# Guidelines for conducting a cognitive modeling study: theory and practice (2/2)

## Valentin Wyart

Lab. de Neurosciences Cognitives et Computationnelles (LNC$^2$)
Institut National de la Santé et de la Recherche Médicale
Ecole Normale Supérieure, Université PSL

valentin.wyart@ens.psl.eu

PSL Data Science Program
https://psl.eu/en/programmes-gradues/programme-data

Paris Artificial Intelligence Research Institute
https://prairie-institute.fr

# Group project

# Ten simple rules for the computational modeling of behavioral data

Robert C Wilson[1,2†*], Anne GE Collins[3,4†*]

[1]Department of Psychology, University of Arizona, Tucson, United States; [2]Cognitive Science Program, University of Arizona, Tucson, United States; [3]Department of Psychology, University of California, Berkeley, Berkeley, United States; [4]Helen Wills Neuroscience Institute, University of California, Berkeley, Berkeley, United States
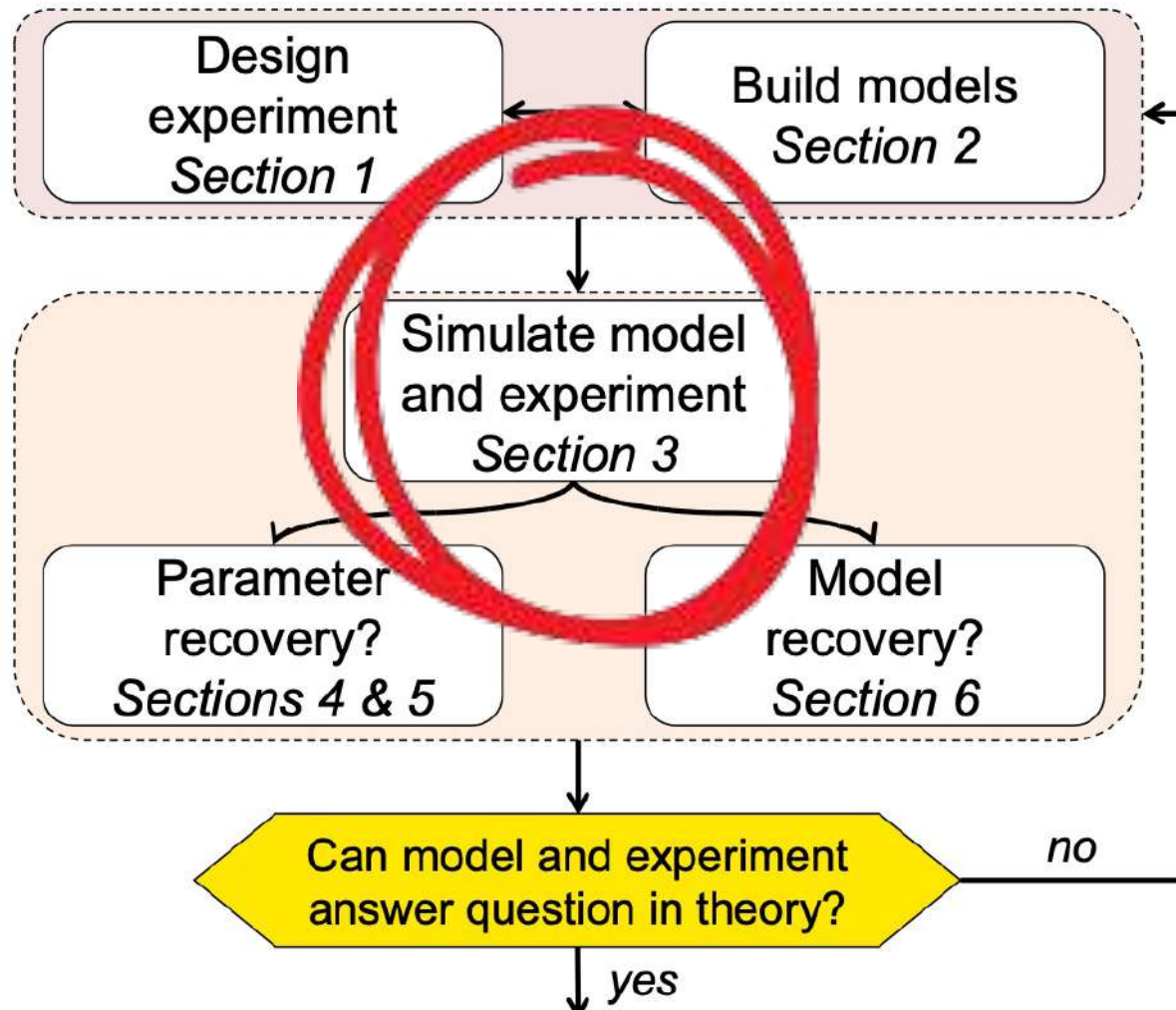
**Abstract** Computational modeling of behavior has revolutionized psychology and neuroscience. By fitting models to experimental data we can probe the algorithms underlying behavior, find neural correlates of computational variables and better understand the effects of drugs, illness and interventions. But with great power comes great responsibility. Here, we offer ten simple rules to ensure that computational modeling is used with care and yields meaningful insights. In particular, we present a beginner-friendly, pragmatic and details-oriented introduction on how to relate models to data. What, exactly, can a model tell us about the mind? To answer this, we apply our rules to the simplest modeling techniques most accessible to beginning modelers and illustrate them with examples and code available online. However, most rules apply to more advanced

# Group project

- Could you open the behavioral dataset?

- Objective: identify the latent cognitive strategy that drives behavior (different for each group)

- Use data mining and modeling approaches:
  - ✓ describe behavior using data mining
  - ✓ identify strategy using data modeling

- Group presentation (15 min/group) on Friday
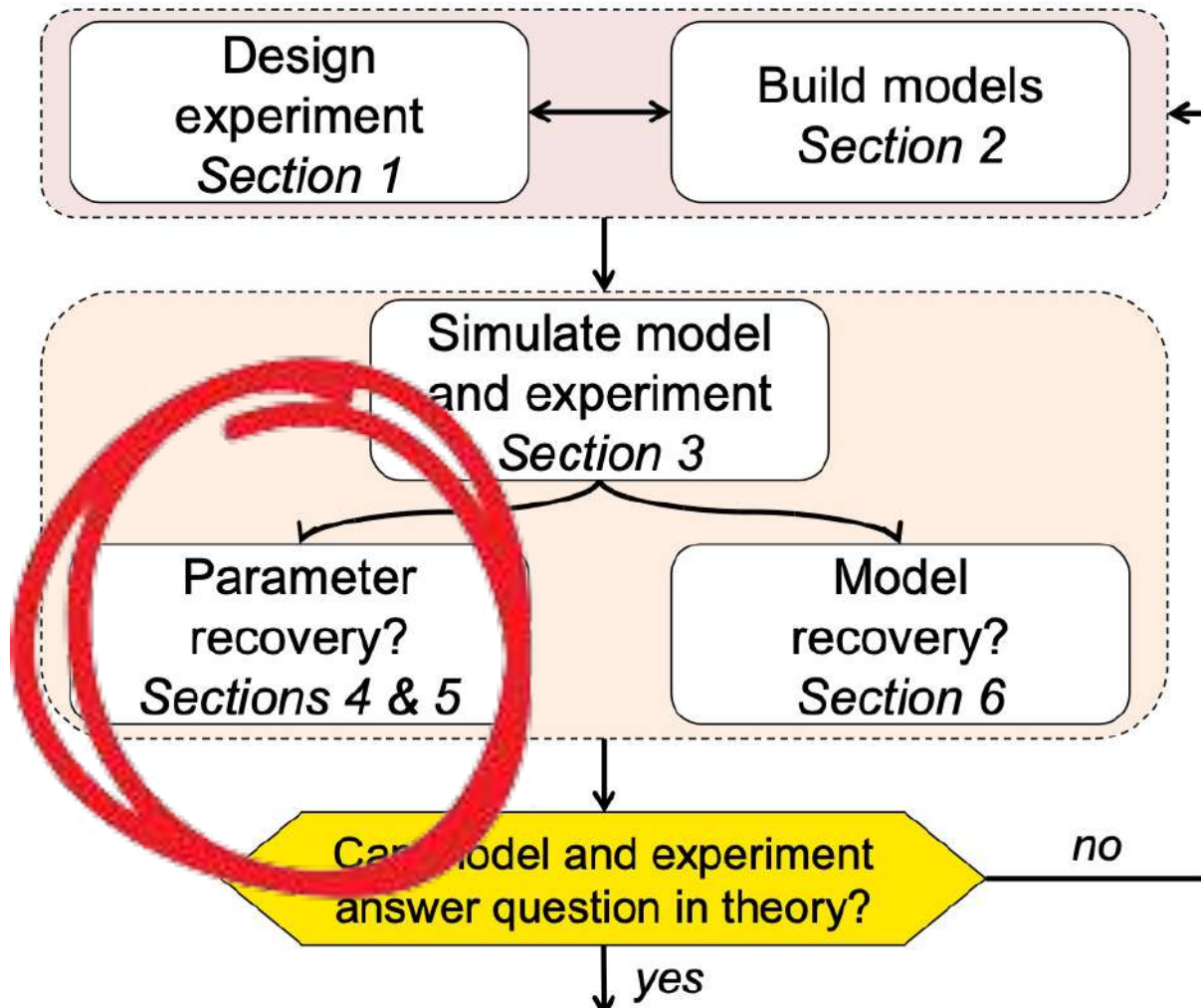
- Don't hesitate to ask for help or advice!

# Modeling guidelines

# Modeling guidelines

- Model simulations are useful to:
  - ✓ check that candidate models make different predictions in the same task
  - ✓ choose task variables (e.g., difficulty)

- What controls difficulty in a <u>stable</u> bandit task?

- Why is it important that <u>all</u> model parameters affect behavioral predictions?

- Why is it important that <u>all</u> candidate models make different behavioral predictions?

# Modeling guidelines

# Modeling guidelines

- Parameter recovery:
  *Before reading too much into fitted parameter values, it is important to check whether the fitting procedure works, by fitting synthetic behavior from a known model whose true parameters are known.*

- Model simulation code is needed:

$$\text{>> behavior} = f(\theta, s)$$

- Model fitting code is needed as well:

$$\text{>> } \hat{\theta}_{\text{MLE}} = \text{argmax}_{\theta} \left( \log\left( p(\text{behavior}|\theta, s) \right) \right)$$

# Modeling guidelines

- Parameter recovery:

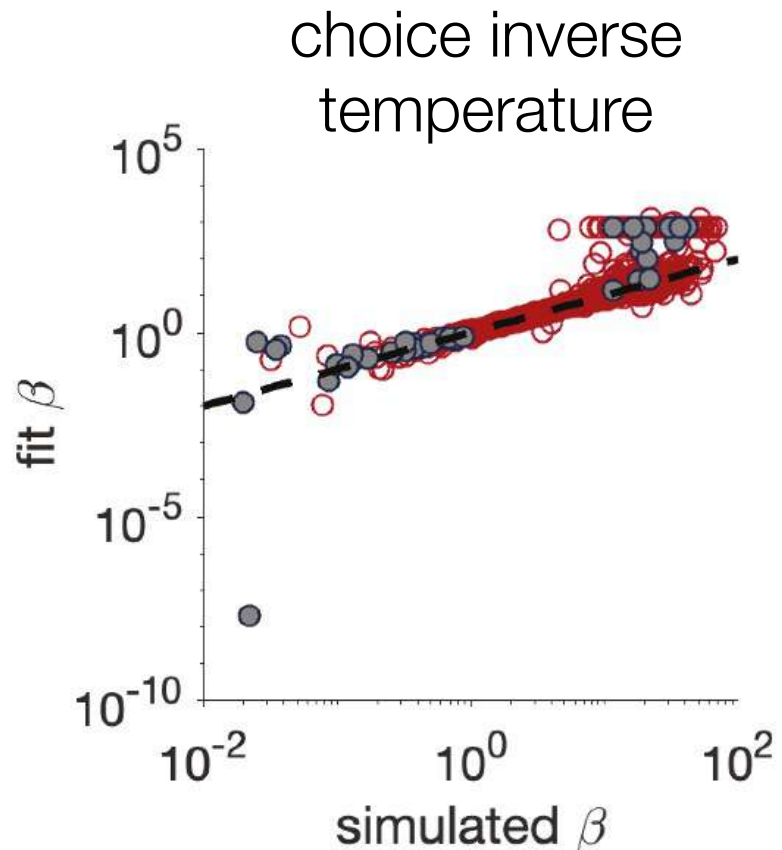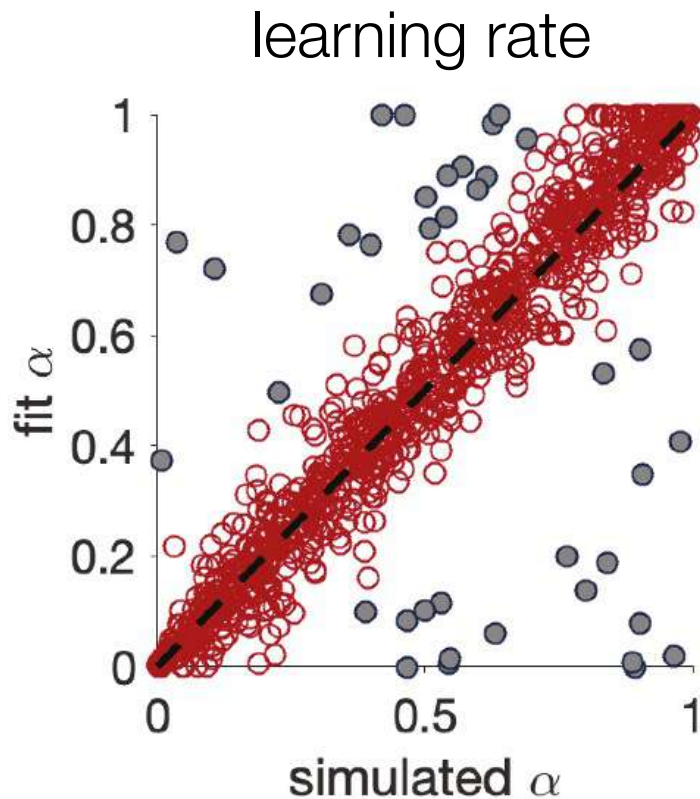## Box 4. Example: parameter recovery in the reinforcement learning model.

We performed parameter recovery with Model 3, the Rescorla Wagner model, on the two-armed bandit task. As before, we set the means of each bandit at $\mu_1 = 0.2$ and $\mu_2 = 0.8$ and the number of trials at $T = 1000$. We then simulated the actions of the model according to *Equations 3 and 4*, with learning rate, $\alpha$, and softmax temperature, $\beta$, set according to

$$\alpha \sim U(0,1) \quad \text{and} \quad \beta \sim \text{Exp}(10) \tag{9}$$

After simulating the model, we fit the parameters using a maximum likelihood approach to get fit values of learning rate, $\alpha$, and softmax parameter, $\beta$. We then repeated this process 1000 times using new values of $\alpha$ and $\beta$ each time.

# Modeling guidelines

- Parameter recovery:



learning rate

choice inverse temperature

# Modeling guidelines

- Parameter recovery:
  output = parameter correlations

- Why is it important that all model parameters affect behavioral predictions?

- Would a model parameter that does not affect behavior in the tested task be recoverable?

- What does parameter confusion mean?

- How can we measure it in practice?

# Modeling guidelines

## Compulsivity is linked to suboptimal choice variability but unaltered reinforcement learning under uncertainty
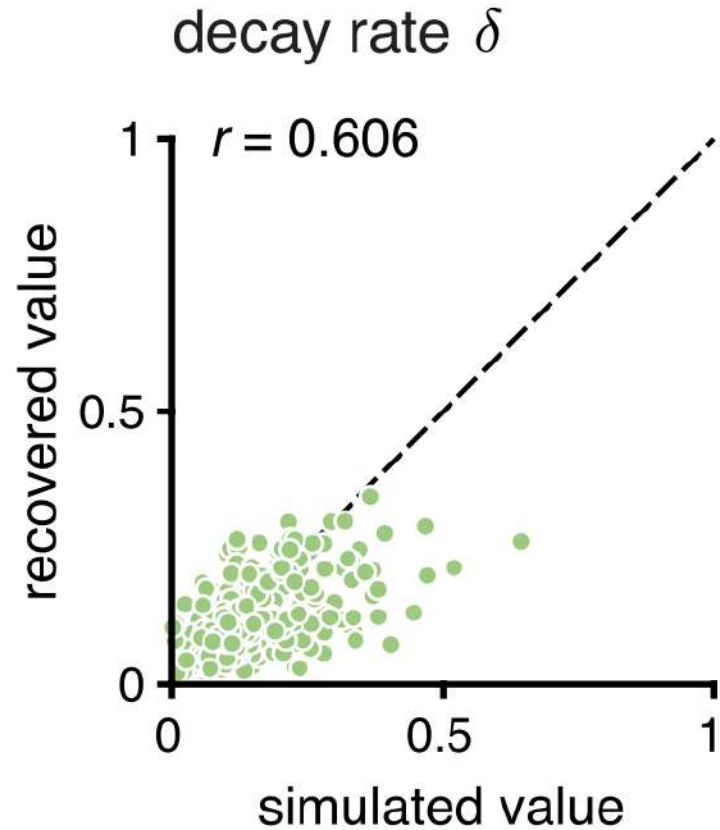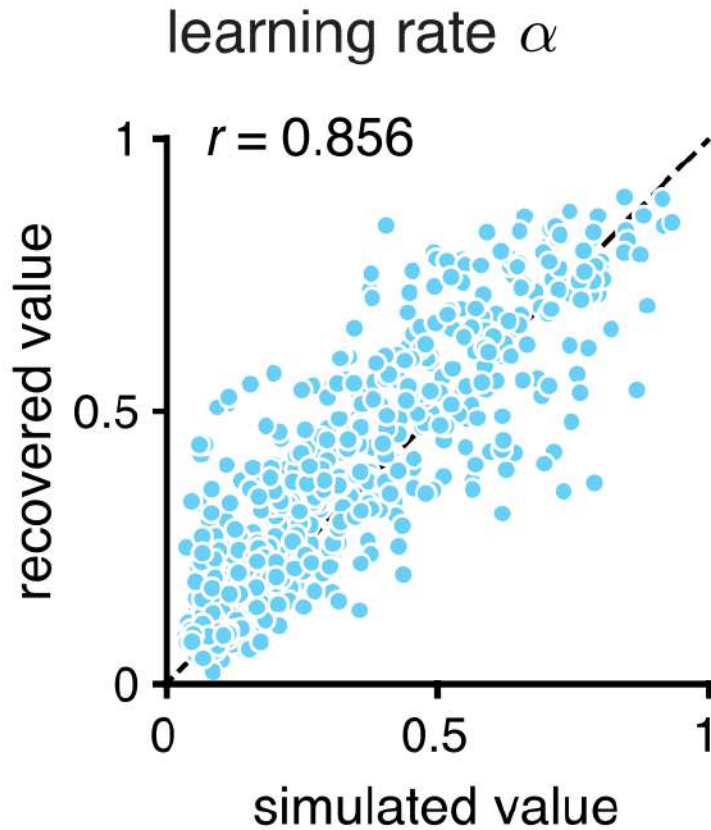
Junseok K. Lee [1,2] ✉, Marion Rouault [1,2,3] & Valentin Wyart [1,2,4] ✉

Compulsivity has been associated with variable behavior under uncertainty. However, previous work has not distinguished between two main sources of behavioral variability: the stochastic selection of choice options that do not maximize expected reward (choice variability) and random noise in the reinforcement learning process that updates option values from choice outcomes (learning variability). Here we study the relation between dimensional compulsivity and behavioral variability using a computational model that dissociates its two sources. Across two independent datasets

# Modeling guidelines

# Modeling guidelines



learning rate $\alpha$

$r = 0.856$

decay rate $\delta$

$r = 0.606$

# Modeling guidelines

# Modeling guidelines

- Model recovery:
  *Before reading too much into model comparison, it is important to check that the comparison procedure works, by comparing models fitted to synthetic behavior whose true model is known.*

- Model simulation code is needed:

$$>> \mathrm{behavior} = f(\theta, s)$$

- Model fitting code is needed as well:

$$>> \mathrm{MLE_M} = \max_\theta \Big( \log\big( p(\mathrm{behavior}|\theta, s) \big) \Big)$$

# Modeling guidelines

- <u>Model recovery:</u>

## Box 5. Example: confusion matrices in the bandit task.

To illustrate model recovery, we simulated the behavior of the five models on the two-armed bandit task. As before, the means were set at $\mu_1 = 0.2$ and $\mu_2 = 0.8$, and the number of trials was set at $T = 1000$. For each simulation, model parameters were sampled randomly for each model. Each simulated data set was then fit to each of the given models to determine which model fit best (according to BIC). This process was repeated 100 times to compute the confusion matrices which are plotted below

# Modeling guidelines

- ## Model recovery:

$p$(fitted | simulated)

fitted model

|  | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | 0.97 | 0.03 | 0 | 0 | 0 |
| 2 | 0.04 | 0.96 | 0 | 0 | 0 |
| 3 | 0.06 | 0 | 0.94 | 0 | 0 |
| 4 | 0.06 | 0 | 0.01 | 0.93 | 0 |
| 5 | 0.03 | 0 | 0.1 | 0.15 | 0.72 |

simulated model

# Modeling guidelines

- <u>Model recovery:</u>
  output = model confusion matrix

- Standard confusion matrix = $p(\text{fitted} \mid \text{simulated})$
  Given behavior from a simulated model, probability of identifying each candidate model as the winning one.

- But what we want is $p(\text{simulated} \mid \text{fitted})$!
  Given a winning model obtained by fitting, probability of each candidate model to have generated behavior.

- Use <u>Bayes rule:</u> $p(\text{simulated} \mid \text{fitted}) \propto$
  $$p(\text{fitted} \mid \text{simulated})\, p(\text{simulated})$$

# Modeling guidelines
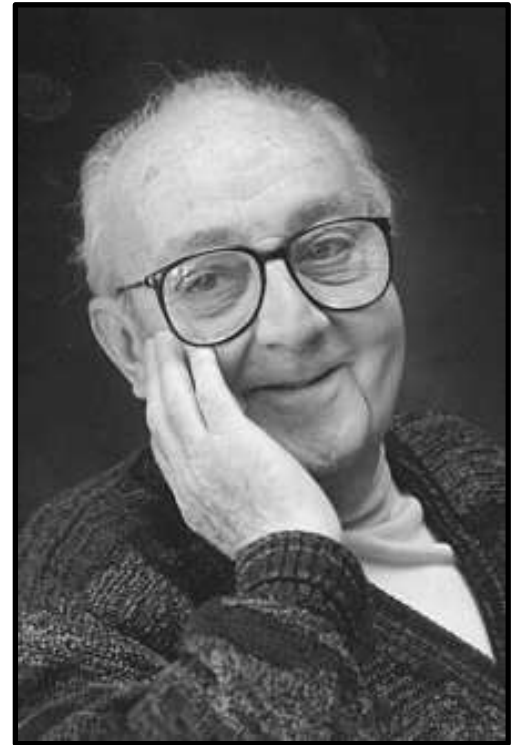
- ## Model recovery:

$p$(fitted | simulated)

$p$(simulated | fitted)

fitted model

fitted model

# Modeling guidelines

- *Essentially, all models are wrong, but some are useful.* (George Box, 1987)

- <u>Scientific worries:</u>
  - ✓ <span style="color:orange">parsimony</span> in theory and model building
  - ✓ wrong but preferably <span style="color:orange">not importantly wrong</span>
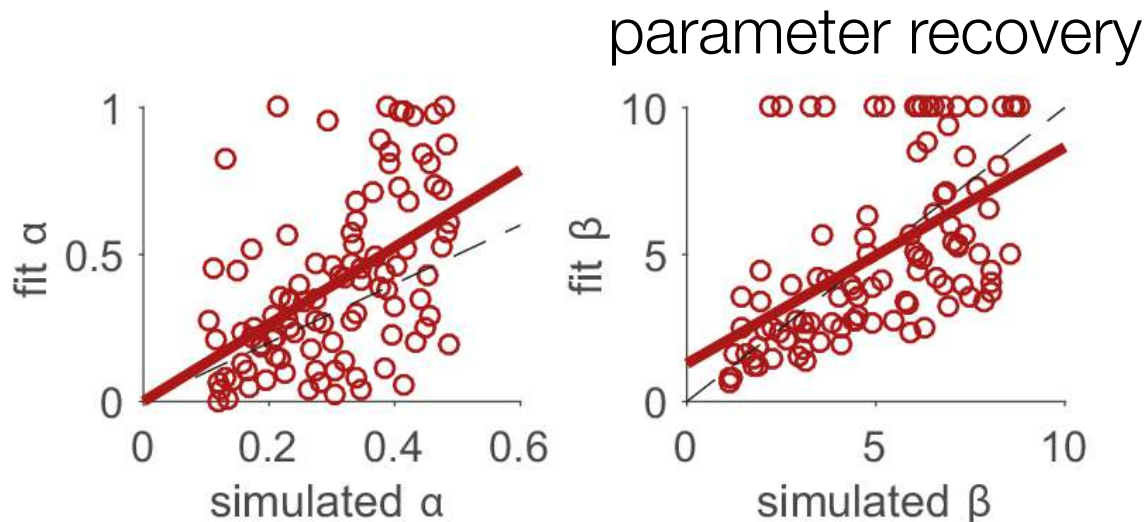
# Modeling guidelines

- *Essentially, all models are wrong, but some are useful.* (George Box, 1987)

- But modeling <u>unimportant</u> model parameters can improve the fitting of important ones!

- Example of choice bias $b$ in <u>TD-based RL</u>:

$$Q_{1,t} = Q_{1,t-1} + \alpha\left(r_t - Q_{1,t-1}\right)$$

$$p_t = 1 \Big/ \left(1 + \exp\left(-\beta\left(Q_{1,t} - Q_{2,t} + b\right)\right)\right)$$

# Modeling guidelines

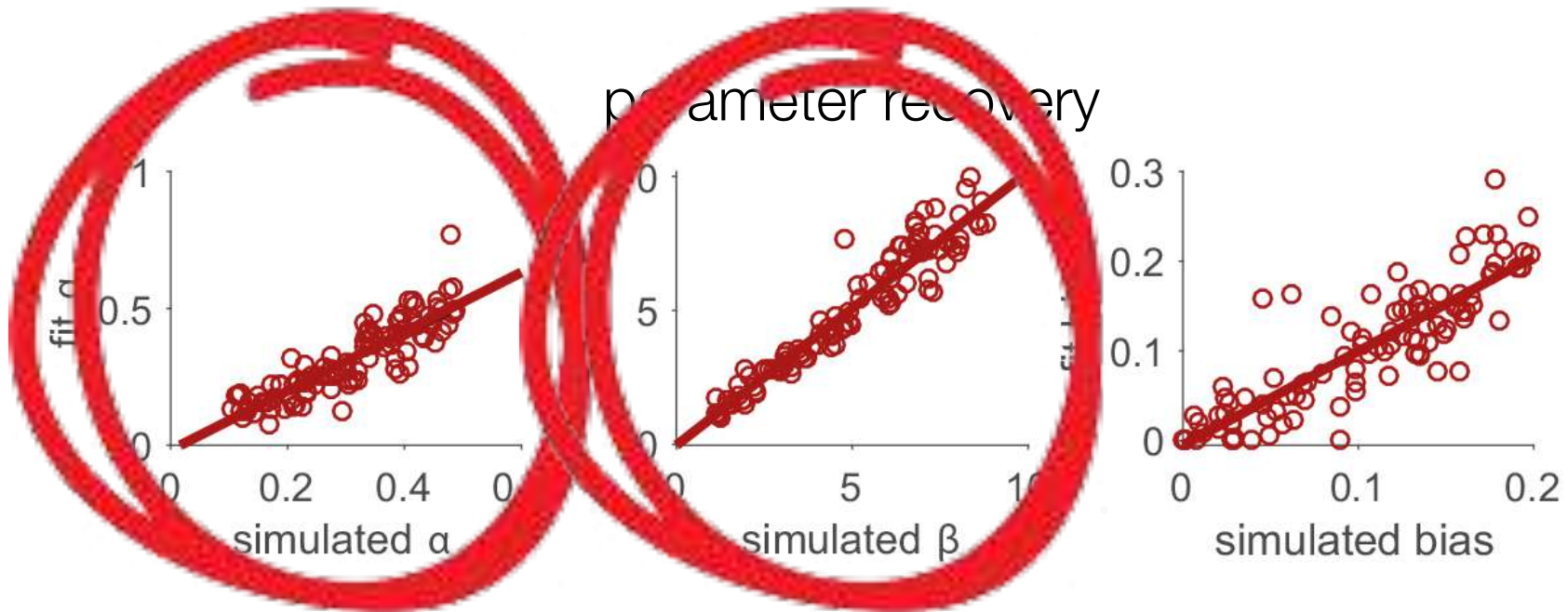- *Essentially, all models are wrong, but some are useful.* (George Box, 1987)

- <u>Simulated model:</u> M3 <span style="color:orange">with choice bias</span>
  <u>Fitted model:</u> M3 <span style="color:orange"><u>without</u> choice bias</span>

parameter recovery



<u>without</u>
choice bias in
fitted model

# Modeling guidelines

- *Essentially, all models are wrong, but some are useful.* (George Box, 1987)

- Simulated model: M3 with choice bias
  Fitted model: M3 with choice bias



parameter recovery

# Modeling guidelines

- *Essentially, all models are wrong, but some are useful.* (George Box, 1987)

- <u>Simulated model:</u> M3 with choice bias
  <u>Fitted model:</u> M3 <u>with</u> choice bias

- What differences between the results of the parameter recovery procedure? Why?

- Do these results conflict with the two worries identified by George Box? Why?

# Modeling guidelines

# When a good fit can be bad

## Mark A. Pitt and In Jae Myung

How should we select among computational models of cognition? Although it is commonplace to measure how well each model fits the data, this is insufficient. Good fits can be misleading because they can result from properties of the model that have nothing to do with it being a close approximation to the cognitive process of interest (e.g. overfitting). Selection methods are introduced that factor in these properties when measuring fit. Their success in outperforming standard goodness-of-fit measures stems from a focus on measuring the generalizability of a model's data-fitting abilities, which should be the goal of model selection.

The explosion of interest in modeling cognitive processes over the past 20 years has fueled the cognitive sciences in many ways. Not only has it opened up new ways of thinking about research problems and possible solutions, but it has also enabled researchers to gain a better understanding of their theories by simulating a computational instantiation of it. Modeling is now sufficiently mainstream that one can get the impression that the

of it. A thorough evaluation of a model requires methods that are sensitive to its quantitative form. Criteria used for evaluating theories [1], such as testing their performance in an experimental setting, do not speak to the quality of the choices that are made in building their quantitative counterparts (i.e. choice of parameters, how they are combined) or their ramifications. The paucity of such model selection methods is surprising given the centrality of the problem itself. What could be more fundamental than deciding between two alternative explanations of a cognitive process?

### How *not* to compare models

Mathematical model are frequently tested against one another by evaluating how well each fits the data generated in an experiment or simulation. Such a test makes sense given that one criterion of model performance is that it reproduce the data. A goodness-of-fit measure (GOF; see Glossary) is invariably used to measure their adequacy in achieving this goal. What is measured is how much a model's predictions deviate from the observed data [2,3]. The model that provides the best fit (i.e. smallest deviation) is favored. The logic of this choice rests on the assumption that the model that provides the best fit to all data must be a closer approximation to the cognitive process under investigation than its competitors [4].

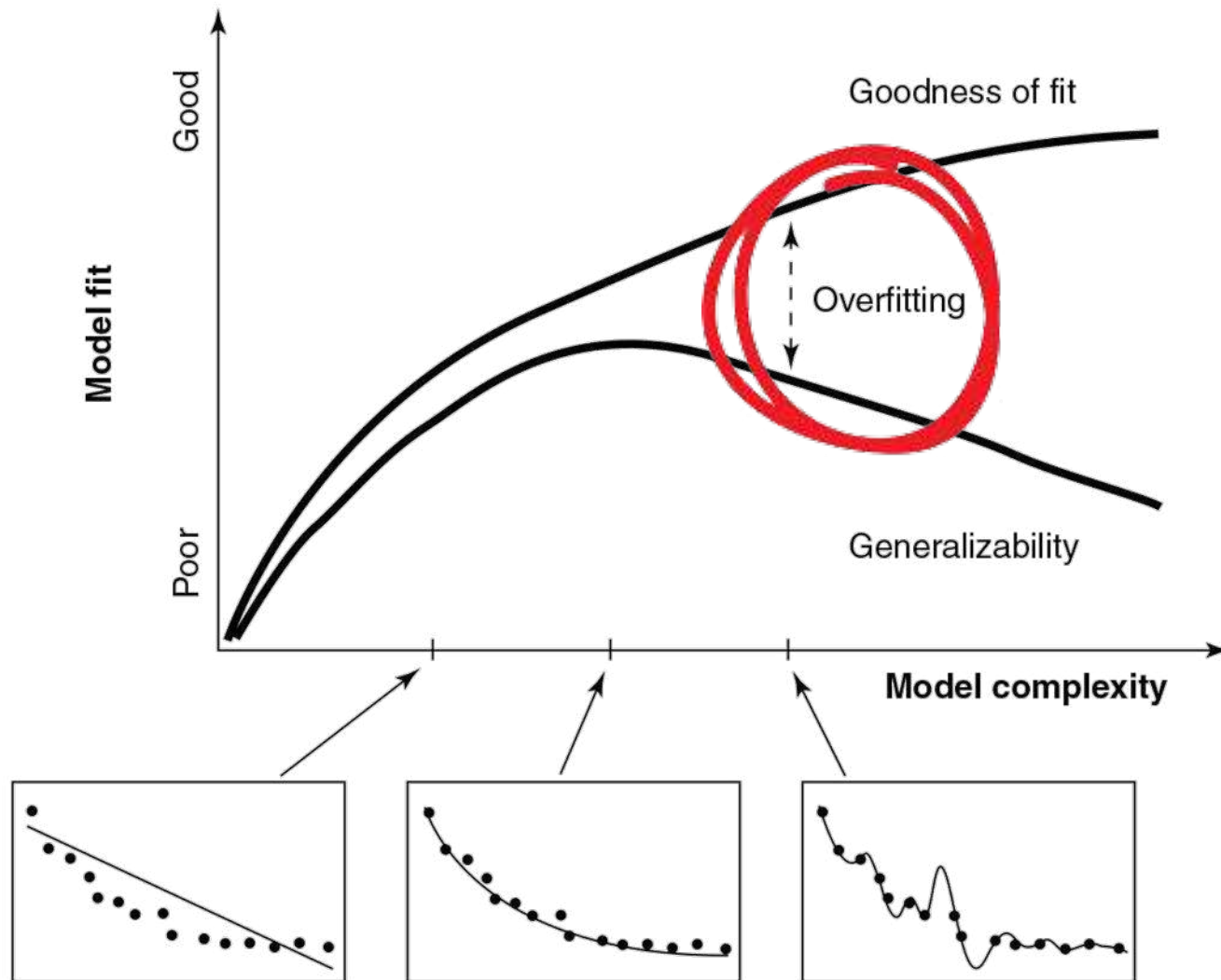Such a conclusion is reasonable if measurements

# Modeling guidelines

- Overfitting issue and Occam's razor

- <u>Law of parsimony:</u> "The simplest explanation is usually the best one."

- Why is this principle important for modeling data?

William of Ockham
(1287–1347)
medieval philosopher

# Modeling guidelines

# Modeling guidelines

- How to deal in practice with overfitting?

- Idea: use a complexity-penalizing metric of fit

- Which of these metrics penalize complexity?
  - ✓ RMSE (root mean squared error)
  - ✓ PVAF (percent variance accounted for)
  - ✓ AIC (Akaike information criterion)
  - ✓ BIC (Bayesian information criterion)

# Modeling guidelines

- How to deal in practice with overfitting?

- Idea: use a complexity-penalizing metric of fit

- Which of these metrics penalize complexity?

**Table II.** Two GOF Measures, four generalizability measures, and the dimensions of complexity to which each is sensitive

| Selection method | Criterion equation | Dimensions of complexity considered |
|---|---|---|
| Root Mean Squared Error | $RMSE = (SSE/N)^{1/2}$ | None |
| Percent Variance Accounted For | $PVAF = 100(1 - SSE/SST)$ | None |
| Akaike Information Criterion | $AIC = -2\ ln(f(y|\theta_0)) + 2k$ | Number of parameters |
| Bayesian Information Criterion | $BIC = -2\ ln(f(y|\theta_0)) + k \cdot ln(n)$ | Number of parameters, sample size |
| Bayesian Model Selection | $BMS = -ln \int f(y|\theta)\pi(\theta)d\theta$ | Number of parameters, sample size, functional form |
| Minimum Description Length | $MDL = -ln(f(y|\theta_0)) + (k/2)ln(n/2\pi) + ln \int \sqrt{det(I(\theta))}d\theta$ | Number of parameters, sample size, functional form |

In the equations above, $y$ denotes observed data, $\theta$ is the model's parameter, $\theta_0$ is the parameter value that maximizes the likelihood function $f(y|\theta)$, $k$ is the number of parameters, n is the sample size, N is the number of data points fitted, SSE is the minimized sum of the squared errors between observations and predictions, SST is the sum of the squares total, $\pi(\theta)$ is the parameter prior density, $I(\theta)$ is the Fisher information matrix in mathematical statistics [a], *det* denotes the determinant of a matrix, and *ln* denotes the natural logarithm of base e.
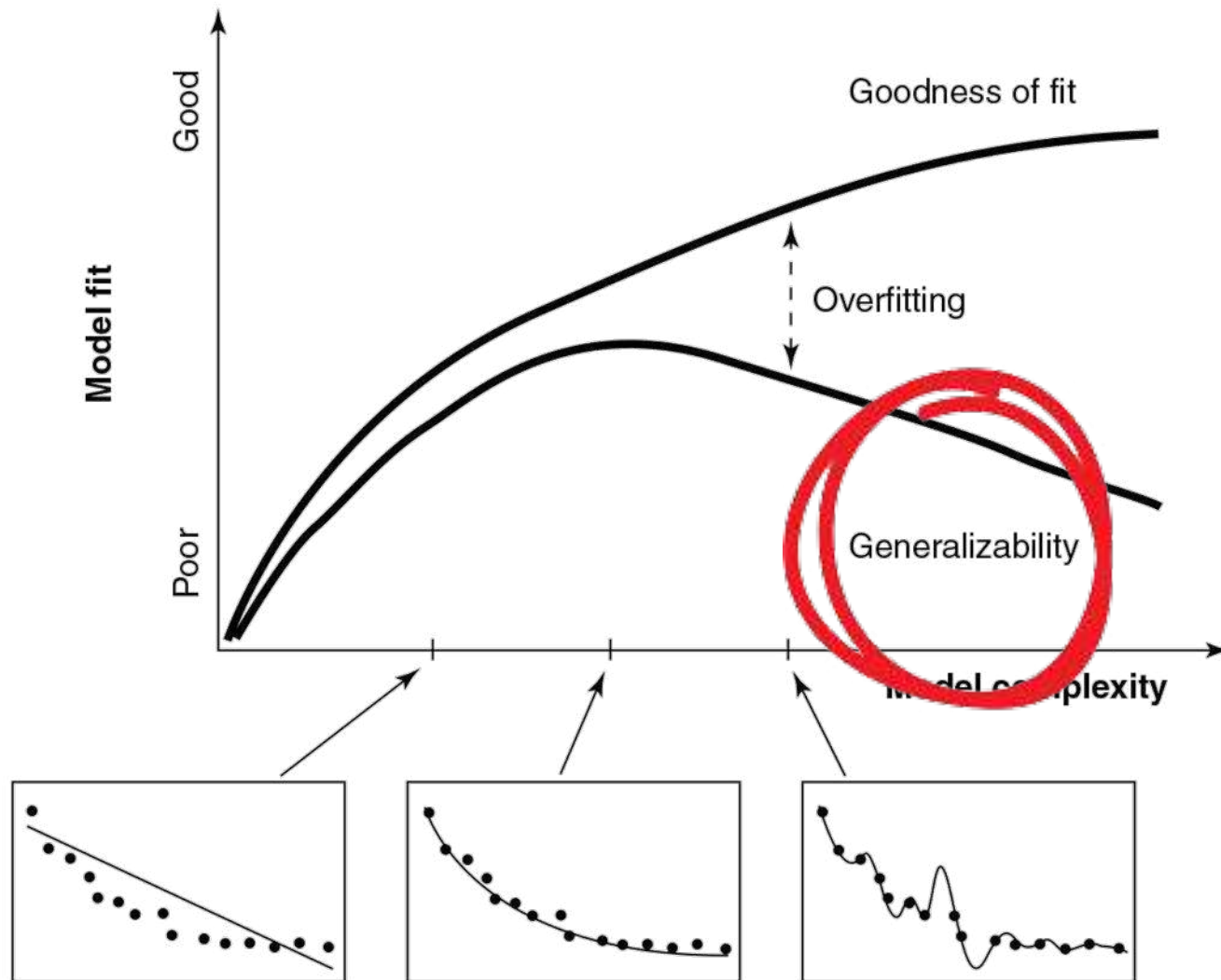
# Modeling guidelines

- How to deal in practice with overfitting?

- Example:
  - ✓ $M_A : y = (1 + x)^{-a}$
  - ✓ $M_B : y = (b + c \cdot x)^{-a}$

Table I. Results of a model recovery simulation in which a GOF measure (RMSE) was used to discriminate models when the source of the error was varied.

| Condition (sources of variation) | Model the data were generated from | | | Model fitted | |
|---|---|---|---|---|---|
| | $M_A$ a = 0.4 | $M_A$ a = 0.6 | $M_B$ | $M_A$ | $M_B$ |
| (1) Sampling error | 100 | – | – | 0.040 (0%) | 0.029 (100%) |
| (2) Sampling error + individual differences | 50 | 50 | – | 0.041 (0%) | 0.029 (100%) |
| (3) Different models | – | 50 | 50 | 0.075 (0%) | 0.029 (100%) |
| (4) Sampling error | – | – | 100 | 0.079 (0%) | 0.029 (100%) |

# Modeling guidelines

- How to deal in practice with overfitting?

- <u>Other example:</u>
  - ✓ $M_1 : y = (1 + x)^{-a}$
  - ✓ $M_2 : y = (b + x)^{-a}$
  - ✓ $M_3 : y = (1 + c \cdot x)^{-a}$

# Modeling guidelines

- How to deal in practice with overfitting?

| Selection method | Model fitted | Model the data were generated from | | |
|---|---|---|---|---|
| | | M₁ | M₂ | M₃ |
| PVAF | M₁ | 0 | 0 | 0 |
| | M₂ | 38 | 97 | 30 |
| | M₃ | 62 | 3 | 70 |
| AIC | M₁ | 79 | 0 | 0 |
| | M₂ | 9 | 97 | 30 |
| | M₃ | 12 | 3 | 70 |
| MDL | M₁ | 86 | 0 | 0 |
| | M₂ | 1 | 92 | 8 |
| | M₃ | 13 | 8 | 92 |

# Modeling guidelines

# Modeling guidelines

- How to deal in practice with overfitting?

- Other idea: use a cross-validation approach

- General procedure:
  - ✓ Fit model on training set
  - ✓ Compute metric of fit on separate test set

- Why does it overcome overfitting?

- Why is it less arbitrary than using a complexity-penalized metric of fit?

# Paper review

# Efficient stabilization of imprecise statistical inference through conditional belief updating

Julie Drevet [1,2] ✉, Jan Drugowitsch [3] and Valentin Wyart [1,2] ✉

Statistical inference is the optimal process for forming and maintaining accurate beliefs about uncertain environments. However, human inference comes with costs due to its associated biases and limited precision. Indeed, biased or imprecise inference can trigger variable beliefs and unwarranted changes in behaviour. Here, by studying decisions in a sequential categorization task based on noisy visual stimuli, we obtained converging evidence that humans reduce the variability of their beliefs by updating them only when the reliability of incoming sensory information is judged as sufficiently strong. Instead of integrating the evidence provided by all stimuli, participants actively discarded as much as a third of stimuli. This conditional belief updating strategy shows good test–retest reliability, correlates with perceptual confidence and explains human behaviour better than previously described strategies. This seemingly suboptimal strategy not only reduces the costs of imprecise computations but also, counterintuitively, increases the accuracy of resulting decisions.

Efficient decision-making about the cause of noisy or ambiguous observations requires the accumulation of multiple pieces of evidence to form accurate beliefs[1,2], a process typically referred to as 'statistical inference'. In stable environments, accu-

bag) were perceived as dark and vice versa (Fig. 1c and Methods). After each marble, participants were asked to identify the bag from which it was drawn (Fig. 1d). Importantly, marbles were not drawn randomly and independently across successive trials, but rather in episodes of multiple draws from the same bag. Decision-making in

# Paper review

- Let's look at this paper and check whether the authors have followed <u>all of the guidelines</u> for modeling behavior…

**Julie Drevet**
Aix-Marseille Université

**Jan Drugowitsch**
Harvard Medical School

# Paper review

- <u>Task:</u> identify the bag (hidden state) from which marbles (observations) are drawn from
  >> hidden-state inference process

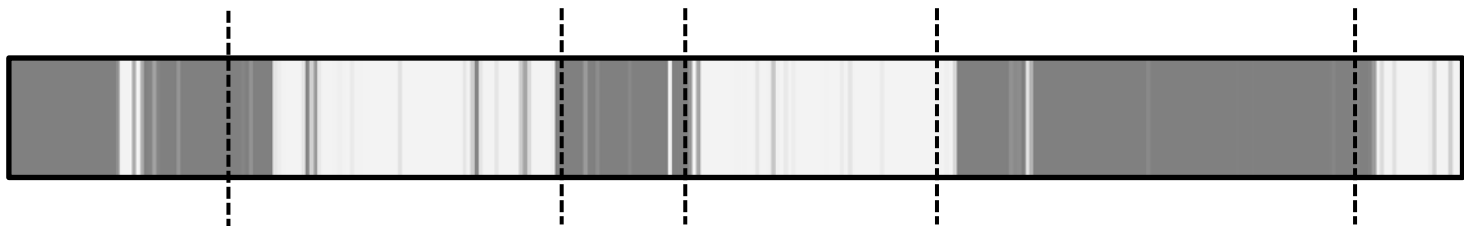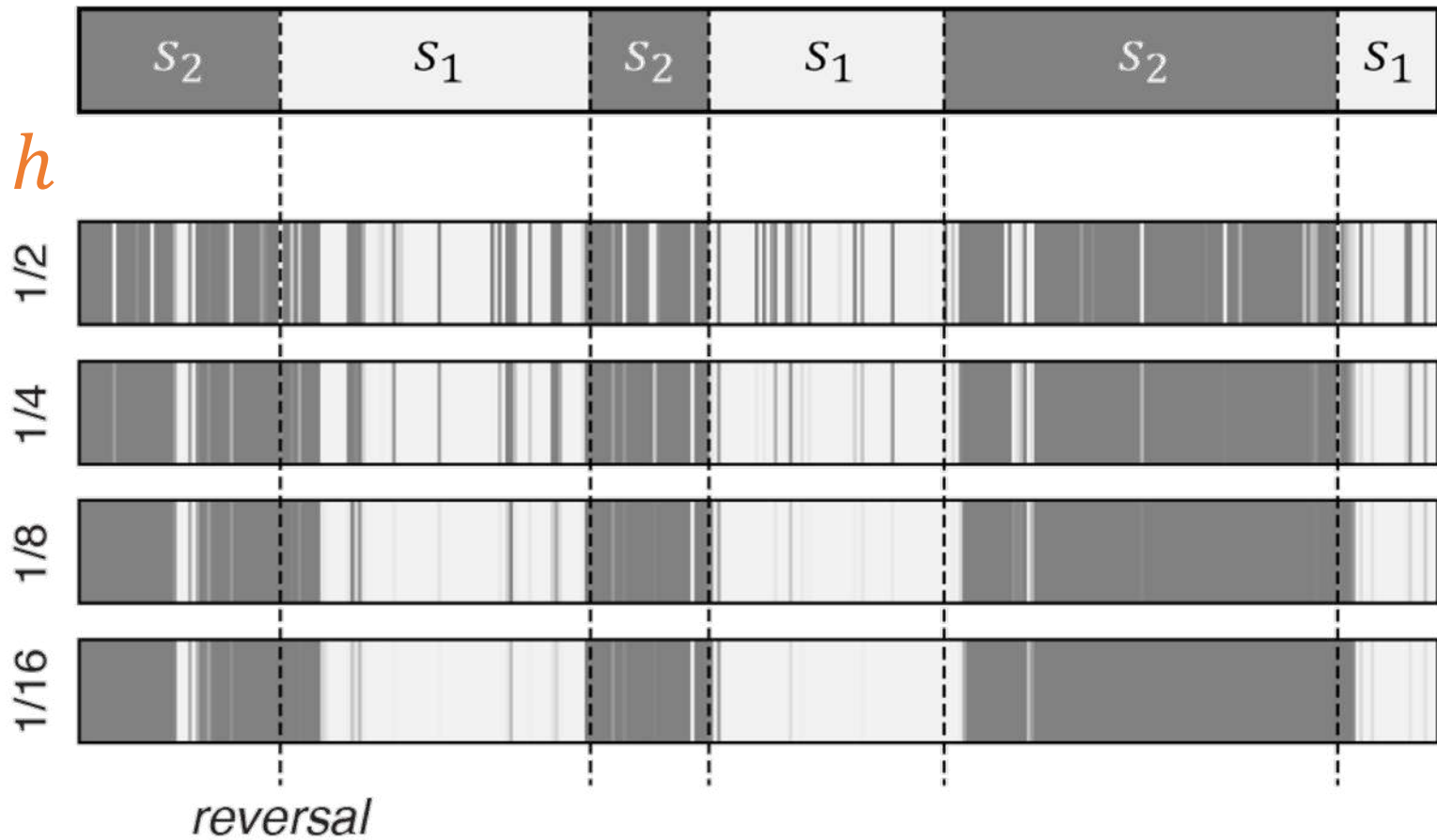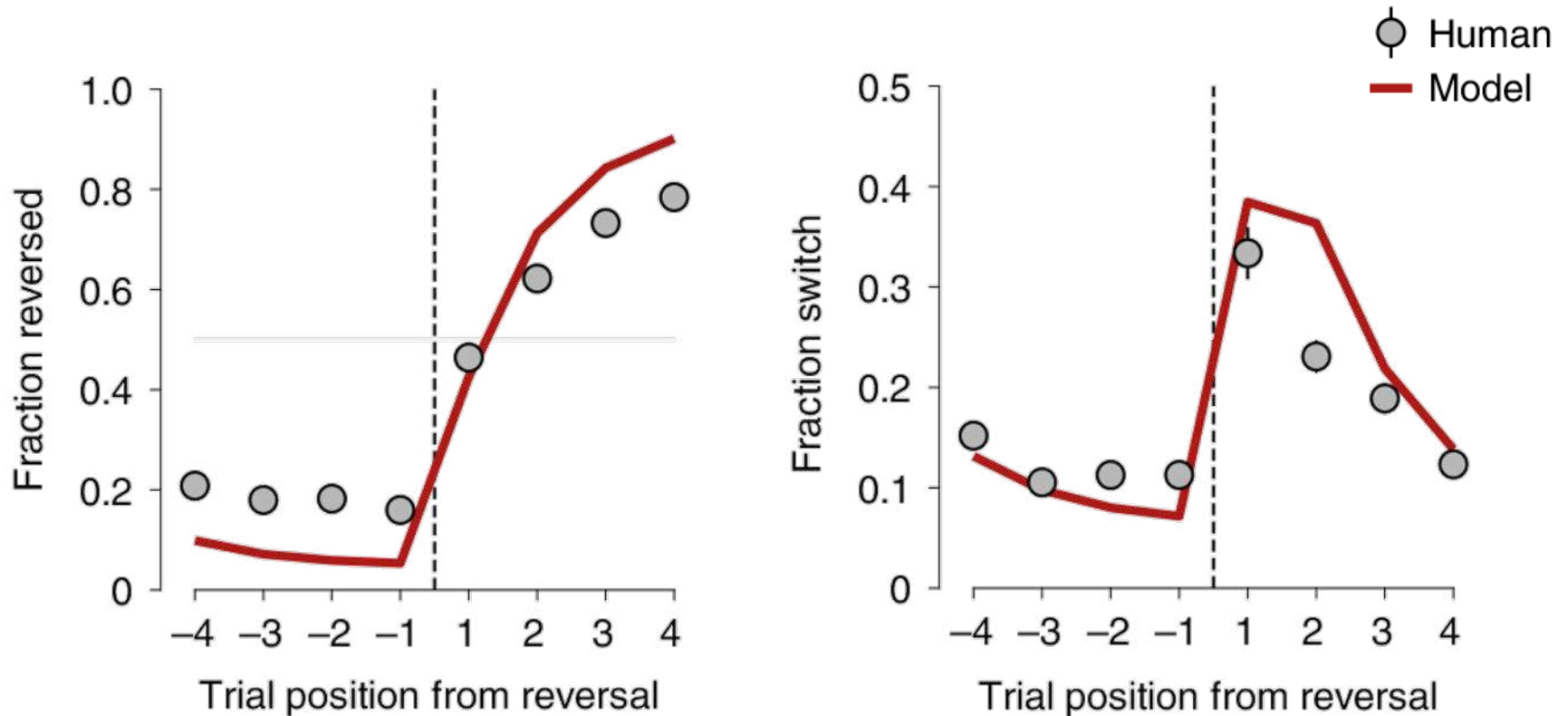bag A

bag B

dark marbles

light marbles

# Paper review

- <u>Task:</u> identify the bag (hidden state) from which marbles (observations) are drawn from
  >> hidden-state inference process



titration procedure
<u>target:</u> 20% errors
(misperceived marbles)

hidden-state $s_t$

$s_2$ $s_1$ $s_2$ $s_1$ $s_2$ $s_1$

sensory evidence $\ell_t$

posterior belief $\mathcal{L}_t$

# Paper review

- <u>Task:</u> identify the bag (hidden state) from which marbles (observations) are drawn from
  >> hidden-state inference process

- Sequential process based on Bayes rule:
  $$>> \mathcal{L}_t = \mathcal{F}(\mathcal{L}_{t-1}, h) + \ell_t$$
  where $h$ = perceived hazard rate (rate of hidden-state reversals)

belief $\mathcal{L}_t = \mathcal{F}(\mathcal{L}_{t-1}, h) + \ell_t$

# Paper review

- Comparison between human behavior and optimal inference…

# Paper review

- Comparison between human behavior and optimal inference… and noisy inference

# Paper review

- Definition of candidate model parameters that could explain human behavior

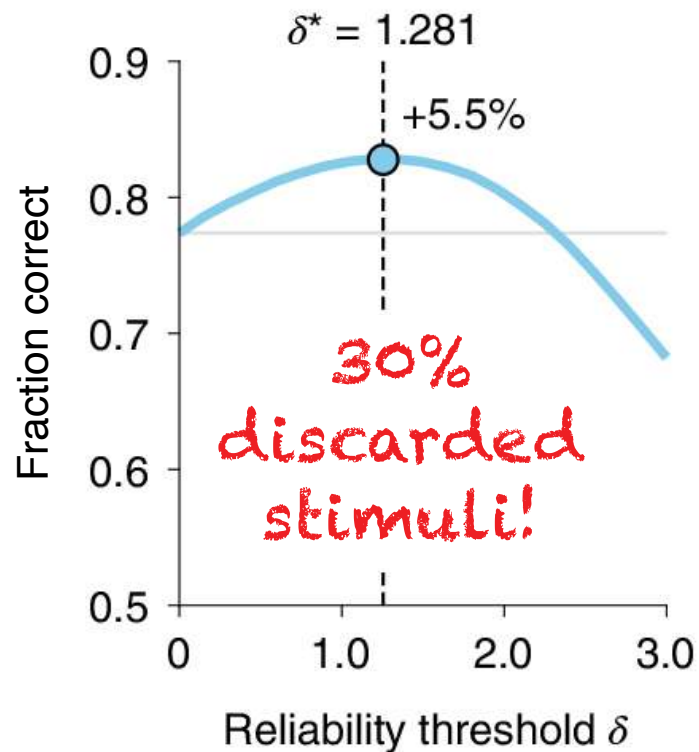hidden-state inference model

# Paper review

- <u>Simulation</u> of candidate model parameters that could explain human behavior

# Paper review

- <u>Simulation</u> of candidate model parameters that could explain human behavior

# Paper review

- <u>Fitting</u> of candidate model parameters that could explain human behavior

# Paper review

- Estimation of candidate model parameters that could explain human behavior

# Paper review

- <u>Estimation</u> of candidate model parameters that could explain human behavior

- <u>Winning model:</u> $M_3$ = conditional inference "ignore marbles whose sensory evidence is less than a reliability threshold $\delta$"

- <u>Prediction:</u> reliability threshold $\delta$ should correlate with how confident participants are at judging marbles as light or dark.

# marble perception task with confidence report

model simulations

human data

$\delta^* = 1.281$

$\delta^* = 1.04$

+5.5%

30% discarded stimuli!

Fraction correct

Reliability threshold $\delta$

Reliability threshold $\delta$

$b_5$   0
      −0.02
      −0.04

# Paper to read

**COGNITIVE SCIENCE**

# Using large-scale experiments and machine learning to discover theories of human decision-making

Joshua C. Peterson[1]\*, David D. Bourgin[1]†, Mayank Agrawal[2,3], Daniel Reichman[4], Thomas L. Griffiths[1,2]

Predicting and understanding how people make decisions has been a long-standing goal in many fields, with quantitative models of human decision-making informing research in both the social sciences and engineering. We show how progress toward this goal can be accelerated by using large datasets to power machine-learning algorithms that are constrained to produce interpretable psychological theories. Conducting the largest experiment on risky choice to date and analyzing the results using gradient-based optimization of differentiable decision theories implemented through artificial neural networks, we were able to recapitulate historical discoveries, establish that there is room to improve on existing theories, and discover a new, more accurate model of human decision-making in a form that preserves the insights from centuries of research.

U nderstanding how people make decisions is a central problem in psychology and economics (*1–3*). Having quantitative models that can predict these decisions has become increasingly important as automated systems interact more closely with people (*4, 5*). The search for such models goes back almost 300 years (*6*) but intensified in the latter half of the 20th century (*7, 8*) as em-pirical findings revealed the limitations of the

narios in which decision-makers face a choice between two gambles, each of which has a set of outcomes that differ in their payoffs and probabilities (Fig. 1A). Researchers studying risky choice seek a theory, which we formal-ize as a function that maps from a pair of gambles, *A* and *B*, to the probability $P(A)$ that a decision-maker chooses gamble *A* over gamble *B*, that is consistent with human decisions for as many choice problems as possible. Dis-covering the best theory is a formidable chal-

This dataset includes >30 times the number of problems in the largest previous dataset (*27*) (Fig. 1B). We then used this dataset to evaluate differentiable decision theories that exploit the flexibility of deep neural networks but use psychologically meaningful constraints to pick out a smooth, searchable landscape of candidate theories with shared assumptions. Differentiable decision theories allow the intu-itions of theorists to be combined with gradient-based optimization methods from machine learning to broadly search the space of theories in a way that yields interpretable scientific explanations.

More formally, we define a hierarchy over decision theories (Fig. 1C) reflecting the addi-tion of an increasing number of constraints on the space of functions. These constraints express psychologically meaningful theoret-ical commitments. For example, one class of theories contains all functions in which the value that people assign to one gamble can be influenced by the contents of the other gamble. If theories in this class are more predictive than those that belong to the simpler classes contained within it (e.g., where the value of gambles are independent), then we know that these simpler theories should be elimi-nated. We enforce each constraint by modify-ing the architecture of artificial neural networks, resulting in differentiable decision theories. This

# Paper to read

- What is the <u>new idea</u>?

- Using <u>interpretable</u> neural network models to fit large-scale behavioral datasets and understand human decisions

- What is the <u>main result</u>?

- The best neural network model, with a context-dependent mixture of theories, predicts human decisions better than the best existing model

# Paper to read

# Paper to read

# Paper to read

# Paper to read

# Paper to read

# Paper to read



context-dependent: decisions affected by previous gambles (context)

# Paper to read

# Paper to read



near-linear utility function (UF)
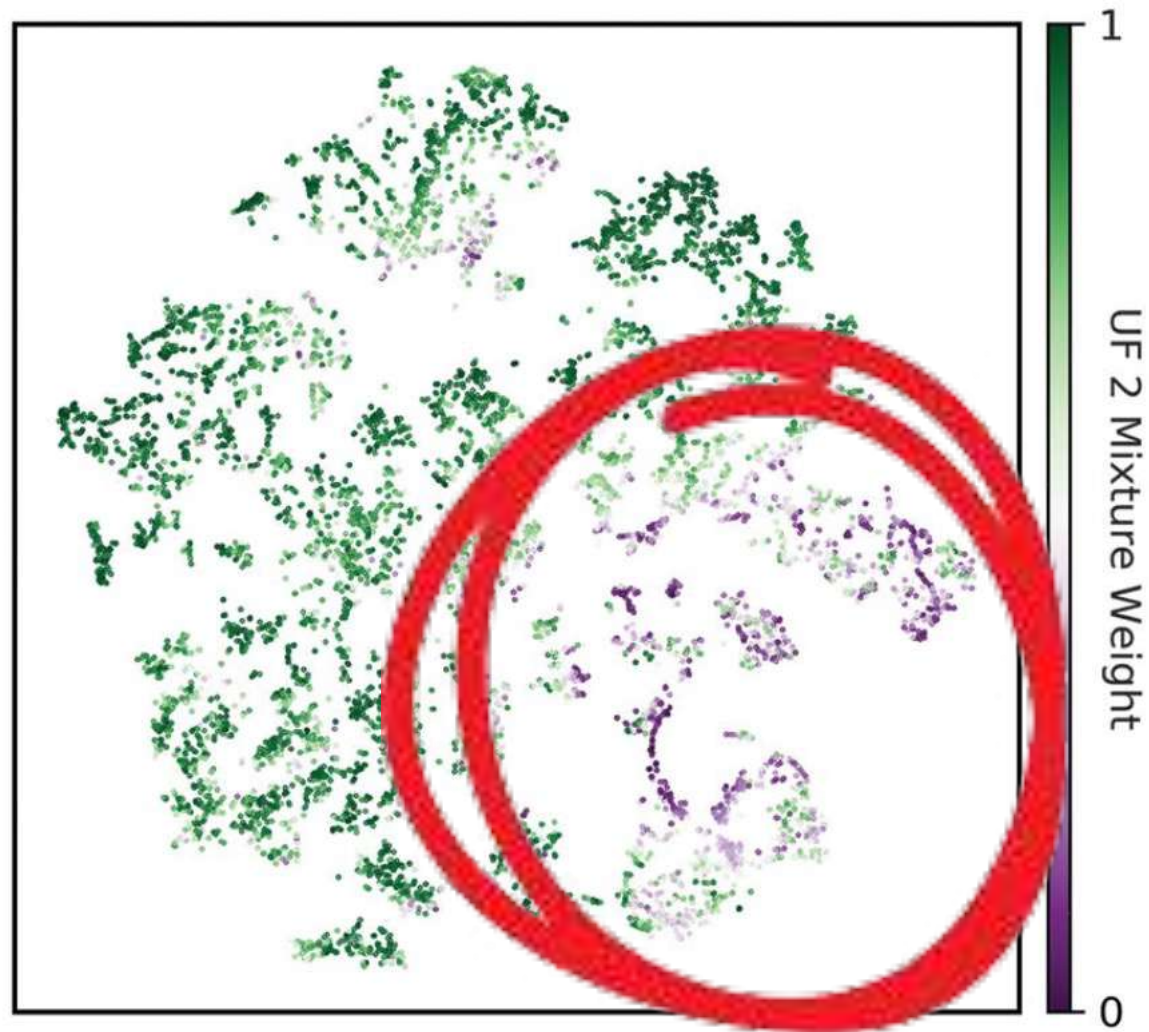linear probability weighting function (PWF)

# Paper to read

loss-averse utility function (UF)
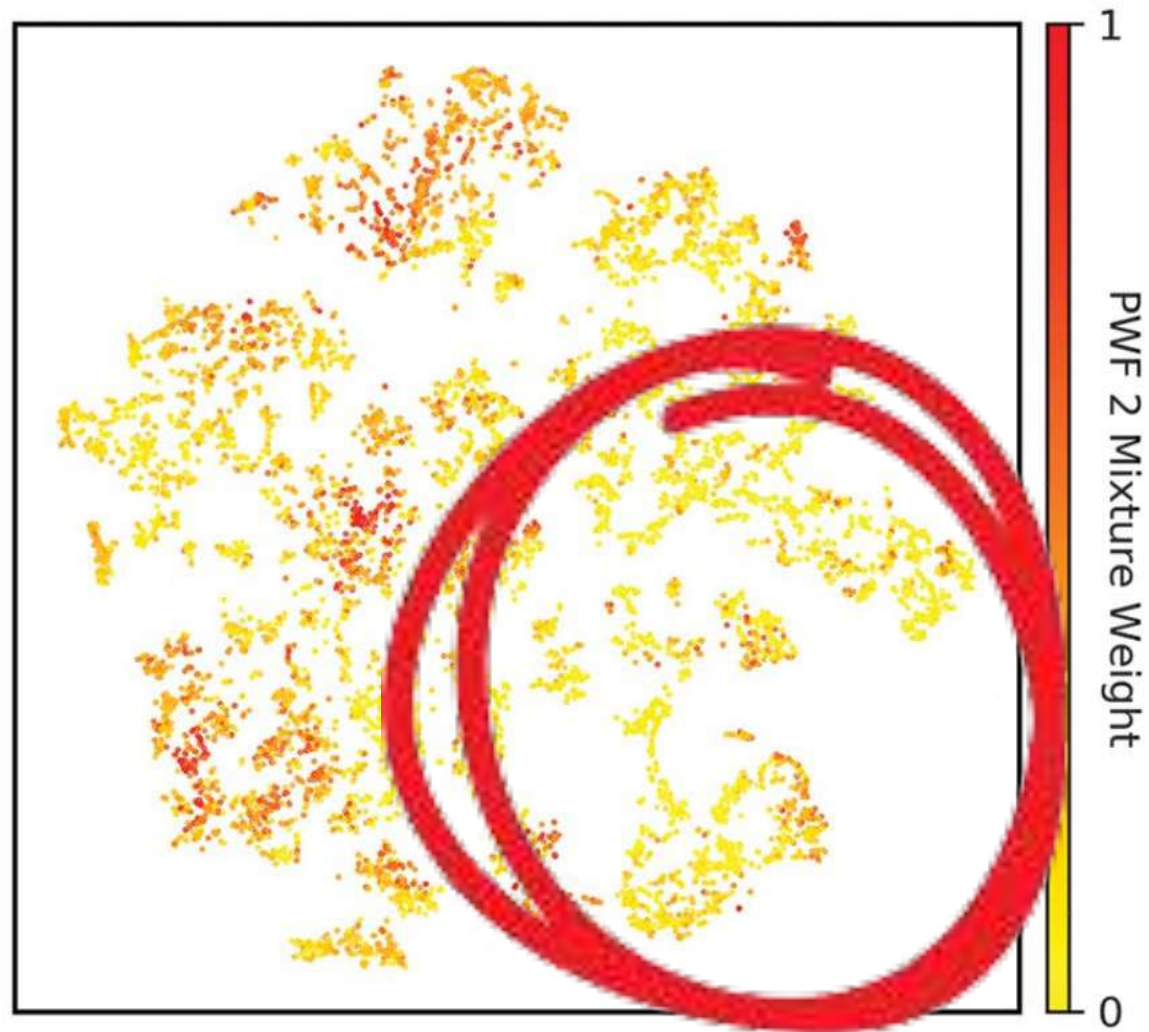distorted probability weighting function (PWF)

# Paper to read

# Paper to read

# Paper to read

# Paper to read

- Have the authors discovered something <u>new</u>?

- They have used machine learning to:
  1/ confirm the architecture of existing theories
  2/ show that there is room for improvement
  3/ provide directions for improving theories

- Can we replace artificial neural networks by computer algorithms?

- Why are computer algorithms <u>more explainable</u> than artificial neural networks?

# Modeling with machine learning

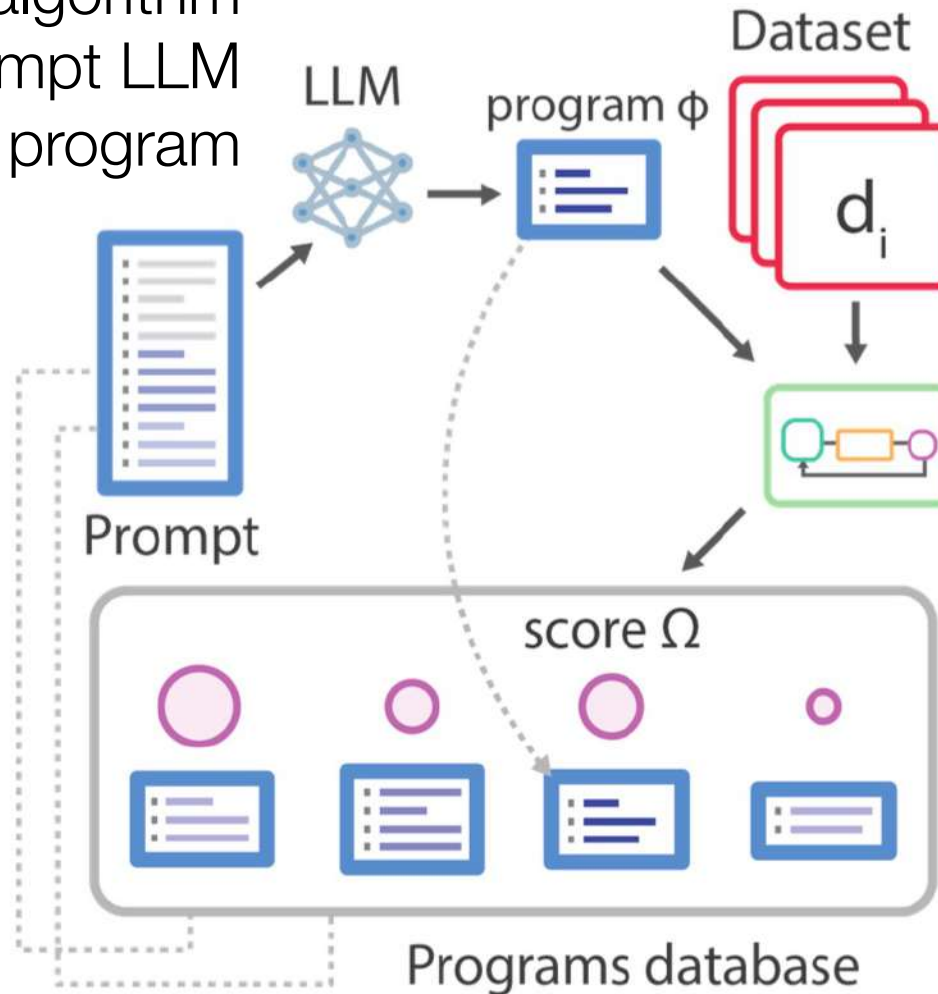# Discovering Symbolic Cognitive Models from Human and Animal Behavior

Pablo Samuel Castro[1], Nenad Tomasev[1], Ankit Anand[1], Navodita Sharma[1], Rishika Mohanta[2,3], Aparna Dev[2], Kuba Perlin[1], Siddhant Jain[1], Kyle Levin[1], Noémi Éltető[1,4], Will Dabney[1], Alexander Novikov[1], Glenn C Turner[2], Maria K Eckstein[1], Nathaniel D Daw[1,5], Kevin J Miller[*,1,6] and Kimberly L Stachenfeld[*,1,7]

[*]Equal contributions, [1]Google DeepMind, [2]Janelia Farm Research Campus, Howard Hughes Medical Institute, Ashburn, VA, USA, [3]The Rockefeller University, New York, NY, USA, [4]Max Planck Institute for Biological Cybernetics, Tübingen, Germany, [5]Princeton Neuroscience Institute, Princeton University, Princeton, NJ, USA, [6]Sainsbury Wellcome Centre, University College London, United Kingdom, [7]Center for Theoretical Neuroscience, Columbia University, New York, NY, USA
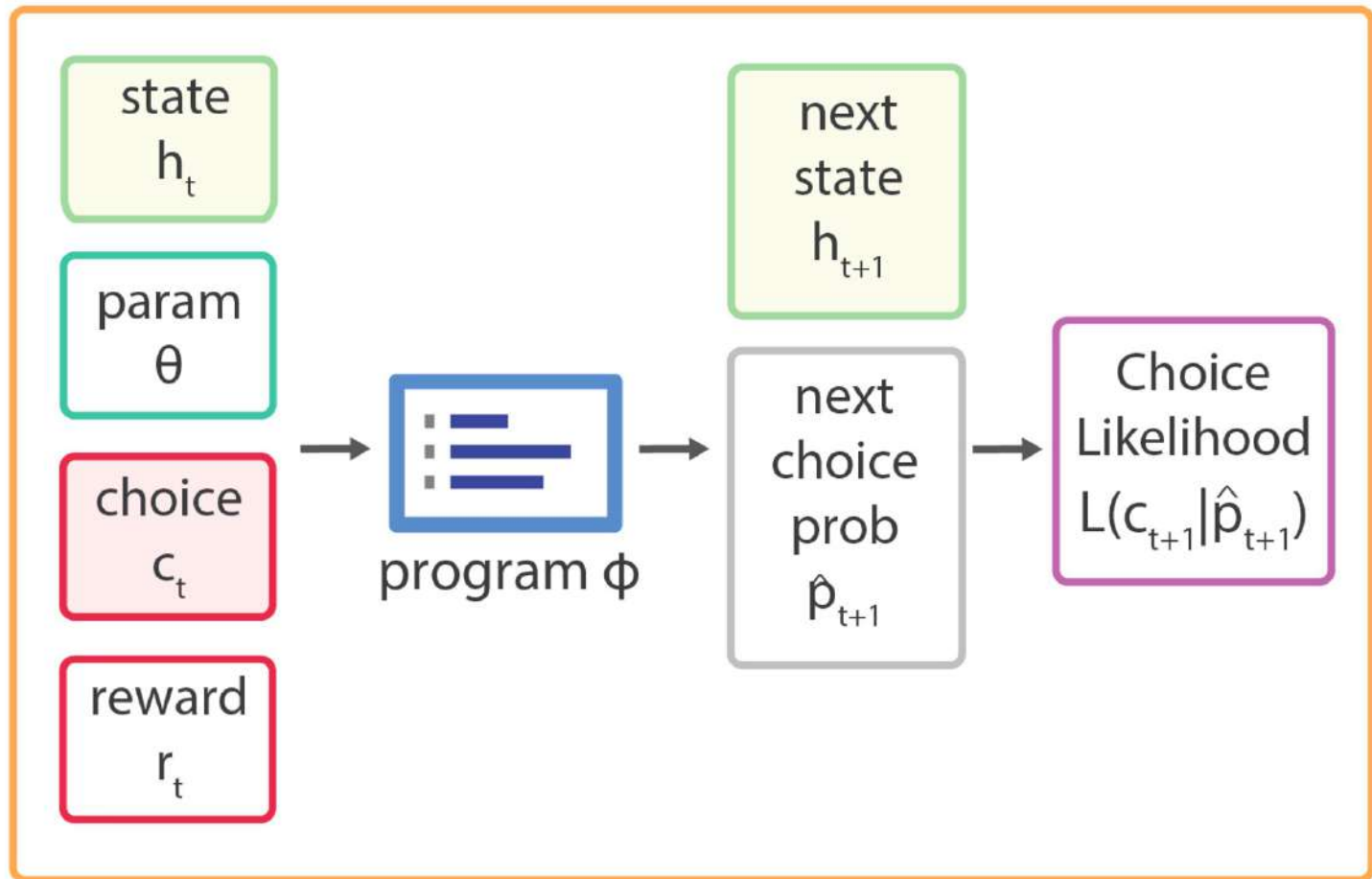
Symbolic models play a key role in cognitive science, expressing computationally precise hypotheses about how the brain implements a cognitive process. Identifying an appropriate model typically requires a great deal of effort and ingenuity on the part of a human scientist. Here, we adapt Romera-Paredes et al. (2024), a recently developed tool that uses Large Language Models (LLMs) in an evolutionary algorithm, to automatically discover symbolic cognitive models that accurately capture human and animal behavior. We consider datasets from three species performing a classic reward-learning task that has been the focus of substantial modeling effort, and find that the discovered programs outperform state-of-the-art cognitive models for each. The discovered programs can readily be interpreted as hypotheses about human and animal cognition, instantiating interpretable symbolic learning and decision-making algorithms. Broadly, these results demonstrate the viability of using LLM-powered program synthesis to propose novel scientific hypotheses regarding mechanisms of human and animal cognition.
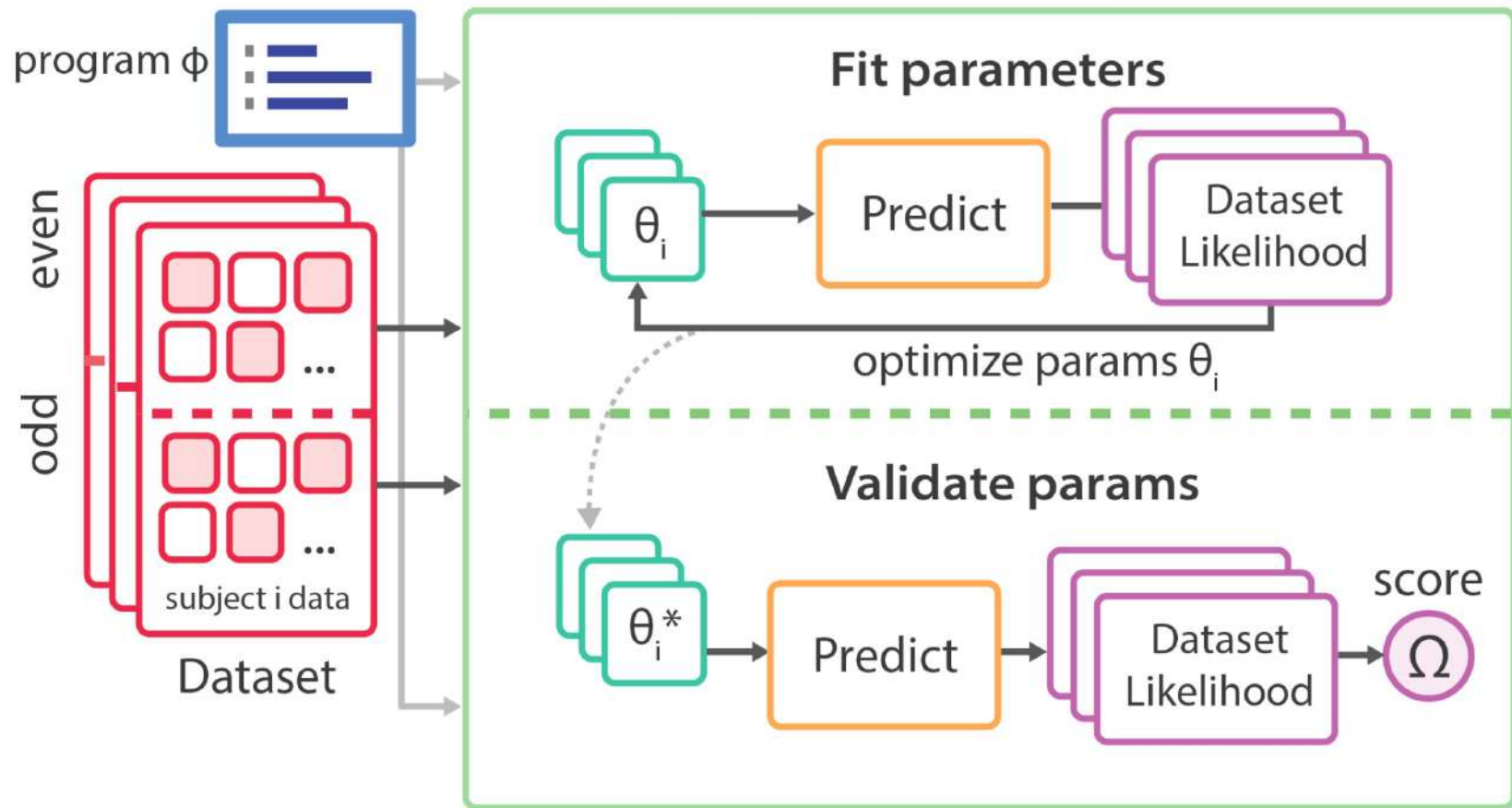
# Modeling with machine learning

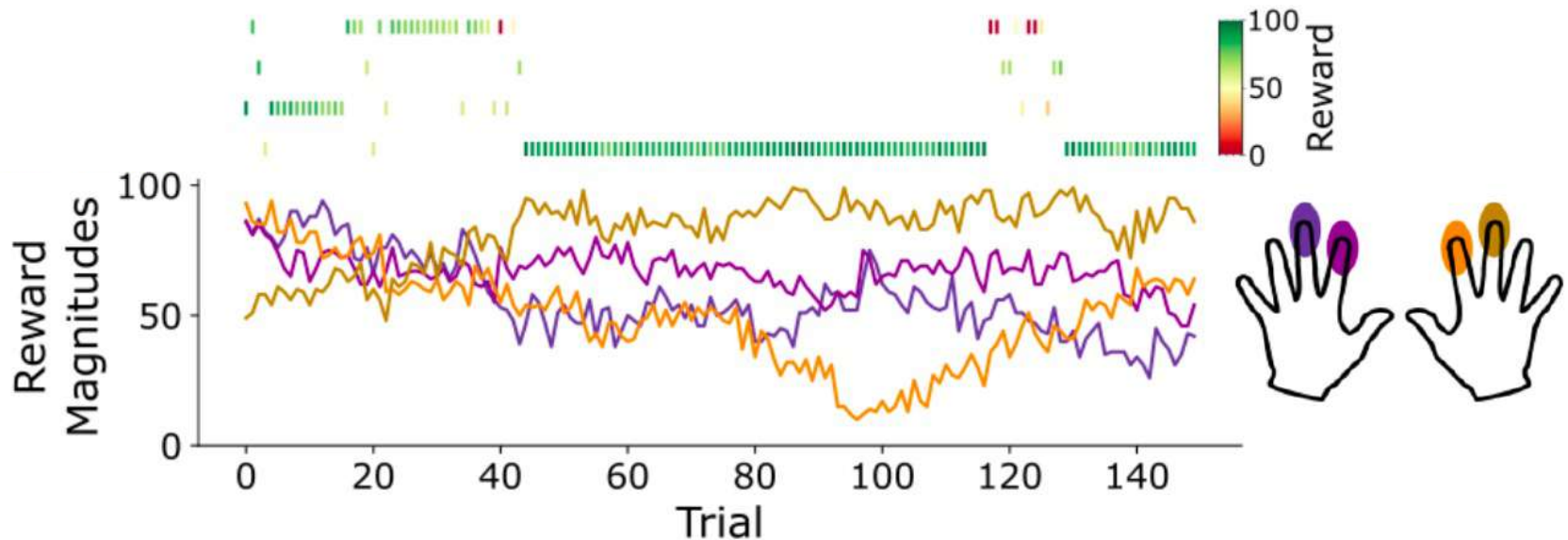evolutionary algorithm used to prompt LLM to improve program
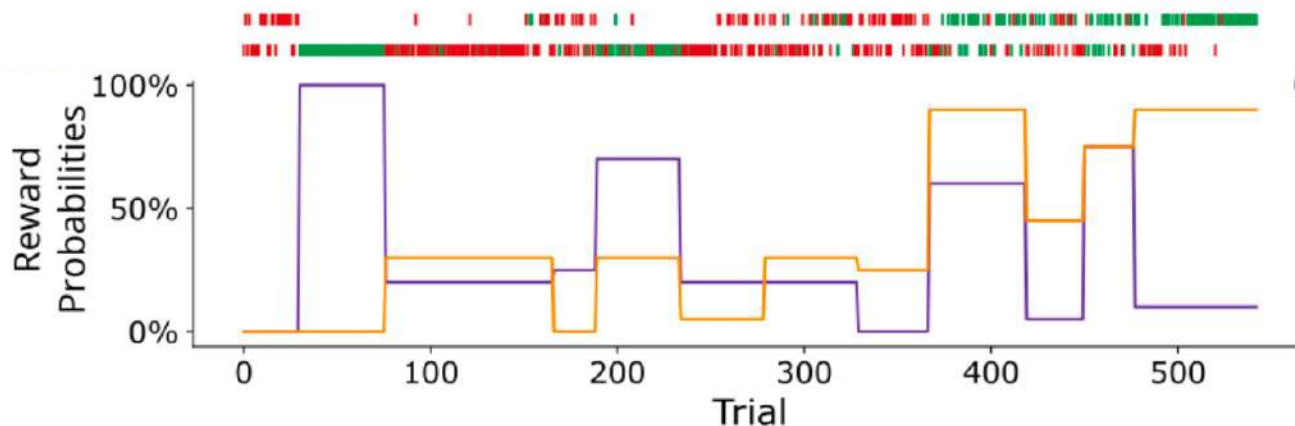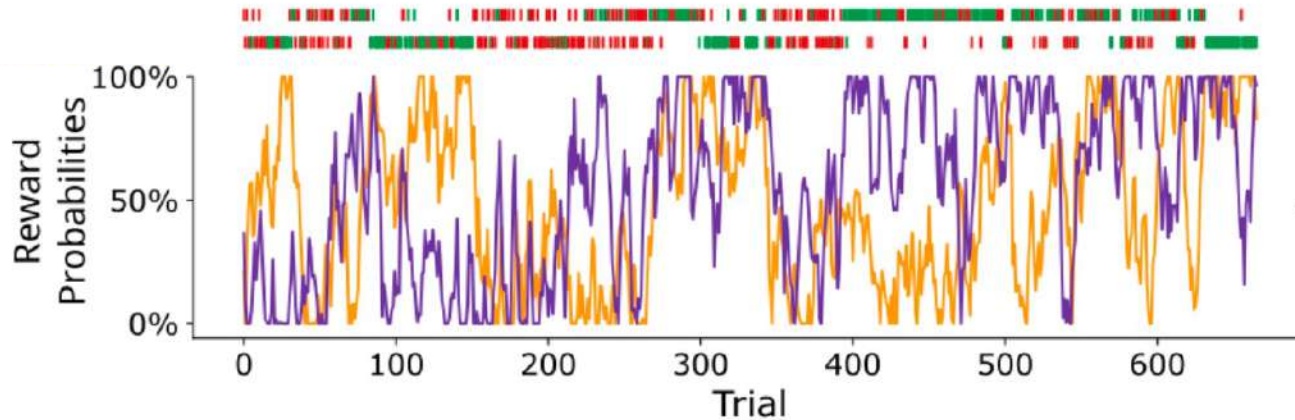
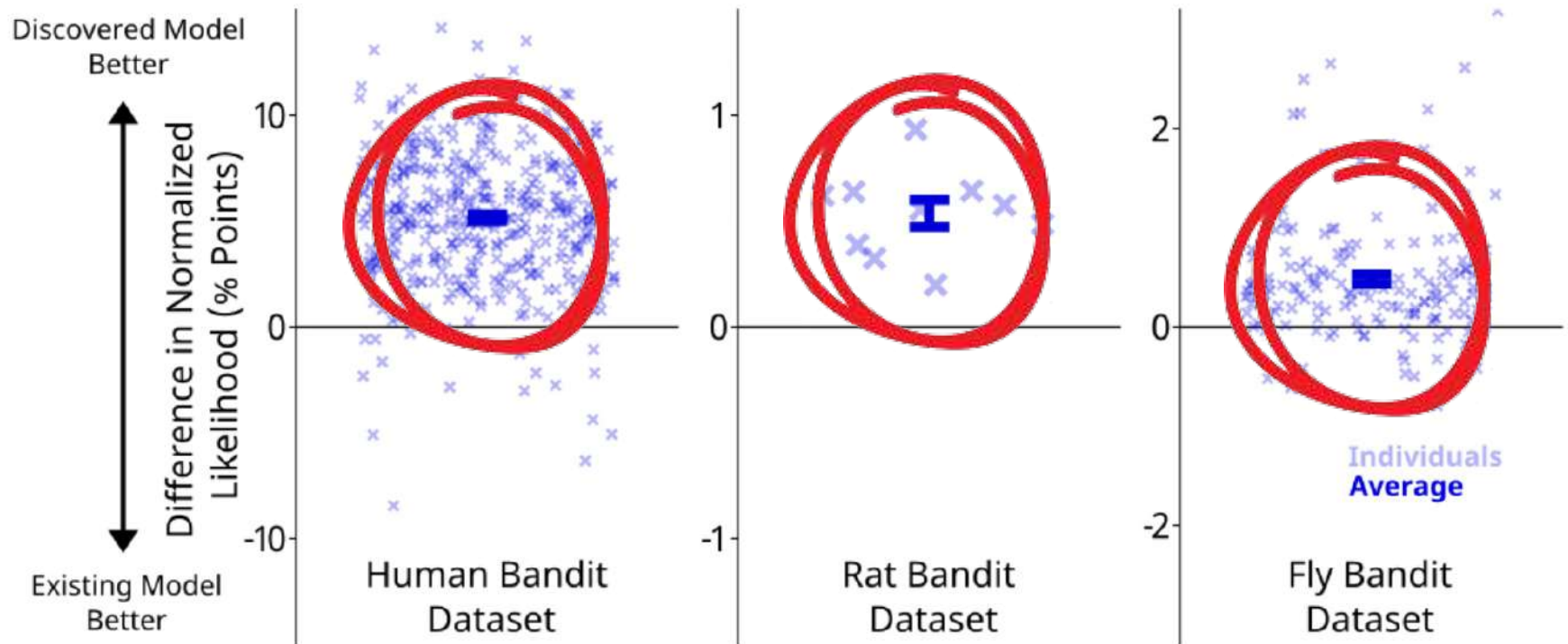# Modeling with machine learning

# Modeling with machine learning

# Modeling with machine learning

# Modeling with machine learning

# Modeling with machine learning

# Modeling with machine learning

**Top Scoring Human Program**    A salient feature of this program was that the bulk of the code—and of the agent's internal state—was devoted to choice history rather than reward tracking. In addition to variables tracking the expected values of the four actions, the best program introduced a number of novel variables that each track different reward-independent statistics of previous choices:

```
q_values = agent_state[:4]
old_choice = agent_state[4]
trial_since_last_switch = agent_state[5]
exploration_rate = agent_state[6]
cumchoice = agent_state[7:11]
```

# Modeling with machine learning

Three independent runs with the Structured2 seed program, (highest scoring program for Human), separately discovered a common (but, to our knowledge, novel) motif whereby the learned values were decayed, at each step, toward their average:

```
# Best Human Bandit Program
q_values = (1 - exploration_rate) * q_values + (
    exploration_rate * jnp.mean(q_values))

# Program 2
q_values = (1-bias) * q_values + bias * jnp.mean(q_values)

# Program 3
updated_q_values = (1 - alpha_choice) * updated_q_values + (
    alpha_choice * jnp.mean(q_values))
```

🙄

# Modeling with machine learning

**Article**

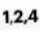# Compulsivity is linked to suboptimal choice variability but unaltered reinforcement learning under uncertainty

Junseok K. Lee [1,2] ✉, Marion Rouault [1,2,3] & Valentin Wyart [1,2,4] ✉

Compulsivity has been associated with variable behavior under uncertainty. However, previous work has not distinguished between two main sources of behavioral variability: the stochastic selection of choice options that do not maximize expected reward (choice variability) and random noise in the reinforcement learning process that updates option values from choice outcomes (learning variability). Here we study the relation between dimensional compulsivity and behavioral variability using a computational model that dissociates its two sources. Across two independent datasets

# Modeling with machine learning

# Coming next

- <u>Practical session:</u> today, 2.00pm, same room

- <u>Guidelines for cognitive modeling:</u>
  Wilson and Collins (2019) Ten simple rules for the computational
  modeling of behavioral data. *eLife*
  https://doi.org/10.7554/eLife.49547 (open-access)

- <u>Contact:</u>
  **Valentin Wyart**     valentin.wyart@ens.psl.eu
  **Lucas Benjamin**     lucas.benjamin78@gmail.com

  Lab. de Neurosciences Cognitives et Computationnelles (LNC[2])
  Institut National de la Santé et de la Recherche Médicale
  Ecole Normale Supérieure, Université PSL