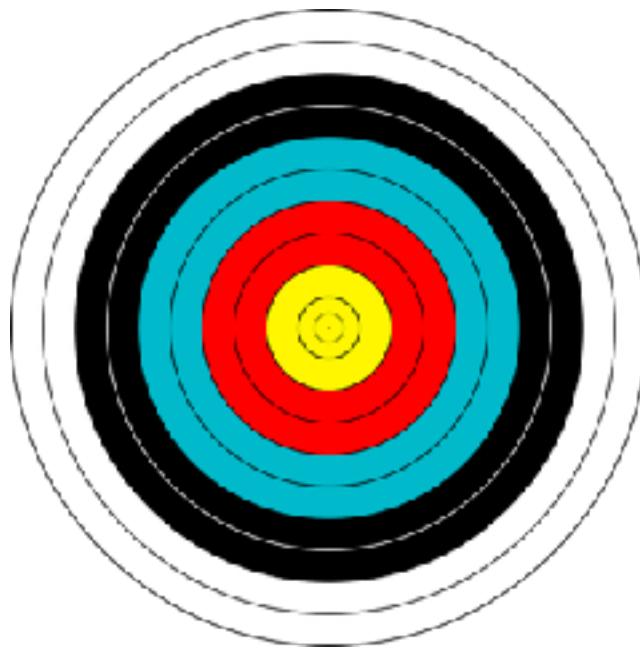
Target



At least 32 trials AND a run of 9 out of 10 correct

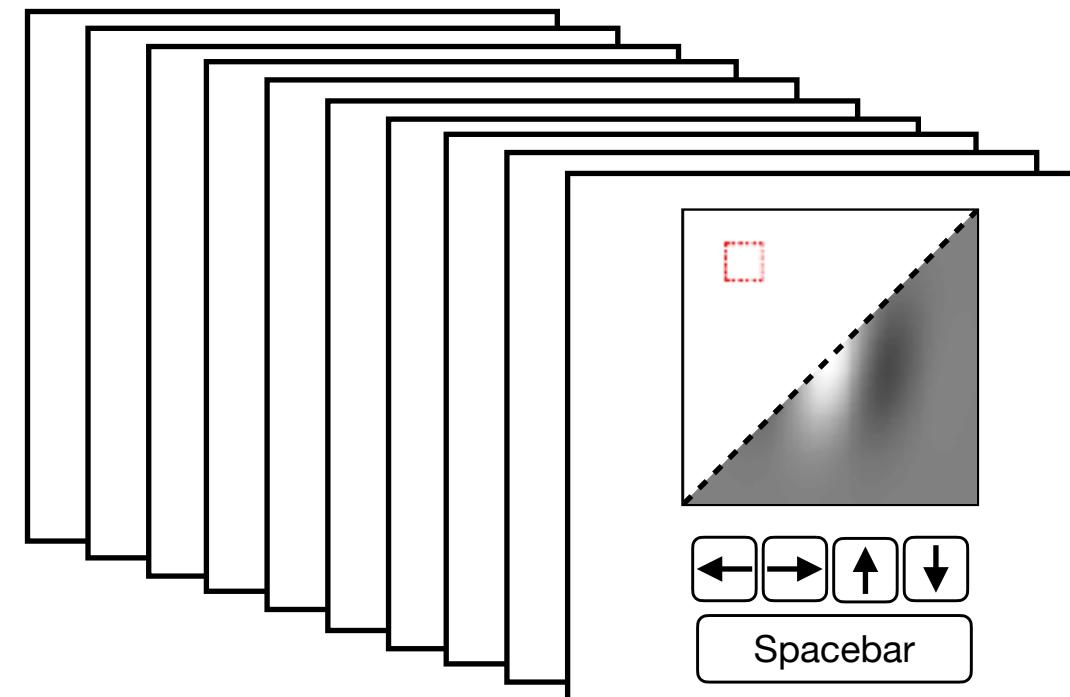
# Task Design

Training Phase



Comprehension  
Questions

Main Search Task



10 rounds  
20 trials in each

Judgment and Confidence

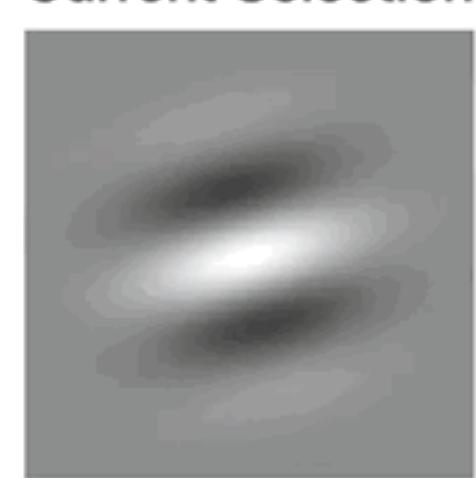
Bonus  
Round

Match target stimuli until learning criterion reached

Current Selection



Current Selection



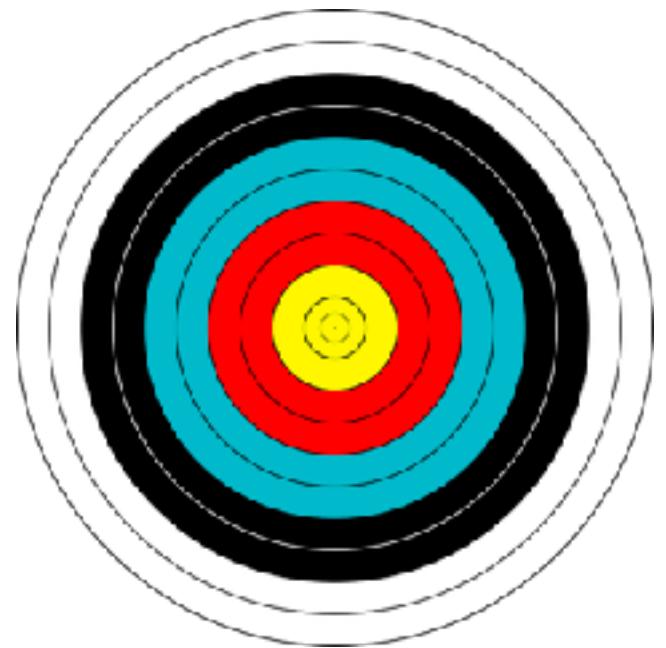
## Bandit Task:

1. Select stimuli using
2. Make selection using Spacebar
3. Reward is displayed and then added to history
4. Start at a random stimuli

At least 32 trials AND a run of 9 out of 10 correct

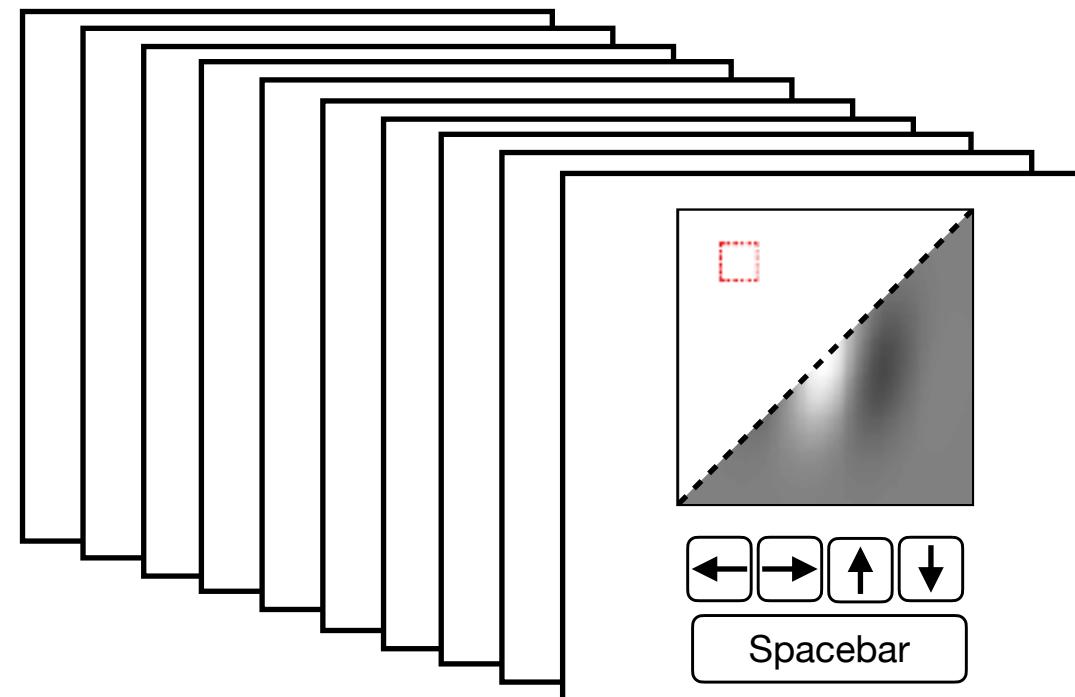
# Task Design

## Training Phase



Comprehension  
Questions

## Main Search Task



10 rounds  
20 trials in each

## Judgment and Confidence

Match target stimuli until learning criterion reached

### Current Selection



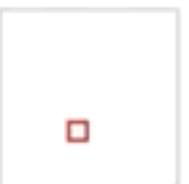
Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.

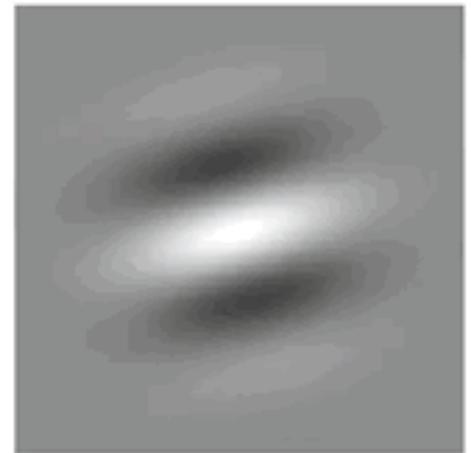
You start from a random location after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target



### Current Selection



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

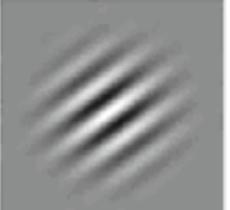
Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.

$\leftarrow$  and  $\rightarrow$  change the tilt while  $\uparrow$  and  $\downarrow$  change the density of stripes.

You start from a random item after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target

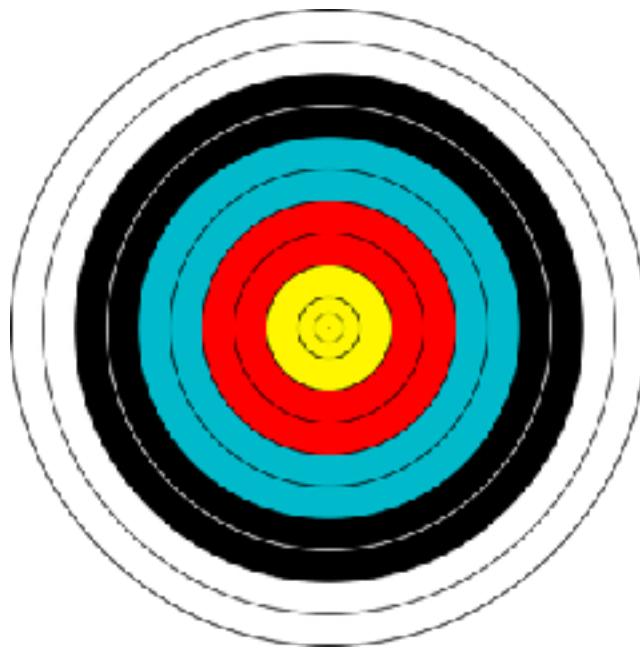


At least 32 trials AND a run of 9 out of 10 correct

Bonus  
Round

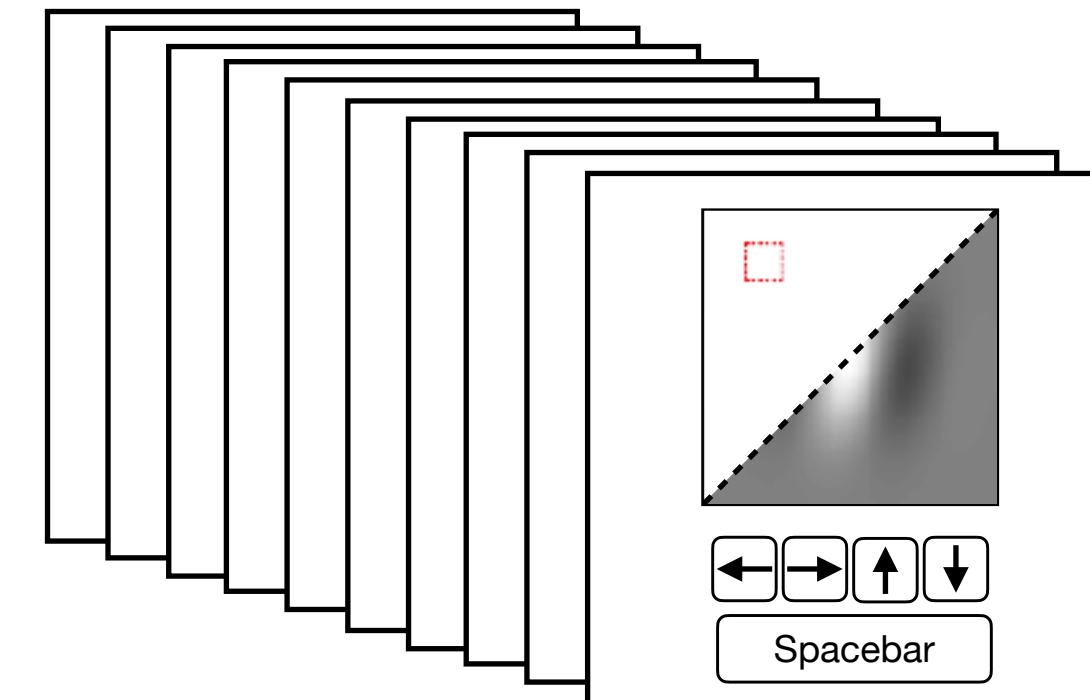
# Task Design

## Training Phase



Comprehension  
Questions

## Main Search Task

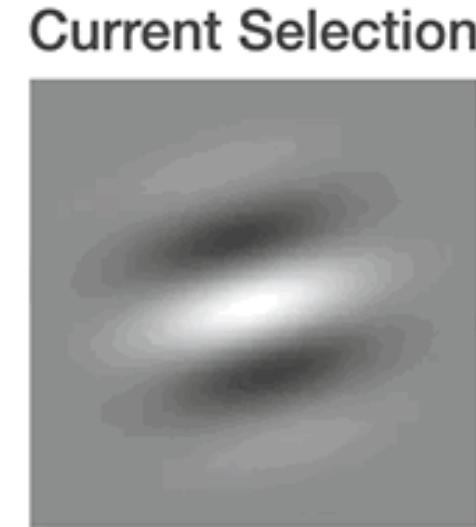
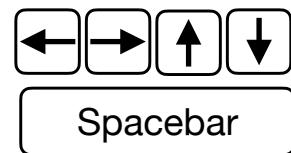


10 rounds  
20 trials in each

## Judgment and Confidence

Match target stimuli until learning criterion reached

### Current Selection



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.

$\leftarrow$  and  $\rightarrow$  change the tilt while  $\uparrow$  and  $\downarrow$  change the density of stripes.

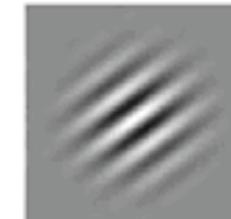
You start from a random location after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

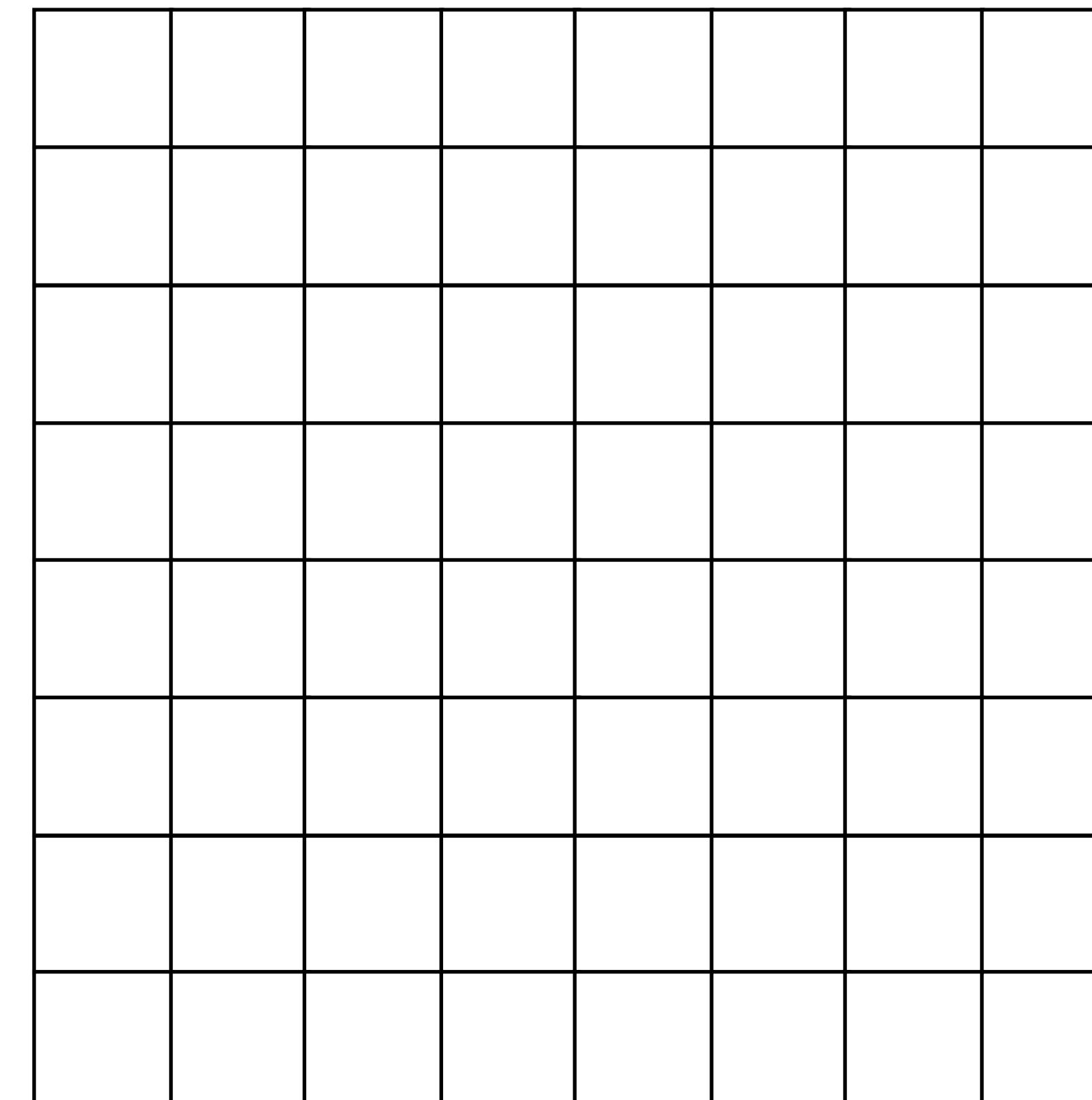
Target



Target



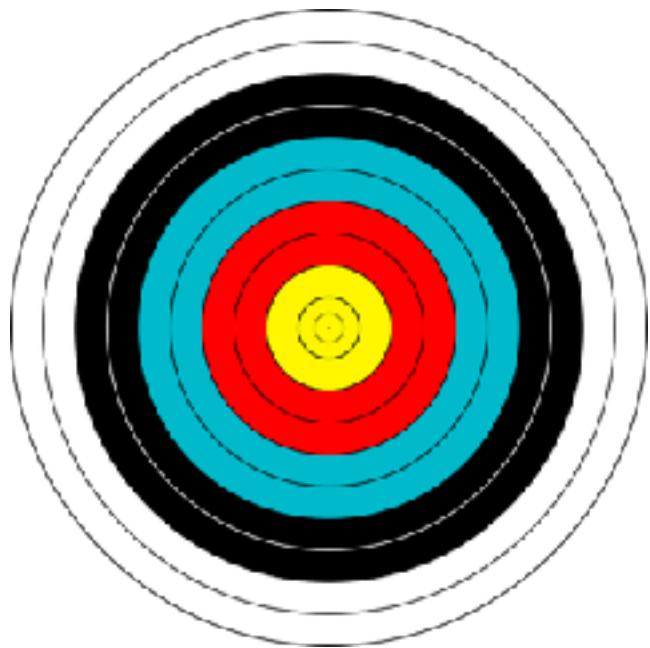
At least 32 trials AND a run of 9 out of 10 correct



Bonus  
Round

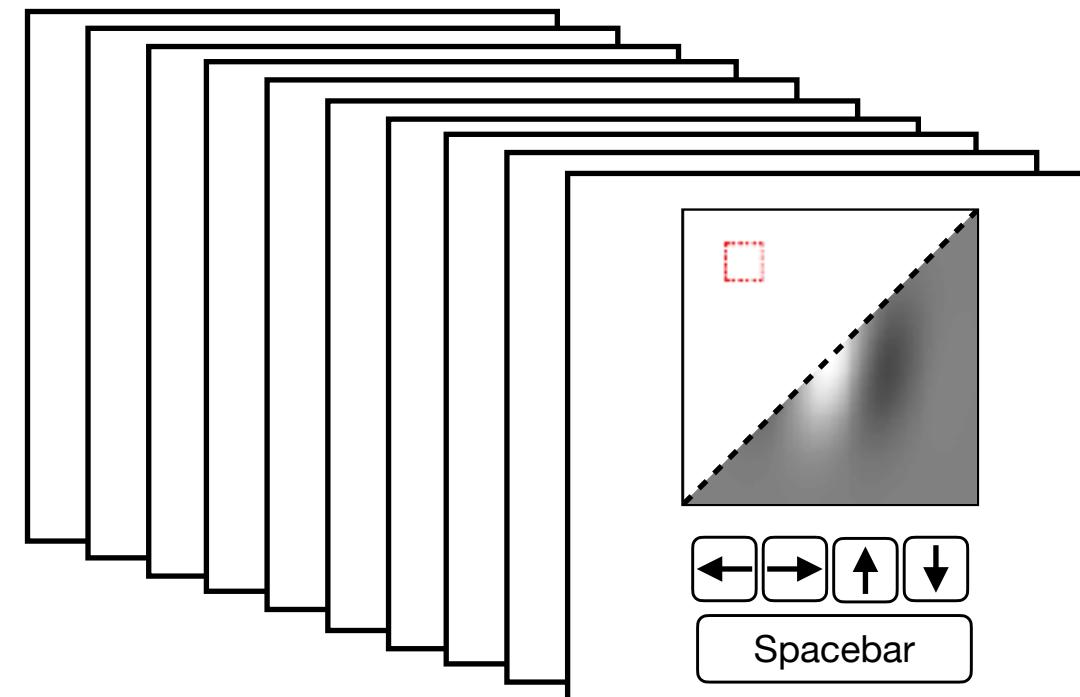
# Task Design

## Training Phase



Comprehension  
Questions

## Main Search Task



10 rounds  
20 trials in each

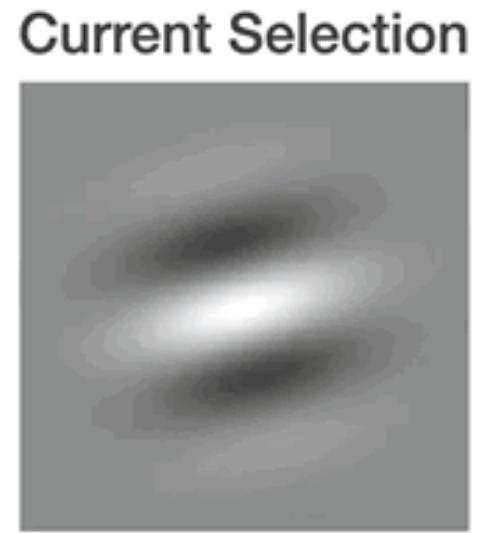
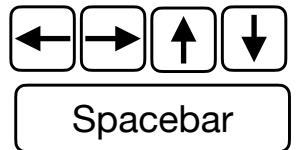
## Judgment and Confidence

Match target stimuli until learning criterion reached

### Current Selection



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

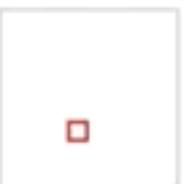
Change selection using **arrow keys** ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press **spacebar** to make a selection.

$\leftarrow$  and  $\rightarrow$  change the tilt while  $\uparrow$  and  $\downarrow$  change the density of stripes.

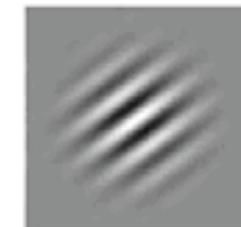
You start from a random location after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target

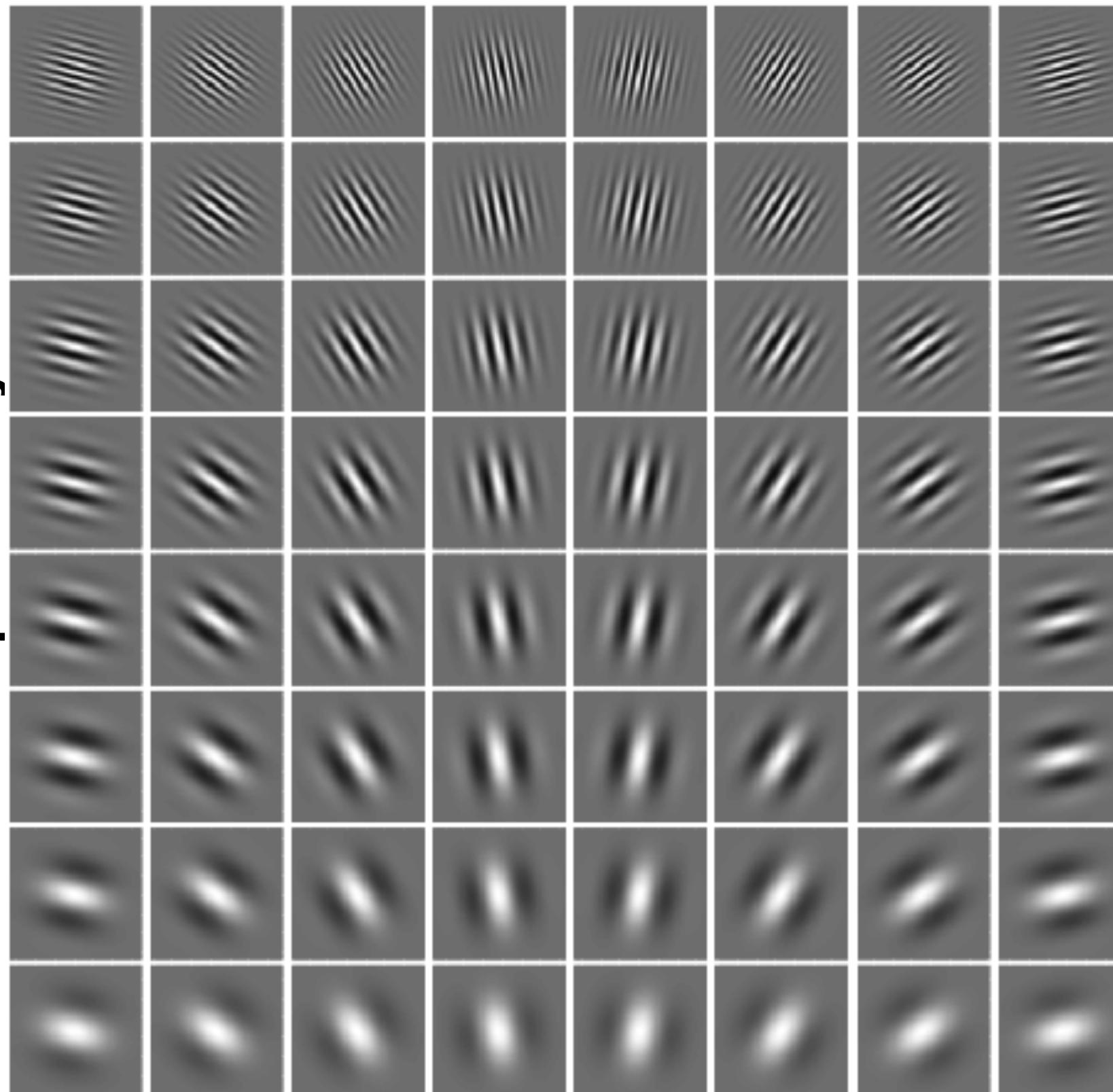


Target



At least 32 trials AND a run of 9 out of 10 correct

Stripe Density

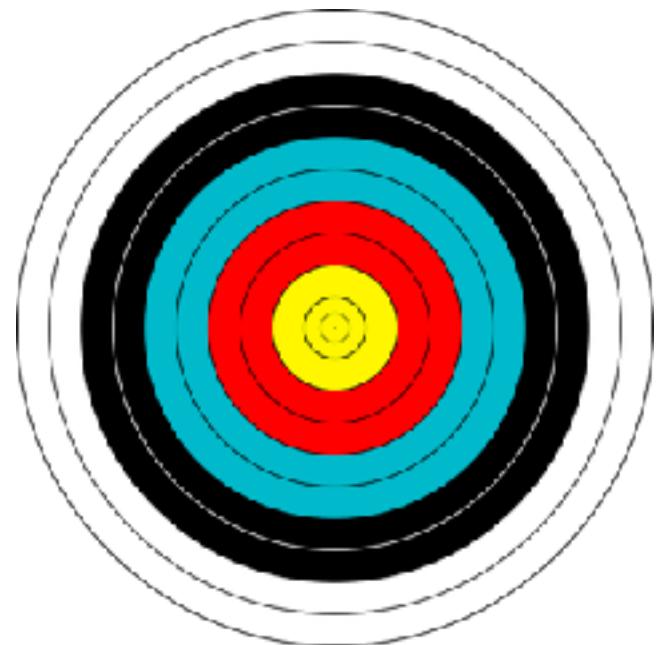


Tilt

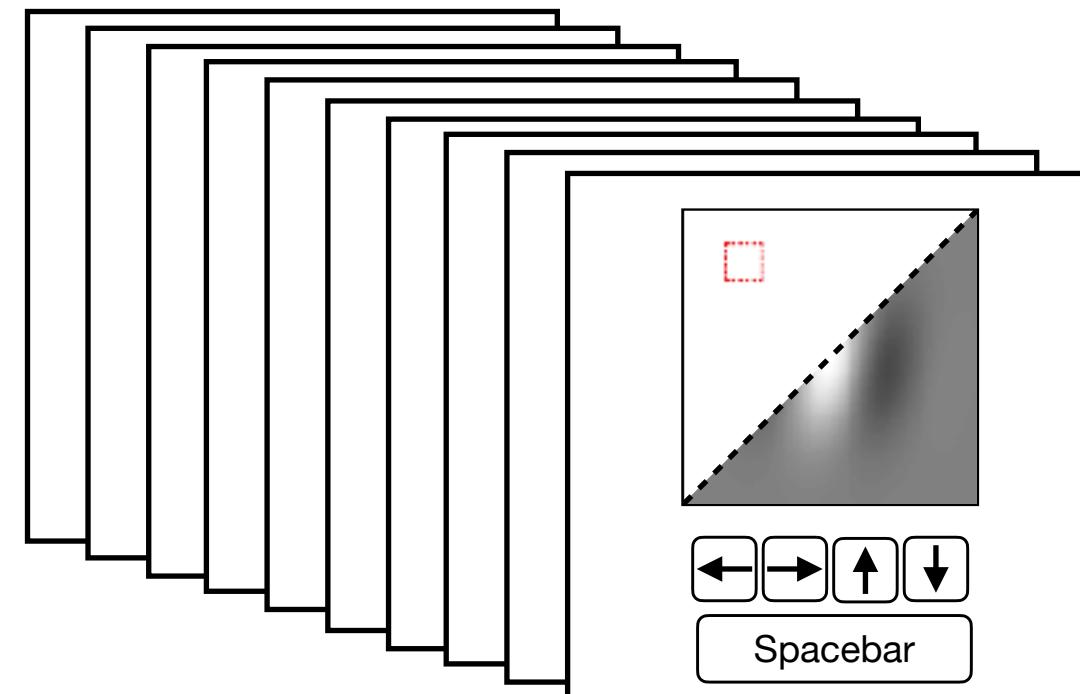
Bonus  
Round

# Task Design

## Training Phase



## Main Search Task



10 rounds  
20 trials in each

## Judgment and Confidence

Match target stimuli until learning criterion reached

### Current Selection



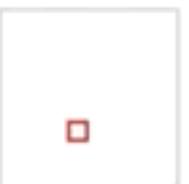
Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.

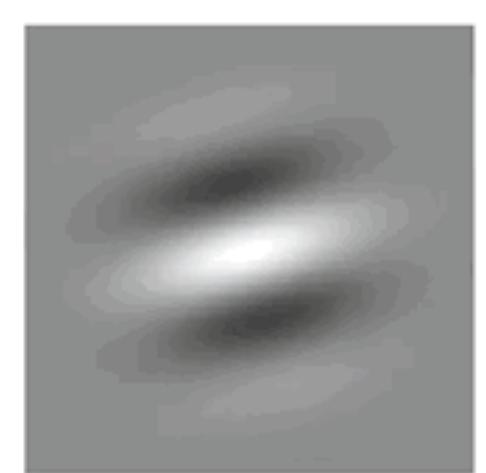
You start from a random location after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target



### Current Selection



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

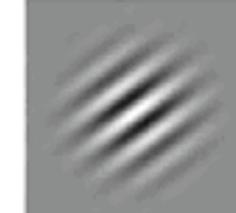
Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.

$\leftarrow$  and  $\rightarrow$  change the tilt while  $\uparrow$  and  $\downarrow$  change the density of stripes.

You start from a random item after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target

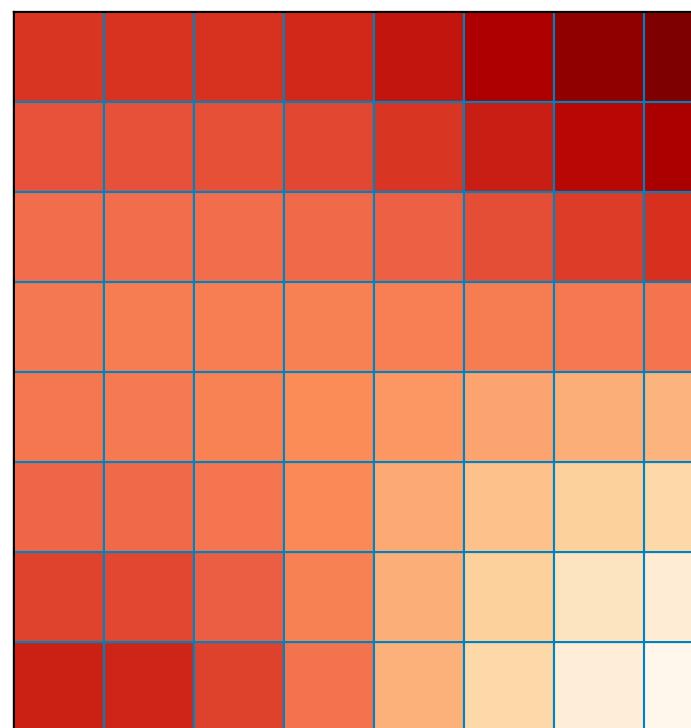
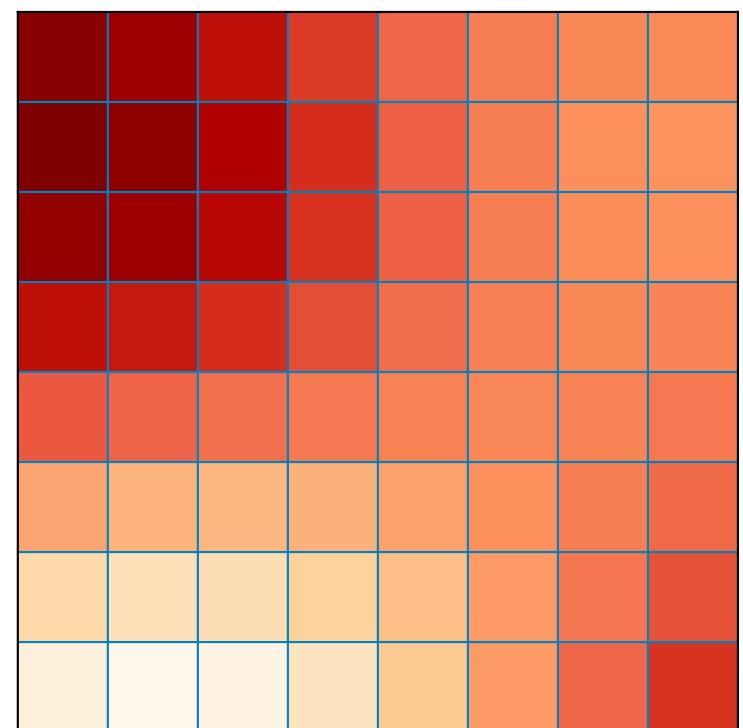


At least 32 trials AND a run of 9 out of 10 correct

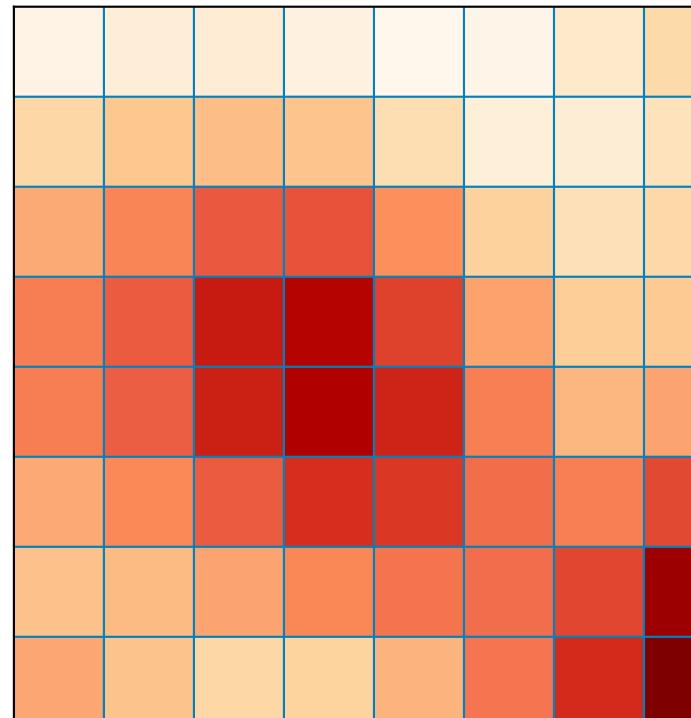
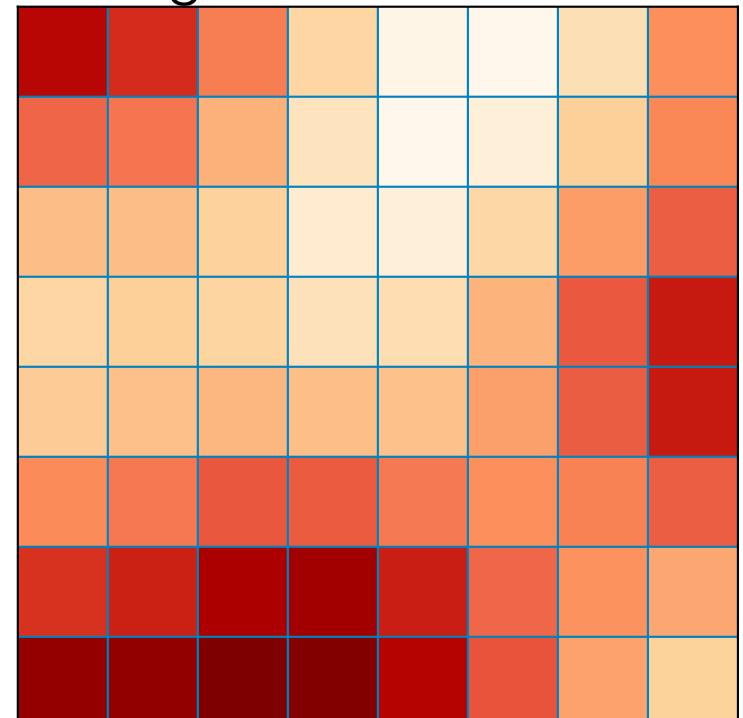
Comprehension Questions

Bonus Round

### Smooth Environments

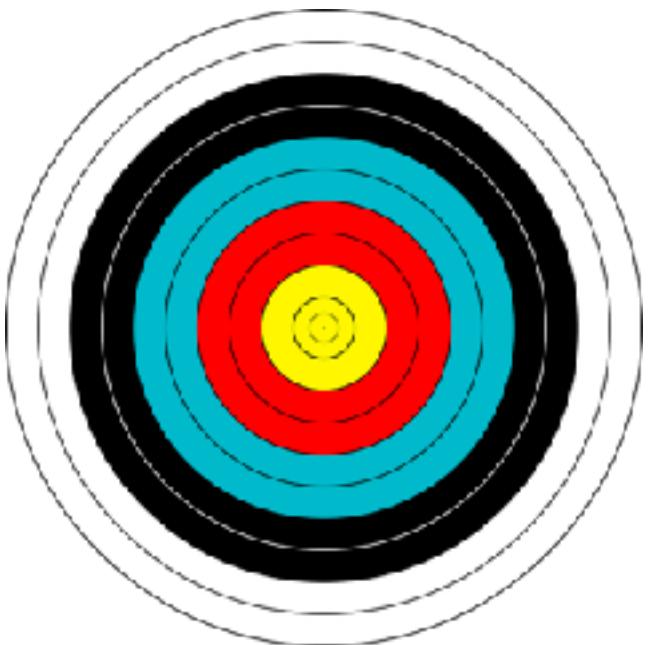


### Rough Environments



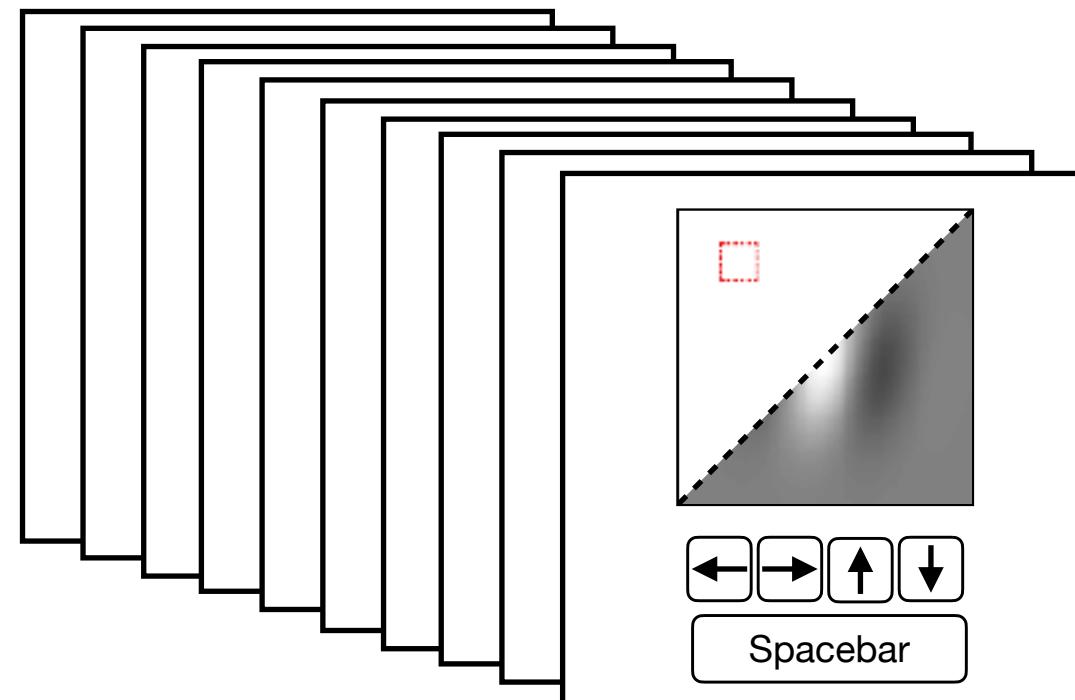
# Task Design

## Training Phase



Comprehension Questions

## Main Search Task



10 rounds  
20 trials in each

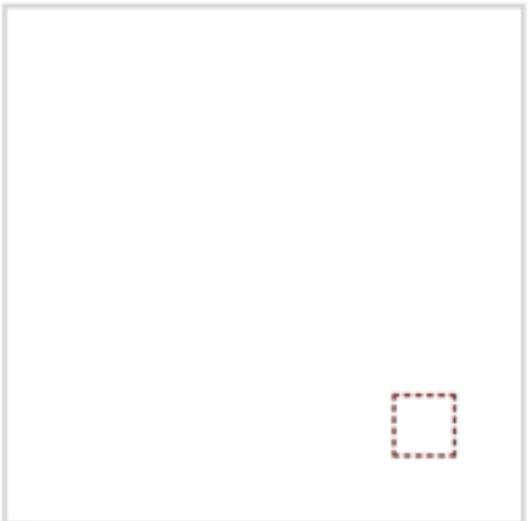
## Judgment and Confidence



After the 15th trial of the 10th round

Match target stimuli until learning criterion reached

### Current Selection



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.

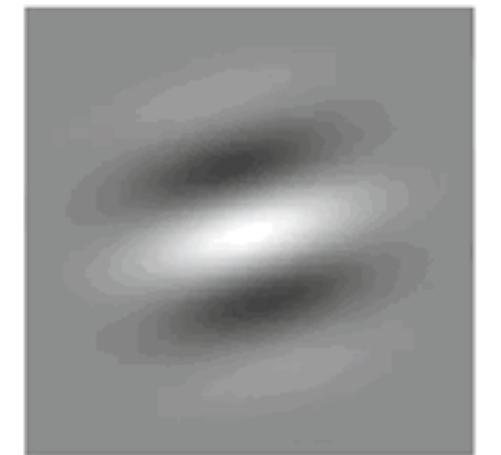
You start from a random location after each choice.

The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target



### Current Selection



Correct Selections: 0 out of 0  
Accuracy (last 10 choices): 0%

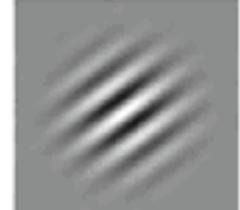
Change selection using arrow keys ( $\leftarrow \rightarrow \uparrow \downarrow$ ) and once you think you've matched the target, press spacebar to make a selection.

$\leftarrow$  and  $\rightarrow$  change the tilt while  $\uparrow$  and  $\downarrow$  change the density of stripes.

You start from a random item after each choice.

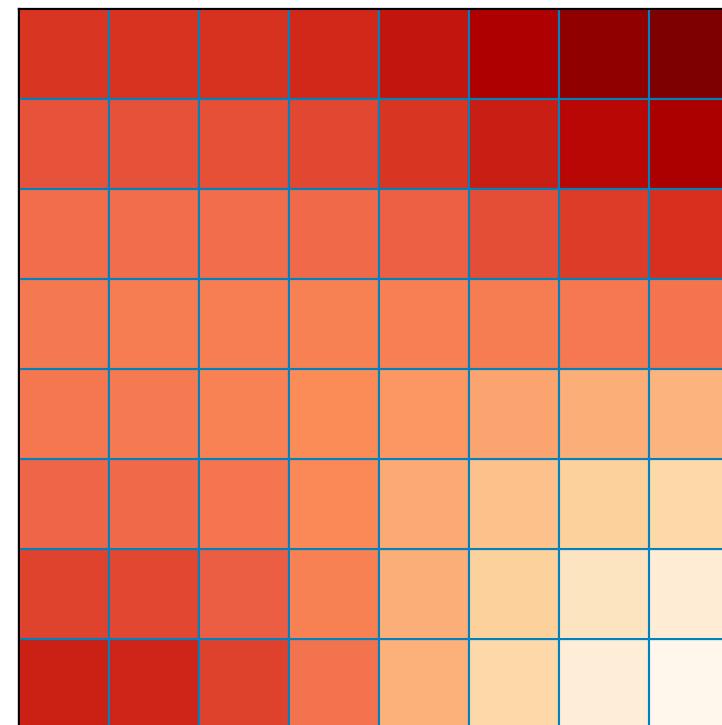
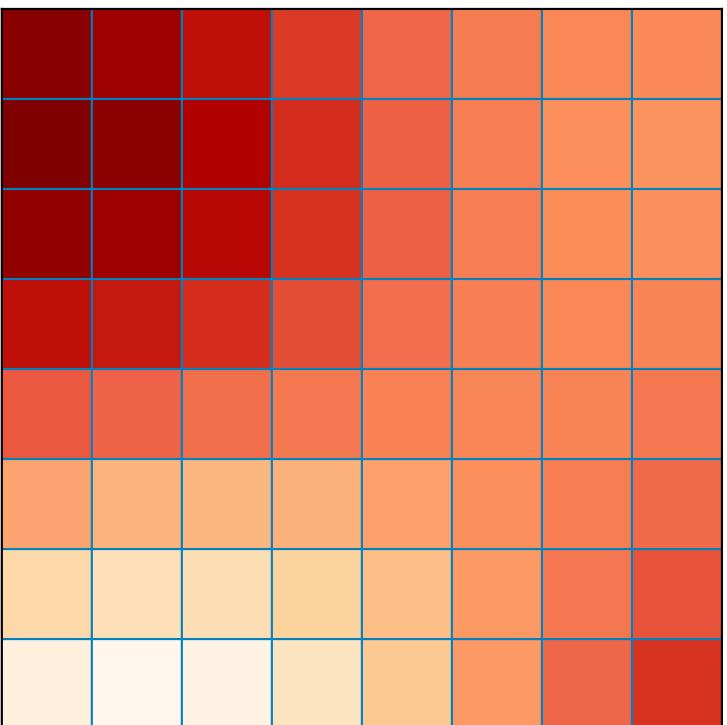
The training session is complete when you have completed at least 32 trials and achieved a run of 9 out of 10 correct responses.

Target

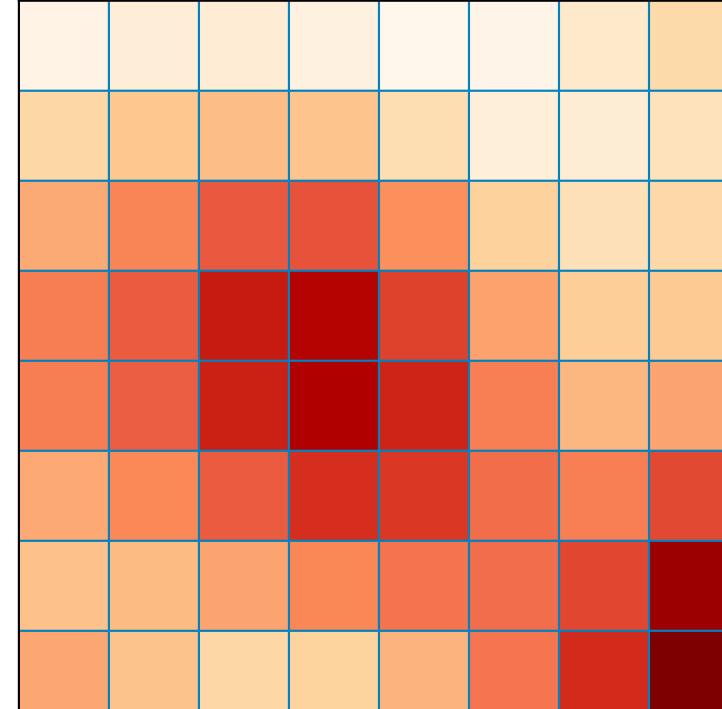
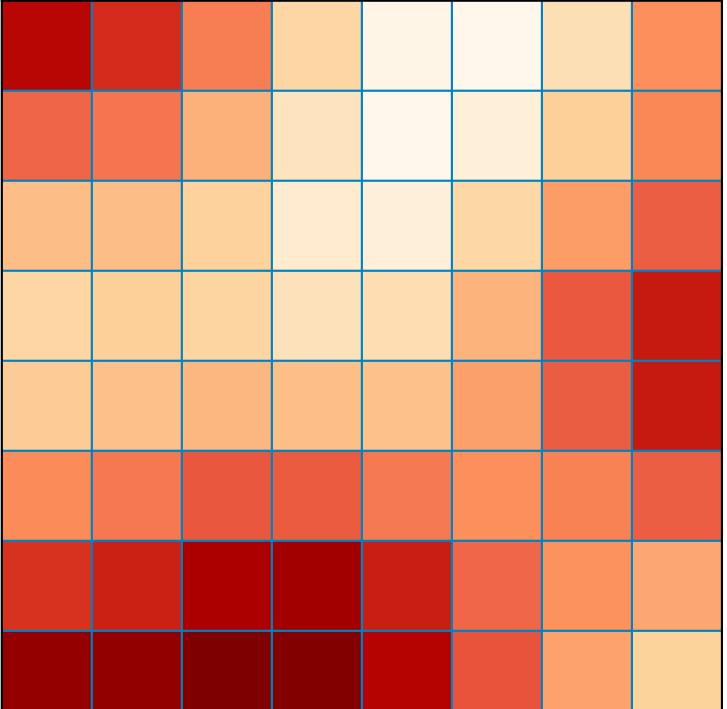


At least 32 trials AND a run of 9 out of 10 correct

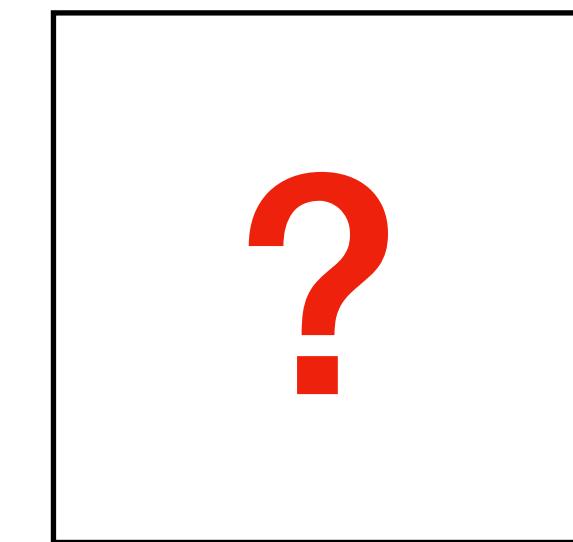
## Smooth Environments



## Rough Environments



Judgments on 10 unobserved stimuli



How many points do you think this item will earn?

45 points

How confident are you?

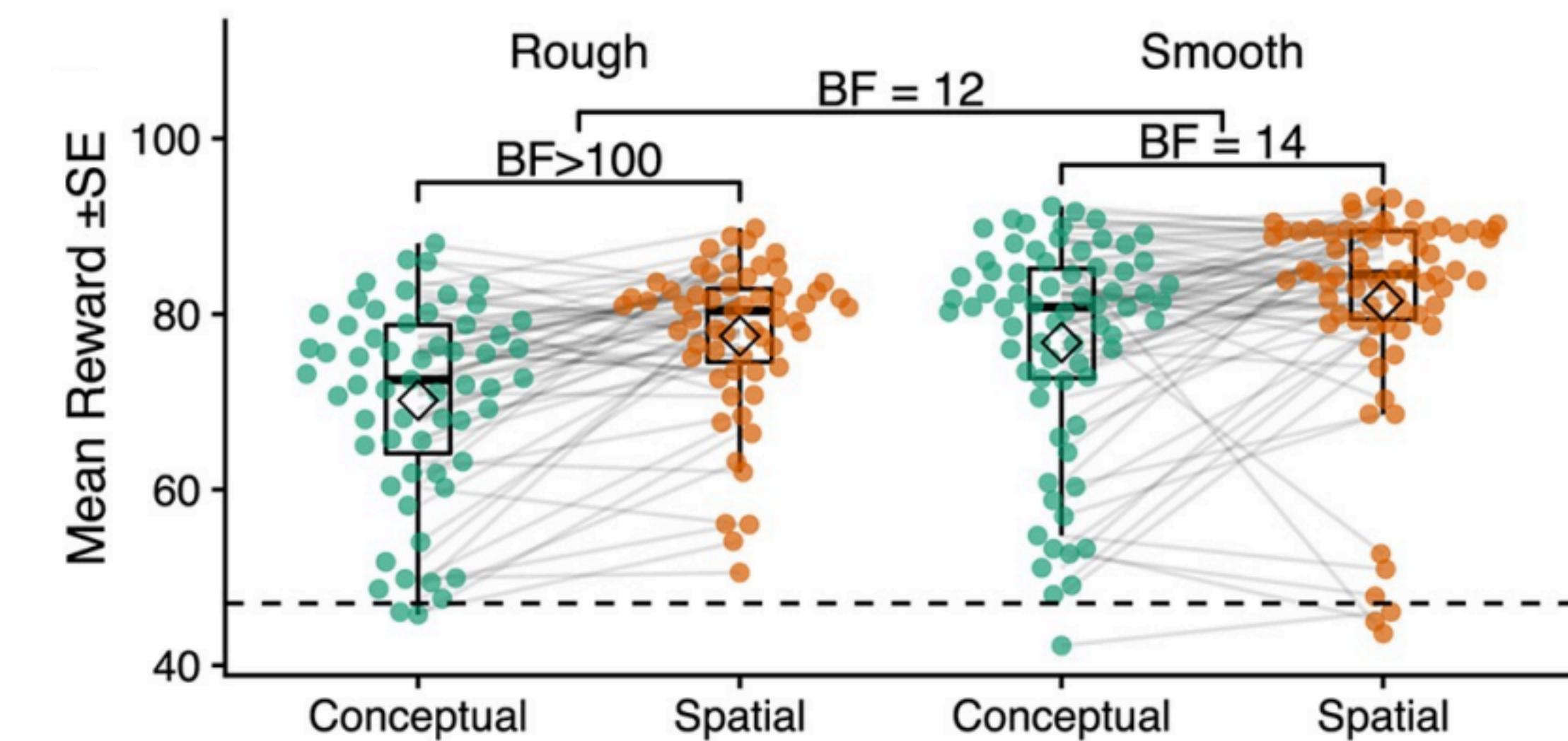
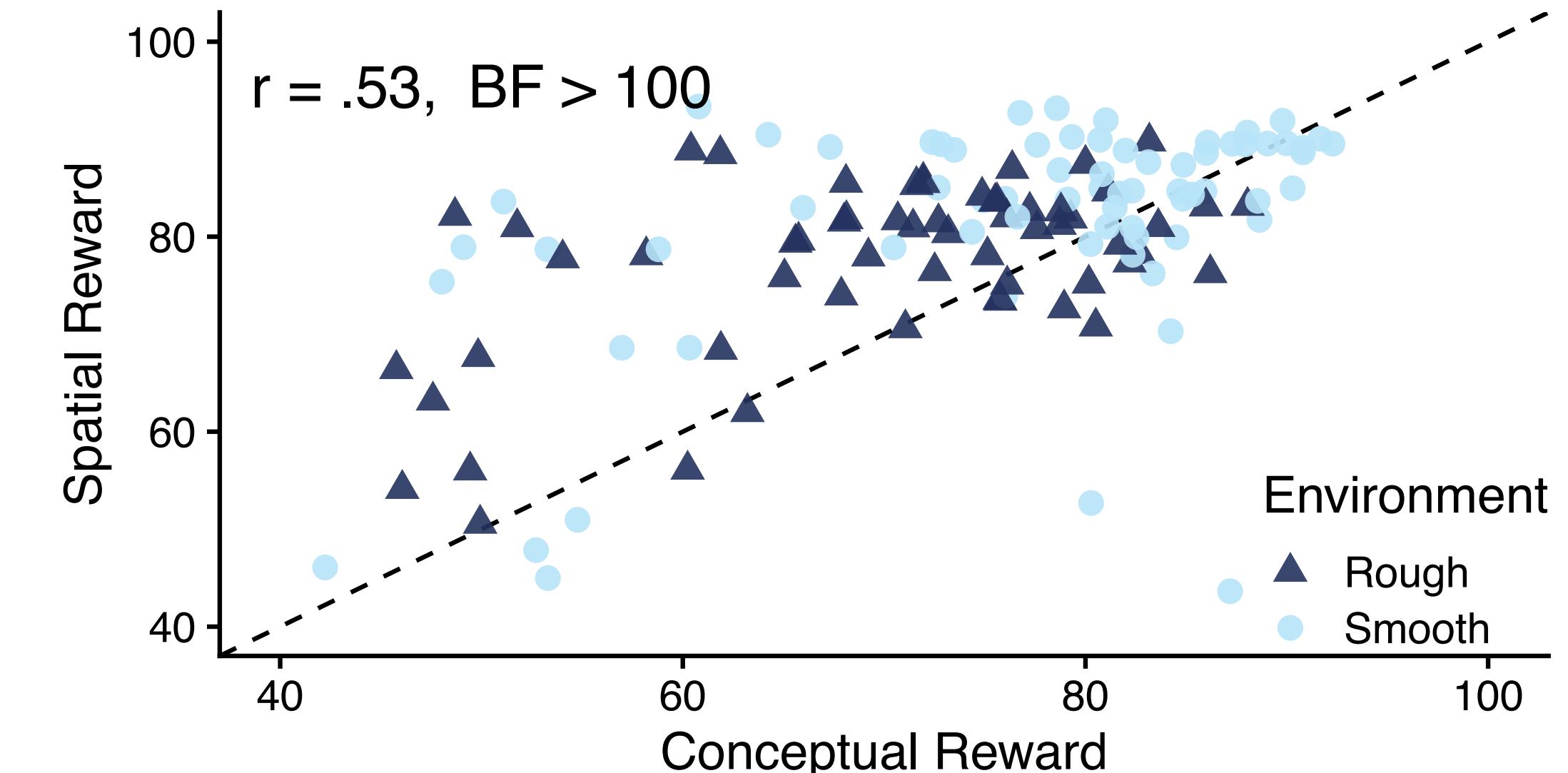
Least confident

Most confident

Submit

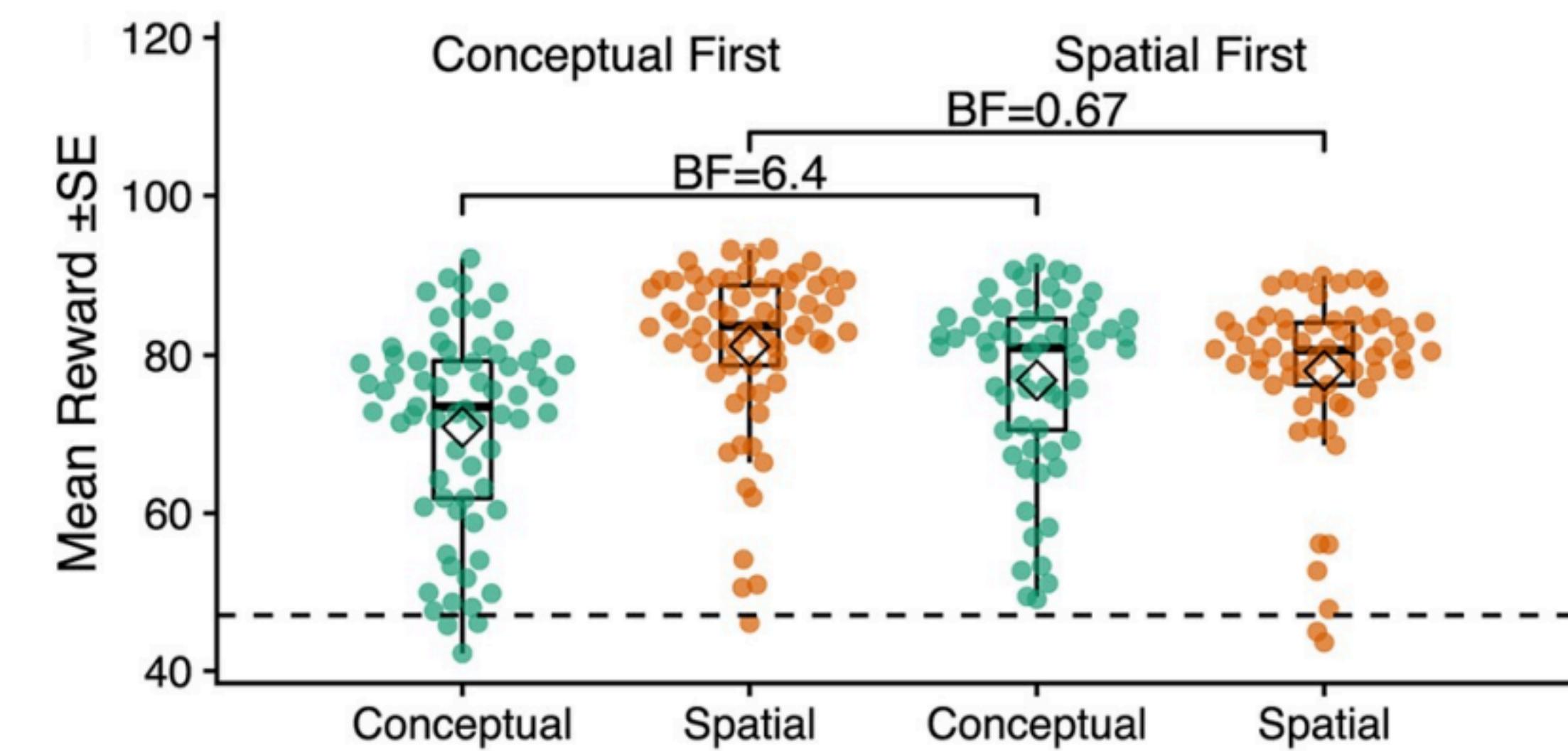
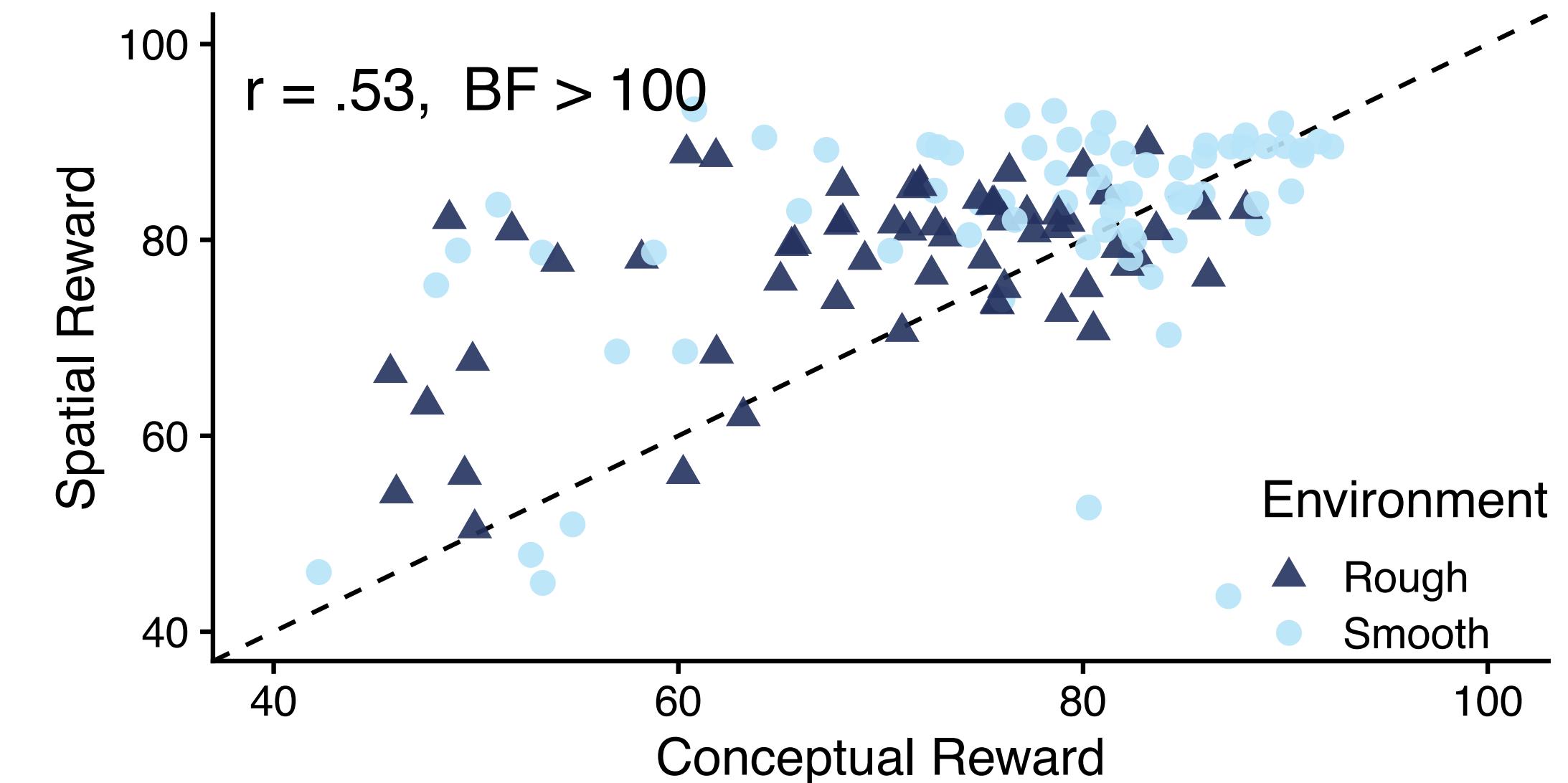
# Behavioral Results

- Correlated performance, but generally better in the spatial task
- This difference can largely be explained by a one-directional transfer effect:
  - Experience with spatial search boosted performance on conceptual search, but not vice versa



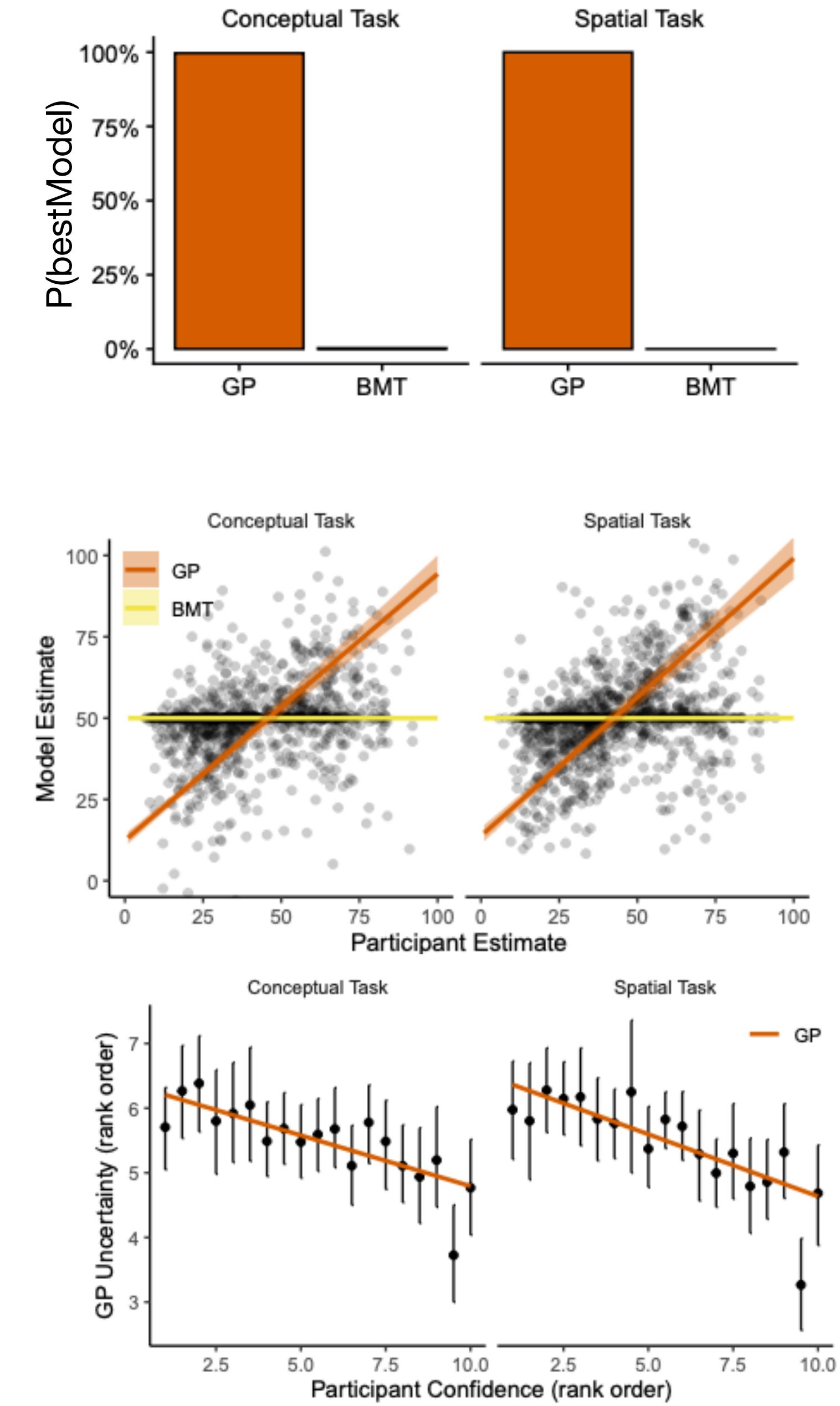
# Behavioral Results

- Correlated performance, but generally better in the spatial task
- This difference can largely be explained by a one-directional transfer effect:
  - Experience with spatial search boosted performance on conceptual search, but not vice versa



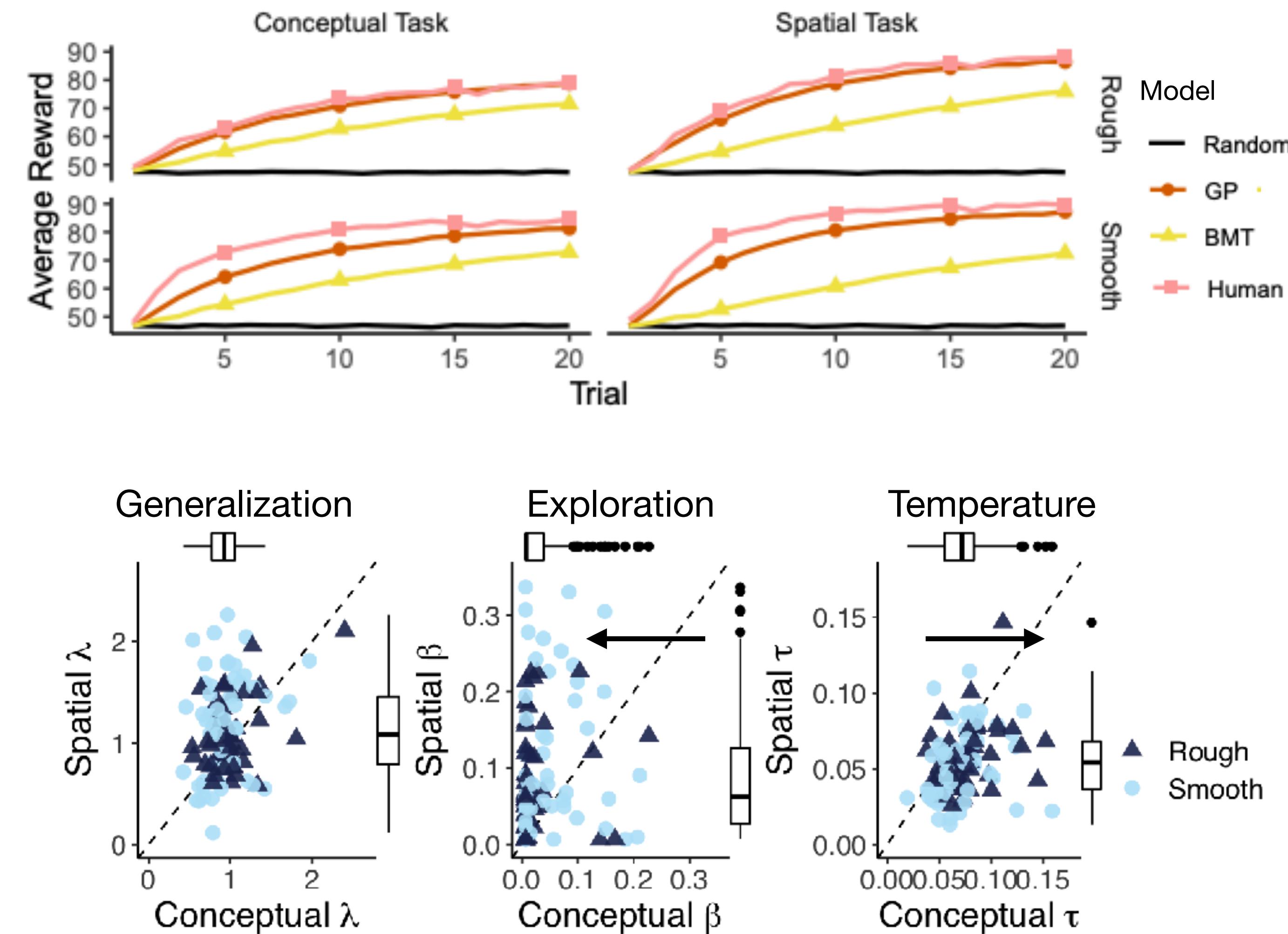
# Modeling Results

- Using group-level Bayesian estimation, we find that **GP-UCB** is the best model of choice behavior
  - $P(\text{bestModel})$  estimates the most prevalent model (corrected for chance); also known as protected exceedance probability
  - **GP-UCB** also predicts bonus round judgments about expected reward and confidence
    - using parameters estimated from rounds 1-9, we can use model simulations to predict participant judgments for unobserved stimuli in round 10
    - BMT makes invariant predictions for novel options, but the GP predictions correspond to participant judgments, where uncertainty is the opposite of confidence



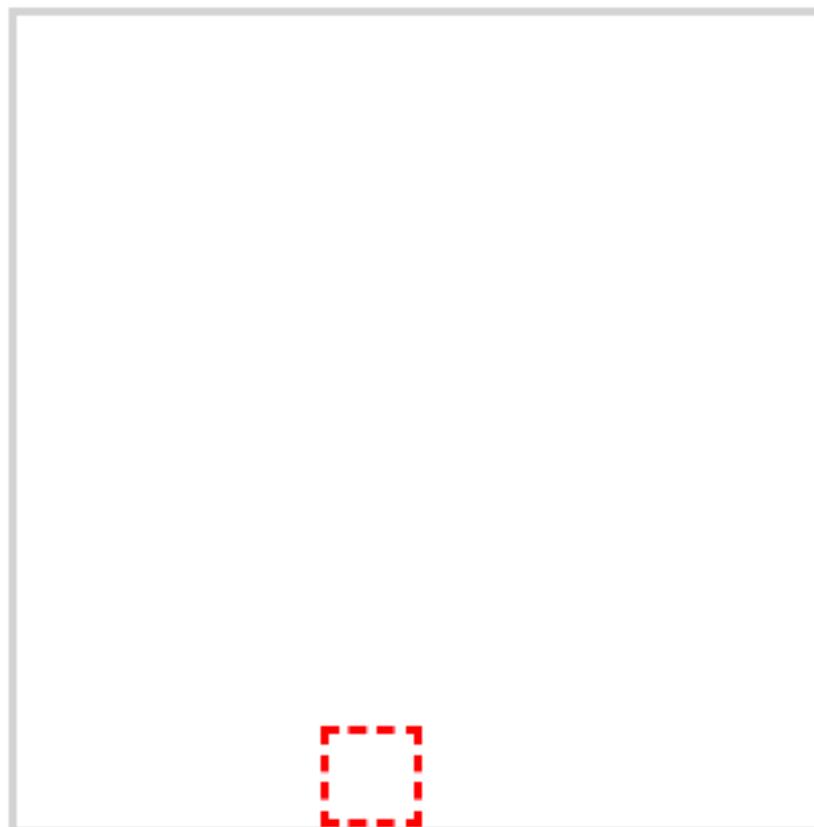
# Modeling Results

- GP-UCB simulated learning curves resemble human performance
- Parameter estimates show similar levels of generalization but change in exploration
- Directed exploration vanishes in the conceptual task, replaced by higher temperature (i.e., random) sampling



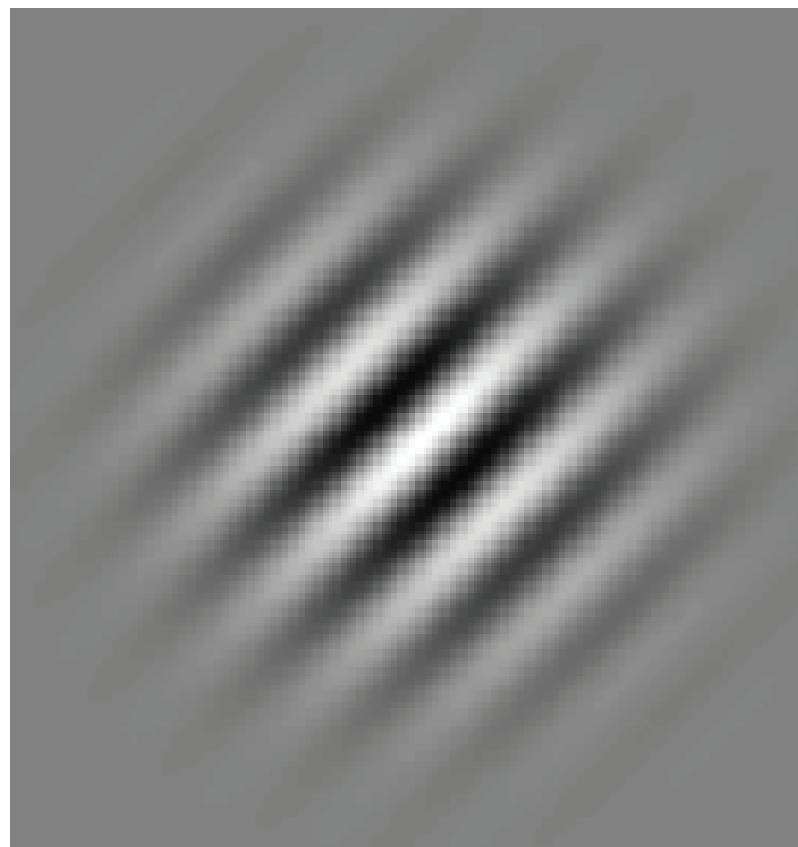
# Summary

Spatial

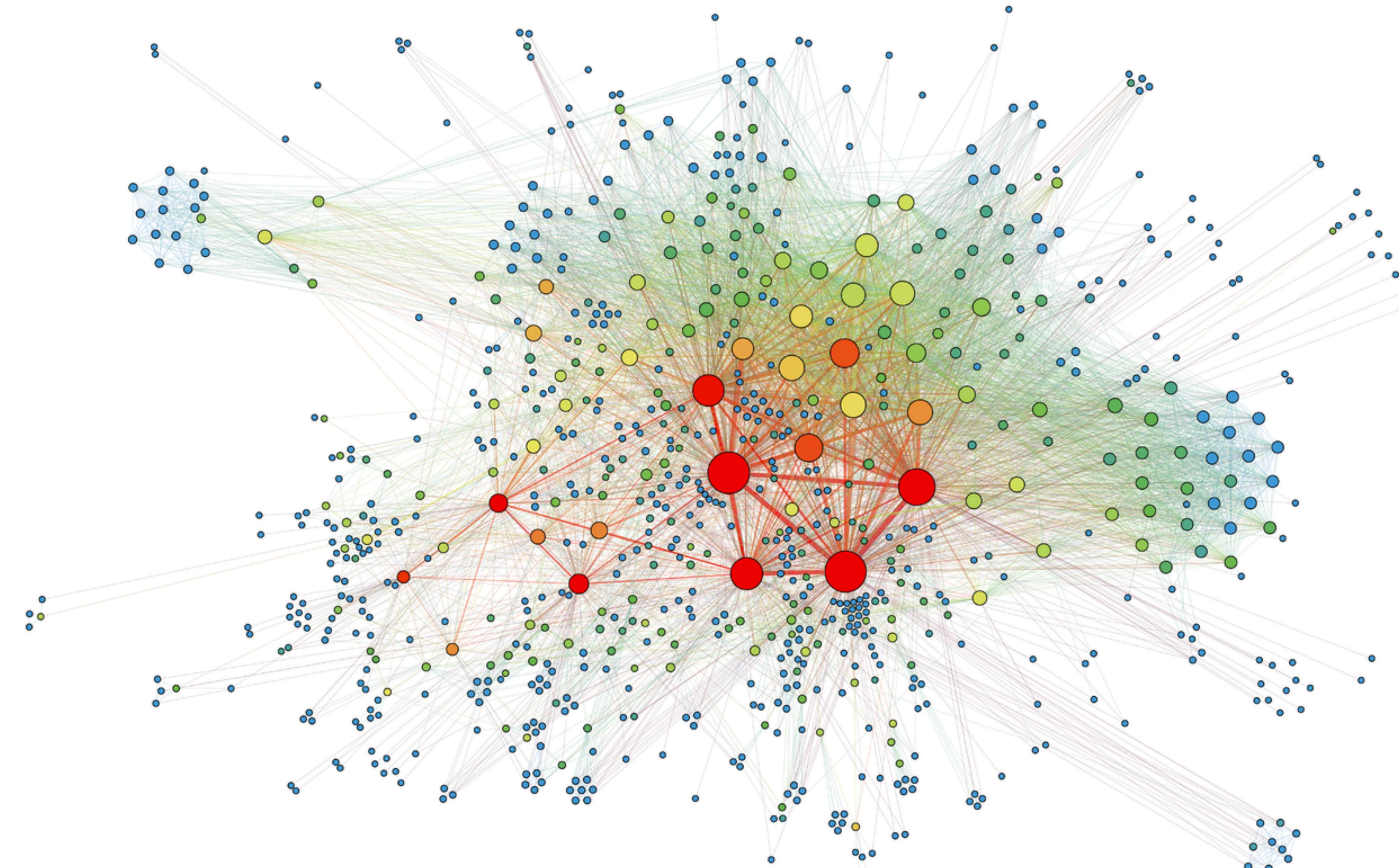


1 day  
↔  
gap

Conceptual

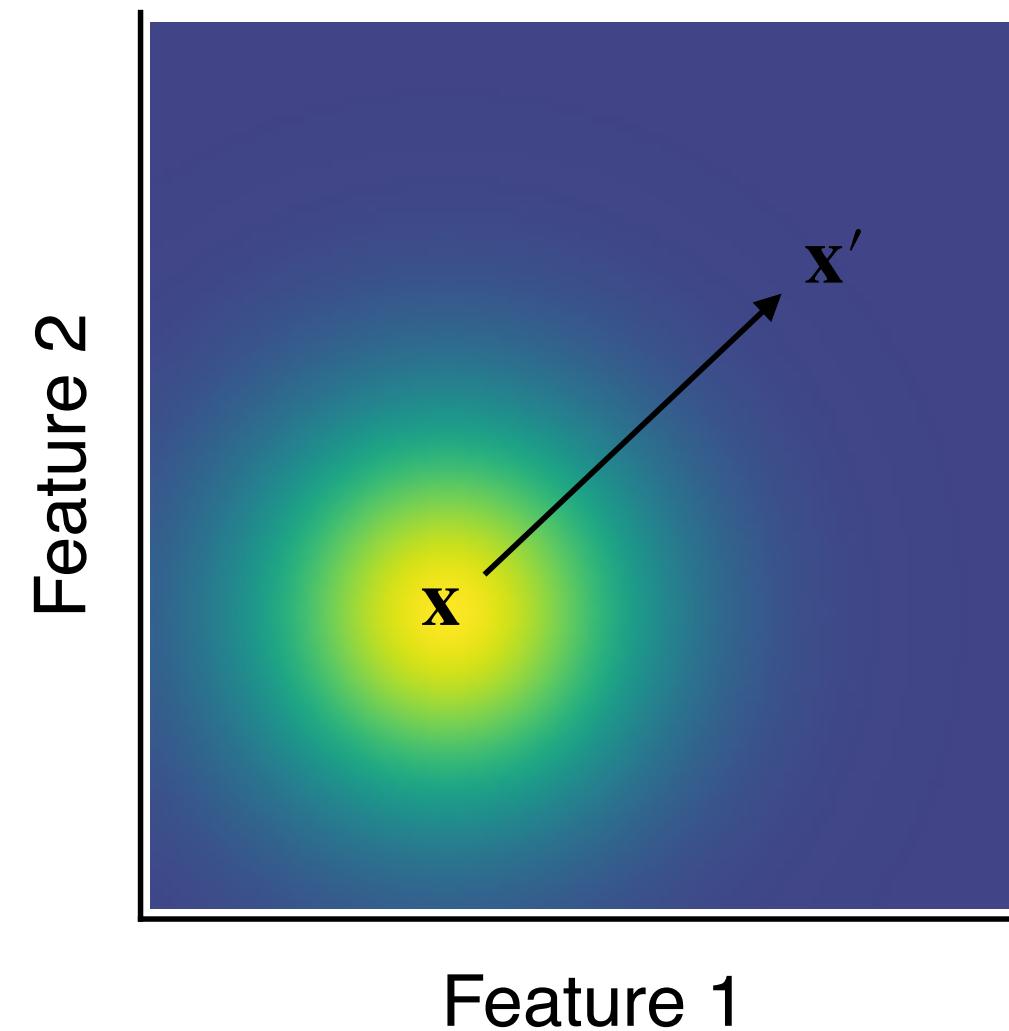


- Similar mechanisms of generalization-guided search in both domains
- But also diagnostic differences:
  - One-directional transfer effect suggest something fundamental about spatial reasoning
  - Switch from directed to random exploration
- Bonus round suggests participants have a good sense of uncertainty (confidence ratings) in both domains, but aren't able to leverage it to direct their exploration in the conceptual task

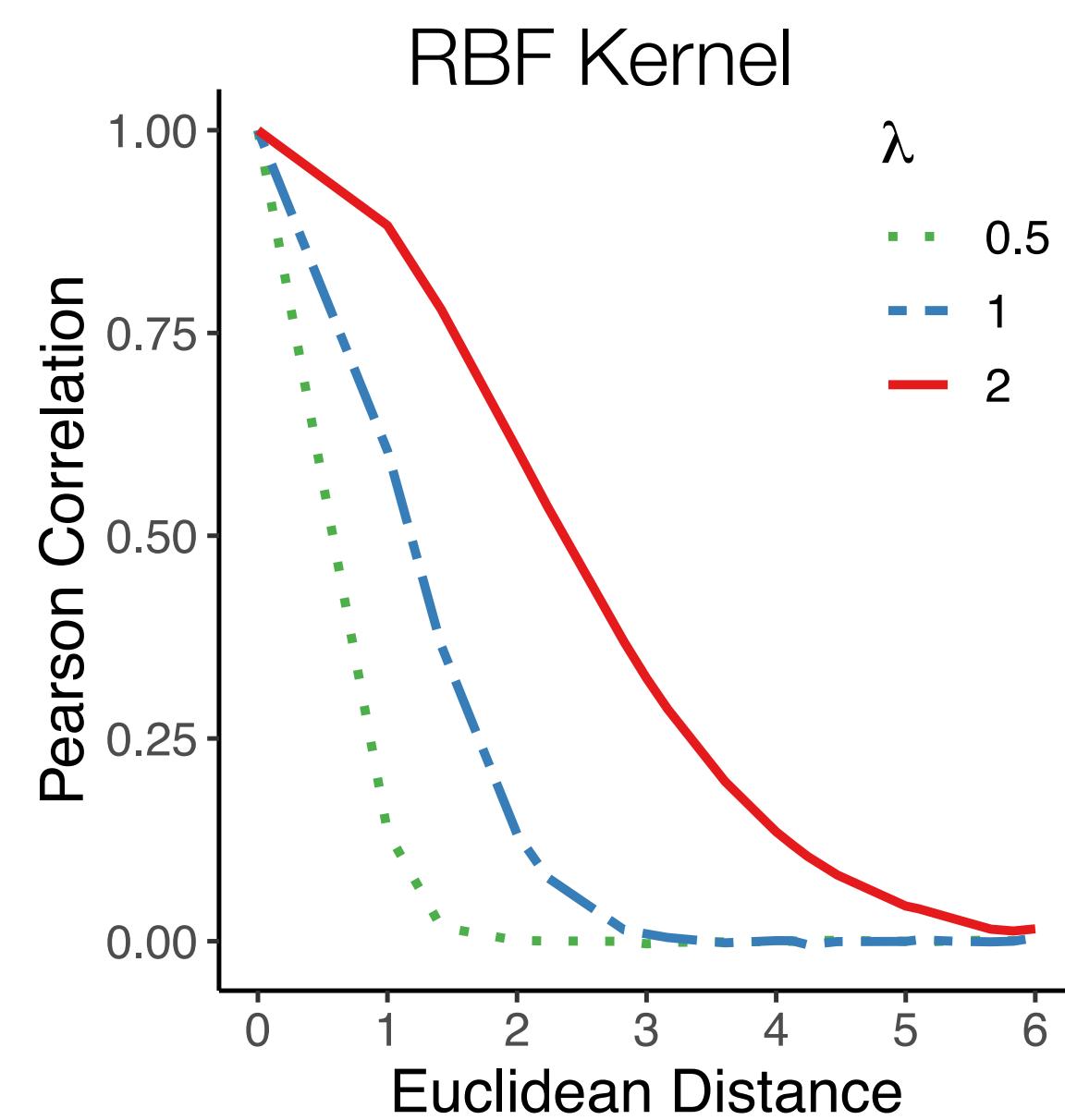
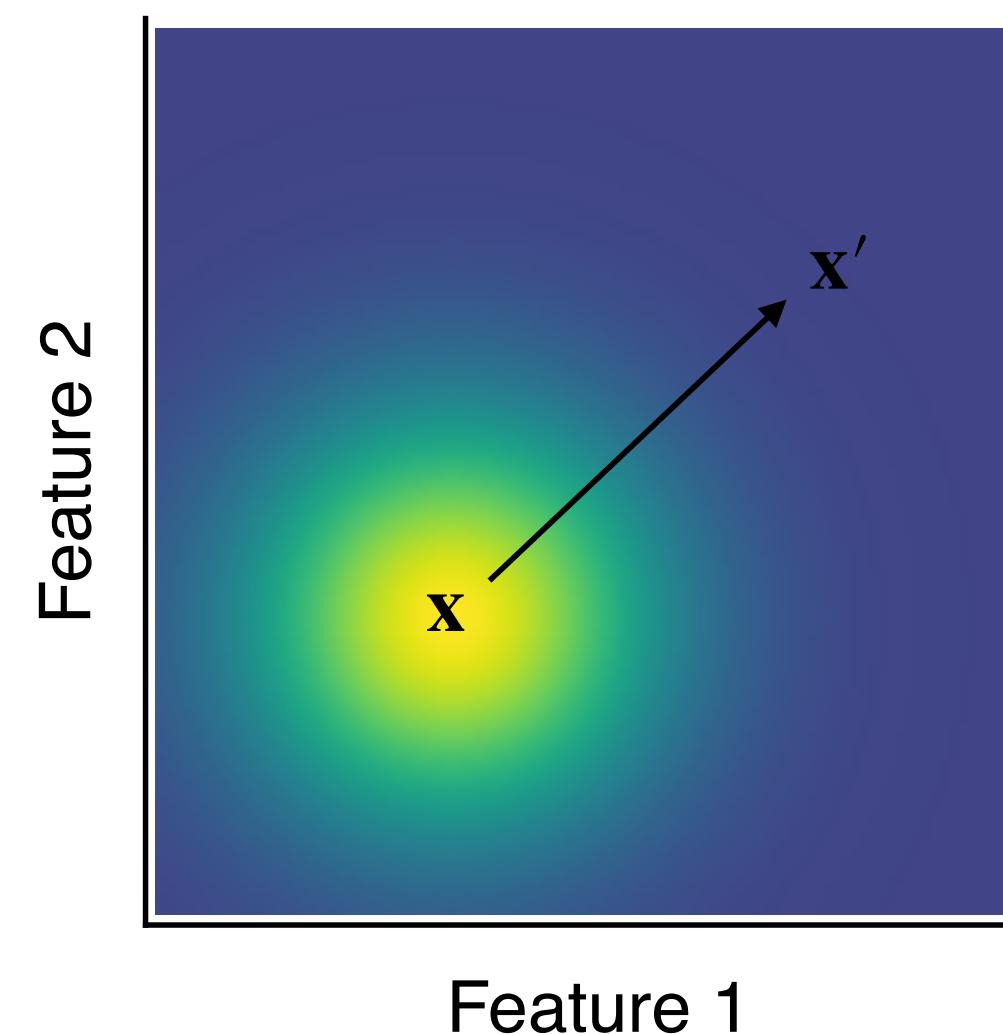


# Exploring Structured Spaces

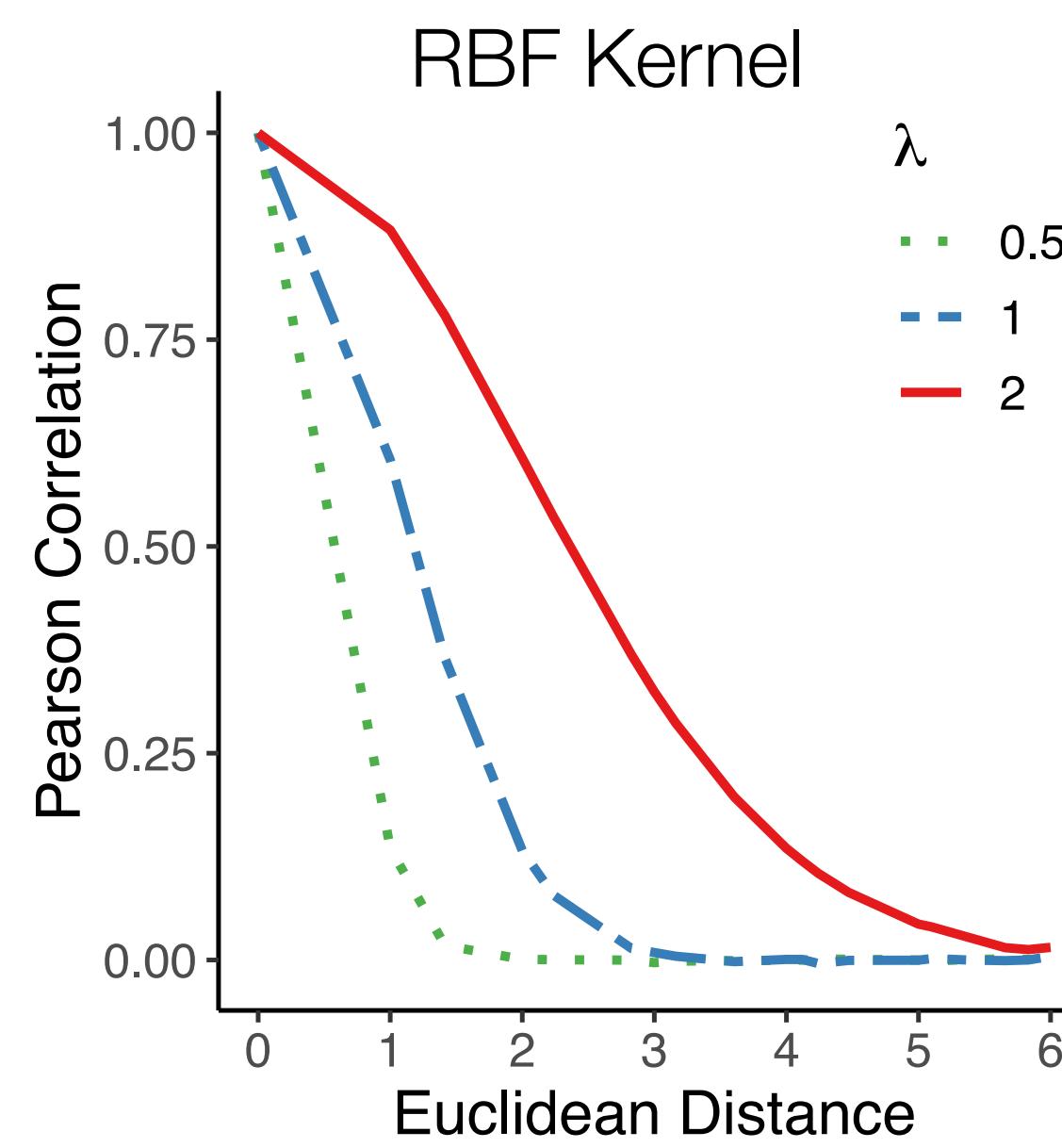
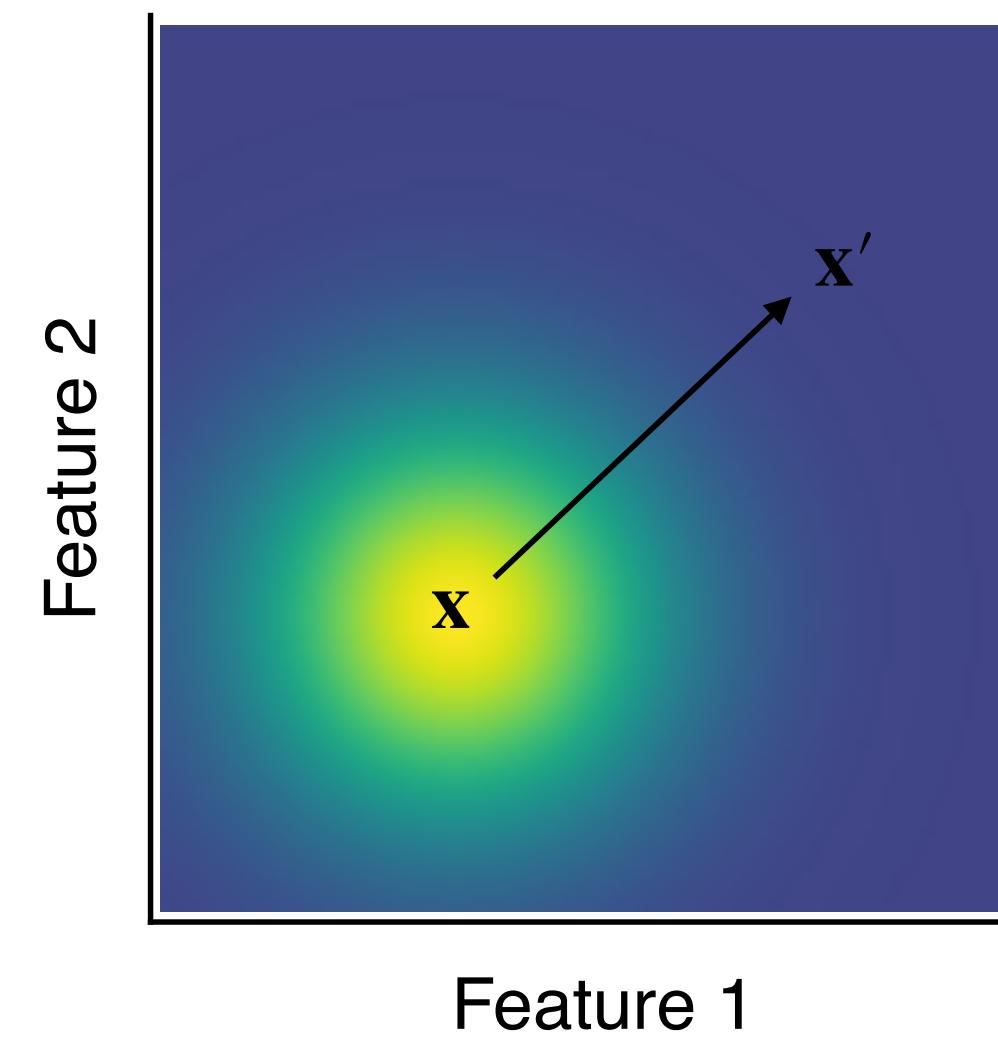
# From continuous to structured spaces



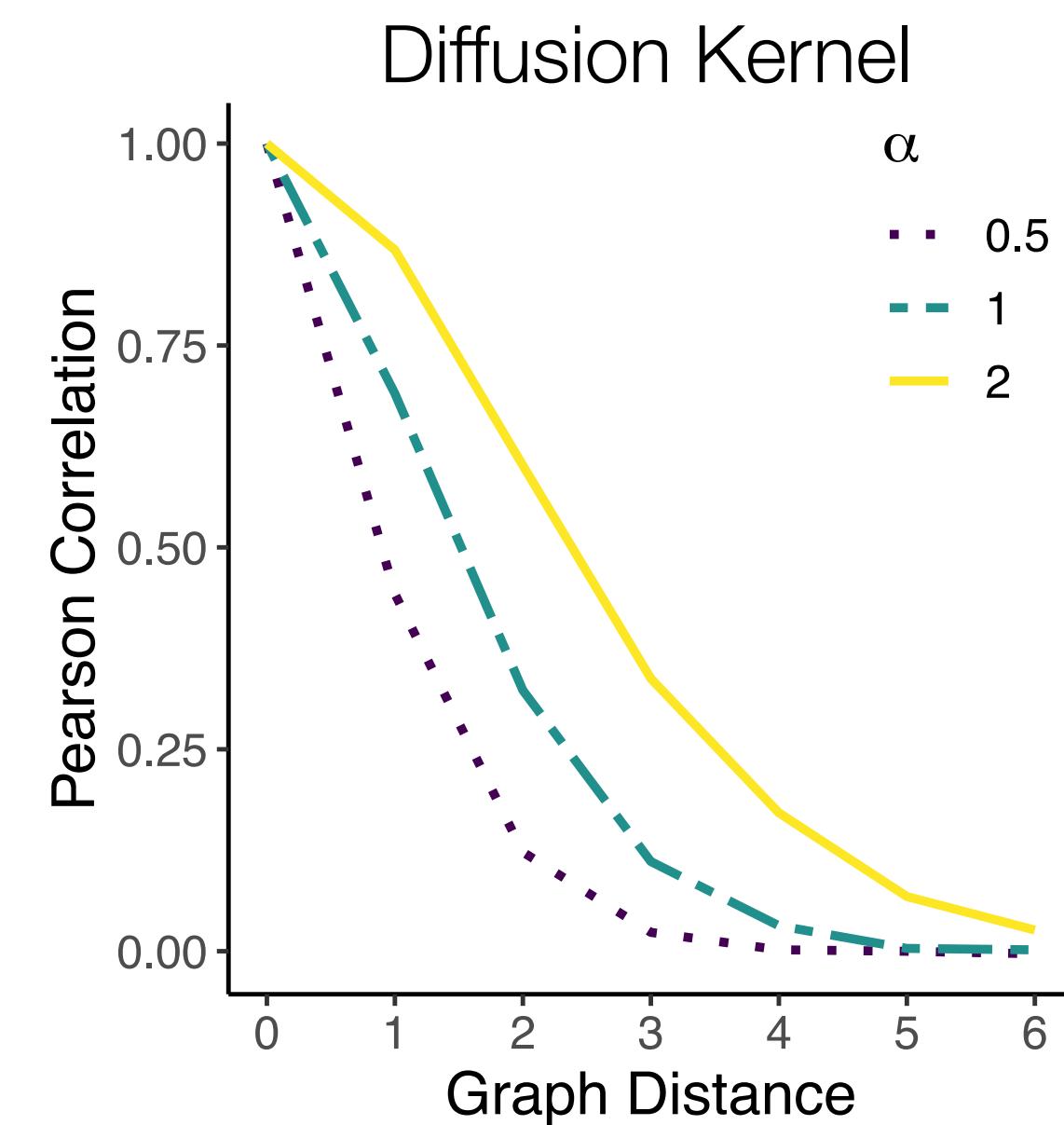
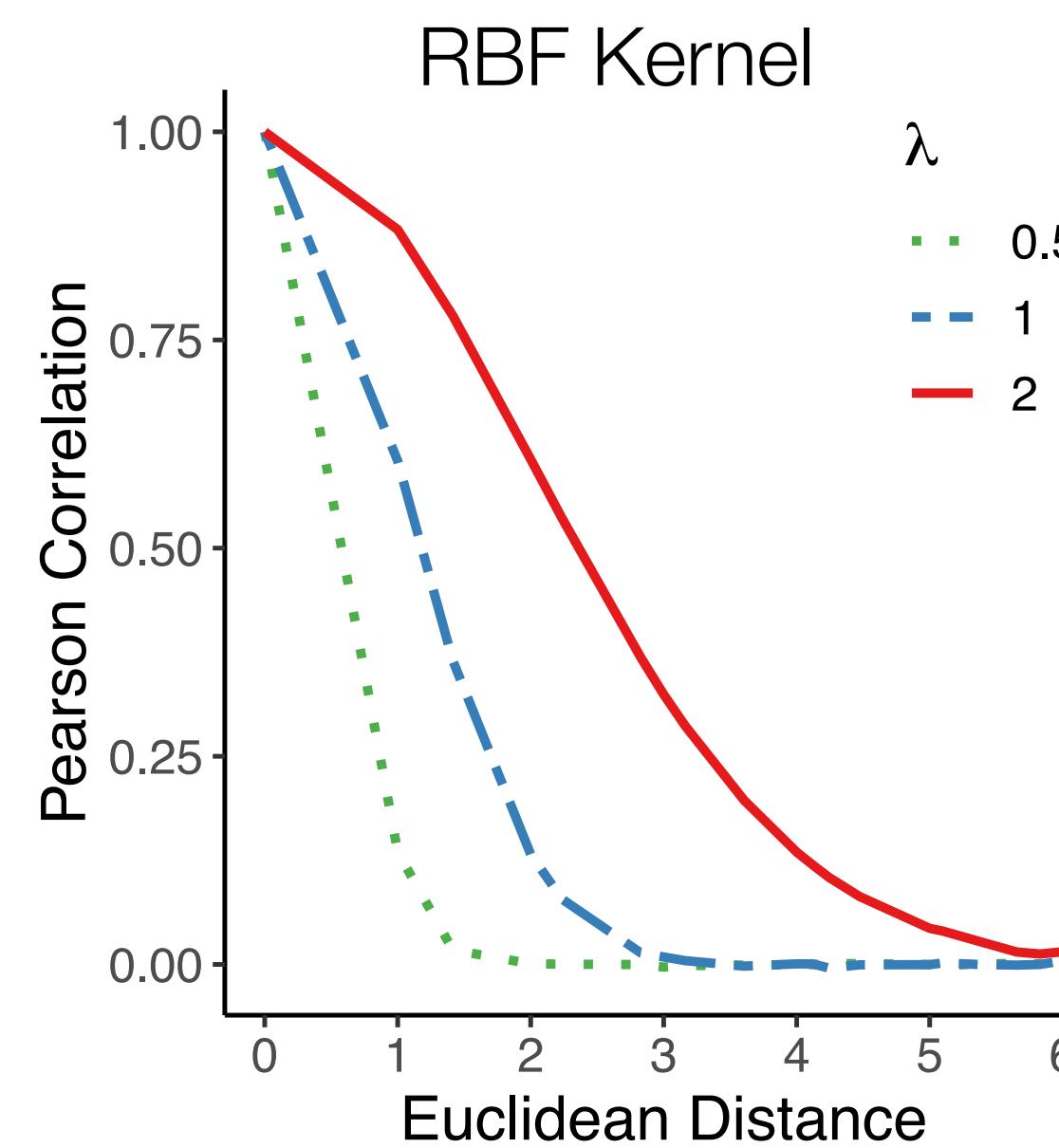
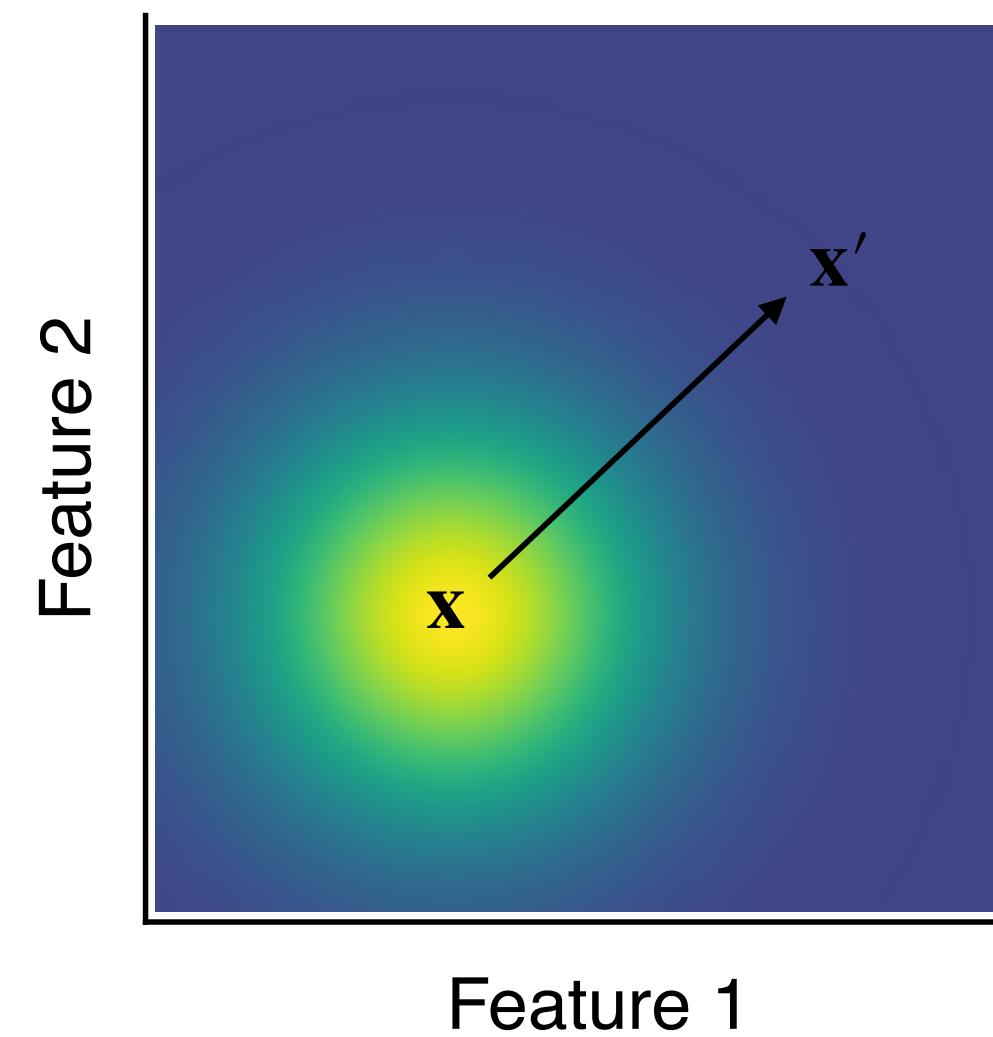
# From continuous to structured spaces



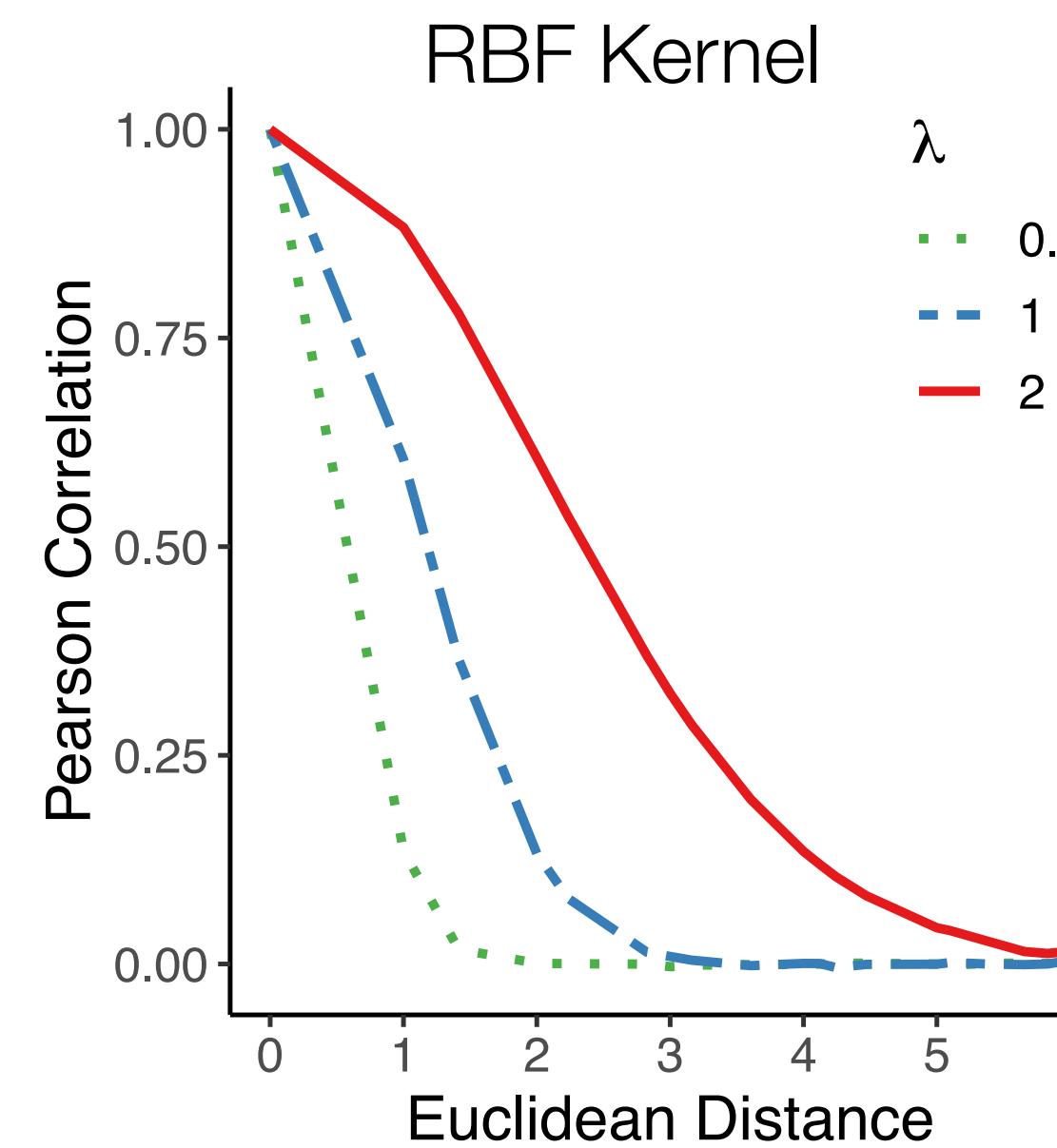
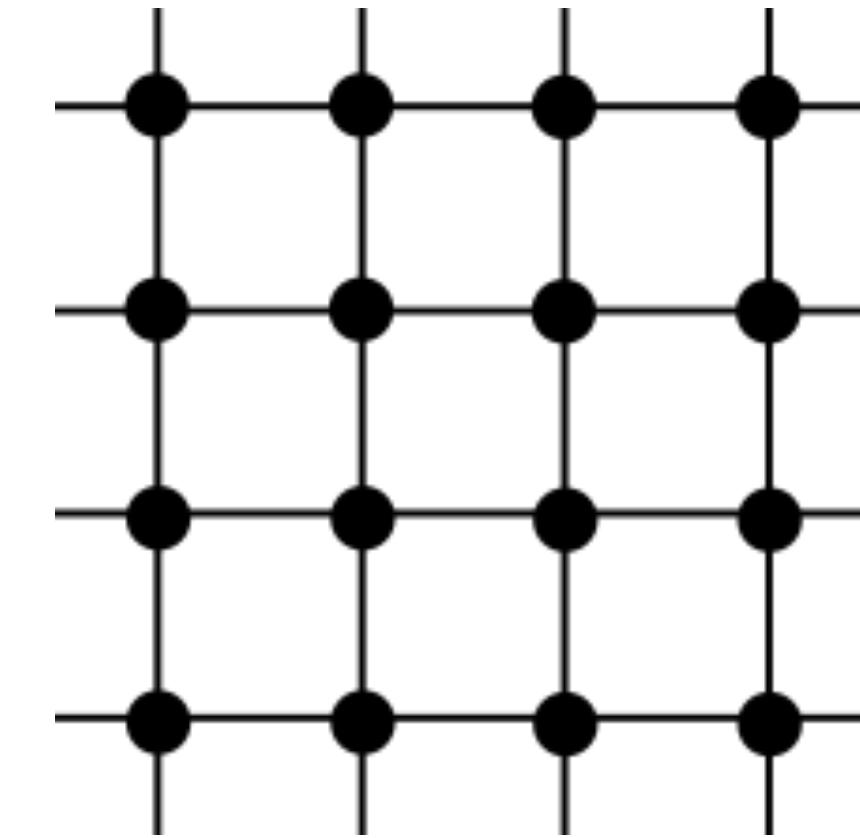
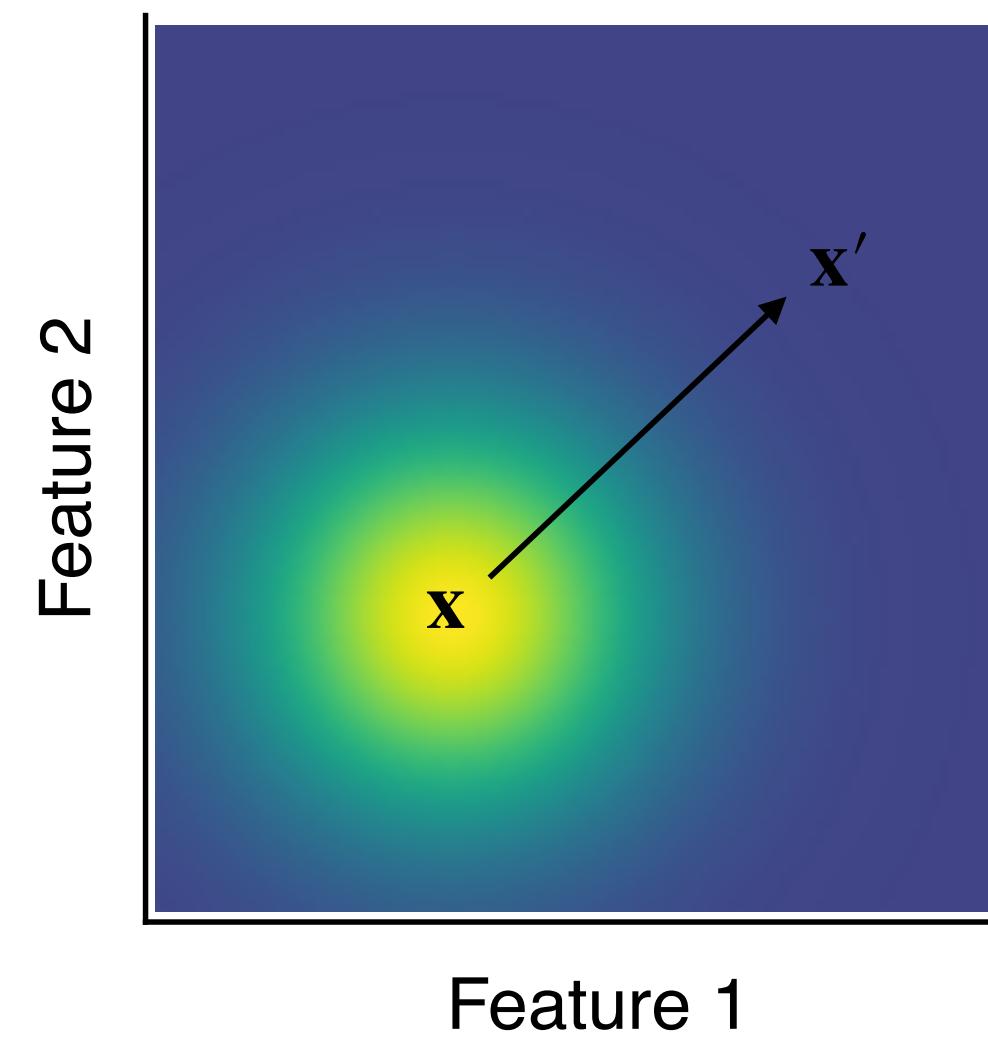
# From continuous to structured spaces



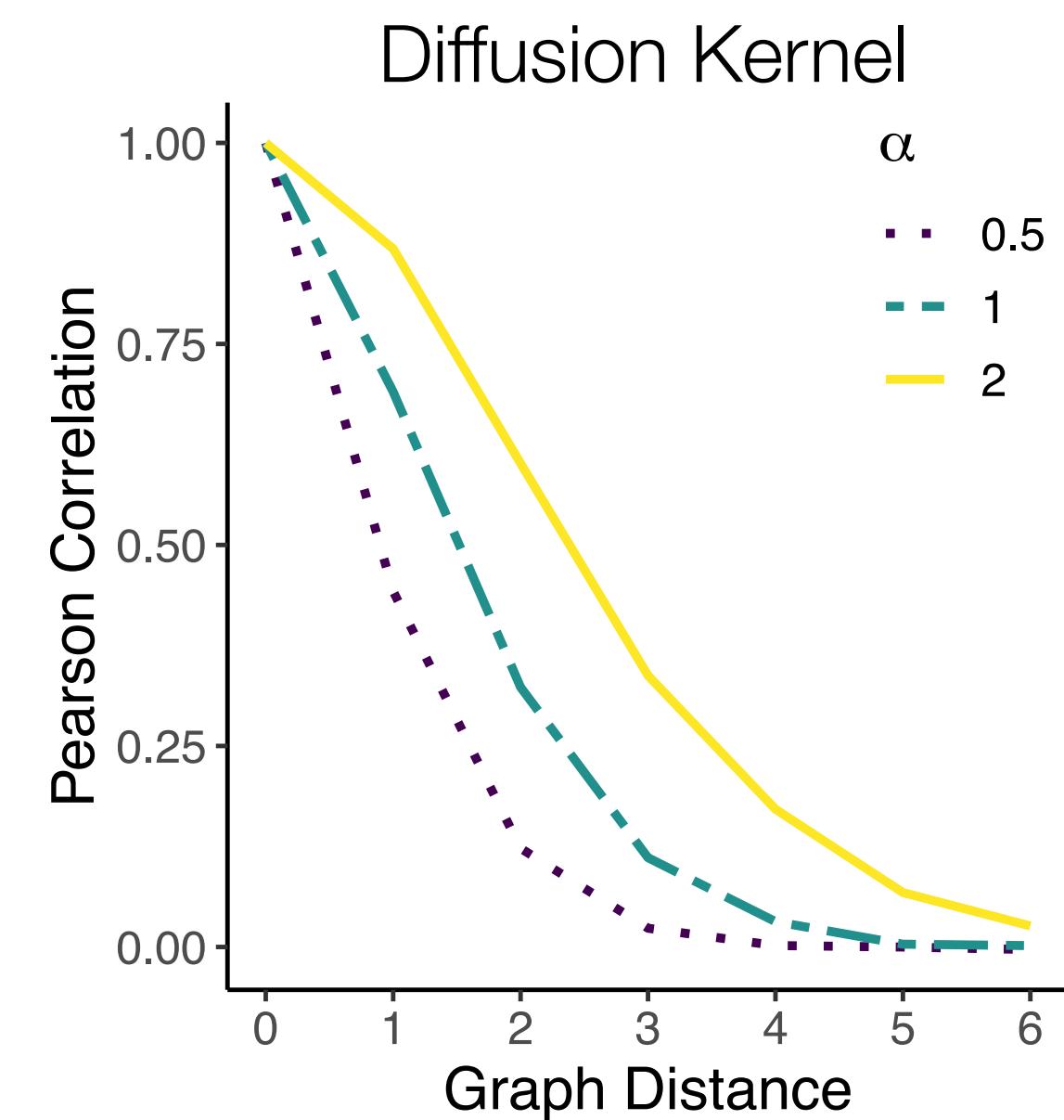
# From continuous to structured spaces



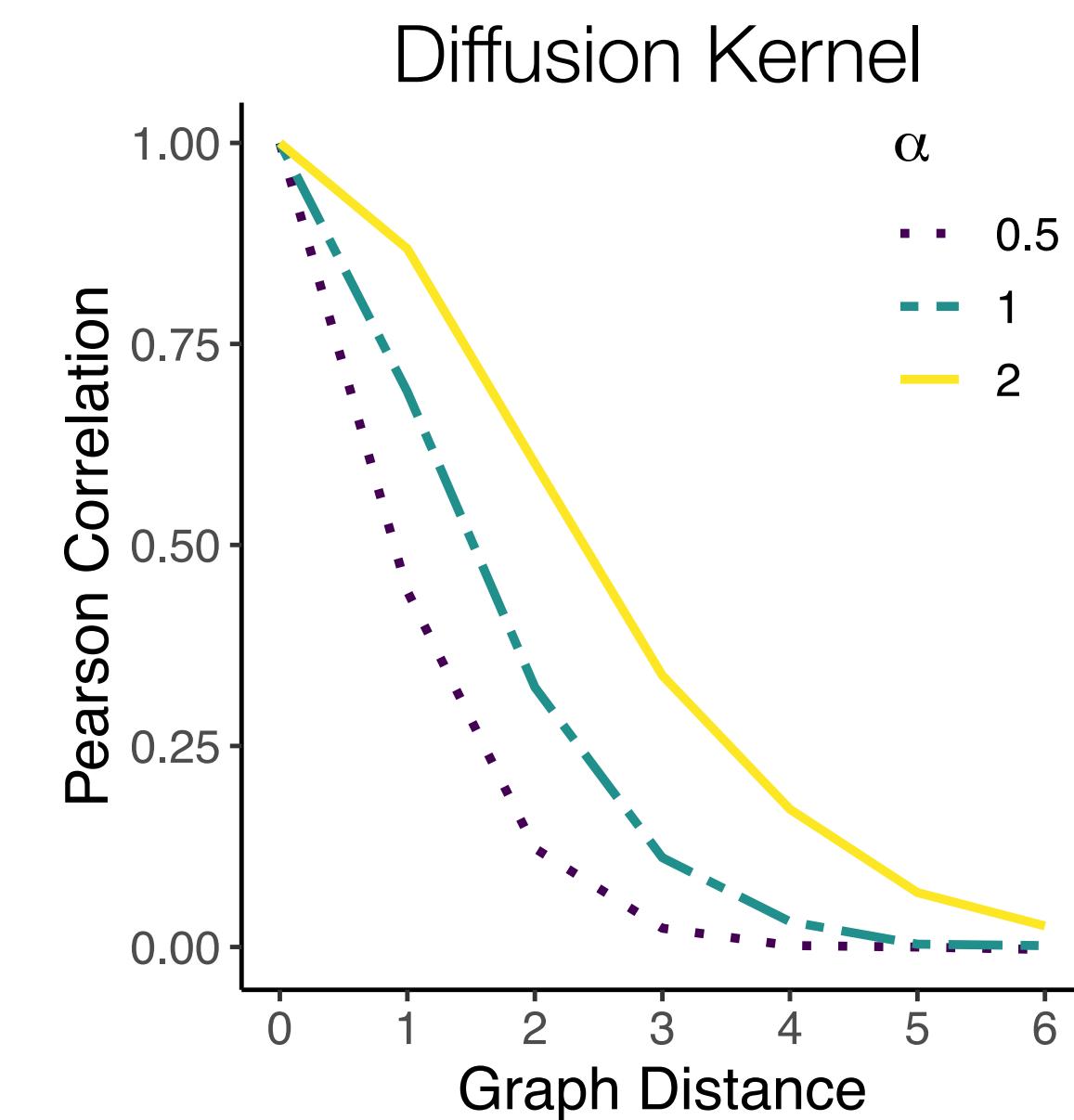
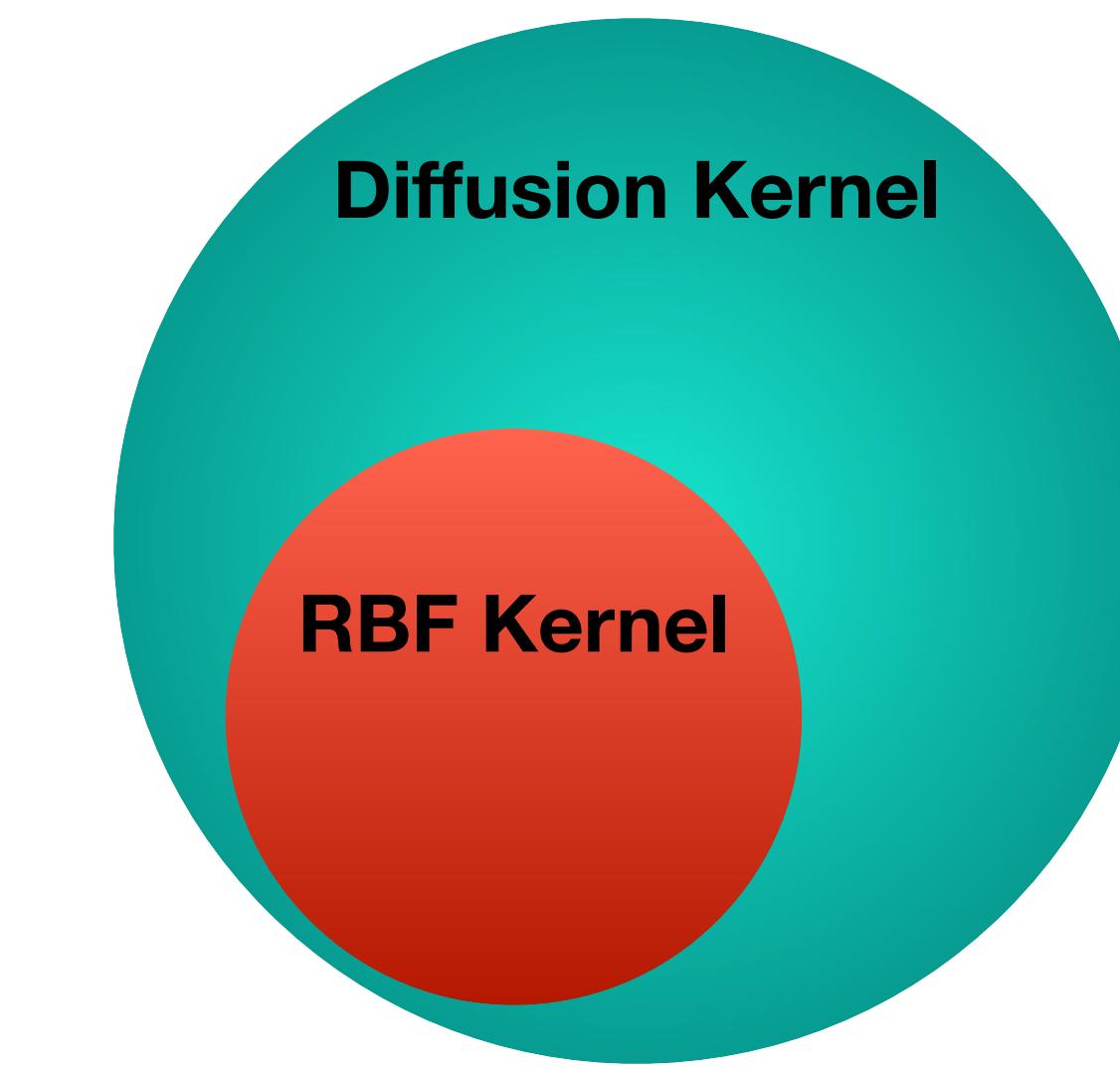
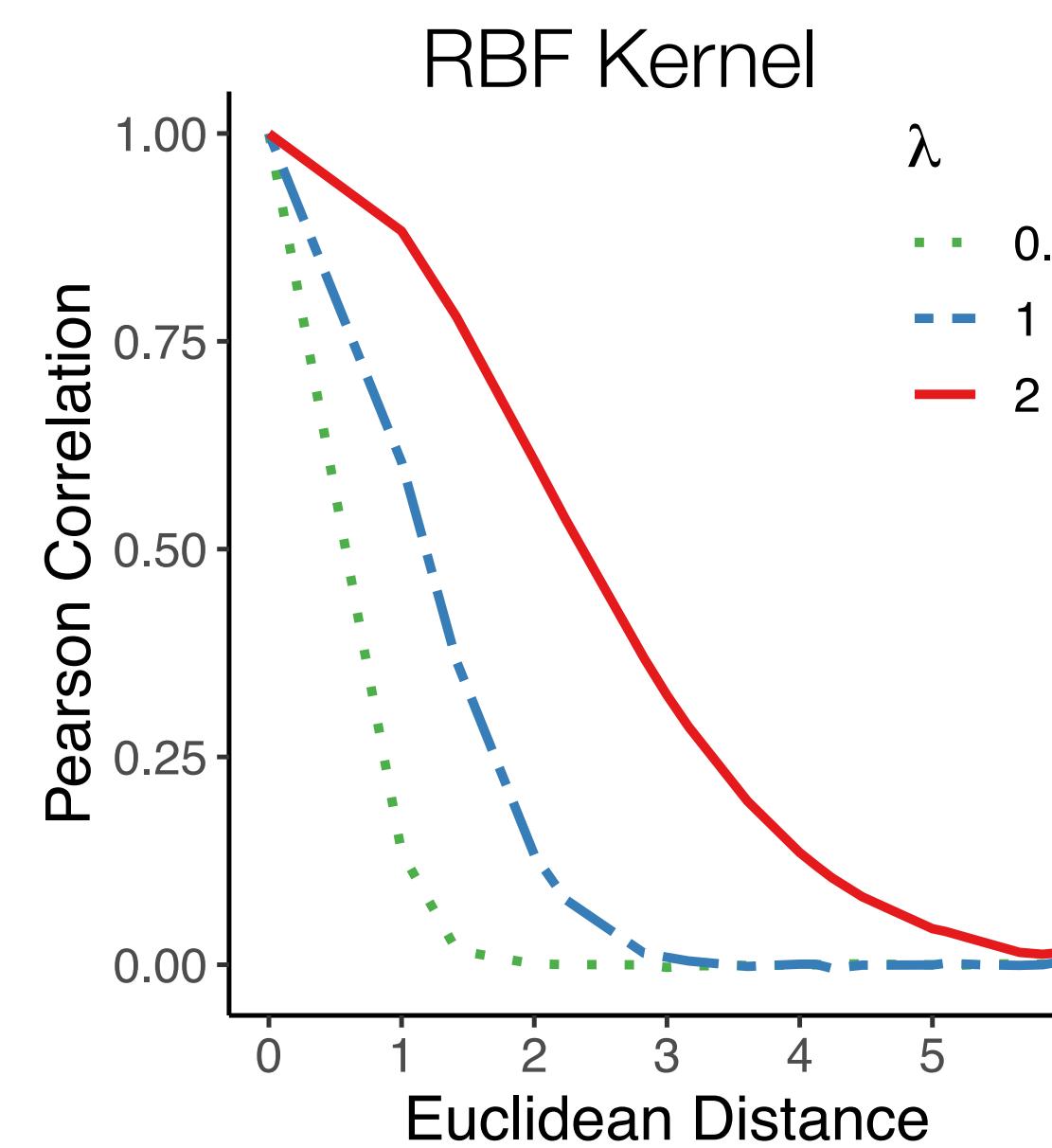
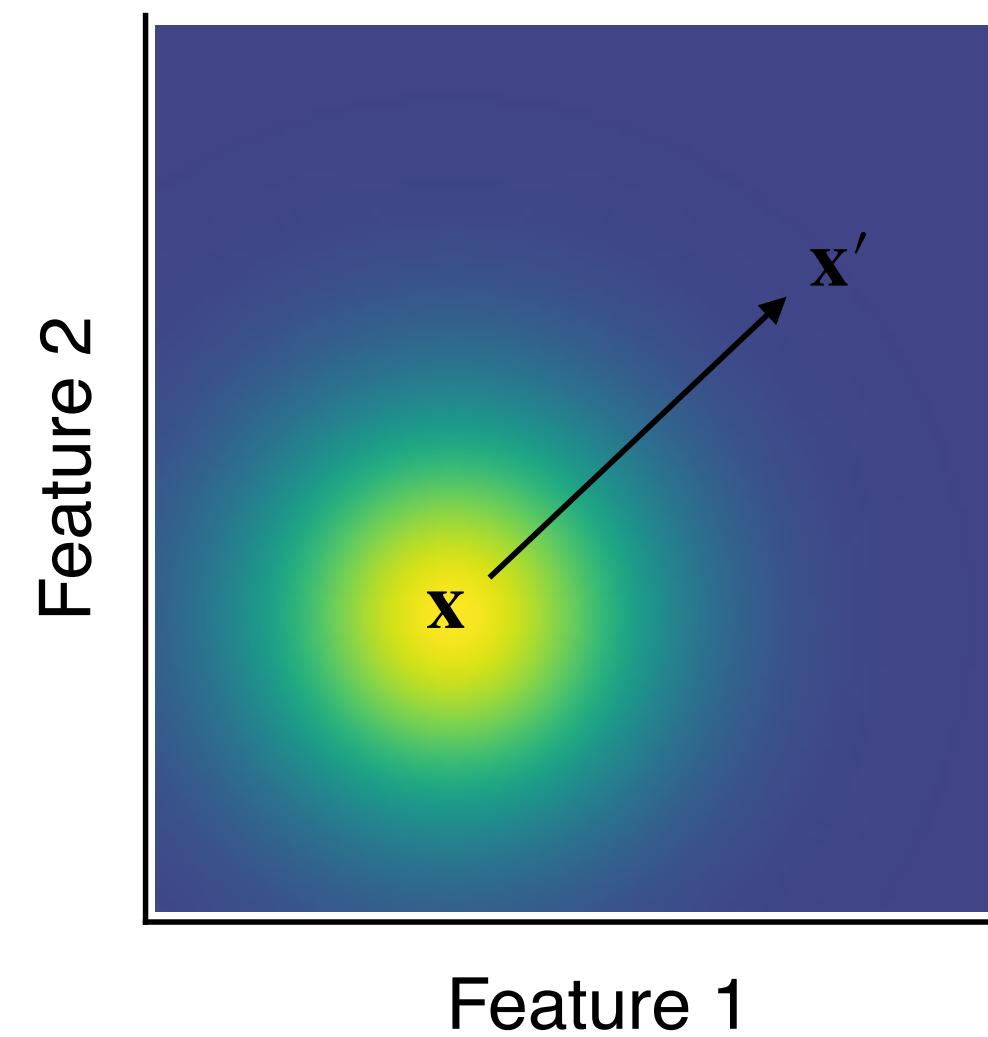
# From continuous to structured spaces

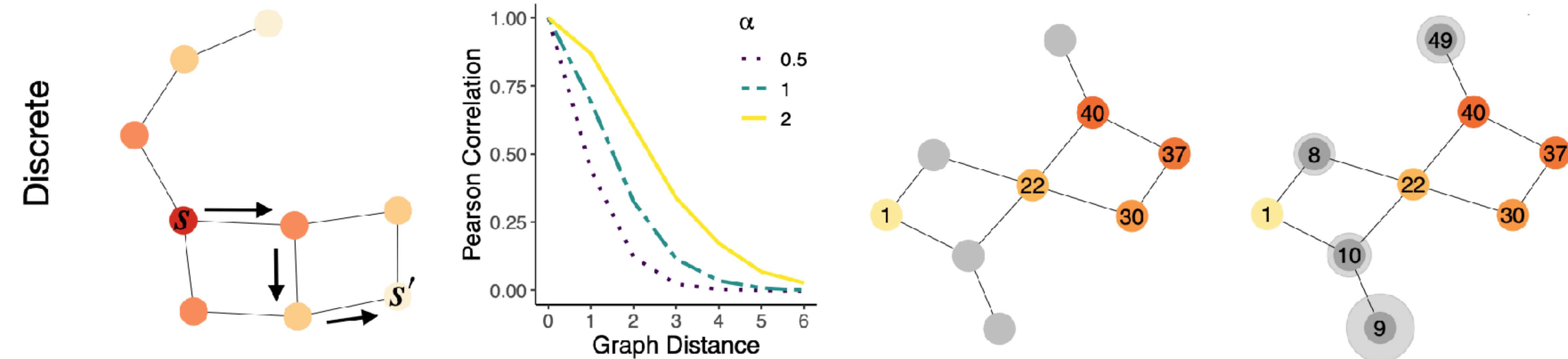
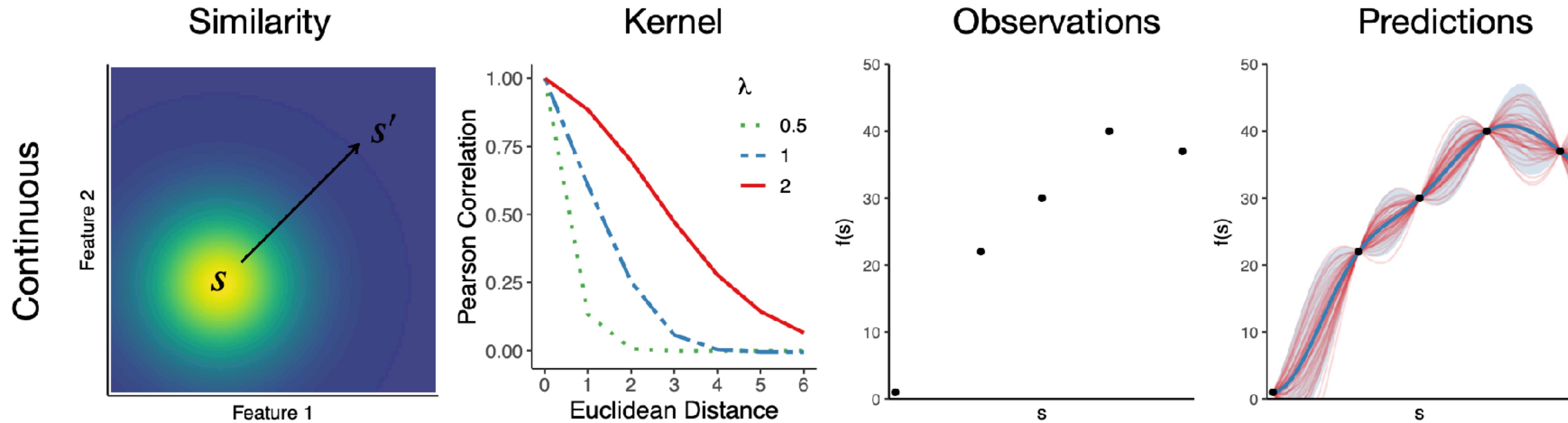


==

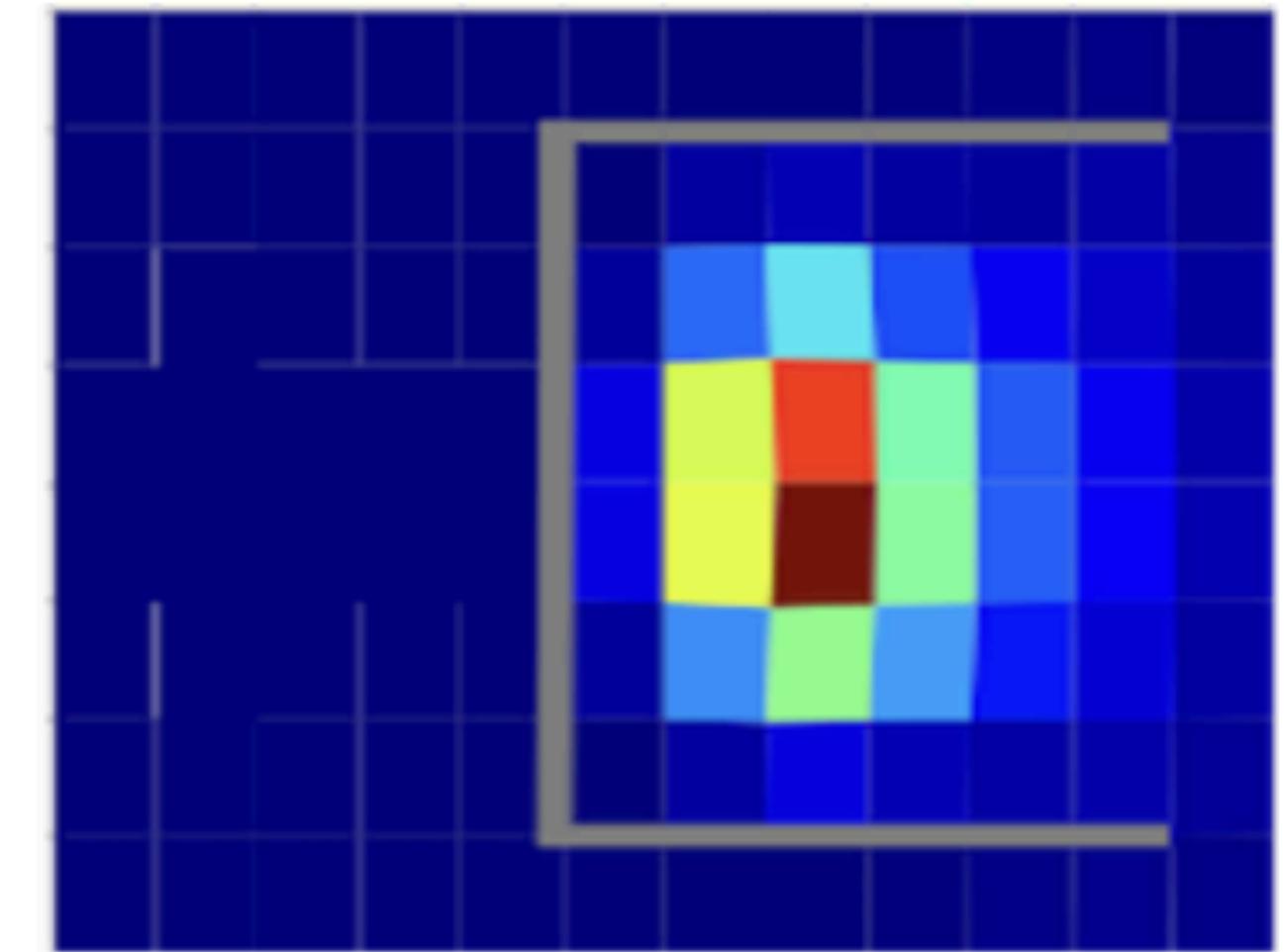
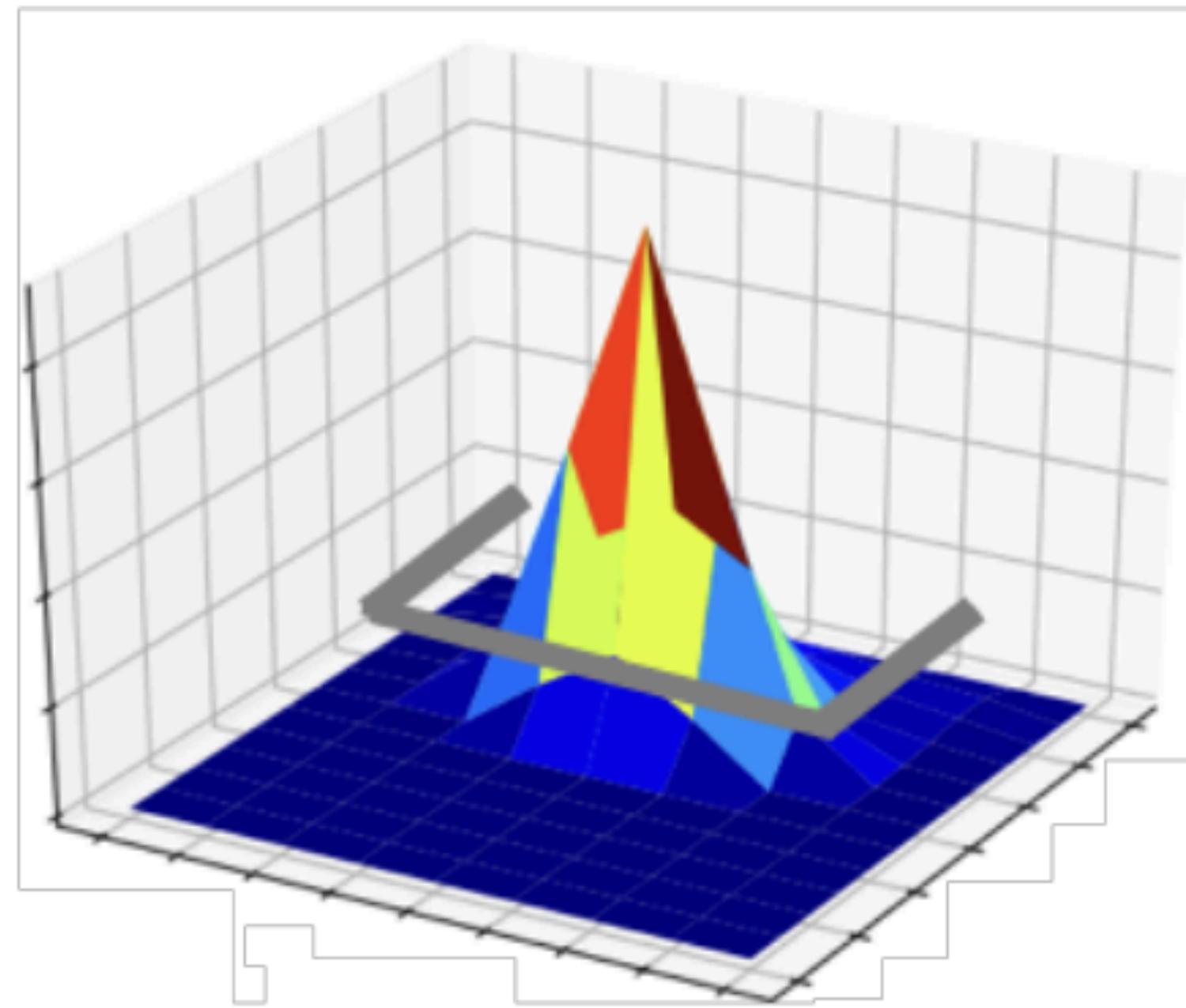
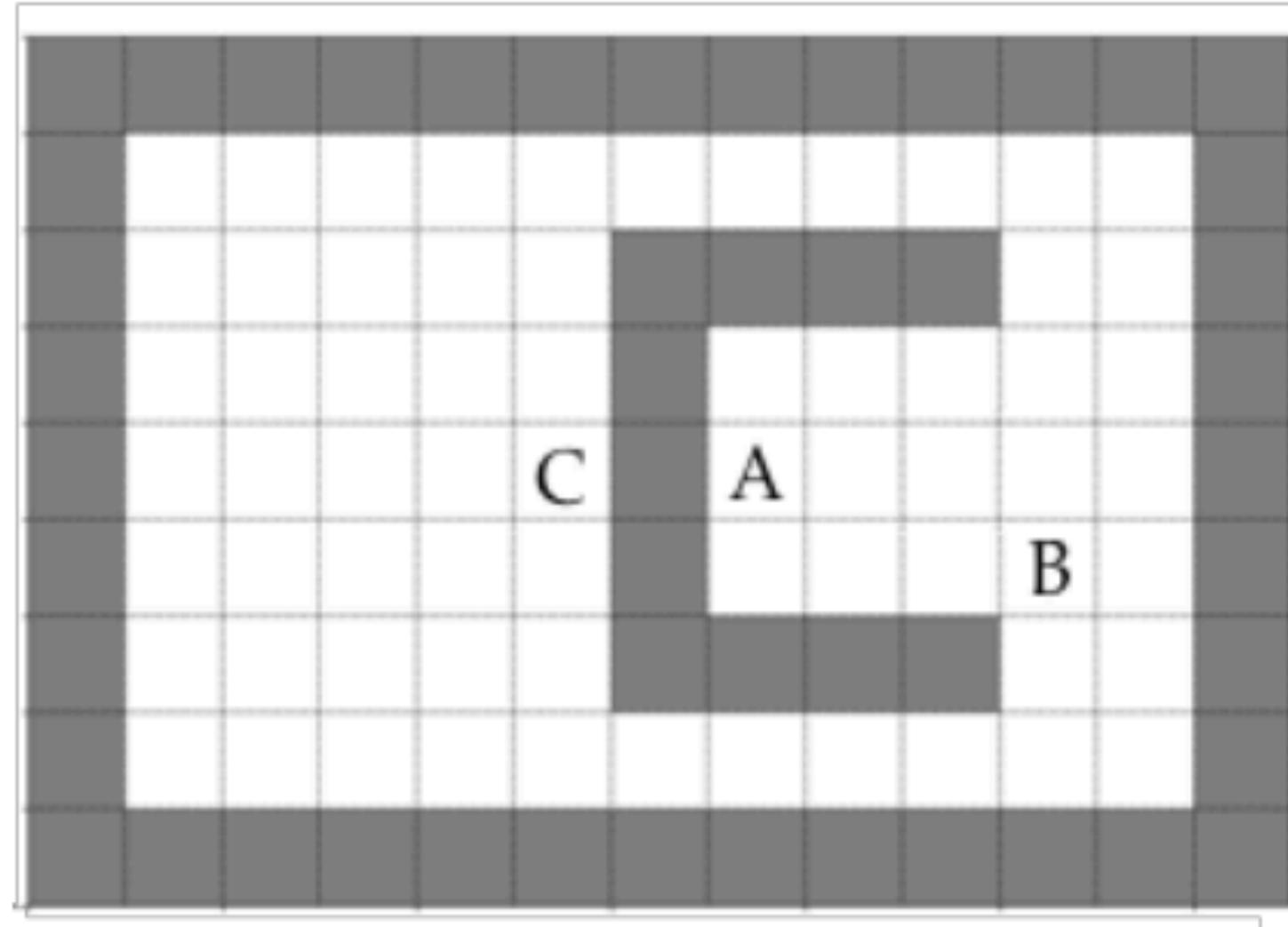


# From continuous to structured spaces





# Generalization based on transition dynamics



Machado et al. (*ICLR* 2018)

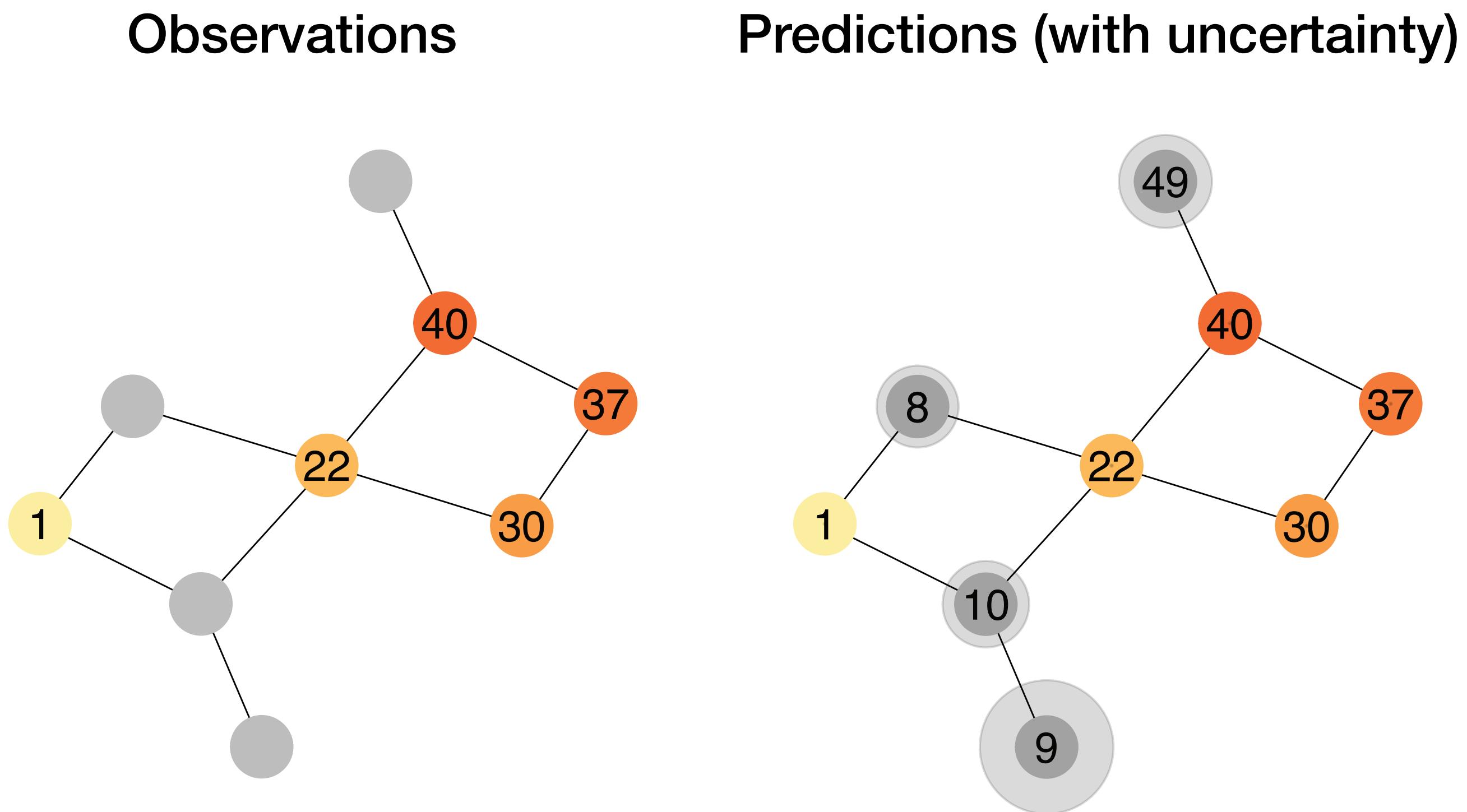
- A indicates a reward
- Even though C is closer than B, the transition dynamics of the environment make it easier for B to reach A

# Diffusion Kernel

- Rather than similarity between features, we use the connectivity structure of the graph to define similarity

$$k_{DF}(s, s') = \exp(-\alpha L)$$

- Where  $L$  is the graph Laplacian
- $\alpha$  is a free parameter (diffusion level)
- The diffusion kernel assumes function values diffuse across the graph according to a random walk



# Experiment 1

**Prediction Task**

Current Network: 4/30  
Current Weighted Error: 10.19

How many passengers do you think will be observed at the selected station?

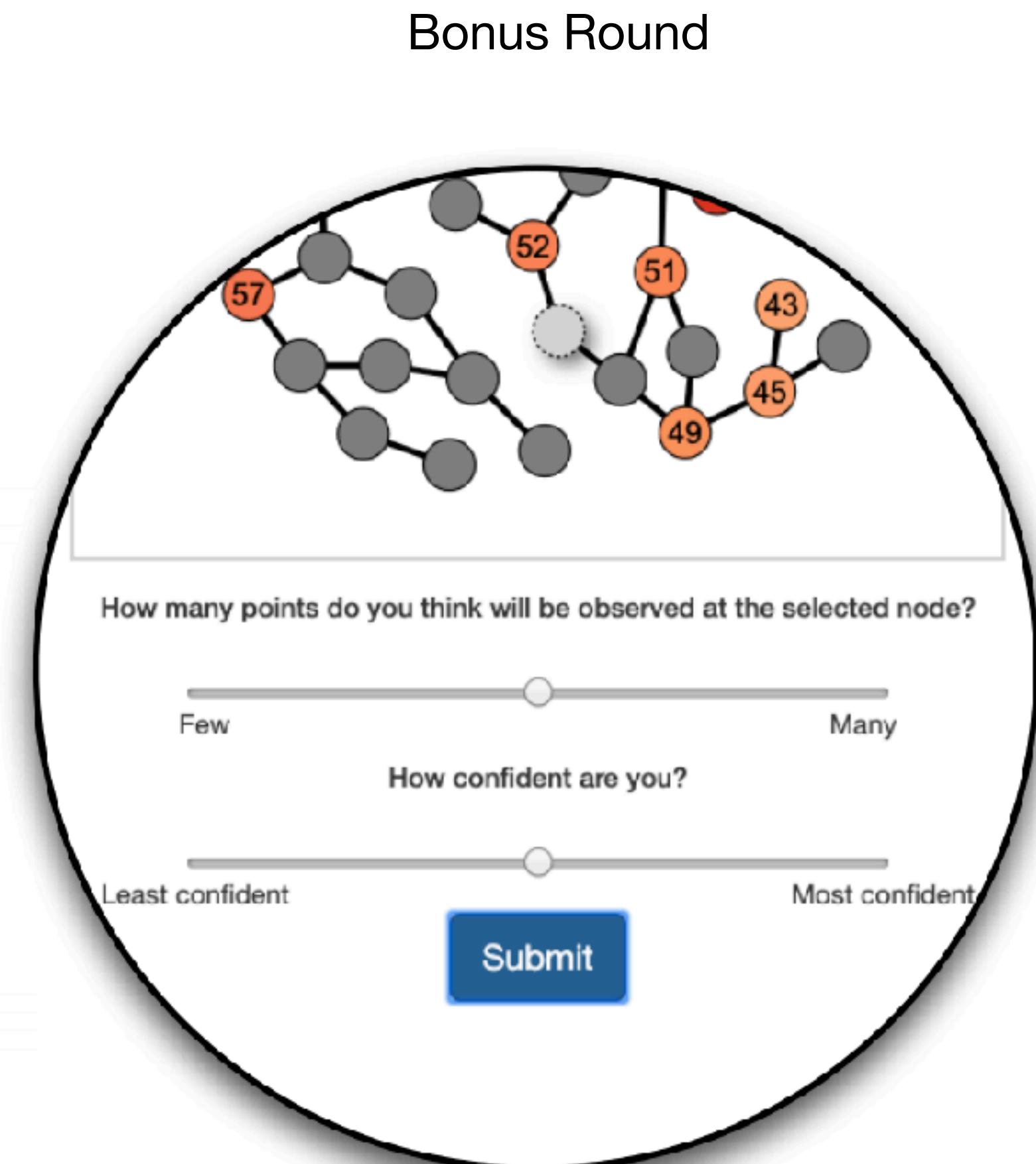
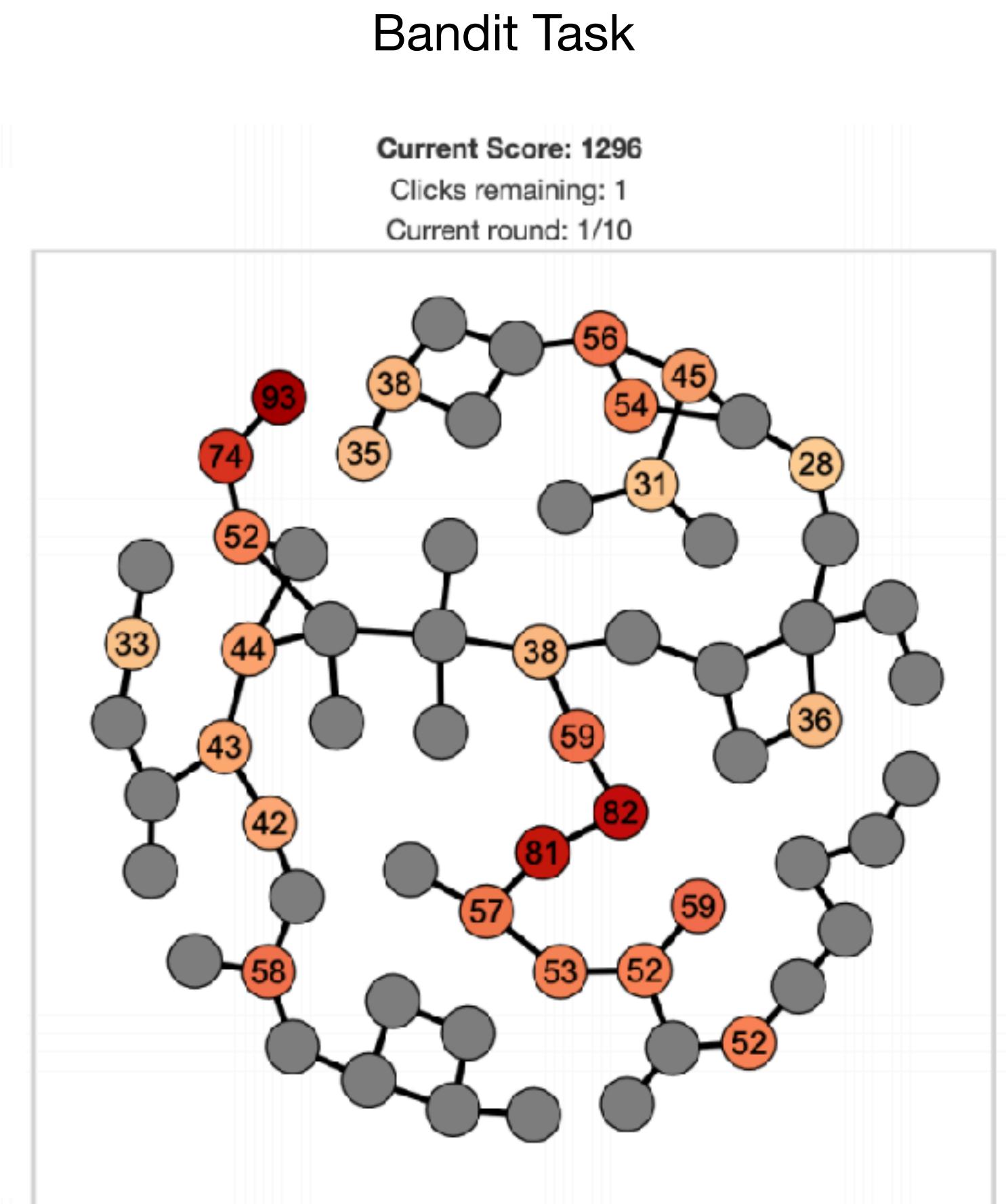
Few                          Many

How confident are you?

Not very confident                  Highly confident

Submit

# Experiment 2



# Experiment 1

**Prediction Task**

Current Network: 4/30  
Current Weighted Error: 10.19

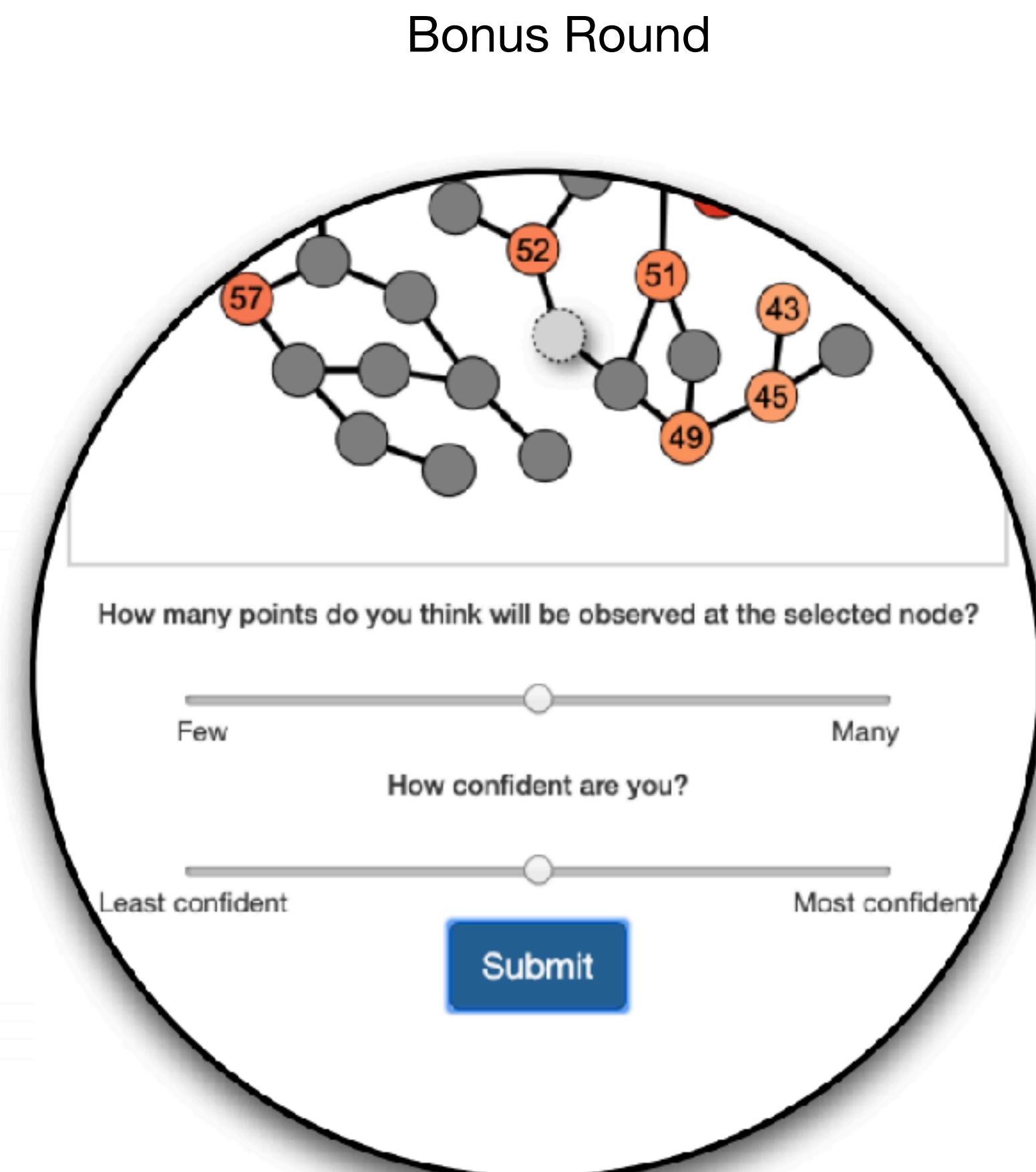
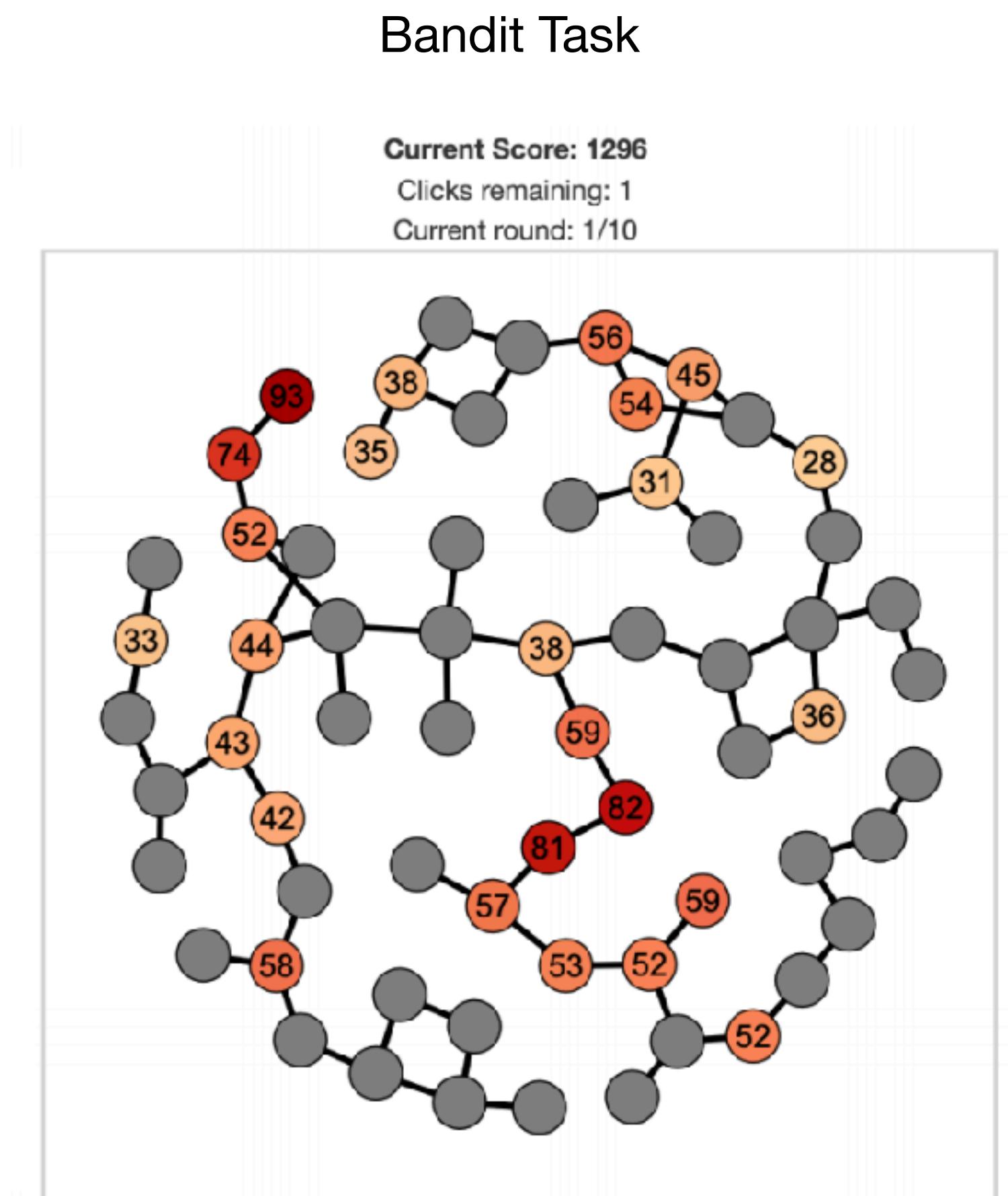
How many passengers do you think will be observed at the selected station?

Few      Many

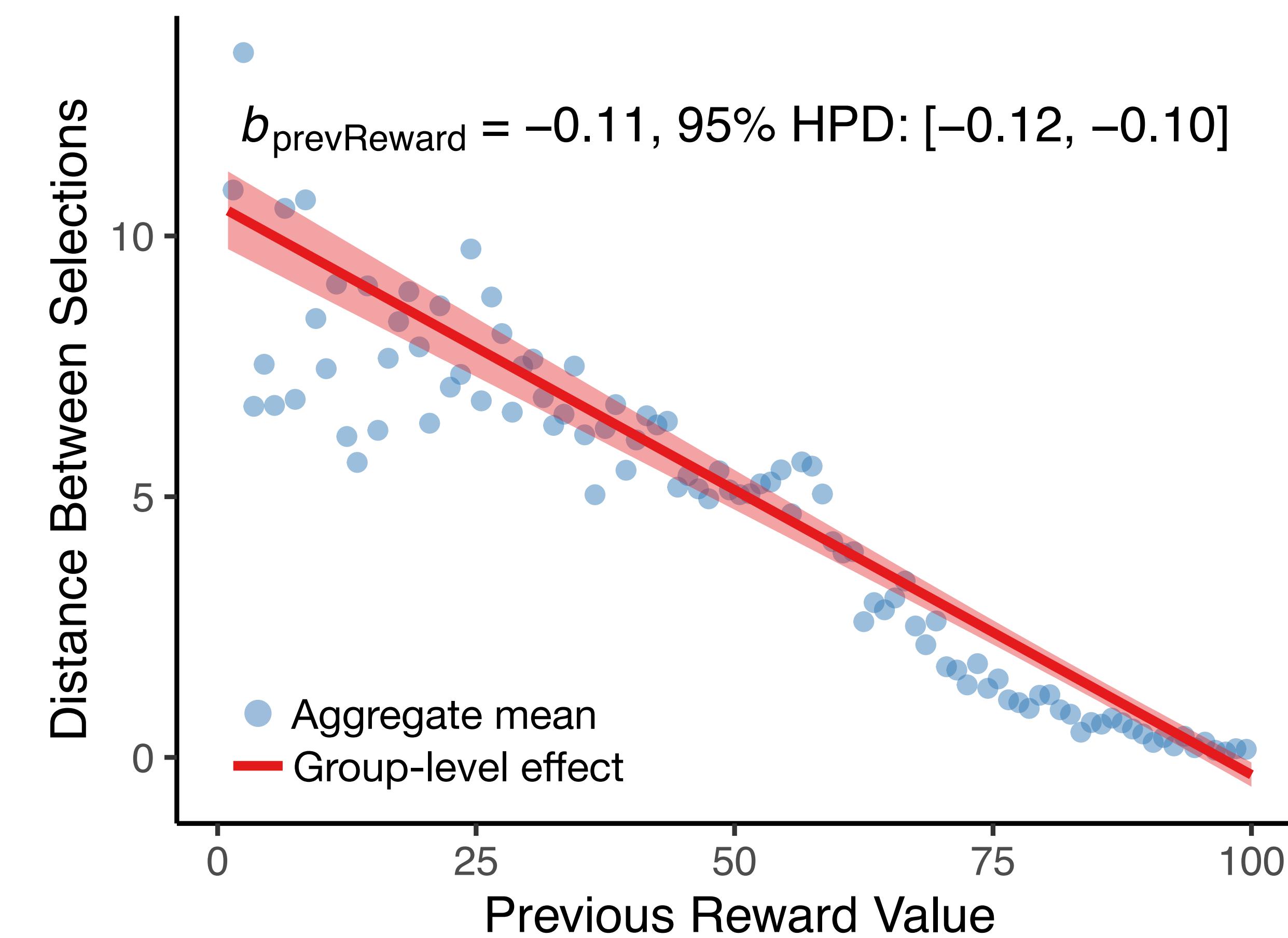
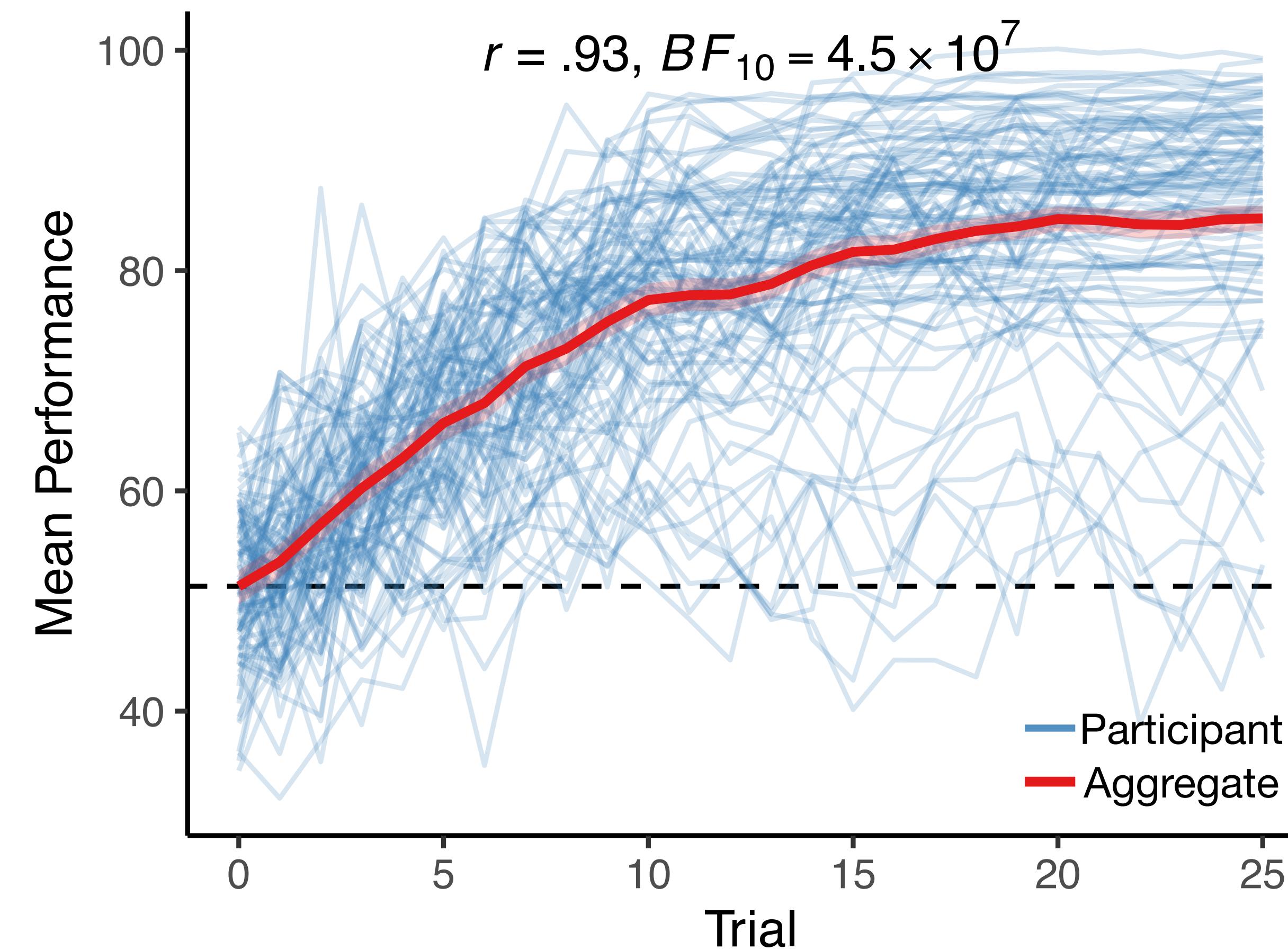
How confident are you?

Not very confident      Highly confident

# Experiment 2

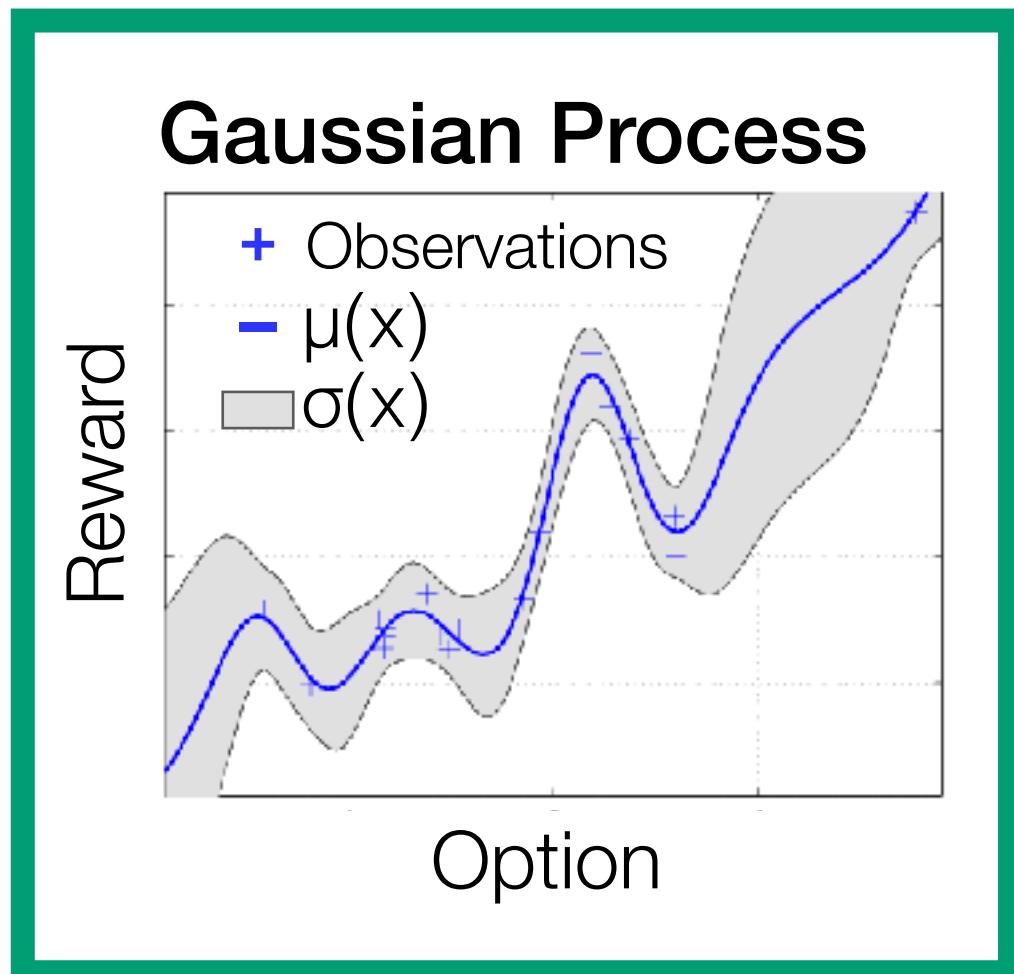


# Behavioral Results

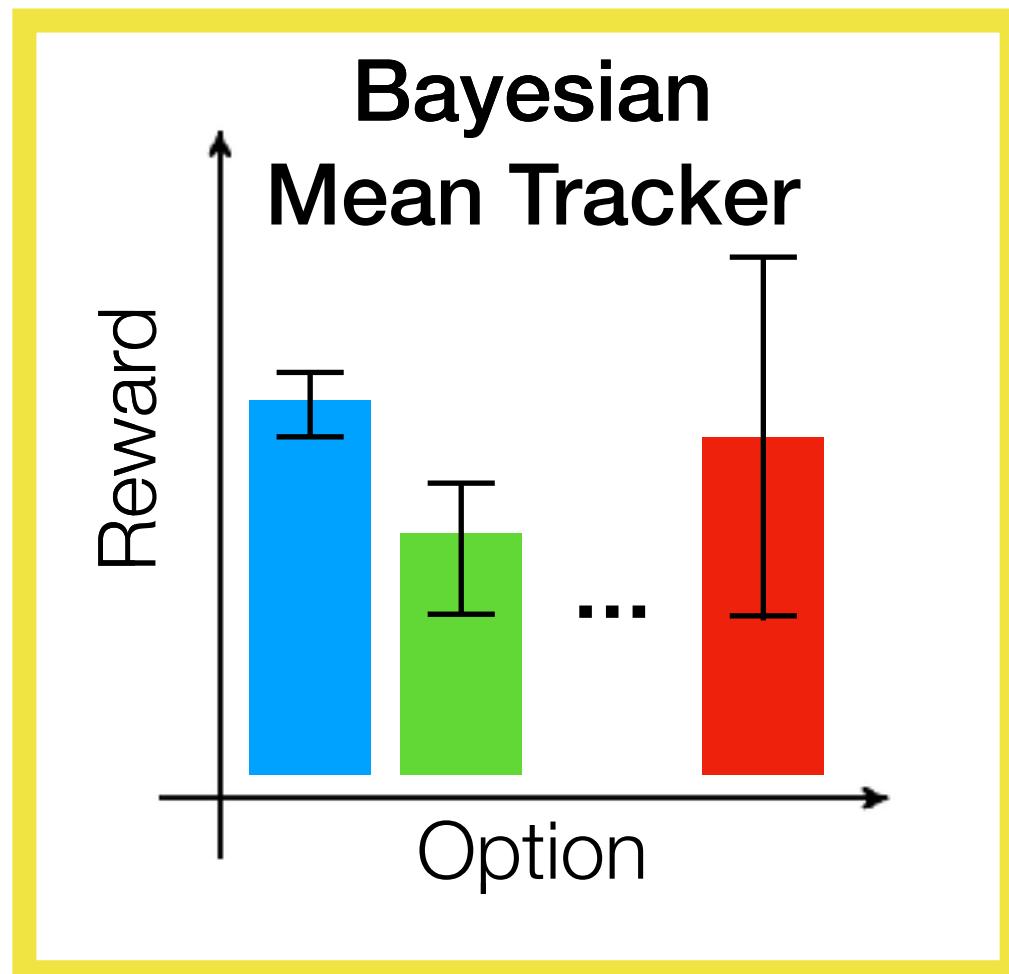


# Model Results

Generalization

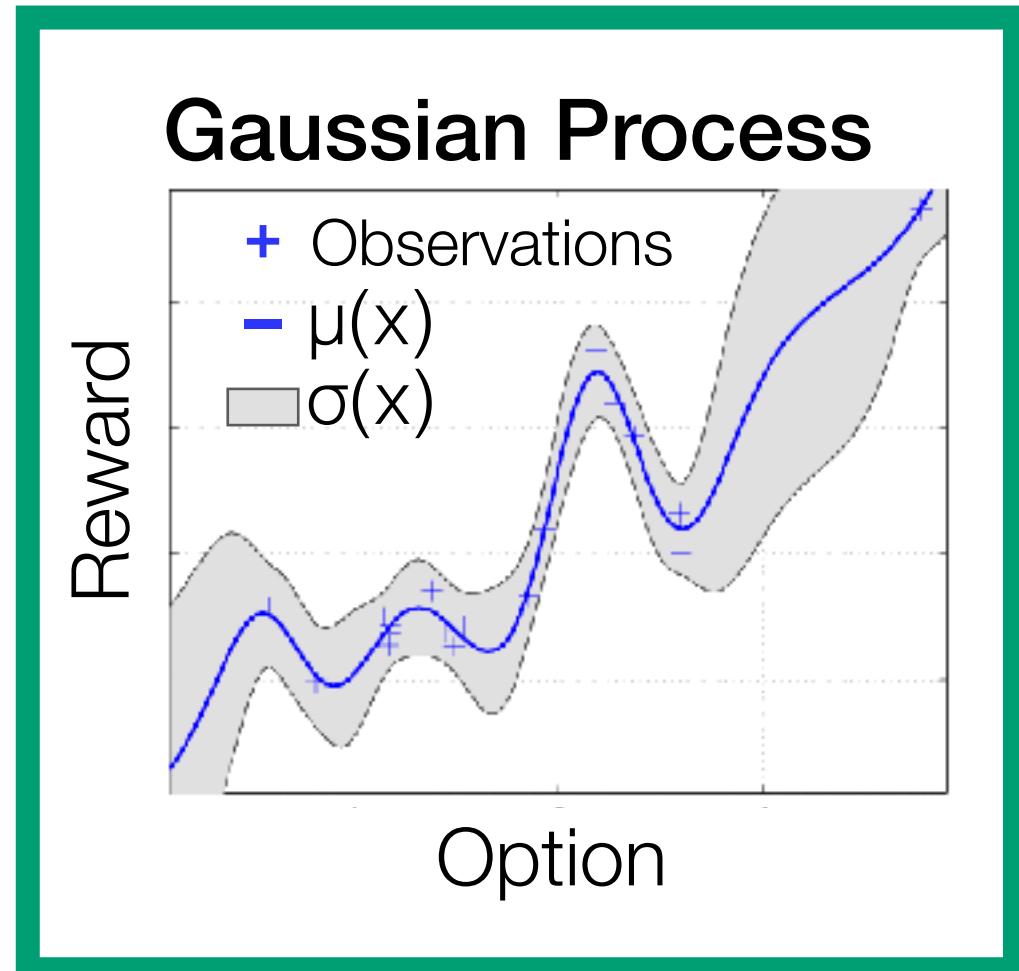


No generalization

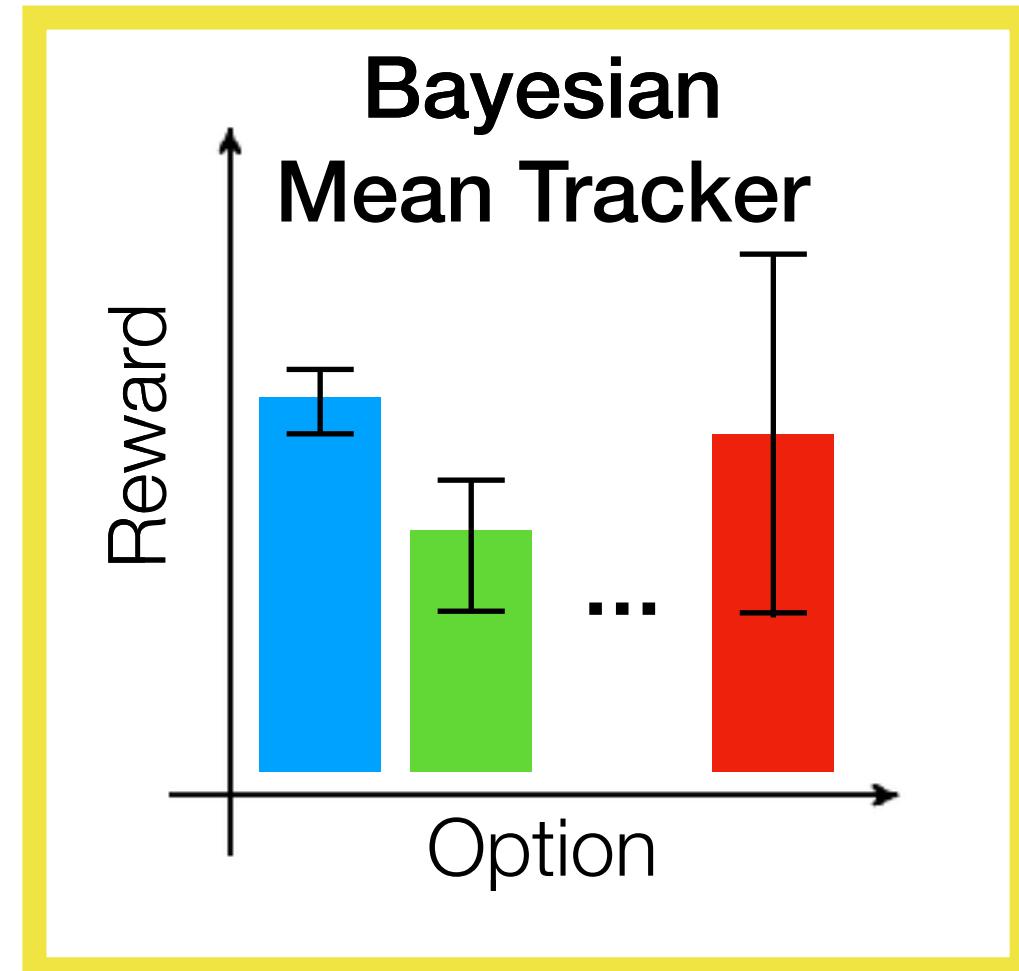


# Model Results

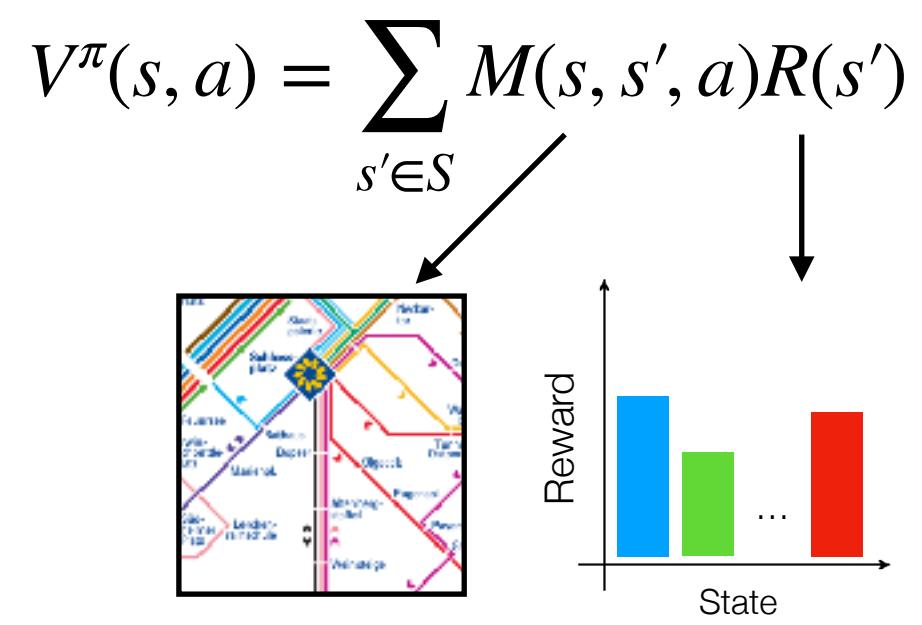
## Generalization



## No generalization

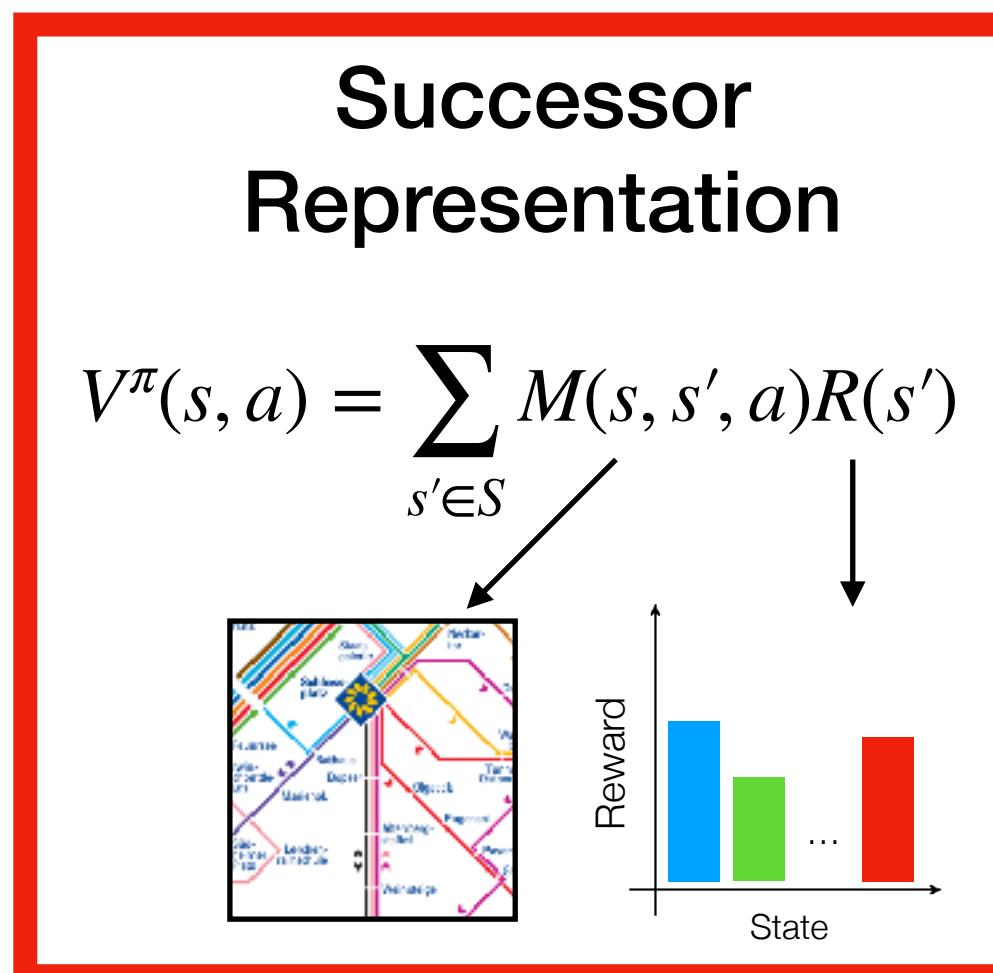
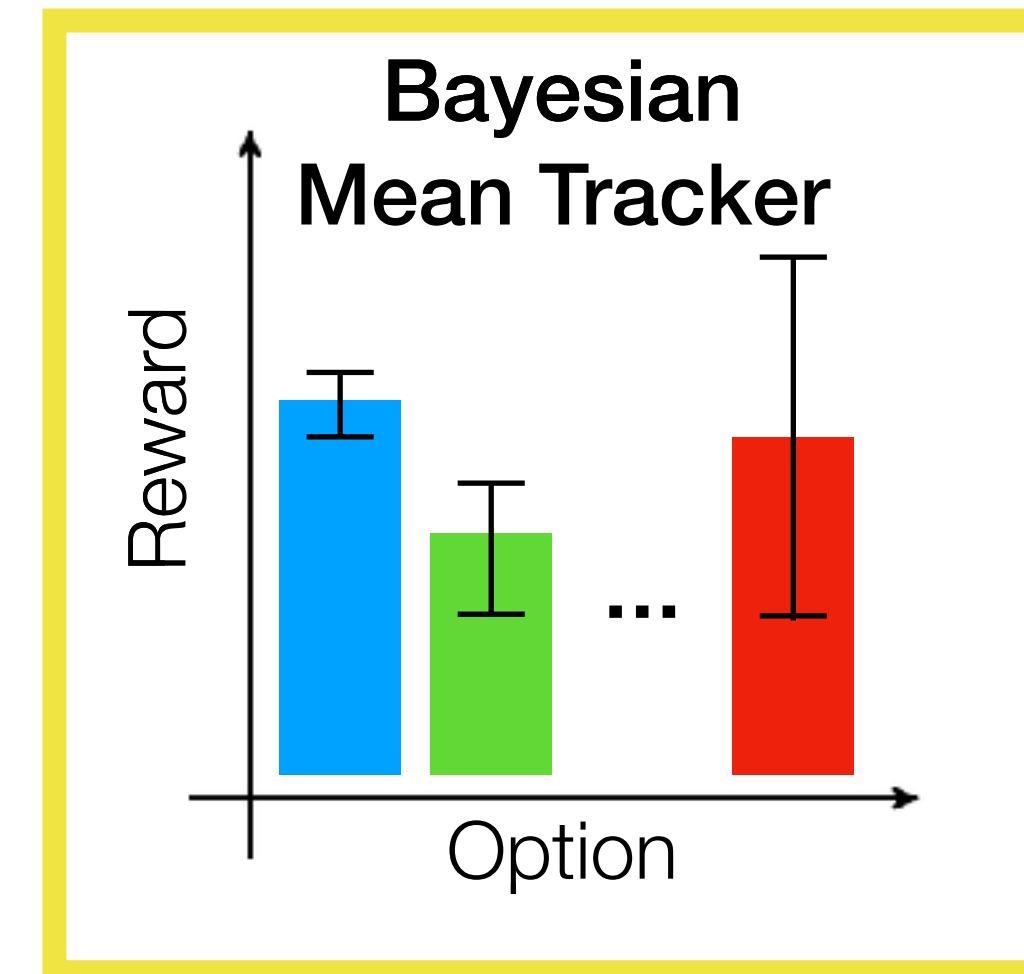
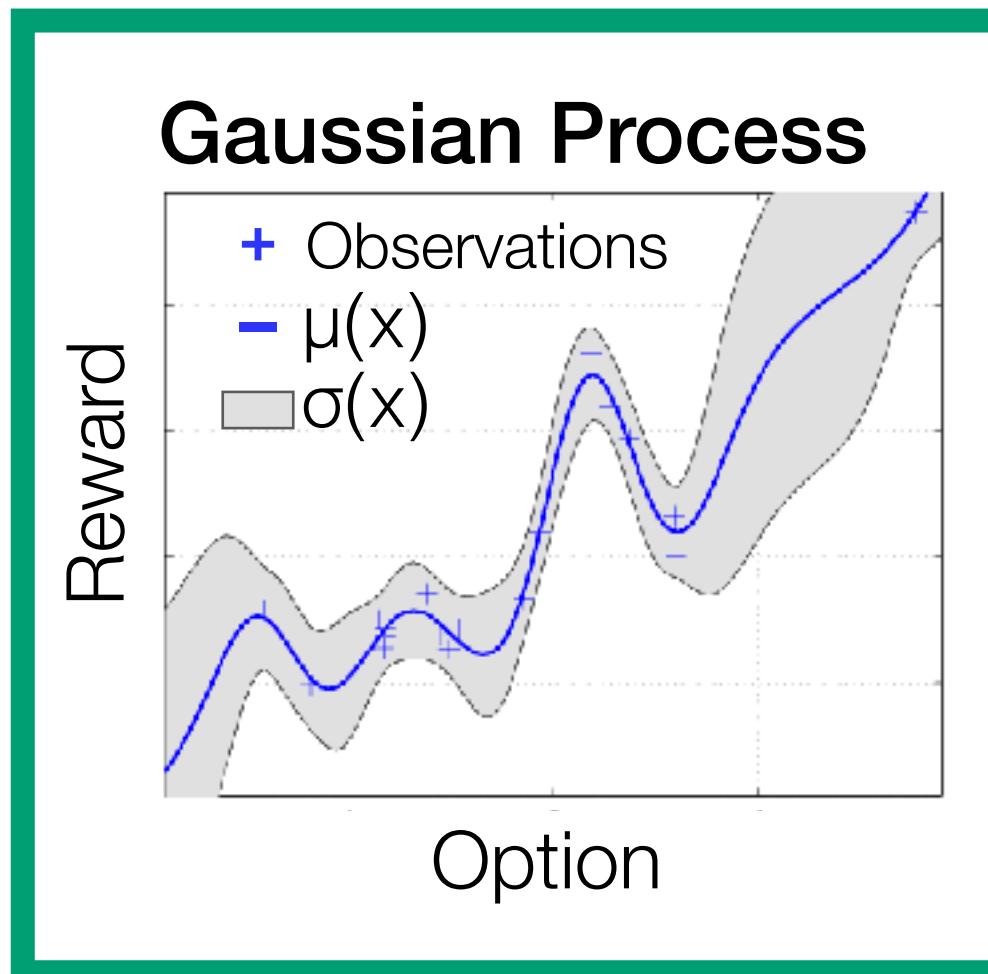


## Successor Representation



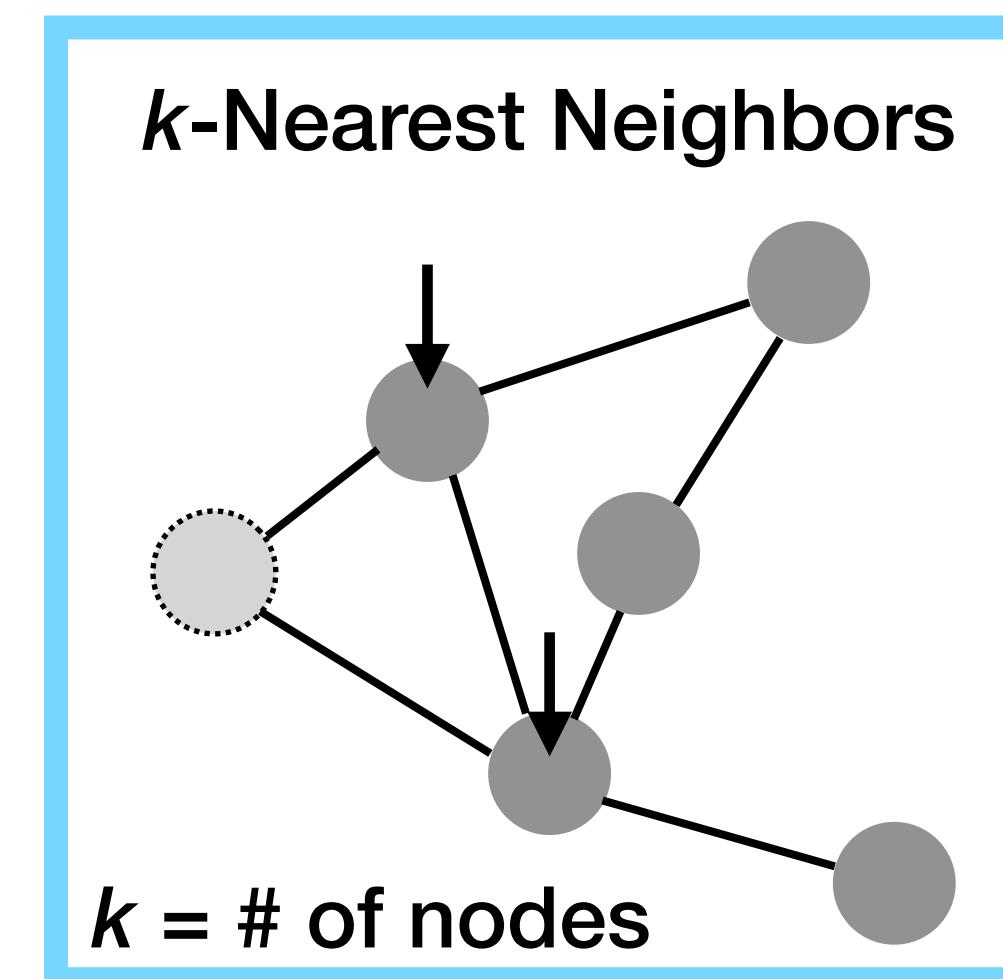
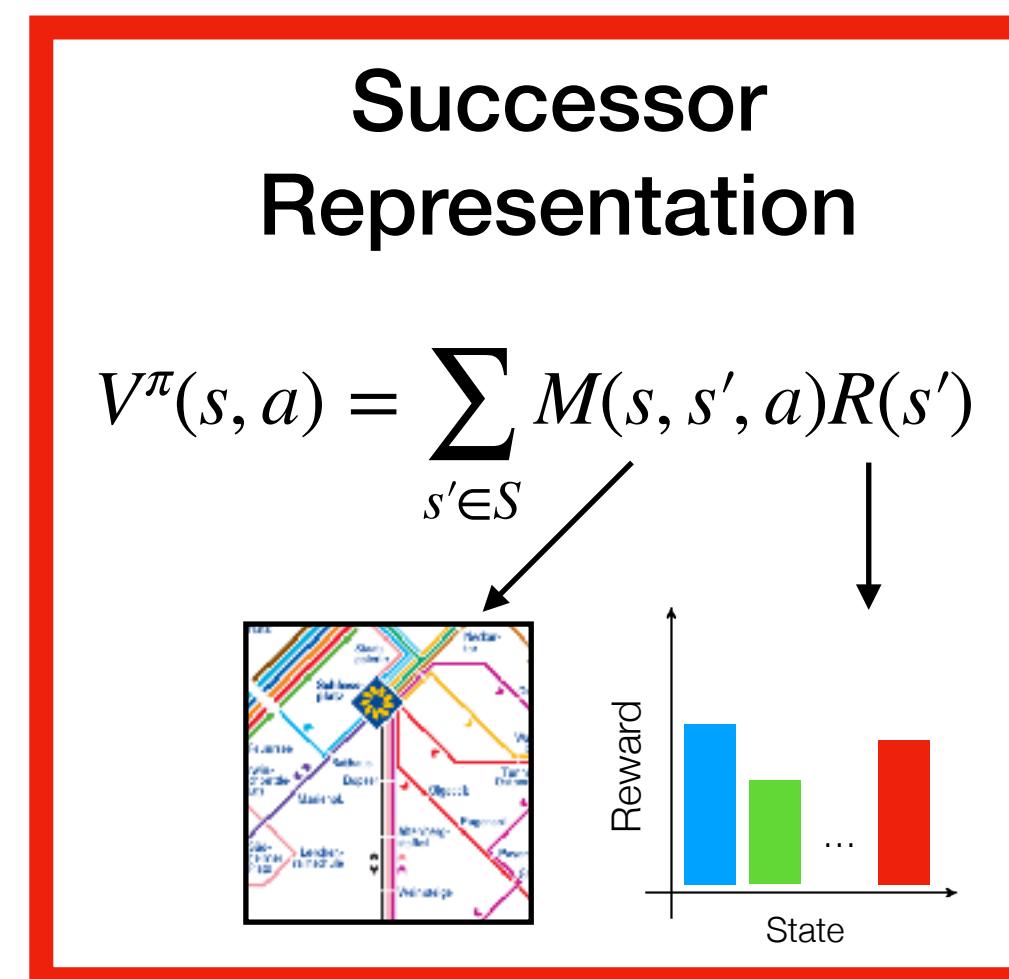
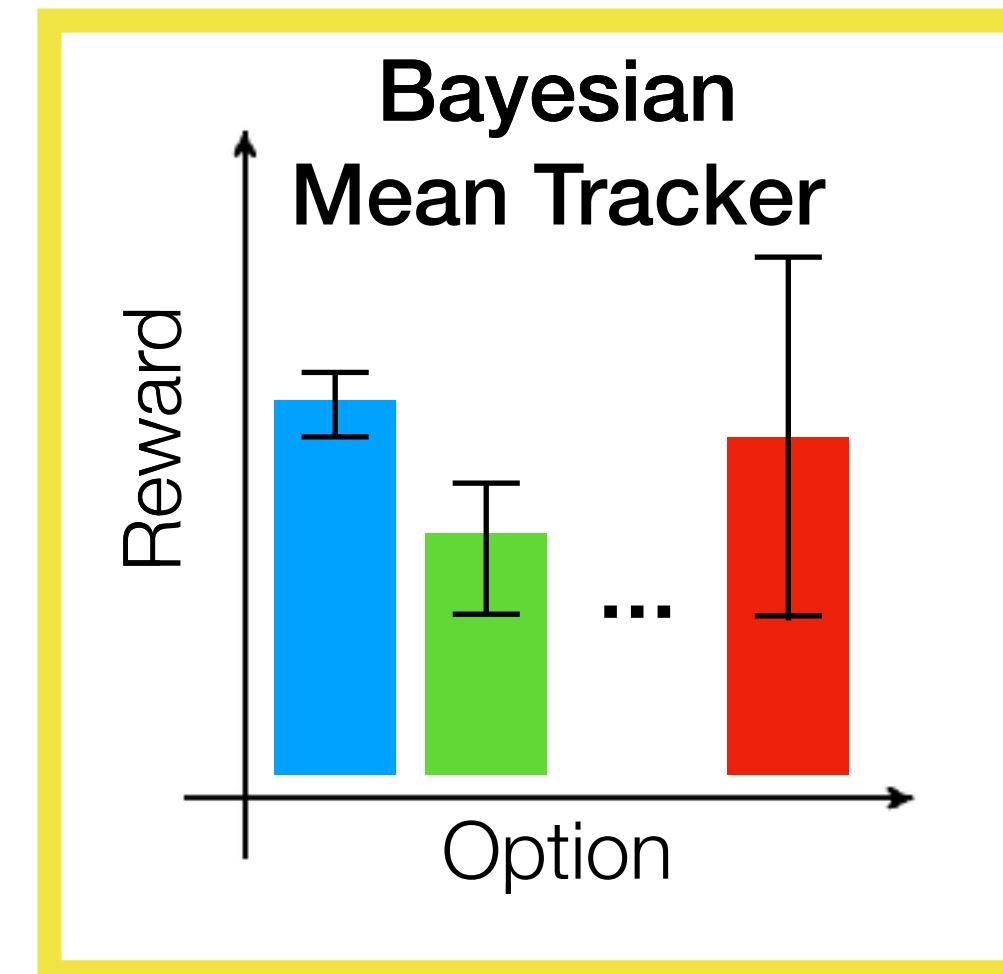
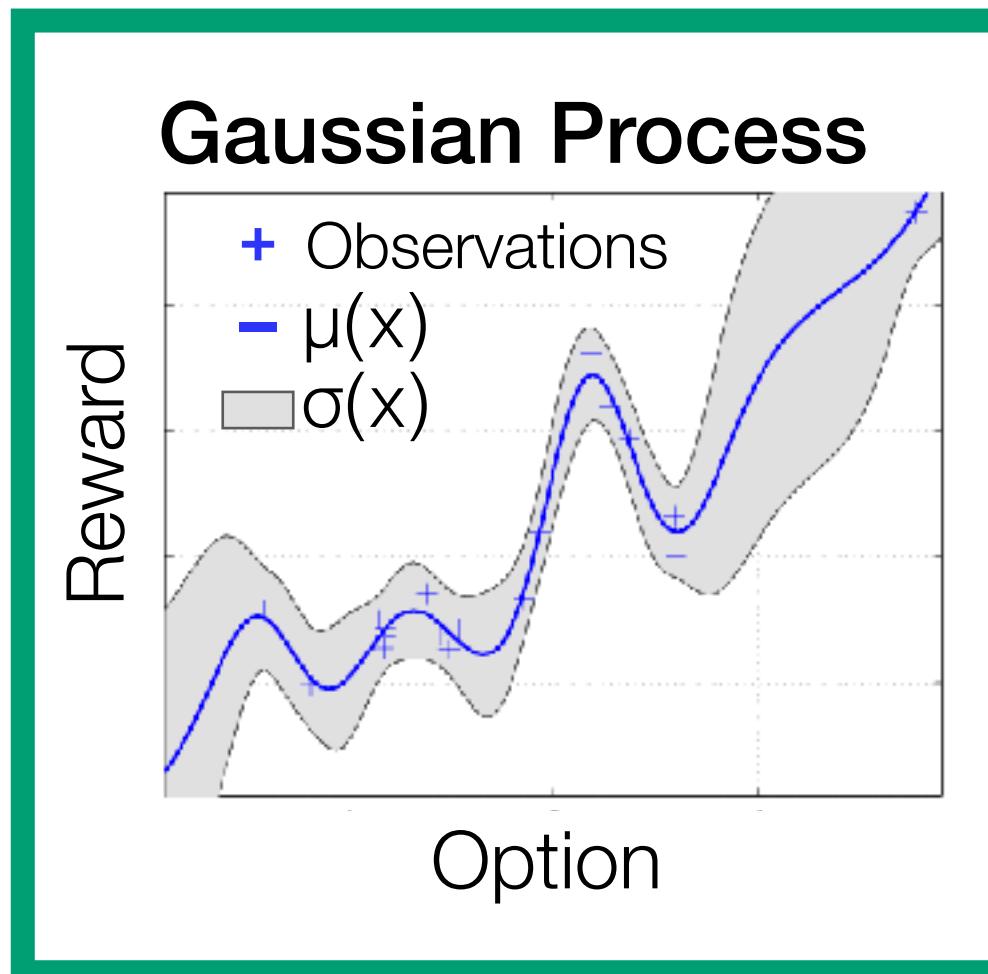
# Model Results

random exploration  
directed + random exploration



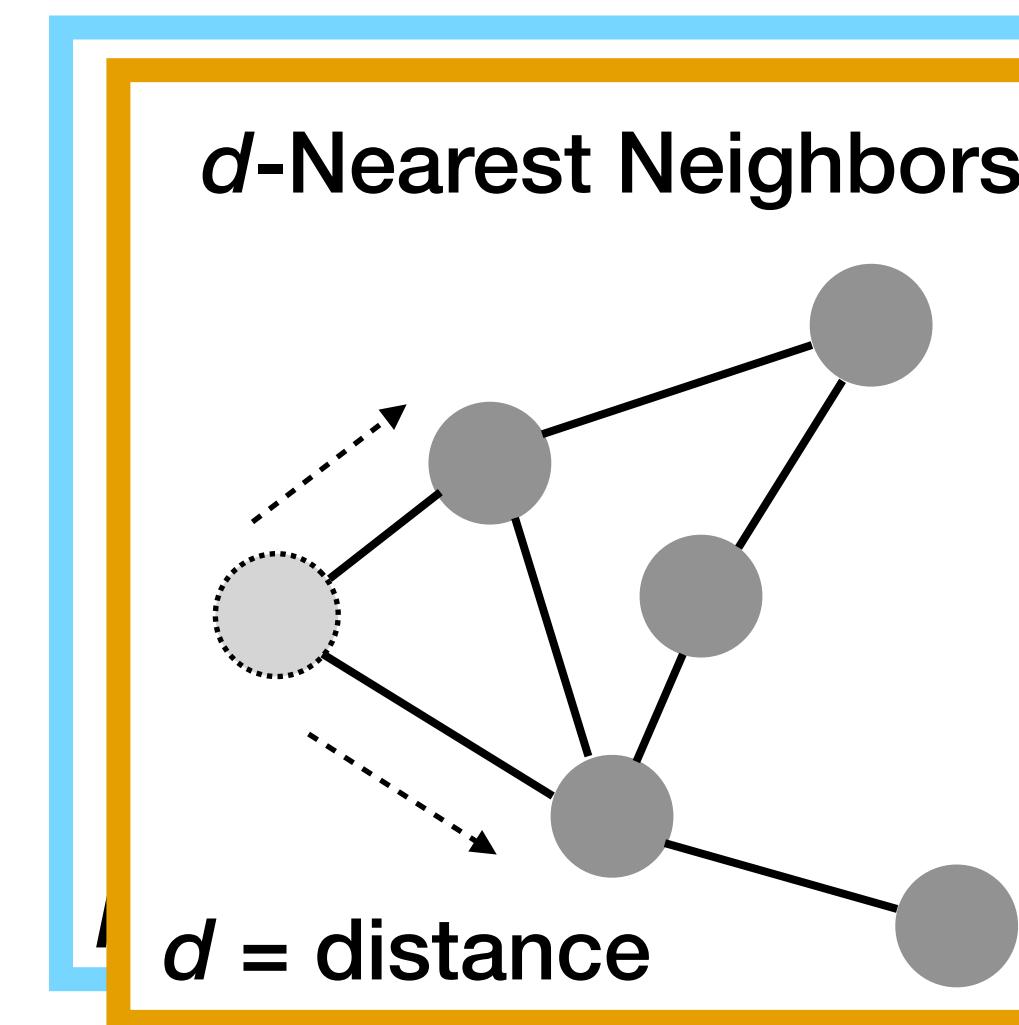
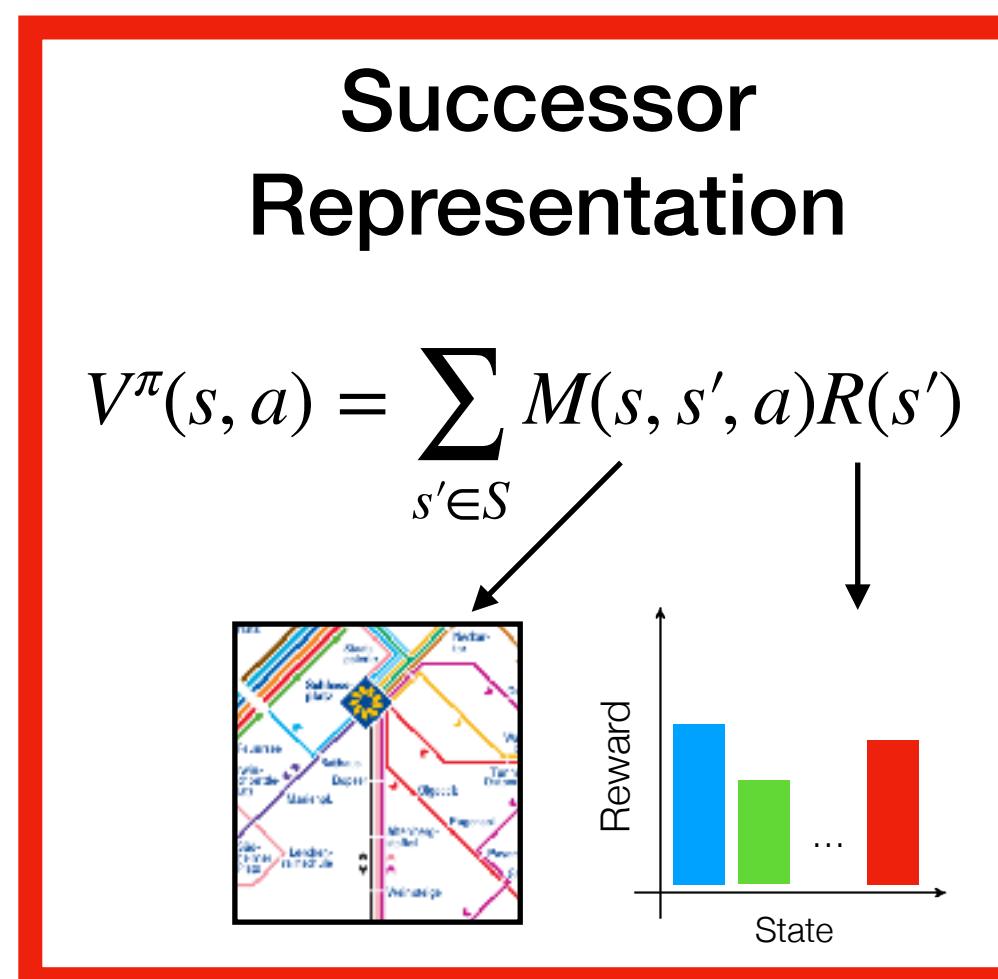
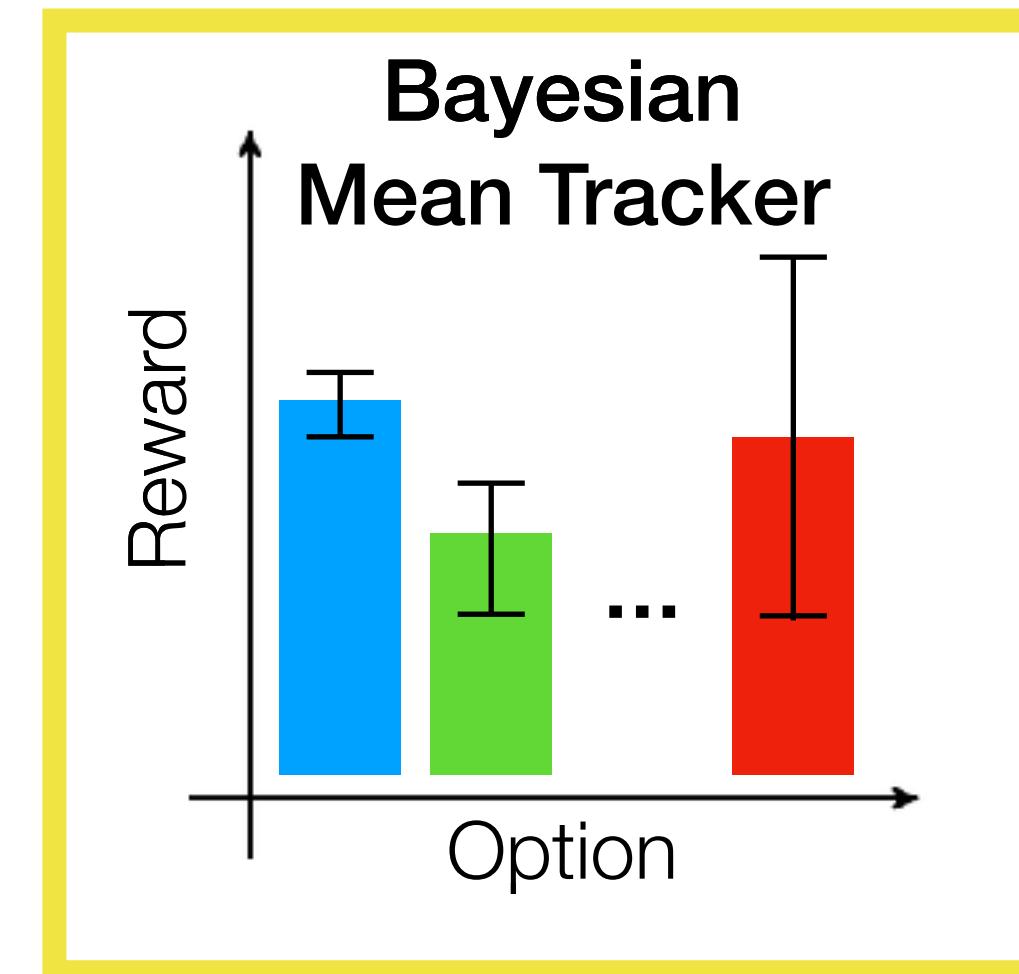
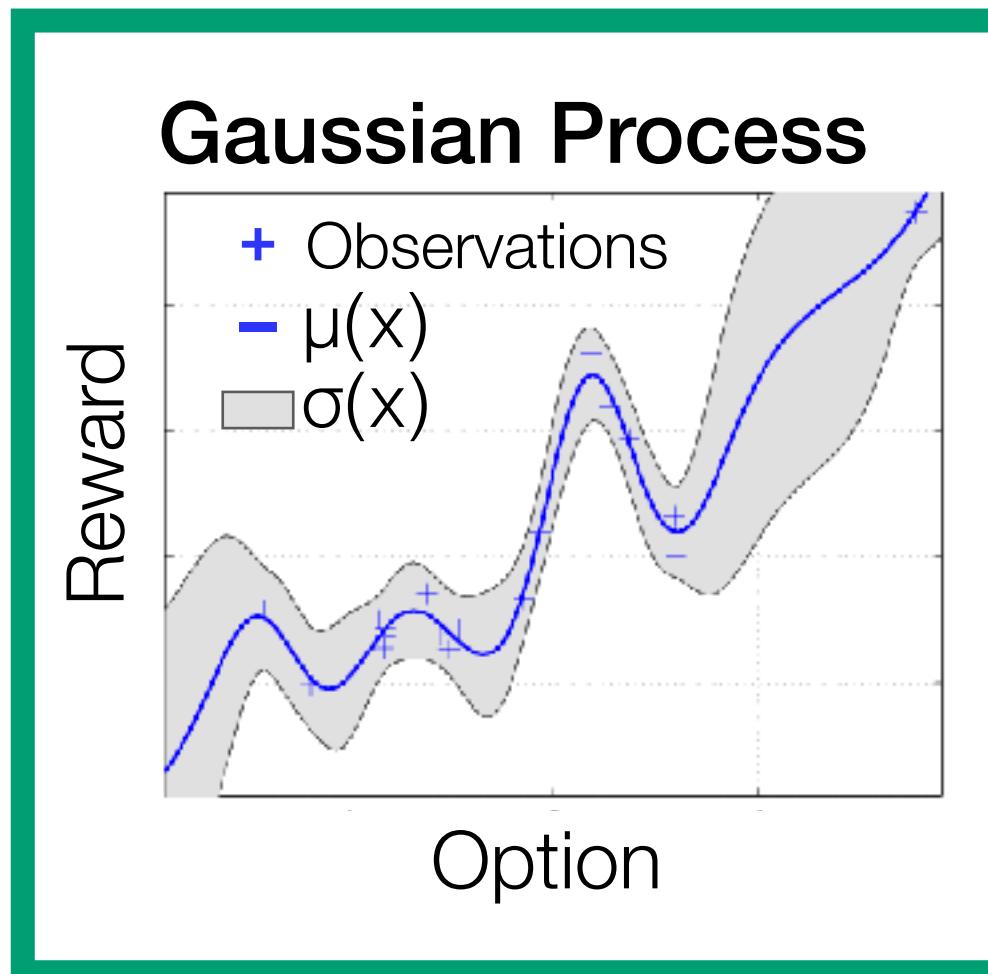
# Model Results

random exploration  
directed + random exploration

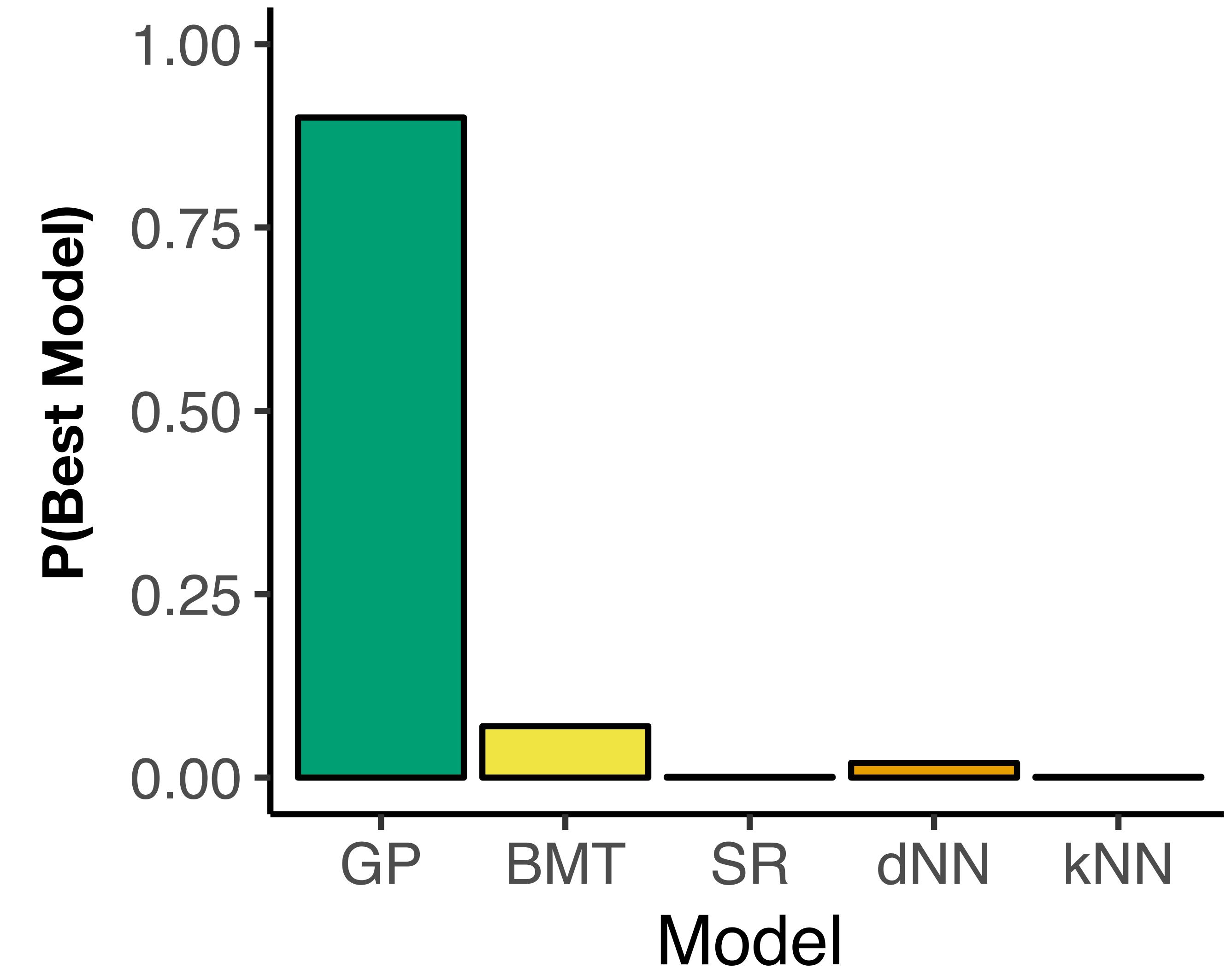
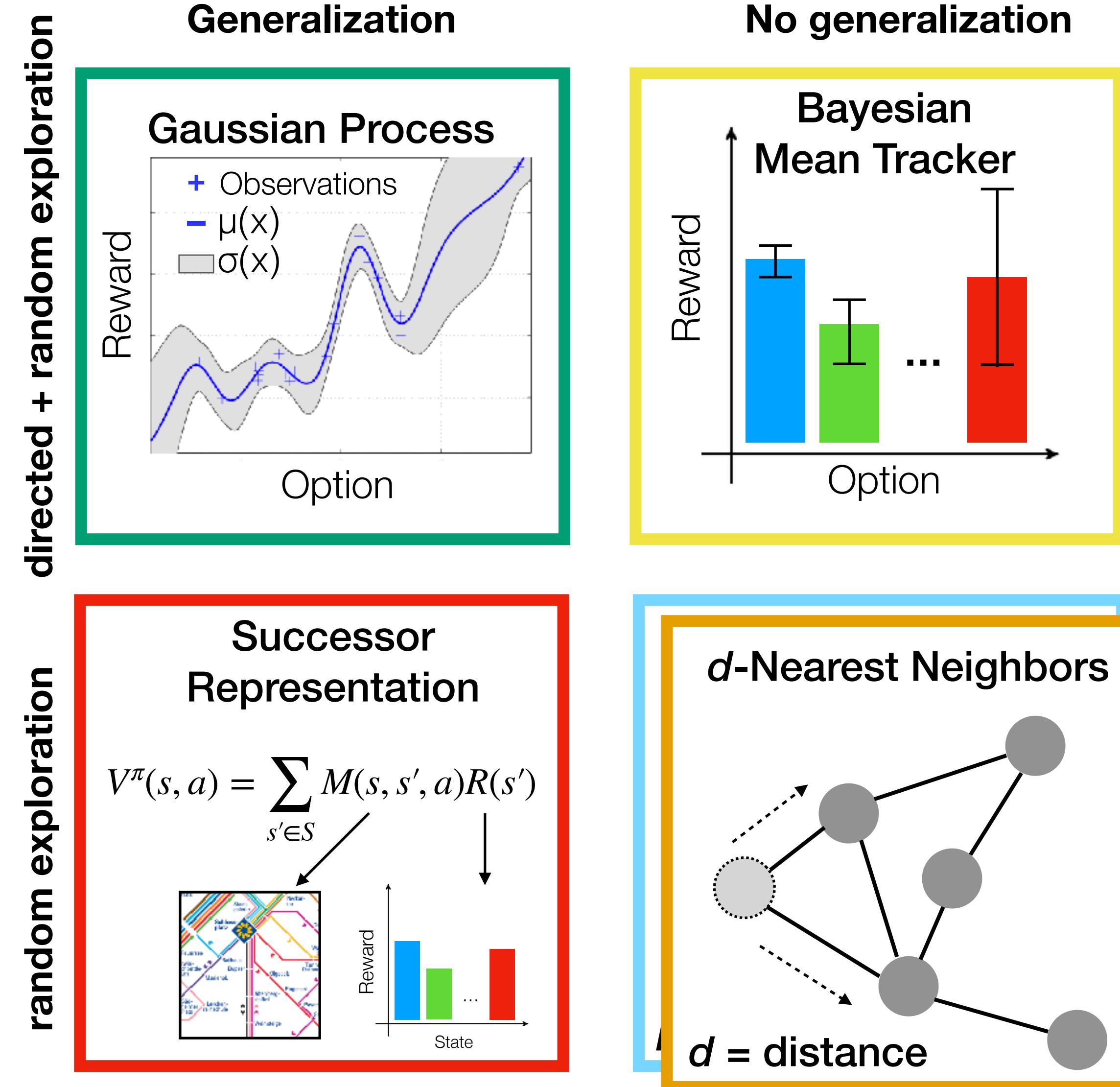


# Model Results

random exploration  
directed + random exploration



# Model Results



# Validation on judgments

The interface is contained within a large black-outlined circle. At the top, there is a network graph with several nodes. Some nodes are orange circles containing numbers (57, 52, 51, 43, 45, 49) and others are grey circles. A dashed grey circle is positioned in the center. Below the graph, a question is displayed:

How many points do you think will be observed at the selected node?

Below the question is a horizontal slider with two endpoints: "Few" on the left and "Many" on the right. A small white circle marks the current position of the slider.

Further down, another question is shown:

How confident are you?

Below this question is another horizontal slider with two endpoints: "Least confident" on the left and "Most confident" on the right. A small white circle marks the current position of the slider.

In the bottom center of the circle is a blue rectangular button with the word "Submit" in white text.

# Validation on judgments

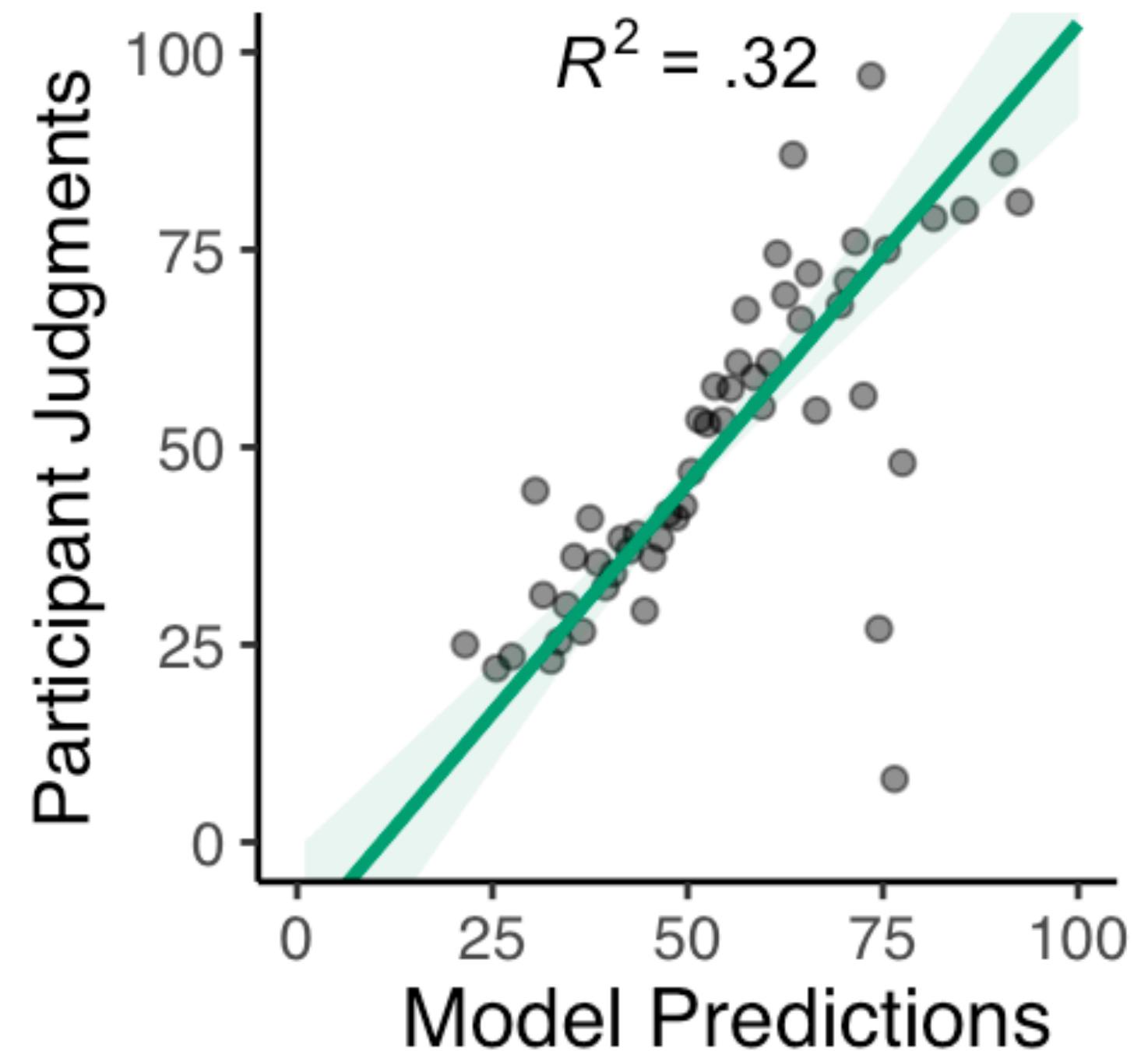
How many points do you think will be observed at the selected node?

Few                          Many

How confident are you?

Least confident                  Most confident

Submit



# Validation on judgments

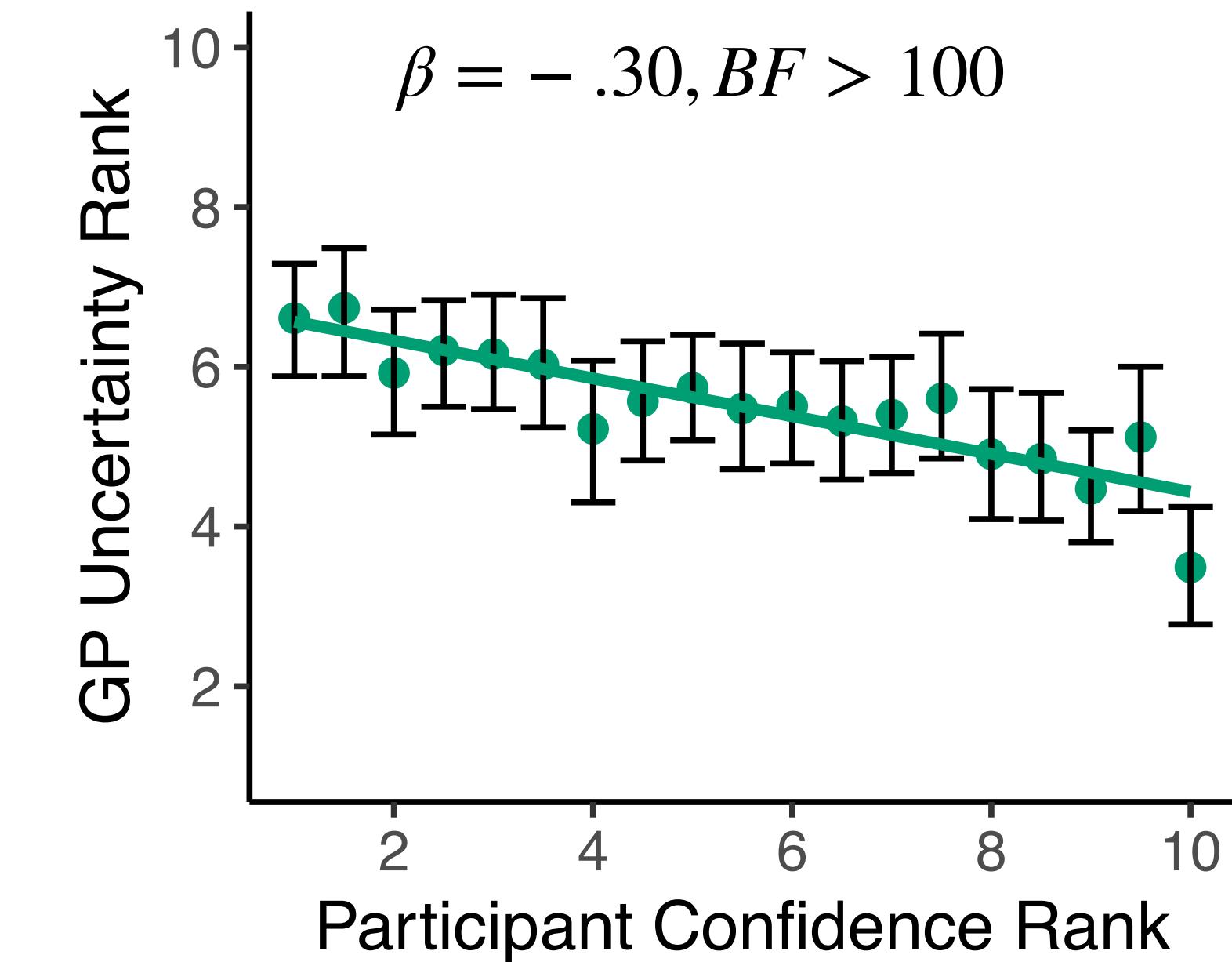
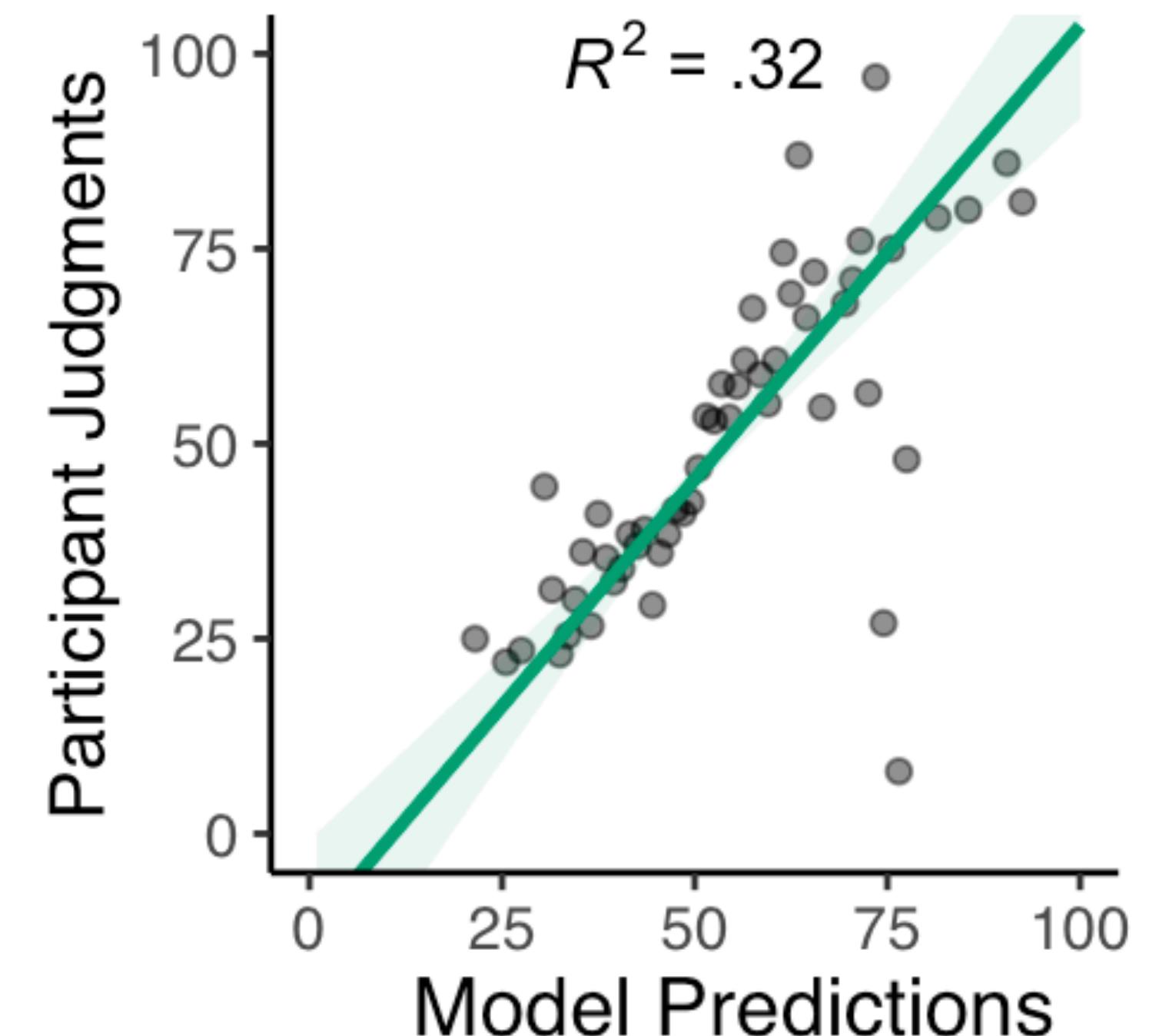
How many points do you think will be observed at the selected node?

Few                          Many

How confident are you?

Least confident                  Most confident

Submit



# Conclusions

- How do we navigate vast problem spaces?
  - Generalization and directed exploration provide a powerful model for efficient learning across many domains
- By modeling generalization as functional inference, we can:
  - Predict search decisions
  - Simulate human-like performance
  - Predict judgments of expected reward and confidence
- Underlying mechanisms of Bayesian inference, kernel similarity, and episodic RL have deep theoretical connections to other models in Neuroscience and Computer Science

# Future directions

- How is the structure of similarity learned?
  - **Successor Representation:** online prediction error about future states? (Dayan, *NeurComp* 1993; Stachenfeld et al., *Nat Neuro* 2017)
  - **Tolman Eichenbaum Machine:** associative learning mechanisms (Whittington et al., *Cell* 2020)
  - **Structure induction:** Bayesian inference about hypothesized structure? (Kemp & Tenenbaum, *PNAS* 2008)
- How do humans keep functional inference tractable as we gain more experience?
  - GPs scale cubically with the size of the data  $\mathcal{O}(n^3)$
  - A large part of this complexity is in computing the underlying uncertainty
  - Perhaps population codes can support flexible generalization with uncertainty (Tano, Dayan, & Pouget, *NeurIPS* 2020)

# Thanks to my collaborators and funders



Eric Schulz  
MPI Tübingen



Björn Meder  
MPI Berlin, Uni. Potsdam



Azzurra Ruggeri  
MPI Berlin



Nico Schuck  
MPI Berlin



Sam Gershman  
Harvard



Simon Ciranka  
MPI Berlin



Anna Giron  
Uni Tübingen



Mona Garvert  
MPI Leipzig



Maarten Speekenbrink  
UCL



Jonathan Nelson  
MPI Berlin, Uni. Surrey



Wouter van den Bos  
Uni Amsterdam



EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN



TÜBINGEN AI CENTER  
BMBF Competence Center for Machine Learning

EBERHARD KARLS  
UNIVERSITÄT  
TÜBINGEN

CyberValley

Intelligent Systems

SPONSORED BY THE  
Federal Ministry  
of Education  
and Research

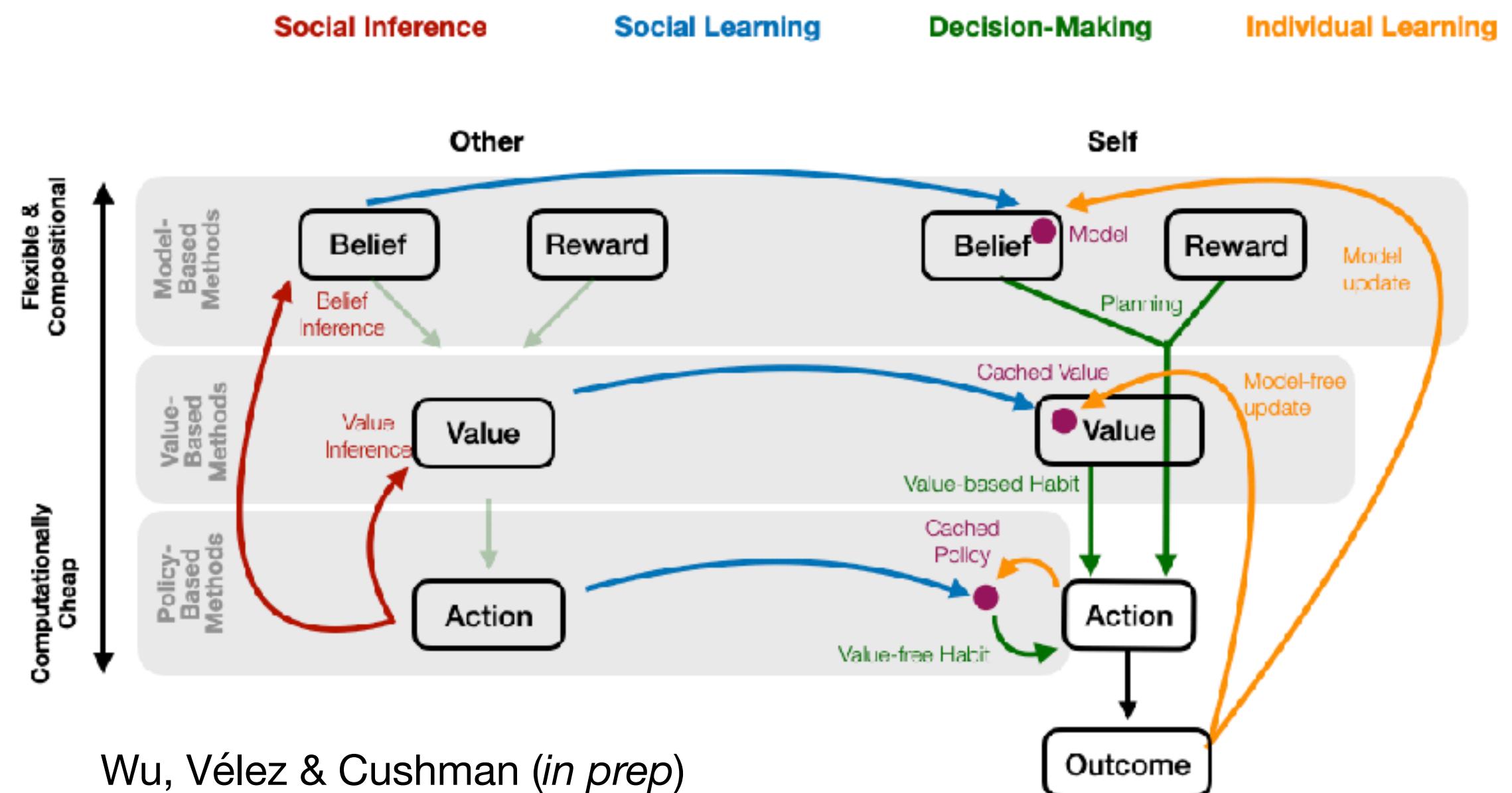
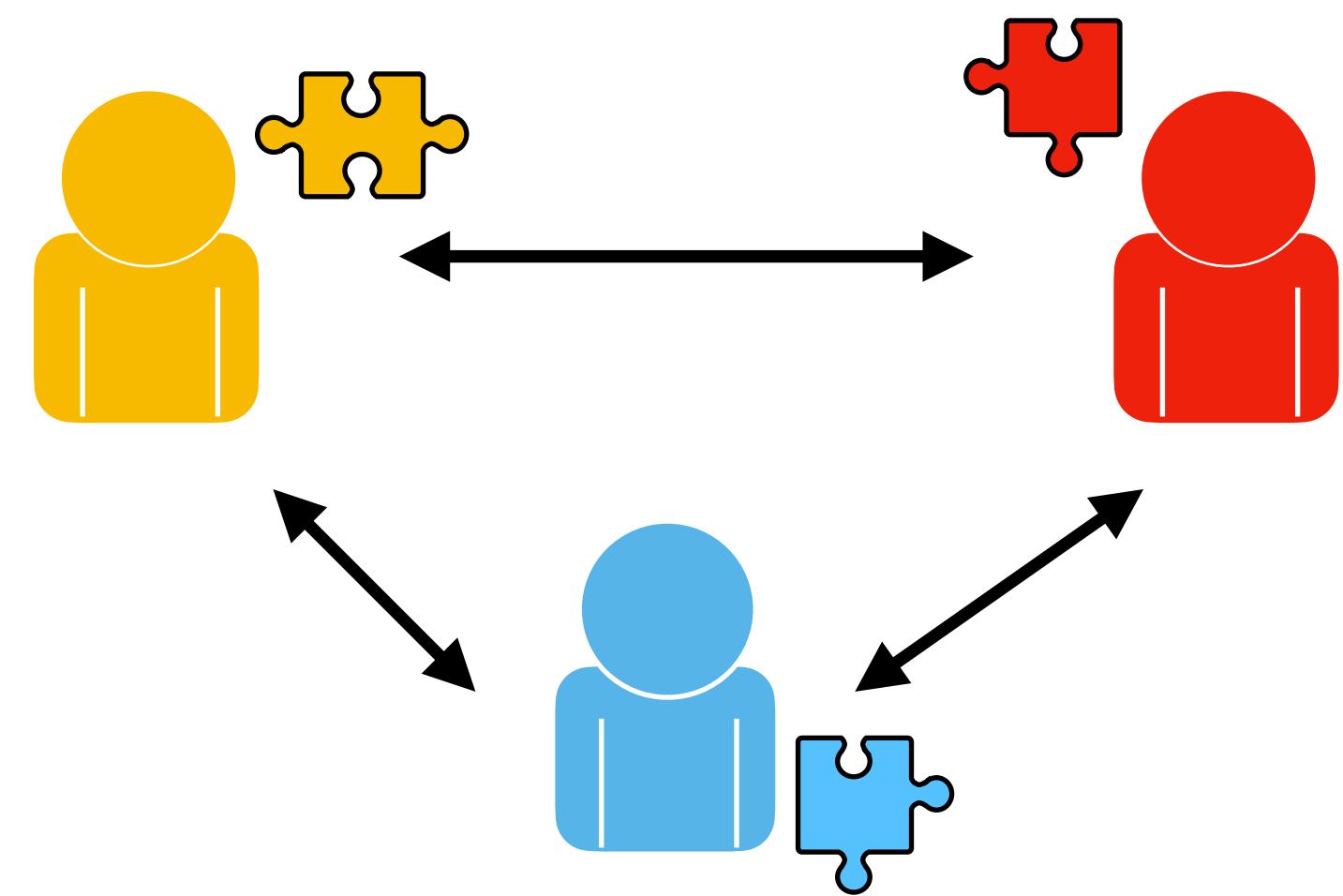
ADVANCING MACHINE INTELLIGENCE WITH ROBUST MACHINE LEARNING

# Fully Funded PhD position: Computational mechanisms of Social Learning

Potential topics include:

- Theory of Mind inference using inverse RL
- Selectivity and specialization in social learning using VR experiments programmed in minecraft
- Cumulative cultural evolution
- Integration of individual and social information

Visit [hmc-lab.com](http://hmc-lab.com) for more info! Apply by **May 15th**



Recreated Visual Field

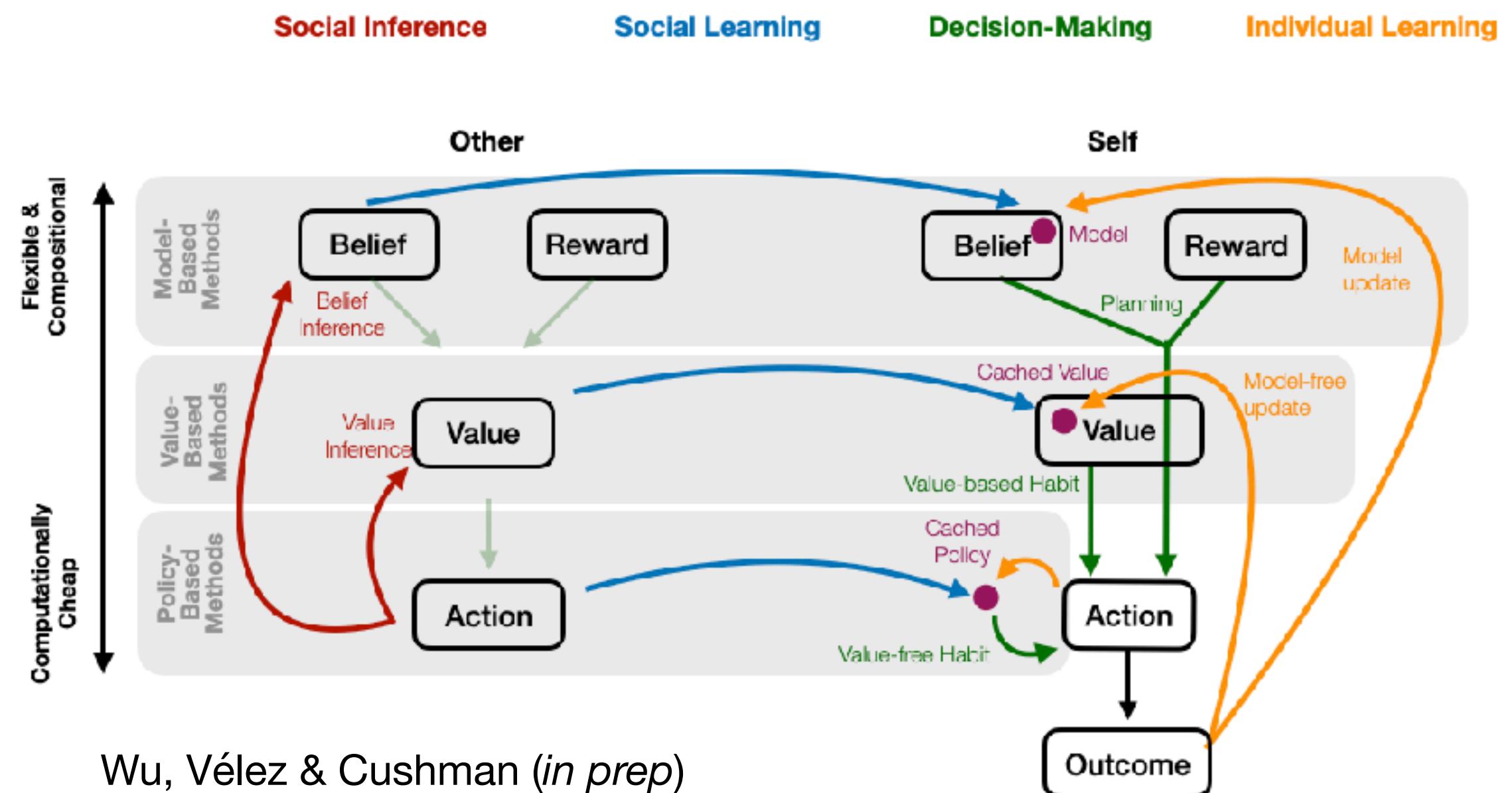
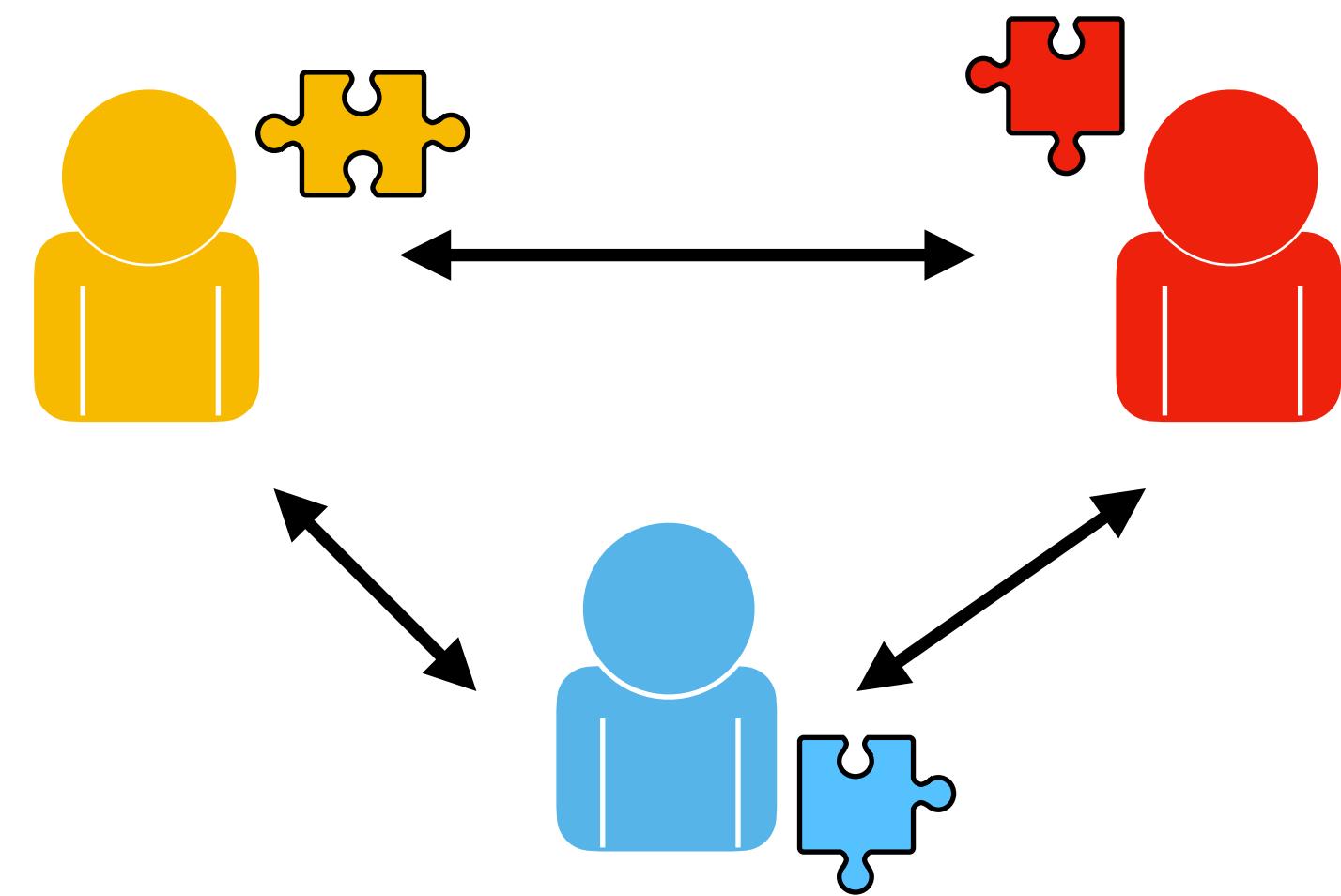


# Fully Funded PhD position: Computational mechanisms of Social Learning

Potential topics include:

- Theory of Mind inference using inverse RL
- Selectivity and specialization in social learning using VR experiments programmed in minecraft
- Cumulative cultural evolution
- Integration of individual and social information

Visit [hmc-lab.com](http://hmc-lab.com) for more info! Apply by **May 15th**

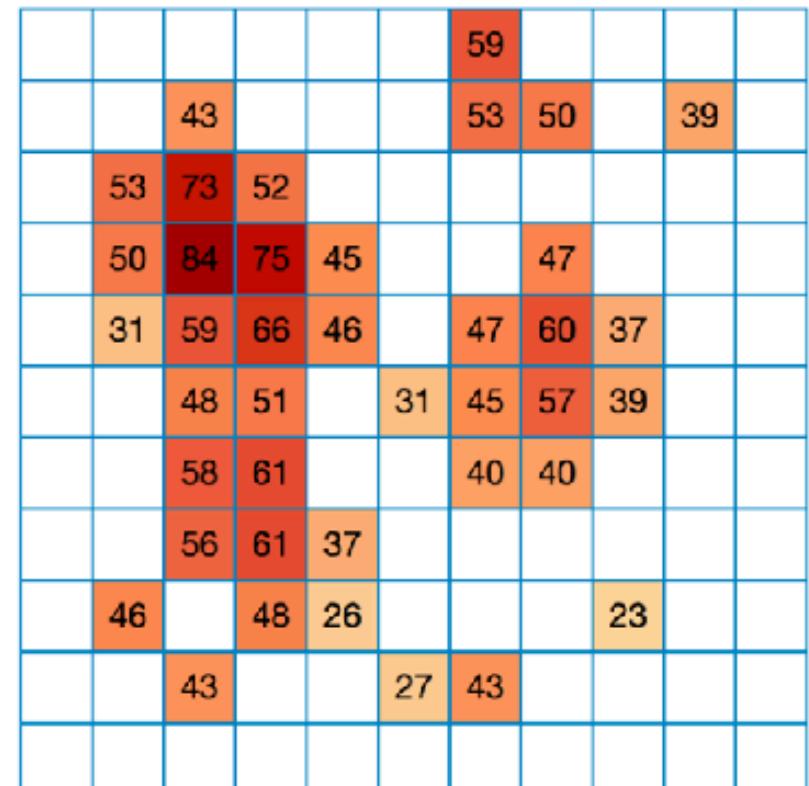


Recreated Visual Field

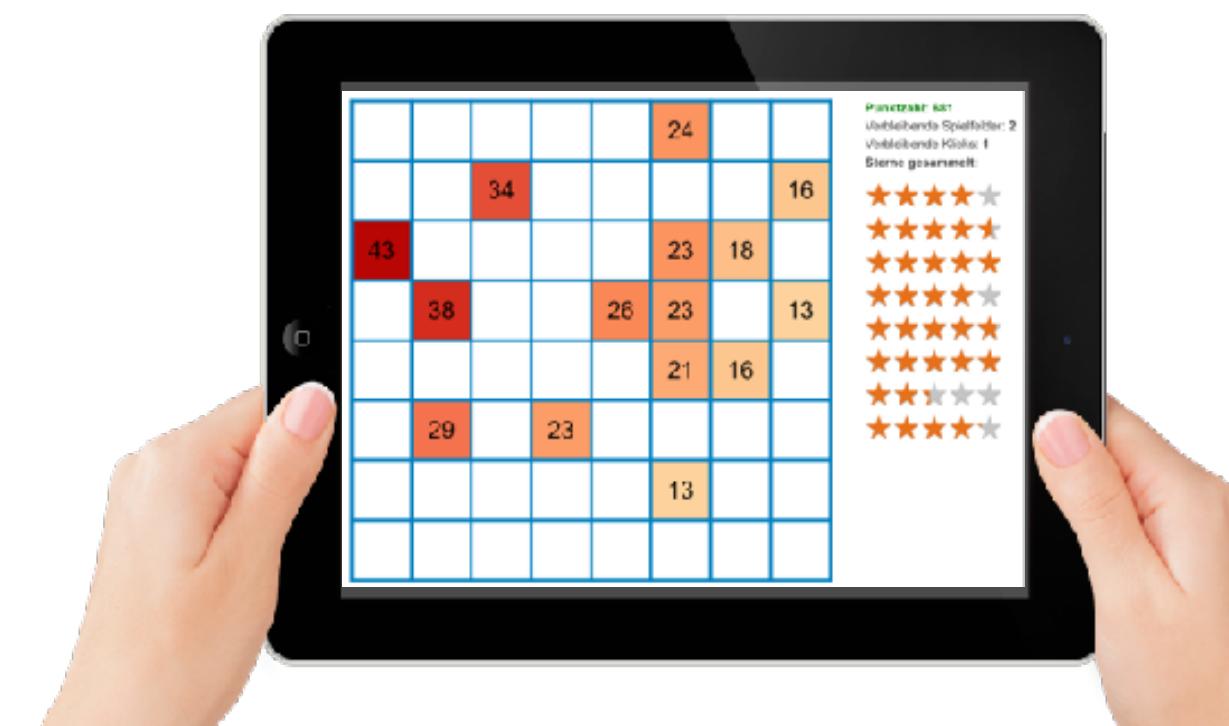


# Questions?

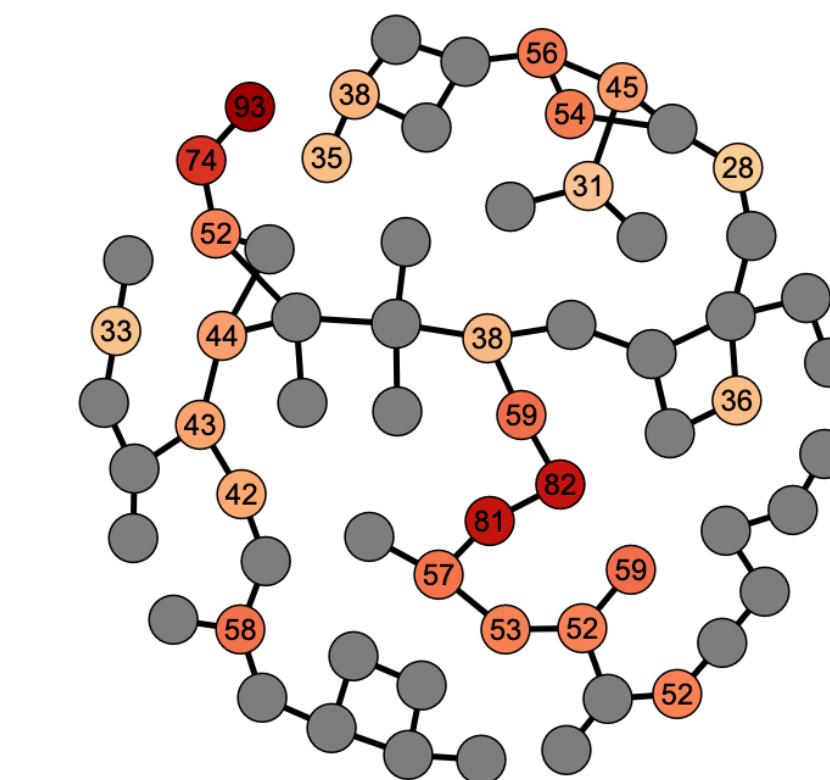
## Grid Search



## Learning Like a Child

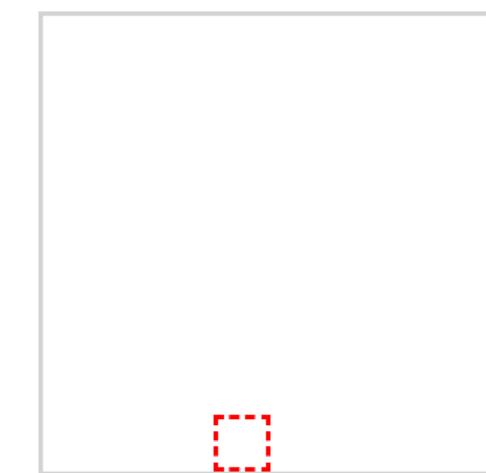


## Graph Search

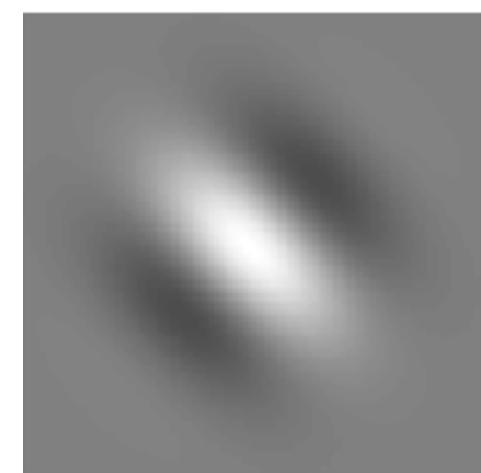


## Conceptual Search

### Spatial



### Conceptual



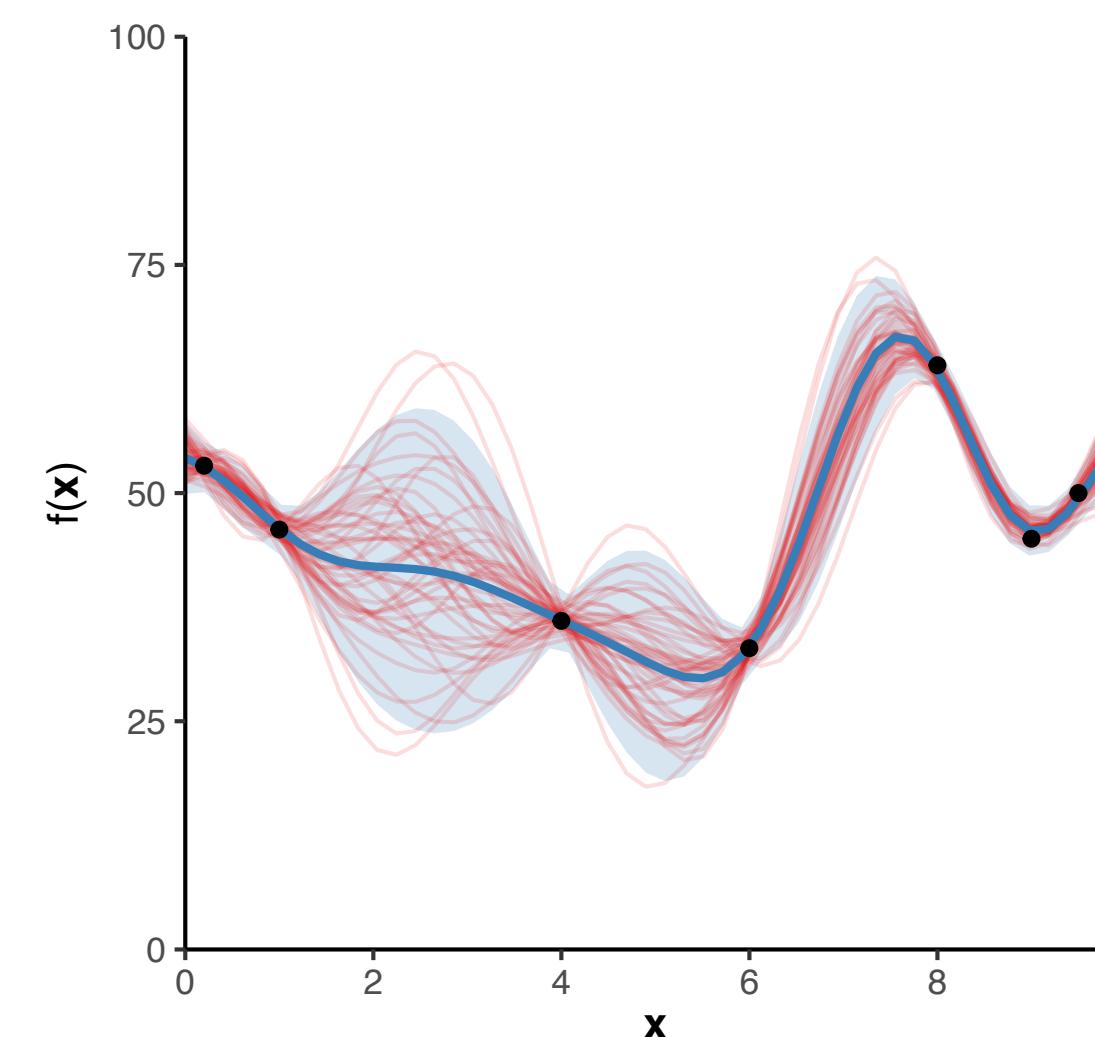
Current Score: 222  
Trials Remaining: 11  
Rounds Remaining: 3

Current Score: 264  
Trials Remaining: 10  
Rounds Remaining: 5

## Safe Search



## Gaussian Processes



## Generalization

