# DeepMind

# REINFORCEMENT LEARNING
## Computational Modeling for Learning and Decision Making

Maria K. Eckstein
Google DeepMind
([mariaeckstein@deepmind.com](mailto:mariaeckstein@deepmind.com))

# Reinforcement Learning (RL)



-> What do both videos have in common?

# What is RL?

Learning from rewards;

and punishment.

# How to Use RL (as a Cognitive Model)?

| Goal | Reward | Ingredients | Algorithm |
|------|--------|-------------|-----------|



+1

action = [→, ←]

state = [           ]

reward = [0, +1]

$$Q(s,a) \leftarrow Q(s,a) + \alpha\ RPE$$

$$RPE = r + \gamma\ Q(s',a') - Q(s,a)$$



action = [jump, stand]

state = [           ]

reward = [0,     ]

**???**

# Questions?

DeepMind

# Lecture Roadmap

# Reinforcement Learning (RL)

1. **Introduction**
2. RL from a psychology perspective
3. RL from an AI perspective
4. RL from a neuroscience perspective
5. Bringing it all together: RL as a cognitive model
6. Conclusion

# Reinforcement Learning (RL)

1. Introduction
2. **RL from a psychology perspective**
3. RL from an AI perspective
4. RL from a neuroscience perspective
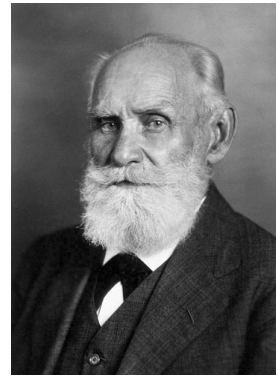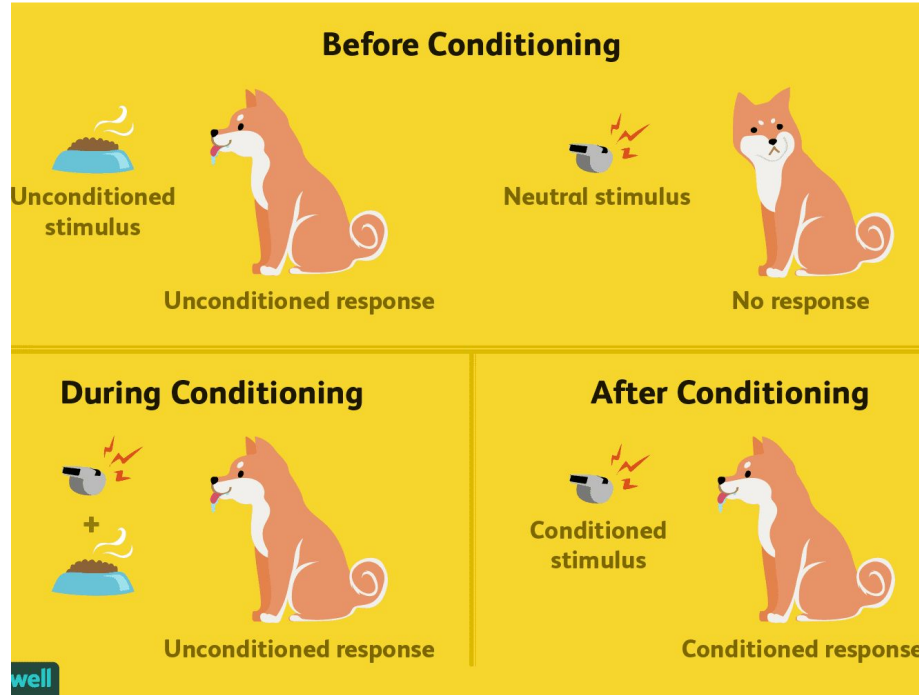5. Bringing it all together: RL as a cognitive model
6. Conclusion

DeepMind
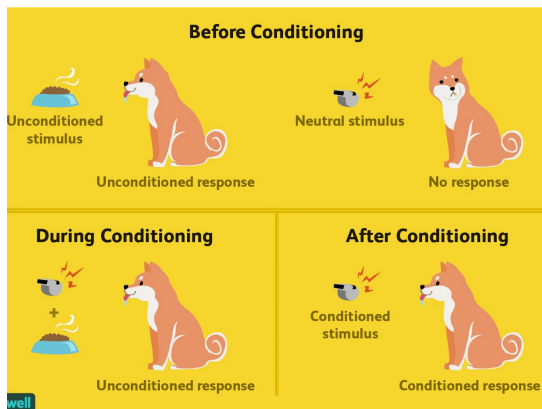
RL from a psychology perspective

# Classical Conditioning



Ivan Pavlov
(1849–1936)

Animals learn associations between US (e.g., food) and neutral CS (e.g., bell) when they reliably co-occur.

# The Rescorla–Wagner Model (1972)



$$\text{RPE} = \lambda - \Sigma[\text{value}(\text{CS})]$$

Combined predictive value of all stimuli

$$\text{value}(\text{CS}) \leftarrow \text{value}(\text{CS}) + \alpha_{CS} * \beta_{US} * \text{RPE}$$

New value (after learning)
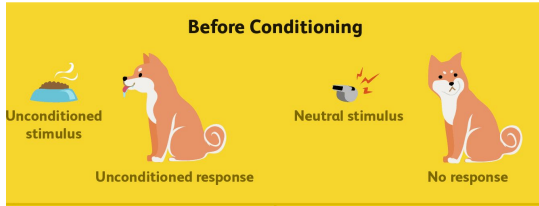
Old value (before learning)

- Stimuli (CS) have "associative strength" (value)
  - Does the stimulus predict a US (reward)?
- When reward arrives, there might a "reward prediction error" (RPE)
  - Was the reward predicted by the present stimuli?
- RPEs trigger learning: update values to predict reward better
  - $\lambda$ is the maximum conditioning possible for the US
  - Learning speed depends on "salience" ($\alpha_{CS}$) and "association value" ($\beta_{US}$)

# Rescorla–Wagner Example

$$\text{RPE} = \lambda - \Sigma[\text{value}(\text{CS})]$$

$$\text{value}(\text{CS}) \leftarrow \text{value}(\text{CS}) + \alpha_{CS} * \beta_{US} * \text{RPE}$$

[[Assume $\alpha_{CS} * \beta_{US} = 0.5$ and $\lambda = 1$]]

**Before Conditioning**

Unconditioned stimulus

Unconditioned response

Neutral stimulus

No response

```
value(bell):        0
λ:                  1
RPE:                1
New value(bell):    0.5
```

**During Conditioning**

Unconditioned response

well

```
value(bell):        0.5
λ:                  1
RPE:                0.5
New value(bell):    0.75
```

**During Conditioning**

Unconditioned response

well

```
value(bell):        0.75
λ:                  1
RPE:                0.25
New value(bell):    0.865
```

???

Conditioned stimulus

Conditioned response

```
value(bell):        1
```

"Conditioned response"

# Blocking Example

$$RPE = \lambda - \Sigma[value(CS)]$$

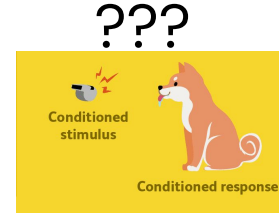$$value(CS) \leftarrow value(CS) + \alpha_{CS} * \beta_{US} * RPE$$

[[Assume $\alpha_{CS} * \beta_{US} = 0.5$ and $\lambda = 1$]]



Unconditioned response

```
value(bell):        1
λ:                  1
RPE:                0
New value(bell):    1  (no change)
```



Unconditioned response

```
value(bell):        1
value(light):       0
Σ[value(CS)]:       1
λ:                  1
RPE:                0
New value(bell):    1  (no change)
New value(light):   0  (no change)
```



???

Unconditioned response

```
value(light):        0
```

No "Conditioned response"

# Operant conditioning

$$\text{RPE} = \text{reward} - \text{value}(\text{action}|\text{state})$$

$$\text{value}(\text{action}|\text{state}) \leftarrow$$

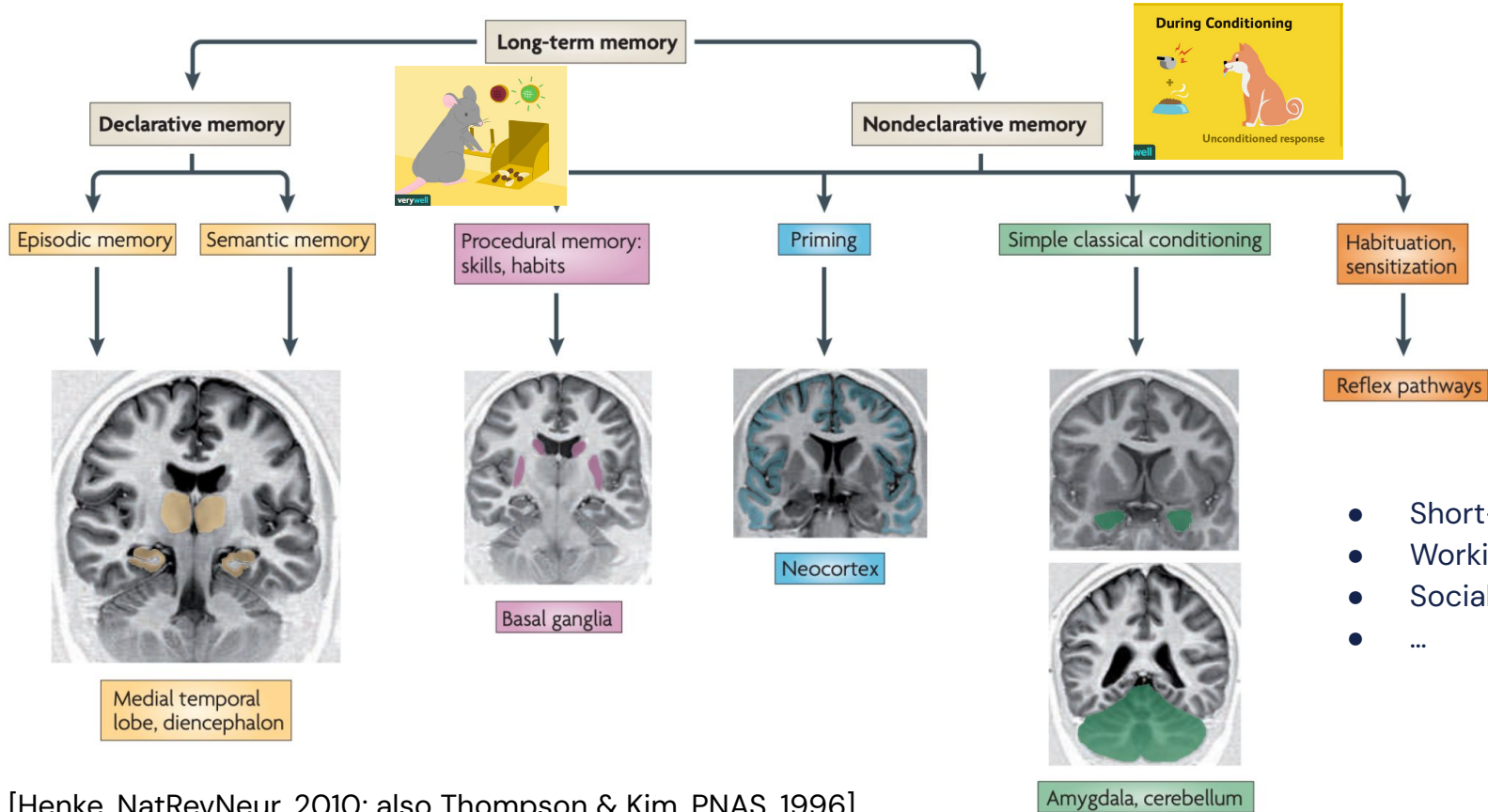$$\text{value}(\text{action}|\text{state}) + \alpha * \text{RPE}$$



```
value(press|lev):         0
reward:                   1
RPE:                      1
New value(press|lev):     0.5
```



```
value(press|lev):         0.5
reward:                   1
RPE:                      0.5
New value(press|lev):     0.75
```

...



```
value(press|lev):         1
```

**Quizz**: According to this theory, what would the trained rat do when it is fully satiated and sees the lever?
A)   Press the lever
B)   Not press the lever

-> Link to "*habitual*" versus "*goal-directed*" behavior.

# Multiple memory systems

[Henke, NatRevNeur, 2010; also Thompson & Kim, PNAS, 1996]

# Questions?

# Reinforcement Learning (RL)

1. Introduction
2. RL from a psychology perspective
3. **RL from an AI perspective**
4. RL from a neuroscience perspective
5. Bringing it all together: RL as a cognitive model
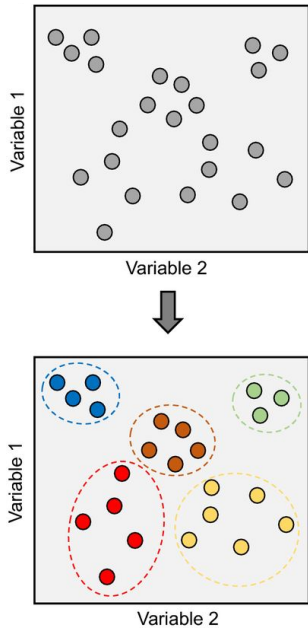6. Conclusion

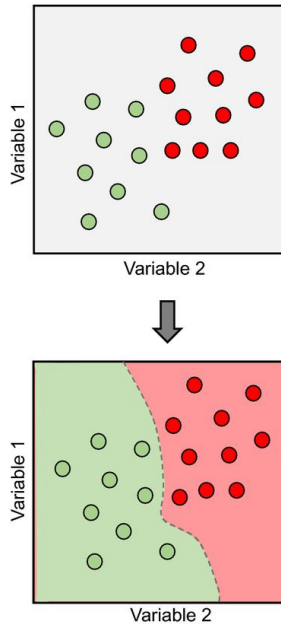# DeepMind

# RL from an AI perspective

Slide credit:
Maria Eckstein
(mariaeckstein@deepmind.com)

# RL in the context of machine learning (ML)

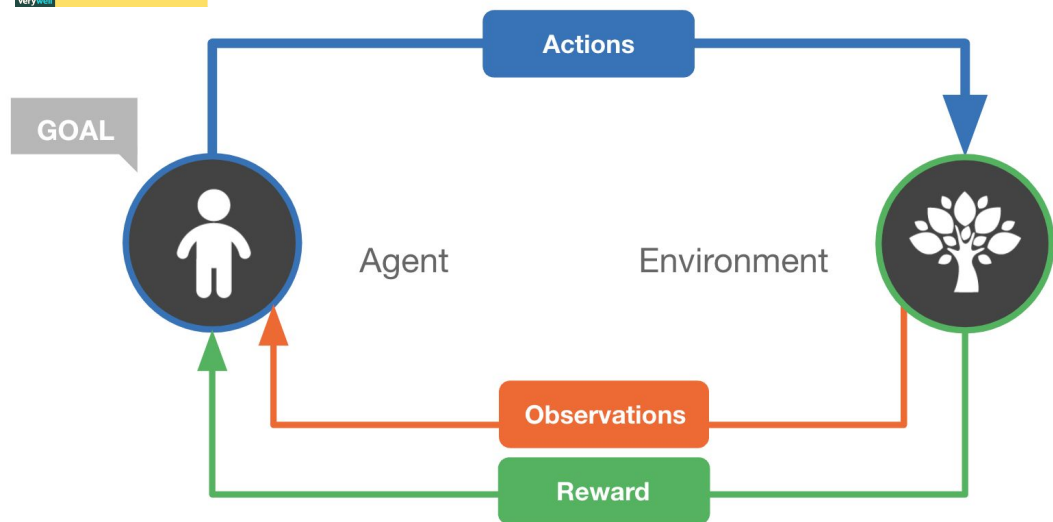**Unsupervised learning:** Learn patterns or structure in data

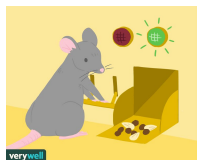(e.g., dimensionality reduction, clustering, ...)

**Supervised learning:** Learn to predict target(s)

(e.g., regression, classification, ...)

**Reinforcement Learning:** Learn from interactions in the world, through a scalar reward signal

*Credit:*
**Maria Eckstein**
(mariaeckstein@deepmind.com)

# RL Ingredients

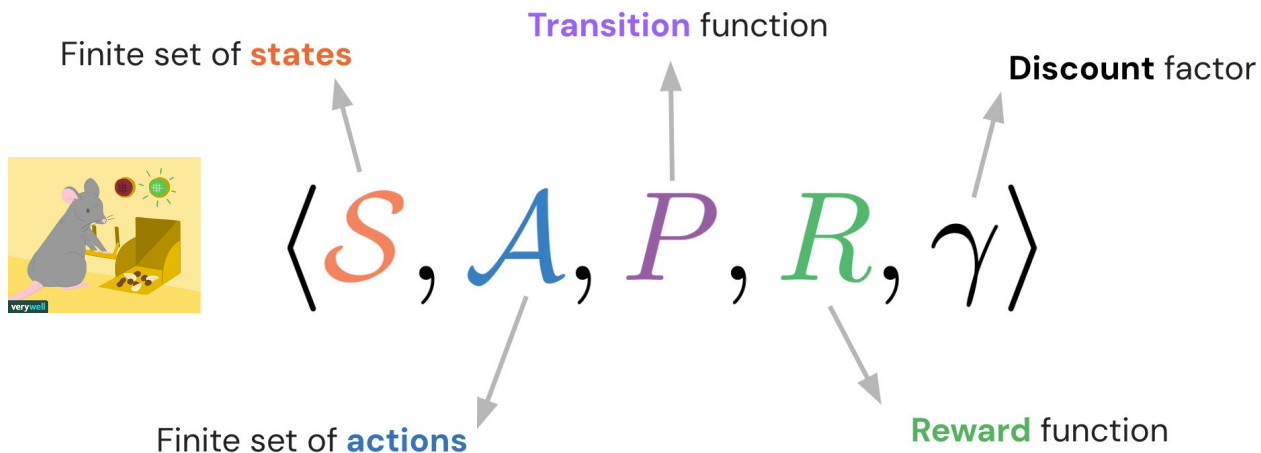**Agent**: Learns a policy π that maps observations to actions, in order to maximize rewards.

**Environment**: E.g., experimental task; game (chess, Starcraft); factory (robotics); fusion reactor; …

**Reward**:

- *Extrinsic* (food, water, hard–coded)
- *Intrinsic* (curiosity, novelty, empowerment, learning progress, compression, explanation, …)

# The Markov Decision Process (MDP)

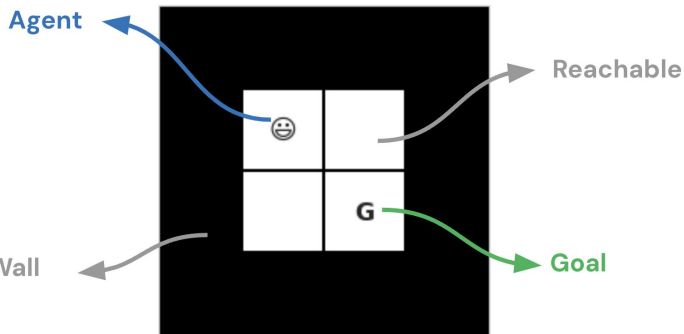**Markov Decision Processes** allow us to *formalize* and *solve* the RL problem.

Finite set of **states**

**Transition** function

**Discount** factor

$$\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$$

Finite set of **actions**

**Reward** function

**Markov Property**: The next state depends only on the current state and action, not on the entire history (e.g., chess).
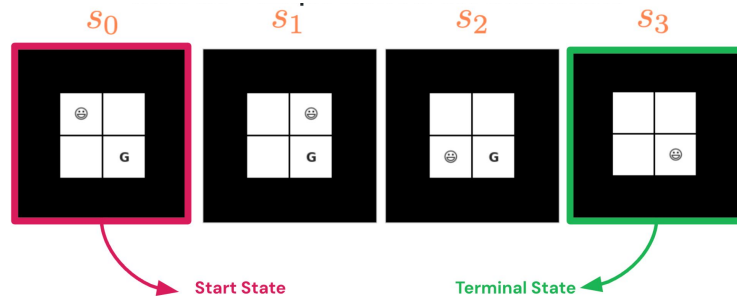
$$\mathbf{P}(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, ..., s_0) = \mathbf{P}(s_{t+1}|s_t, a_t)$$

Future　Present　　　　Past　　　　　Future　Present

*Credit:*
**Maria Eckstein, Jane Wang, Feryal Behbahani (mariaeckstein@deepmind.com)**

# Grid Worlds

State space $\mathcal{S}$



Start State

Terminal State

Size of the world: **[ 2 X 2 ]**

Agent

Reachable

Wall

Goal

Action Space $\mathcal{A}$



Up

Left       Right

Down

Transition model $P$



$s_0$  $s_1$

$s_2$  $s_3$

up / left

$s_0$  right  $s_1$

down

$s_2$  $s_3$

up / left   left   up / right

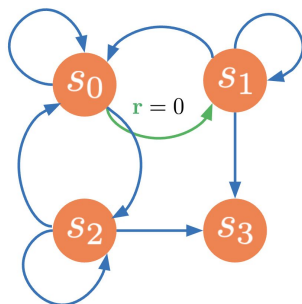$s_0$  right  $s_1$

up      down       down

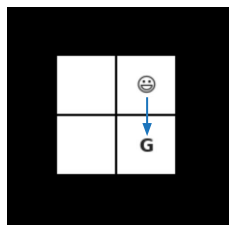$s_2$  right  $s_3$

left/down

Rewards $R$

Empty cell: 0
Wall: –5
Goal: +10

# Policy and Values

In MDP terms:



[assume γ = 0.9]

$$Q^{\pi^*}(s_0, a_{right}) = 0 + 0.9 * 10 = 9$$



$$Q^{\pi^*}(s_1, a_{down}) = 10 + 0.9 * 0 = 10$$

**Agent's goal**: Maximize (γ-discounted) sum of future rewards:

$$\mathbf{G_t} = \mathbf{r_t} + \gamma \mathbf{r_{t+1}} + \gamma^2 \mathbf{r_{t+2}} + \dots$$

$$\underbrace{\phantom{\mathbf{G_t} = \mathbf{r_t} + \gamma \mathbf{r_{t+1}} + \gamma^2 \mathbf{r_{t+2}}}}_{\textbf{Return}}$$

To achieve this, the agent learns an action **policy π**:

$$\mathbf{a_t} \sim \boldsymbol{\pi}(\mathbf{a_t} | \mathbf{s_t})$$

***How do we find this policy?***

Using values! Once we have (optimal) values, executing the optimal policy is easy:

$$\boldsymbol{\pi}^*(\mathbf{s}) = \max_{\mathbf{a}} \mathbf{Q}^*(\mathbf{s}, \mathbf{a})$$

This works because values are defined as:

$$\mathbf{Q}^{\boldsymbol{\pi}}(\mathbf{s_t}, \mathbf{a_t}) = \mathbb{E}_{\boldsymbol{\pi}} \left[ \mathbf{r_t} + \gamma \mathbf{r_{t+1}} + \gamma^2 \mathbf{r_{t+2}} + \dots | \mathbf{s_t}, \mathbf{a_t} \right]$$
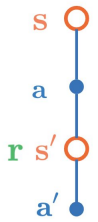
# Learning the value function

$$Q^\pi(s_t, a_t) = \mathbb{E}_\pi \left[ r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + ... | s_t, a_t \right]$$

**Practically**: We can't predict the future! (And we don't want to...)

$$P(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, ..., s_0) = P(s_{t+1}|s_t, a_t)$$

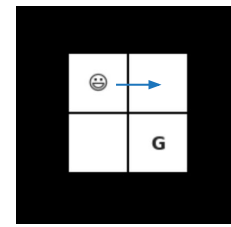**SARSA (on-policy control)**

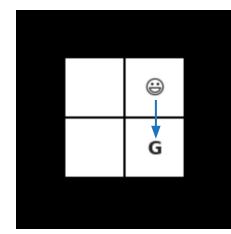- Bootstrapping value updates based on "on-policy" (actual) experience

s

a

r s′

a′

reward   **Value next state**   Old value

$$\underbrace{Q(s,a)}_{\text{new value}} \leftarrow \underbrace{Q(s,a)}_{\text{old value}} + \alpha(\underbrace{r + \gamma Q(s', a') - Q(s,a)}_{\text{RPE}})$$
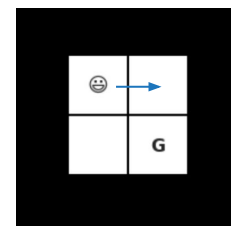
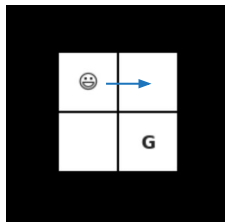**Q-learning (off-policy control)**

- Bootstrapping value updates based on "off-policy" (hypothetical) experience

s

a

r s′

a′

Best avail. action
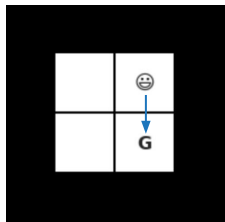
$$Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \underbrace{\max_{a'} Q(s', a')} - Q(s,a))$$

$$\pi^*(s) = \max_a Q^*(s,a)$$

All Q's=0

Q(s0,right) <- 0+α
(0+γ*0−0) = 0

Q(s1,down) <- 0+α
(10+γ*0−0) = 5

Q(s0,right) <- 0+α
(0+γ*5−0) = 2.25

# Temporal Difference (TD) Learning

```
Value(s) += α * RPE
RPE = r - Value(s)
```

**Learning rate**

$$V_{t+1}\left(s_t\right) \leftarrow V_t\left(s_t\right) + \eta \delta_t$$

**Reward prediction error**

**New estimate of value of current state**

**Old estimate of value of current state**

$$\delta_t = R_t + \gamma V_t\left(s_{t+1}\right) - V_t\left(s_t\right)$$

**Reward prediction error**

**Actual observed value of current state (written the recursive way)**

**Old estimate of value of current state**

# Learning the value function

**Problem**: We can't predict the future! (And we don't want to...)

### SARSA (on-policy control)

- Bootstrapping value updates based on "on-policy" (actual) experience

s O
a |
r s′ O
a′ |

$$\underbrace{Q(s,a)}_{\text{new value}} \leftarrow \underbrace{Q(s,a)}_{\text{old value}} + \alpha(\underbrace{\overbrace{r}^{\text{reward}} + \overbrace{\gamma Q(s',a')}^{\substack{\text{Value next}\\\text{state}}} - \overbrace{Q(s,a)}^{\text{Old value}}}_{\text{RPE}})$$

### Q-learning (off-policy control)

- Bootstrapping value updates based on "off-policy" (hypothetical) experience

s O
a |
r s′ O
a′

$$Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma \overbrace{\max_{a'} Q(s',a')}^{\text{Best avail. action}} - Q(s,a))$$

$$\pi^*(s) = \max_a Q^*(s,a)$$

All Q's=0

Q(s0,right) <- 0+α
(0+γ*0−0) = 0

Q(s1,down) <- 0+α
(10+γ*0−0) = 5
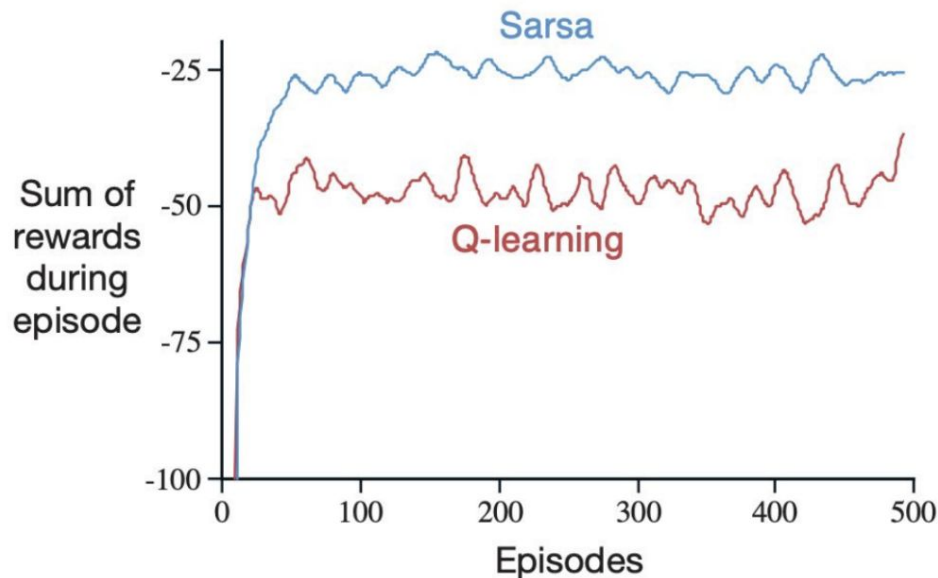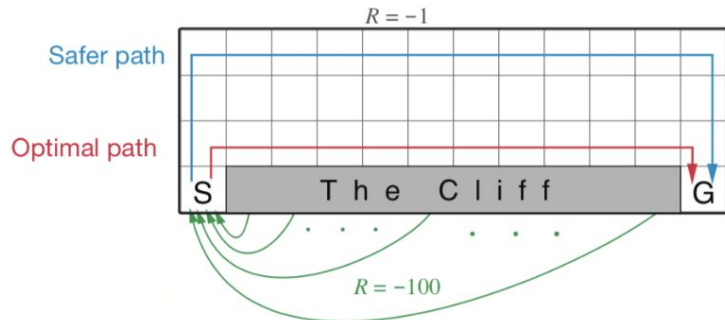
Q(s0,right) <- 0+α
(0+γ*5−0) = 2.25

# SARSA vs Q-Learning: The Cliff-walking Example



Exploration

Sutton & Barto. Reinforcement Learning: An Introduction. (Chapter 6)

# SARSA vs Q-Learning: The Cliff-walking Example

- **Q–learning** learns the **optimal path** while its online performance is worse than **SARSA**.

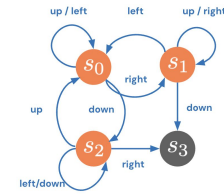- **SARSA** learns the **safer path**.

# Cheat Sheet

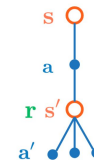**Rescorla Wagner**:    keep track of reward expectations



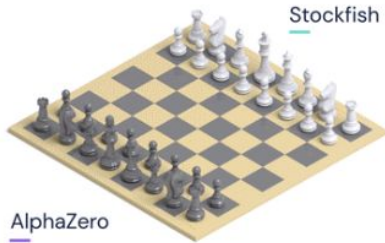**TD Learning**:    +over time



**SARSA**:    +control (on-policy)



**Q-Learning**:    +control (off-policy)

# Real-world Reinforcement Learning: **Examples**

**Game playing**

Stockfish

AlphaZero

**Robotics / Manipulation**

**User personalization**

**Self-driving cars**

**Managing energy usage**

# Questions?

# Reinforcement Learning (RL)

1. Introduction
2. RL from a psychology perspective
3. RL from an AI perspective
4. **RL from a neuroscience perspective**
5. Bringing it all together: RL as a cognitive model
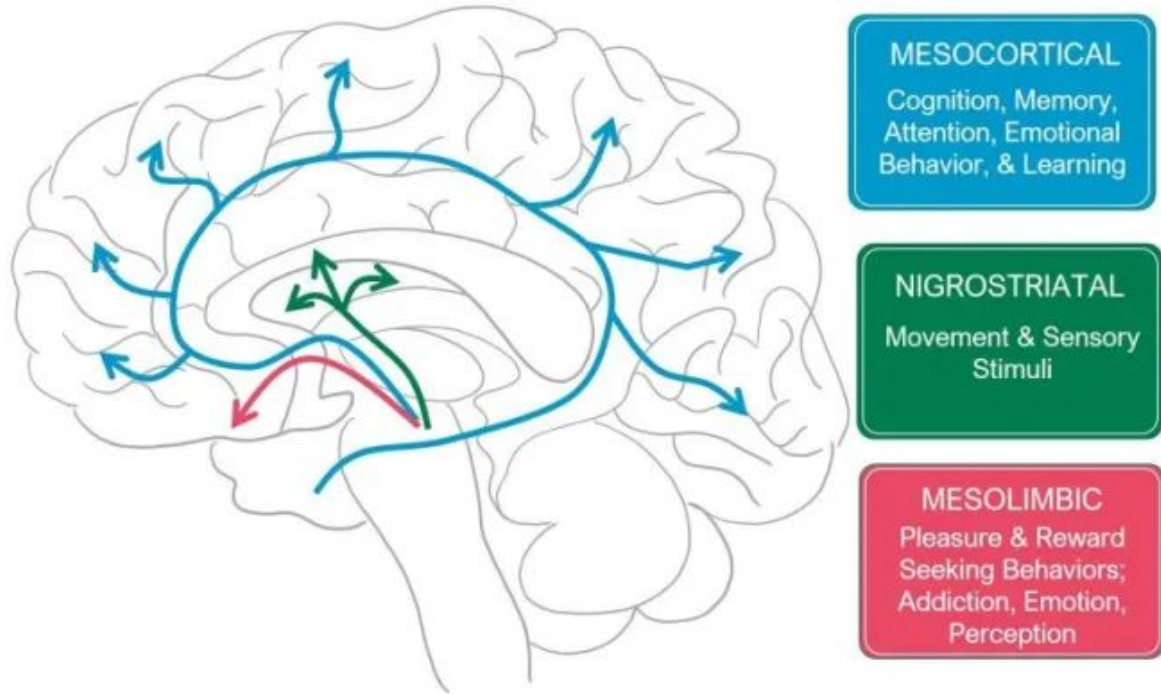6. Conclusion

# DeepMind

# RL in neuroscience

# The Neurotransmitter Dopamine

**MESOCORTICAL**
Cognition, Memory, Attention, Emotional Behavior, & Learning

**NIGROSTRIATAL**
Movement & Sensory Stimuli

**MESOLIMBIC**
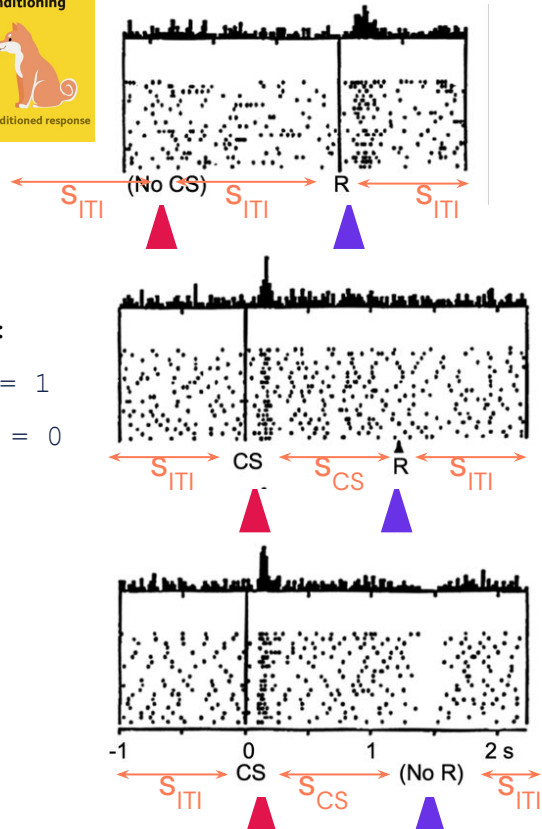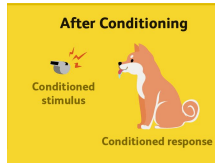Pleasure & Reward Seeking Behaviors; Addiction, Emotion, Perception

Essential for theory of reinforcement learning!

# Dopamine Reward Prediction Errors

**Quizz**: What does dopamine firing represent?

$$\text{TD RPE} = r + \gamma\, V(s') - V(s)$$



After Conditioning

Values:
$$V(s_{CS}) = 1$$
$$V(s_{ITI}) = 0$$

$$\text{RPE} = r + \gamma\, V(s_{ITI}) - V(s_{ITI})$$
$$= 0 + \gamma\, 0 \qquad - 0 = \mathbf{0}$$

$$\text{RPE} = r + \gamma\, V(s_{ITI}) - V(s_{ITI})$$
$$= 1 + \gamma\, 0 \qquad - 0 = \mathbf{1}$$

$$\text{RPE} = r + \gamma\, V(s_{CS}) - V(s_{ITI})$$
$$= 0 + \gamma\, 1 \qquad - 0 = \mathbf{0.9}$$

$$\text{RPE} = r + \gamma\, V(s_{ITI}) - V(s_{CS})$$
$$= 1 + \gamma\, 0 \qquad - 1 = \mathbf{0}$$

$$\text{RPE} = r + \gamma\, V(s_{CS}) - V(s_{ITI})$$
$$= 0 + \gamma\, 1 \qquad - 0 = \mathbf{0.9}$$

$$\text{RPE} = r + \gamma\, V(s_{ITI}) - V(s_{CS})$$
$$= 0 + \gamma\, 0 \qquad - 1 = \mathbf{-1}$$

MESOCORTICAL
Cognition, Memory, Attention, Emotional Behavior, & Learning

NIGROSTRIATAL
Movement & Sensory Stimuli

MESOLIMBIC
Pleasure & Reward Seeking Behaviors; Addiction, Emotion, Perception
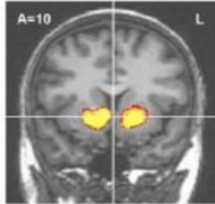
- Converging evidence across studies and species
- Mostly in simple conditioning paradigms

[Niv, 2009]

[Montague et al., 1996; Schultz et al., 1997]

Credit:
Maria Eckstein
(mariaeckstein@deepmind.com)

# Human fMRI



money
value predicted
(Daw et al 2006)

faces
attractiveness
(O'Doherty et al 2003)

Coke or Pepsi
degree favored
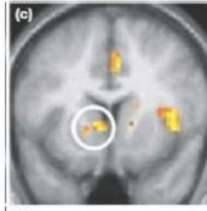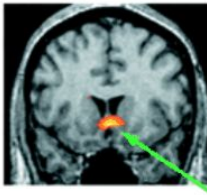(McClure et al. 2004)

money
gain vs loss
(Kuhnen & Knutson
2005)

food odors
valued vs devalued
(Gottfreid et al 2003)

juice
unpredictable vs
predictable
(Berns et al 2001)

Rewards / reward
anticipation activate:

- Ventromedial
  prefrontal cortex
- Orbitofrontal cortex
- Striatum

➢ *Generalized
  appetitive
  function?*

# Questions?

# Reinforcement Learning (RL)

1.  Introduction
2.  RL from a psychology perspective
3.  RL from an AI perspective
4.  RL from a neuroscience perspective
5.  **Bringing it all together: RL as a cognitive model**
6.  Conclusion

# DeepMind

# RL for Cognitive Modeling

Slide credit:
Maria Eckstein
(mariaeckstein@deepmind.com)

# What is Cognitive Modeling?

**Goal**: Understand behavior, cognitive process

**Method**:

- Find model (e.g., RL, Regression, DDM, ...)
- "Fit" model (find best parameters, using cross-entropy loss / negative log likelihood)
- Expand model
  - e.g., forgetting; reward vs punishment [Frank et al., 2004]; WM [Collins & Frank, 2012]; counterfactuals [Boorman et al., 2011]; ...
- Model comparison (AIC, BIC, WAIC, ...)

**Result**:

- "Cognitive process"
- Fitted parameters (individual differences)
- Normative understanding (optimality)
- Quantitative methods, statistics
- Complex, multi-step processes
- Precise prediction

$$RPE = r + \gamma\, Q(s',a') - Q(s,a)$$

$$Q(s,a) \leftarrow Q(s,a) + \alpha * RPE$$

RL

Slide credit:
**Maria Eckstein**
**(mariaeckstein@deepmind.com)**

# What is RL Modeling?

| Goal | Reward | Ingredients | Algorithm |
|------|--------|-------------|-----------|



+1

**a**ction = [→, ←]

**s**tate = [  ]

**r**eward = [0, +1]

$$\text{RPE} = r + \gamma \, Q(s',a') - Q(s,a)$$
$$Q(s,a) \leftarrow Q(s,a) + \alpha * \text{RPE}$$
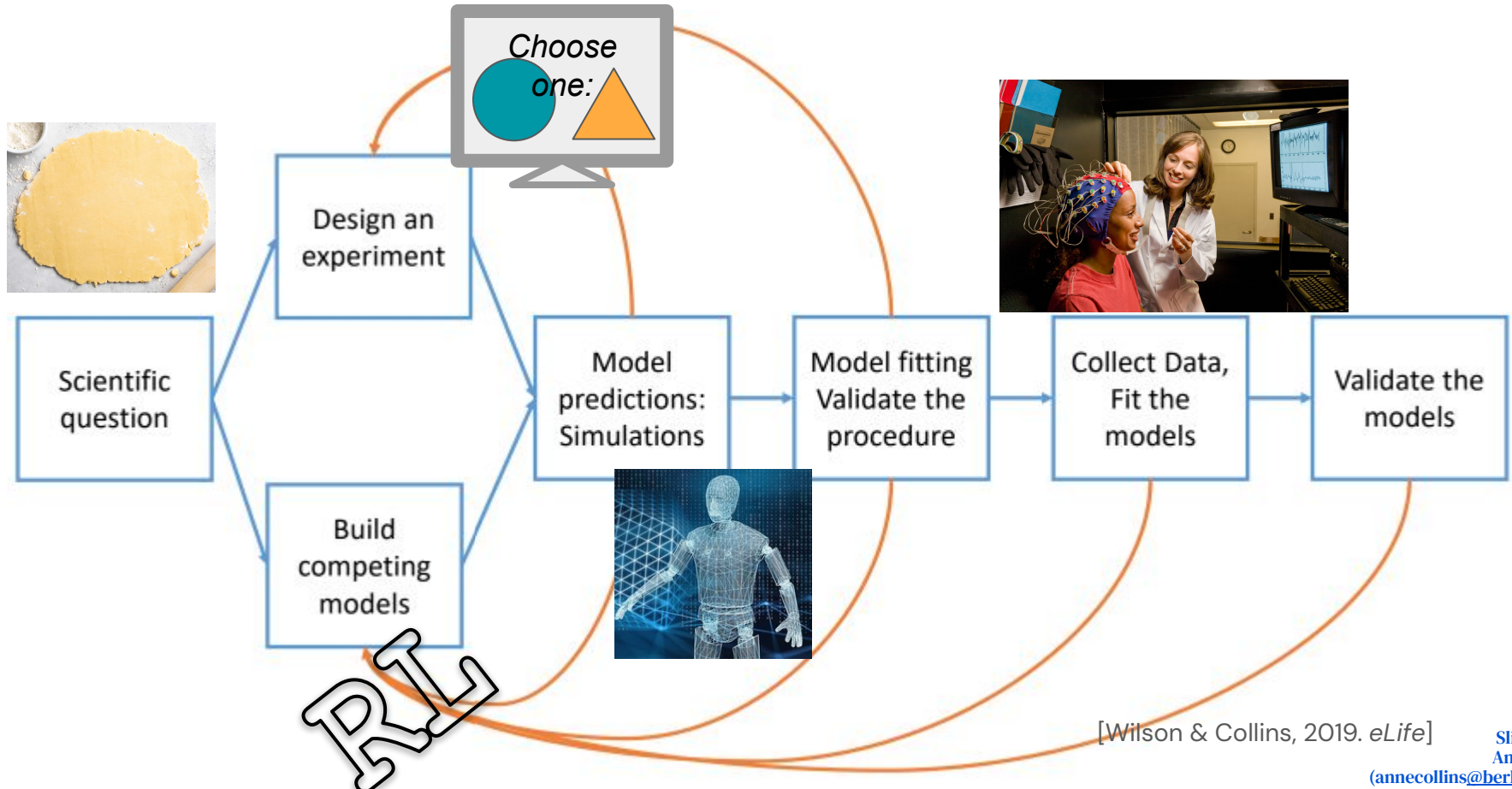


*Choose one:*

+1

**a**ction = [ F  H ]
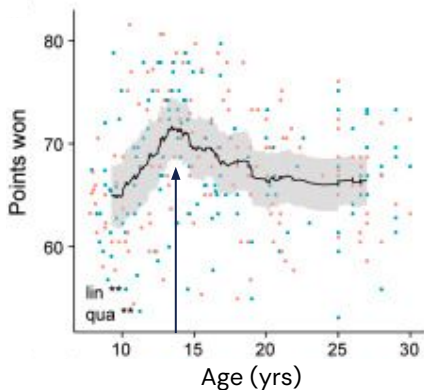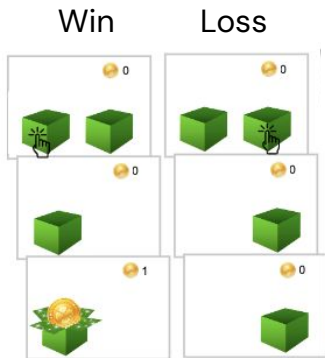
**s**tate = [ ●▲,  ▲● ]

**r**eward = [0, +1]

$$\text{RPE} = r + \gamma \, Q(s',a') - Q(s,a)$$
$$Q(s,a) \leftarrow Q(s,a) + \alpha * \text{RPE}$$

# A Recipe for Cognitive Modeling

[Wilson & Collins, 2019. *eLife*]

# Learning to Reversal Learn

**Goal**: Understand age trajectory of reversal learning

Win    Loss





Points won / Age (yrs)

- Best performance at ~13–15

**Why**? Cognitive mechanism?

[Eckstein, Master, Dahl, Wilbrecht & Collins, 2022. *DCN*]

$$RPE = r - Q(s,a)$$
$$Q(s,a) \leftarrow Q(s,a) + \alpha * RPE$$



perseverance

$\alpha$ positive

Age (yrs)

*Slide credit*:
Maria Eckstein
mariaeckstein@deepmind.com

# Learning to Reversal Learn

$$p(s_t | a_t, r_t) \propto$$
$$p(a_t, r_t | s_t) * p(s_t)$$

**Goal**: Understand age trajectory of reversal learning

Win     Loss



- Best performance at ~13–15

**Why**? Cognitive mechanism?

[Eckstein, Master, Dahl, Wilbrecht & Collins, 2022. *DCN*]



*Slide credit*:
Maria Eckstein
(mariaeckstein@deepmind.com)

# Model-based or model-free RL?



(a)

0s
+ <2s
+ 3s
+ <2s
+ 3s
+ 1.5s

(b)

Stage 1: s

Stage 2: s'



stay probability

common
rare

rewarded   unrewarded   rewarded   unrewarded   rewarded   unrewarded

**Model-free: SARSA**

*At both stages:*

RPE = r – Q(s,a) + Q(s',a')

$Q_{MF}$(s,a) <– Q(s,a) + α * RPE

**Model-based**

*Stage 2:*

RPE = r – Q(s',a')

Q(s',a') <– Q(s',a') + **α** * RPE

*Stage 1:*

$Q_{MB}$(s,a) = p($s_A$'|s,a) * $max_a$ Q($s_A$',a') + p($s_B$'|s,a) * $max_a$ Q($s_B$',a')

**Hybrid**
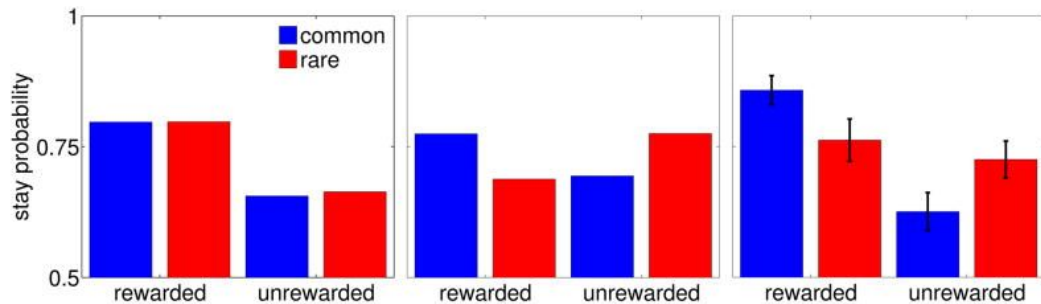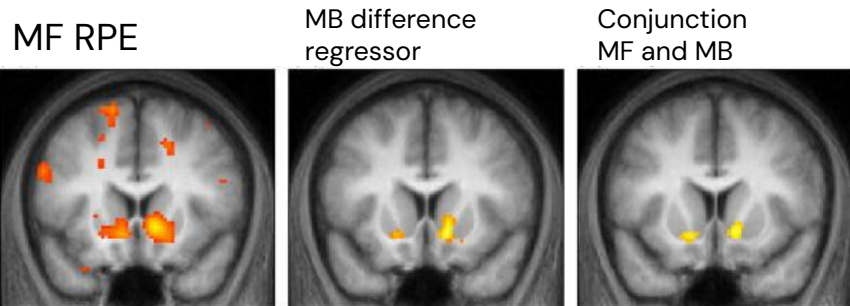
Q(s,a) = w * $Q_{MF}$(s,a) + (1 – w) * $Q_{MB}$(s,a)

# Model-based or model-free RL?

**Results**:

- Hybrid models (LL=3.364) fits data better than MF alone (LL=3.418) or MB alone (LL=3.501)
- Fitted value of w (median across subjects): 0.39



MF RPE      MB difference regressor      Conjunction MF and MB

---

**Model-free: SARSA**

*At both stages:*

$$RPE = r - Q(s,a) + Q(s',a')$$

$$Q_{MF}(s,a) \leftarrow Q(s,a) + \alpha * RPE$$

---

**Model-based**

*Stage 2:*

$$RPE = r - Q(s',a')$$

$$Q(s',a') \leftarrow Q(s',a') + \alpha * RPE$$

*Stage 1:*

$$Q_{MB}(s,a) = p(s_A'|s,a) * \max_a Q(s_A',a') + p(s_B'|s,a) * \max_a Q(s_B',a')$$

---

**Hybrid**

$$Q(s,a) = w * Q_{MF}(s,a) + (1 - w) * Q_{MB}(s,a)$$

# Questions?

*Slide credit:*
**Maria Eckstein**
**(mariaeckstein@deepmind.com)**

# Reinforcement Learning (RL)

1. Introduction
2. RL from a psychology perspective
3. RL from an AI perspective
4. RL from a neuroscience perspective
5. Bringing it all together: RL as a cognitive model
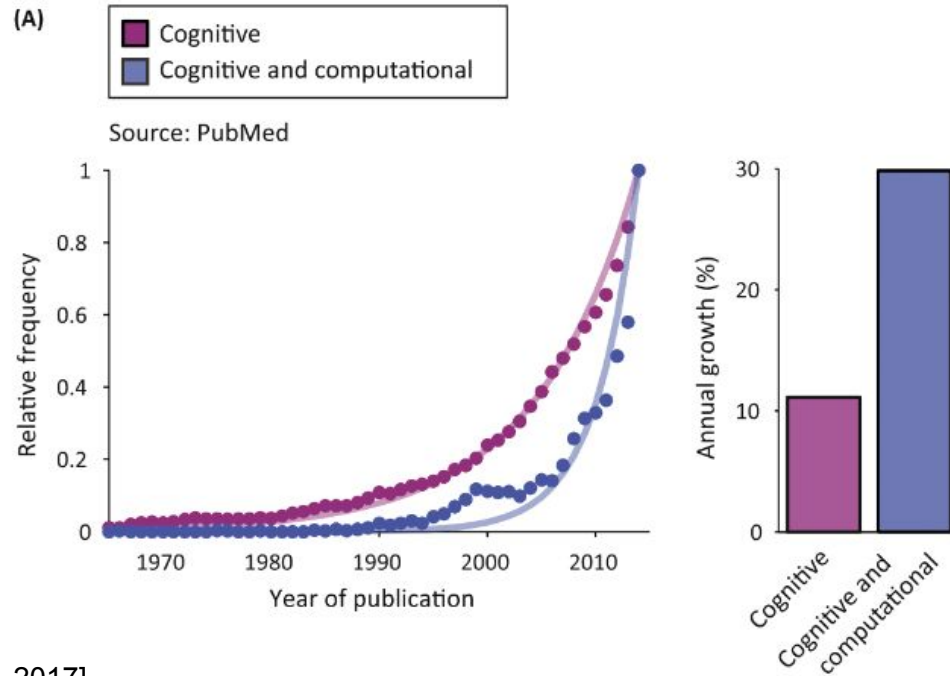6. **Conclusion**

# DeepMind

# Conclusion

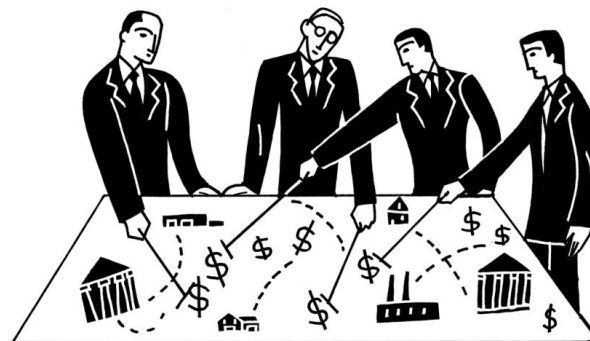# Computational modeling is on the rise!



(A)

Cognitive
Cognitive and computational

Source: PubMed

[Palminter et al., 2017]

# Where do rewards come from?



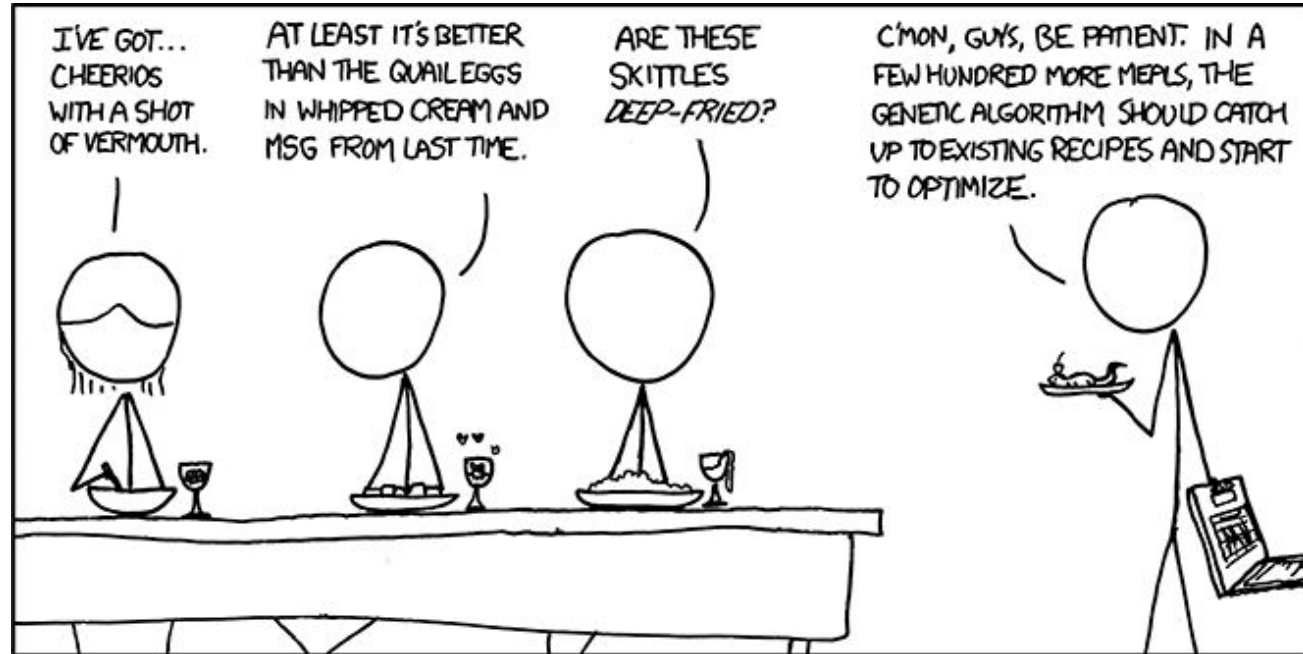Evolution?



Economists?

- Intrinsic / extrinsic?
- Innate / learned?
- Context-dependent?
- Individual differences?

# Exploration

- Epsilon–greedy / softmax?
- Structured exploration?
- Intrinsic goals?
- Sparse rewards

# Credit Assignment

Beginning



End



How to link distal outcomes to earlier causes despite many intervening events?

How to generalize over similar + different instances?

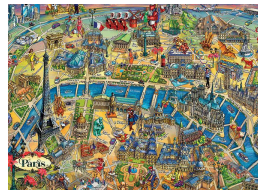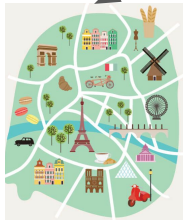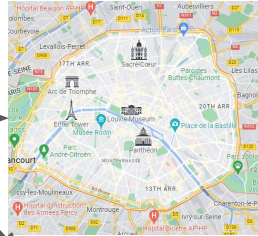How to use knowledge of structure inform credit assignment?

Slide Credit:
Kim Stachenfeld
(stachenfeld@deepmind.com)

# Models as Maps

### Original



### Model



- Cognitive model = map
  - Smaller, more abstract
  - Loose information

- Different maps
  - Depending on the purpose
  - No one "true" map

# Questions?

*Slide credit:*
**Maria Eckstein**
**(mariaeckstein@deepmind.com)**

# Want to Learn More?

**Books**

- [Reinforcement Learning: an Introduction by Sutton & Barto](#)
- [Algorithms for Reinforcement Learning by Csaba Szepesvari](#)

**Lectures and course**

- [Neuromatch Lecture on RL by Jane Wang and Feryal Behbahani](#)
- [RL Course by David Silver](#)
- [Reinforcement Learning Course | UCL & DeepMind](#)
- [Emma Brunskill Stanford RL Course](#)
- [RL Course on Coursera by Martha White & Adam White](#)

**More practical**

- [Spinning Up in Deep RL by Josh Achiam](#)
- [Acme white paper](#) & [Colab tutorial](#)
- [OpenAI Gym](#)

Reinforcement
Learning

An Introduction
second edition

Richard S. Sutton and Andrew G. Barto

# Acknowledgements

**Kim Stachenfeld, Anne Collins, Jane Wang, Feryal Behbahani**, Nathaniel Daw, Chris Knutsen, Kevin Miller, Zeb Kurth-Nelson, Matt Botvinick, Chris Summerfield
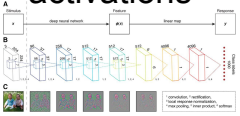
Dog
tricks
by
Justy

# Bonus

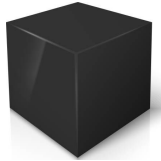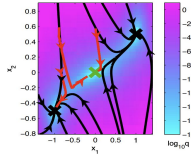# Theory–driven vs Data–driven Models

**Analyze activations**



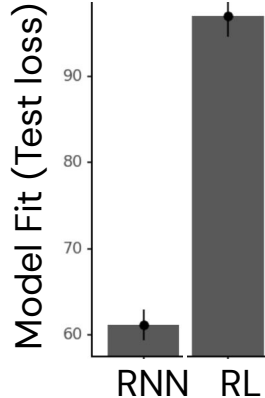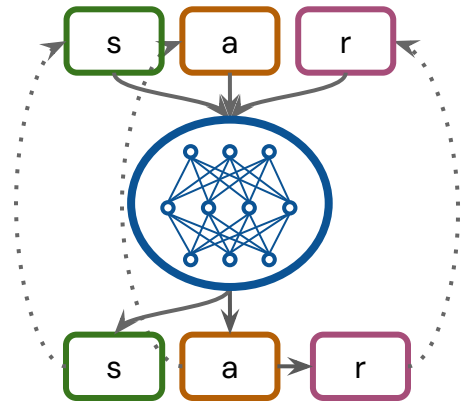**Explainability…**

**Analyze dynamics**



"You won 81 points!"



**Cognitive development**

[van den Bos et al., 2012; Lefebvre et al., 2017; Nussenbaum & Hartley, 2019; Master et al., 2020; Eckstein et al., 2022]

**Brain function**

[Daw et al., 2006; O'Doherty et al., 2007; Dayan & Niv, 2008; Miller et al., 2017; Starkweather et al., 2018]

y=10

**Psychiatry**

RPE

[Maia & Frank, 2011; Montague et al., 2012; Huys et al., 2016; Redish & Gordon, 2016; Hauser et al., 2019]

r = -0.376
p = 0.049

z = -4

Abstraction, Model-based, Habits, Exploration, Sequences, …

## Vanilla RNN
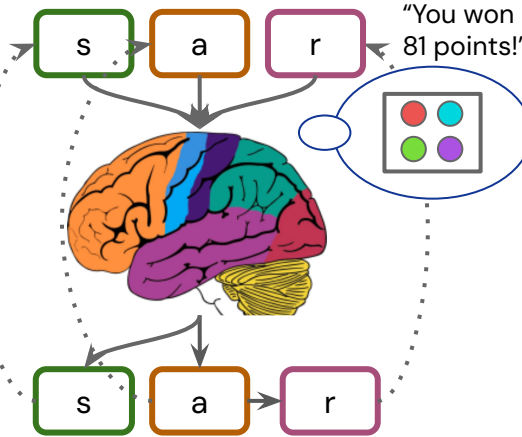


**Model Fit (Test loss)**


RNN    RL

**Trade-offs**

- Predictive power (RNN) vs Interpretability (RL)
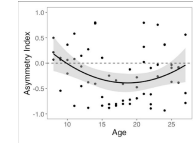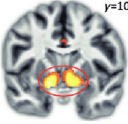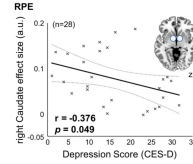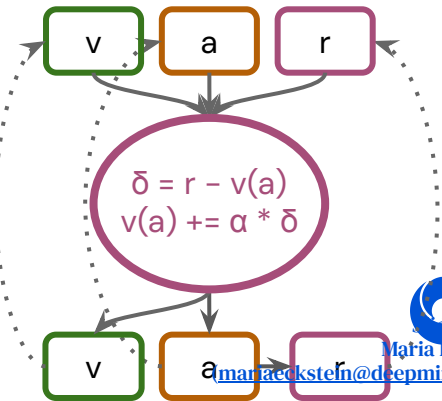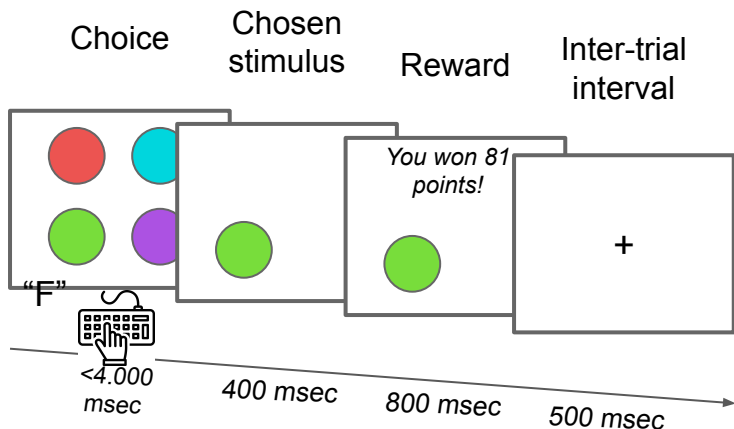- What makes a good model?
  [Navarro, 2019; Box, 1979; Eckstein et al., 2021]

**Uncover the cognitive process**

- *Why* is RL underperforming?
- *Which* cognitive processes are missing?
- *Which* assumptions are wrong?

## Classic RL

$\delta = r - v(a)$
$v(a) += \alpha * \delta$

Credit: Maria Eckstein
mariaeckstein@deepmind.com

# Dataset



Choice — Chosen stimulus — Reward — Inter-trial interval

*You won 81 points!*

"F"

<4.000 msec    400 msec    800 msec    500 msec

**Key points**

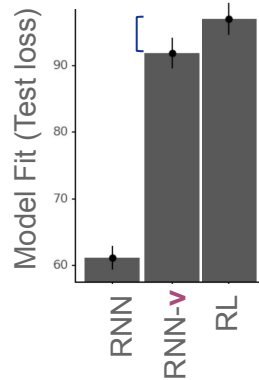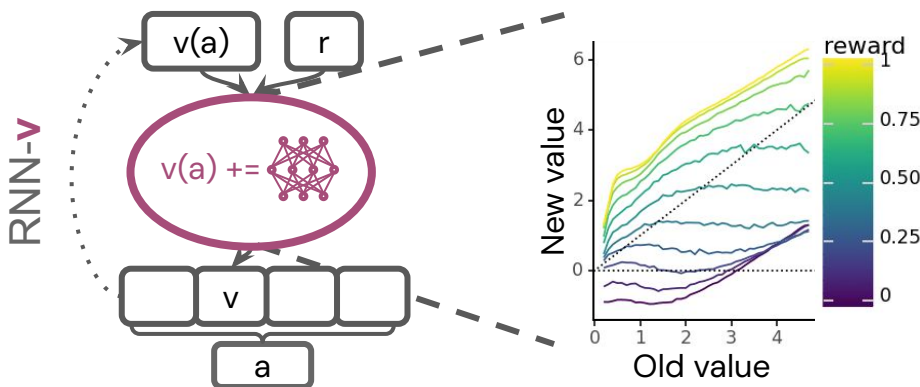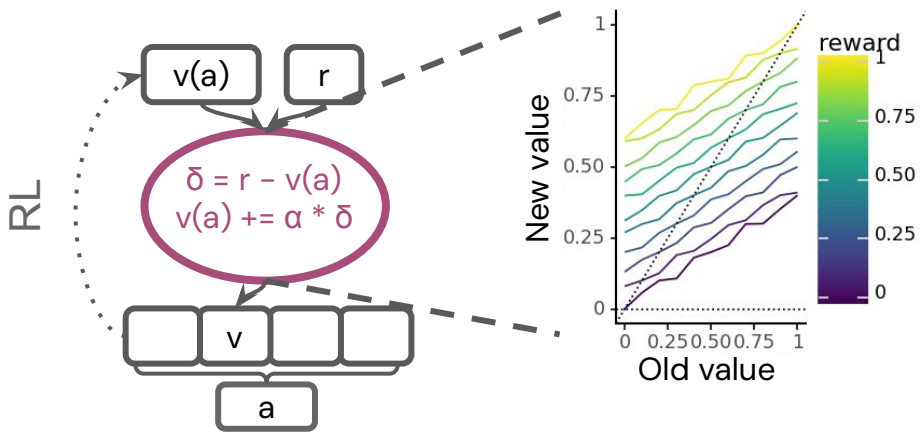- 4–armed bandit
- Arms drift independently

**Original task**

- 14 participants, 150 trials, fMRI [Daw et al., 2006]

**Our version**

- 880 participants
- Several blocks (1 training block, several testing blocks à 150 trials)
- Online (prolific)
- Exclusion: 2% of participants, 0.6% of blocks

Slide credit:
**Maria Eckstein
(mariaeckstein@deepmind.com)**

# A Different Value Update, Learned from Data



RL

$$\delta = r - v(a)$$
$$v(a) += \alpha * \delta$$

RNN-**v**

$$v(a) += \text{[neural network]}$$

**Conclusion**

- RNN–**v** fits better than RL
- Human learning is different from pure RL theory
- But still a big gap in model fit
- Test other assumptions of RL

*Slide credit*:
**Maria Eckstein**
(mariaeckstein@deepmind.com)

[Eckstein, Daw, Summerfield, & Miller, 2023, *CogSci*]

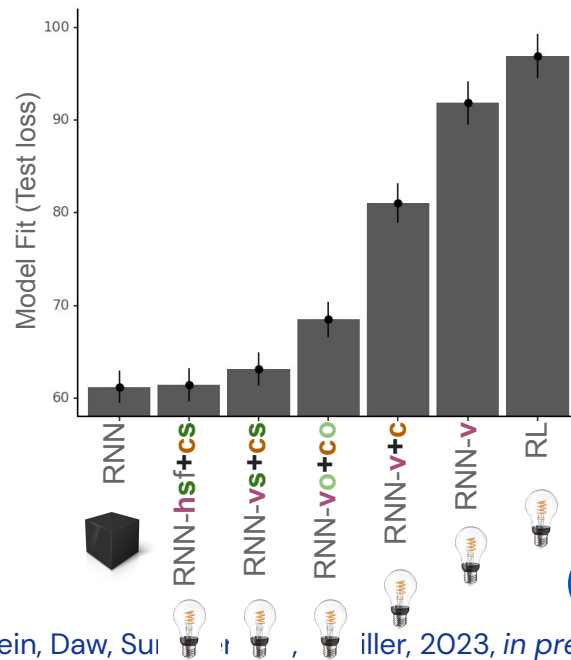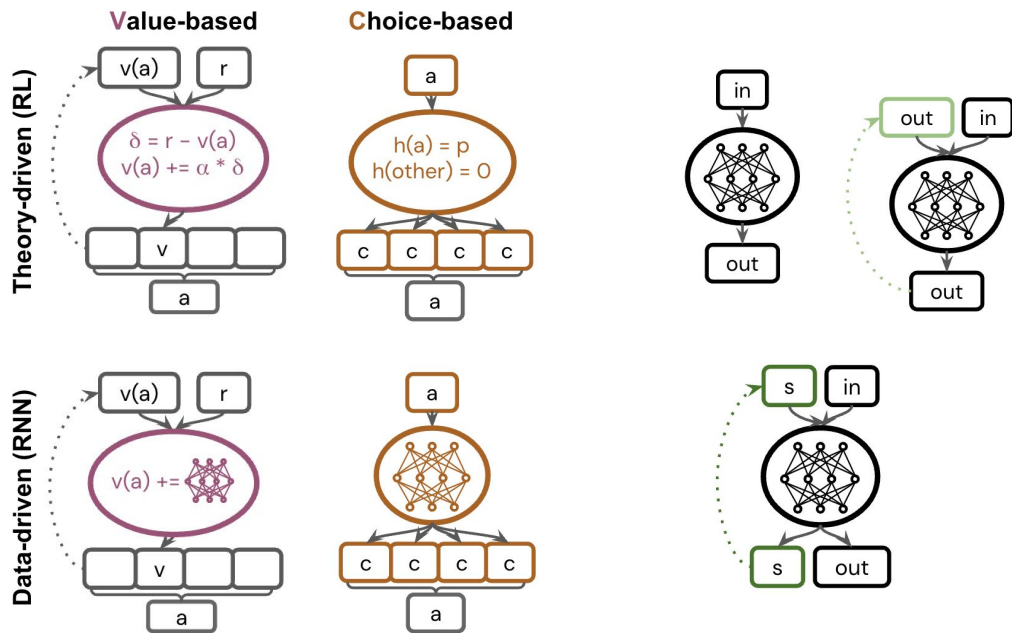# Testing other assumptions of RL

*Slide credit:*
**Maria Eckstein**
(mariaeckstein@deepmind.com)

## Reward–independent processes

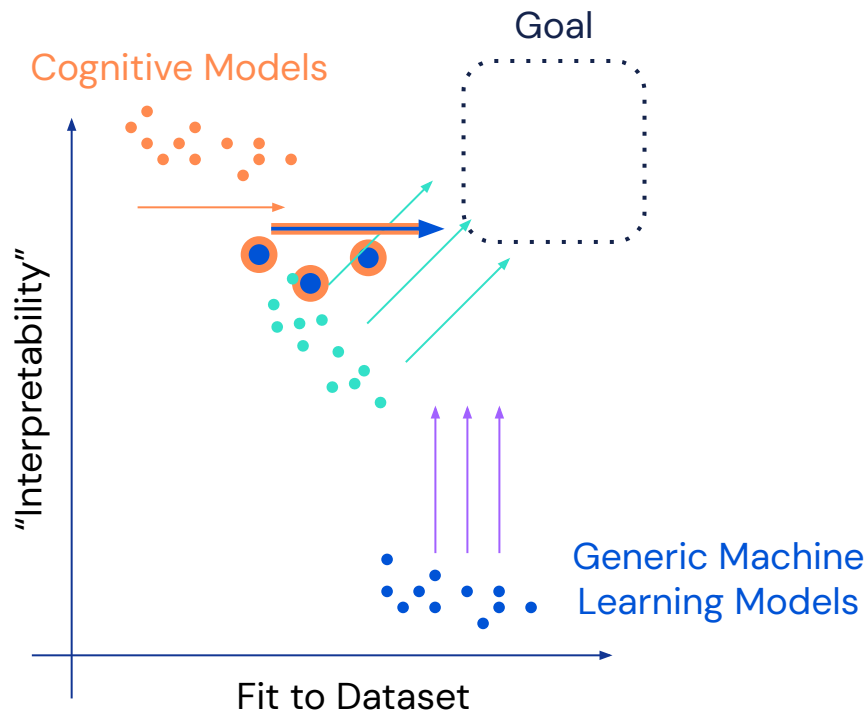[e.g., Gillan et al., 2015; Miller et al., 2019; Sugawara & Katahira, 2021; …]

## Memory / Context

[e.g., Collins & Frank, 2012; Palminteri et al., 2015; Davidow et al., 2016; Gershman & Daw, 2017; Wang et al., 2018; Ramani, 2019; …]



[Eckstein, Daw, Su———  ———iller, 2023, *in prep*]

# Conclusions: A Landscape of Possibilities

- Quantitative models of behavior: A key tool for Comp. Cog. Neuro.

- Classic Cognitive modeling
- ML models as benchmarks
- ML models for post-hoc interpretability
- Interpretability-encouraging architectures
- Hybrid models

- Combining them!
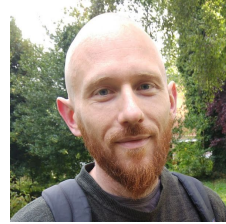
- Your ideas?

# Acknowledgements

Slides:

Collaborators at GDM:



Anne Collins

Kim Stachenfeld

Zeb Kurth–Nelson

Kevin Miller

Nathaniel Daw

Chris Summerfield