

Competition between Model-based and Model-free Learning

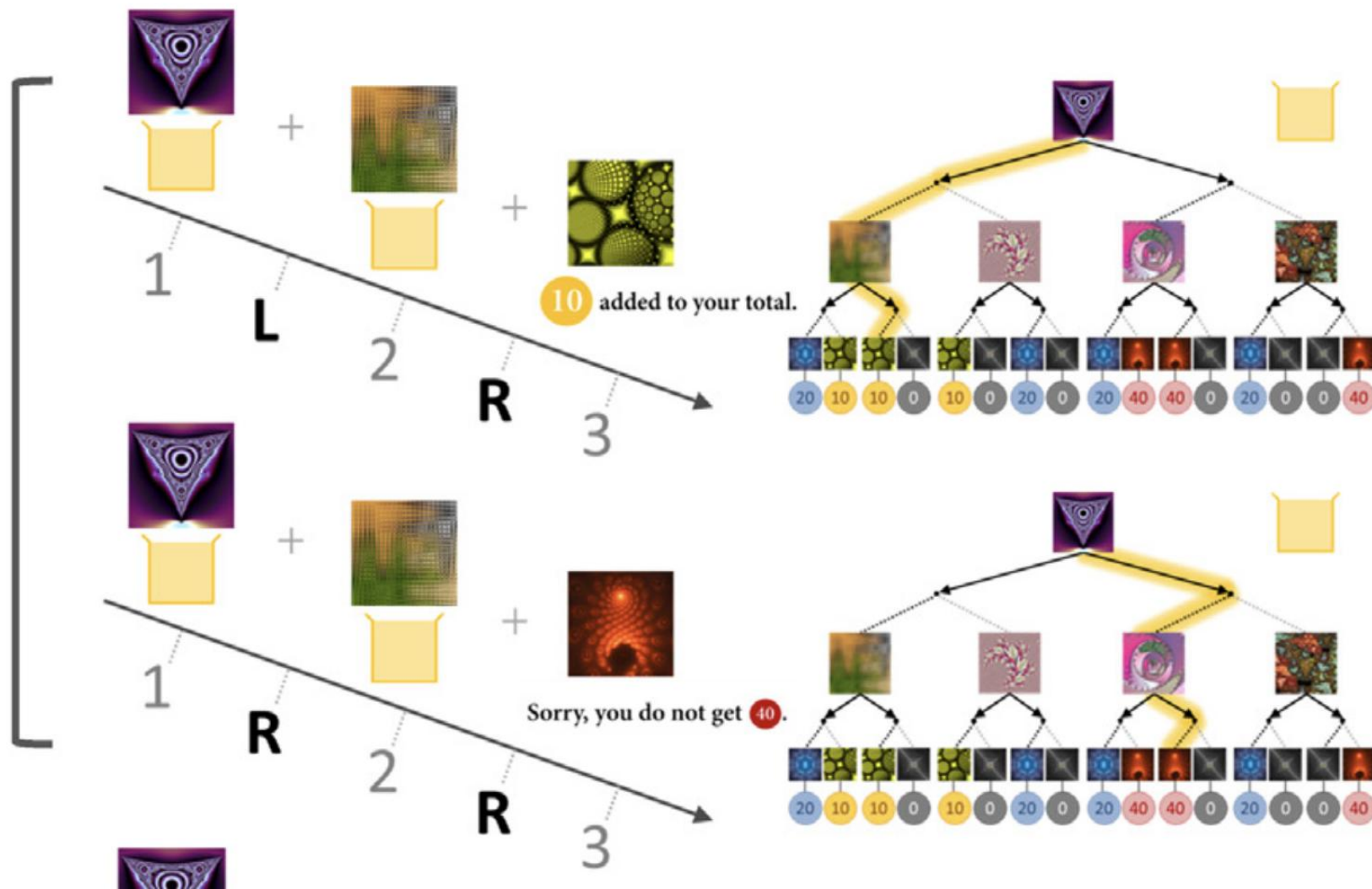
— — Arbitration model

张博涛 进展汇报

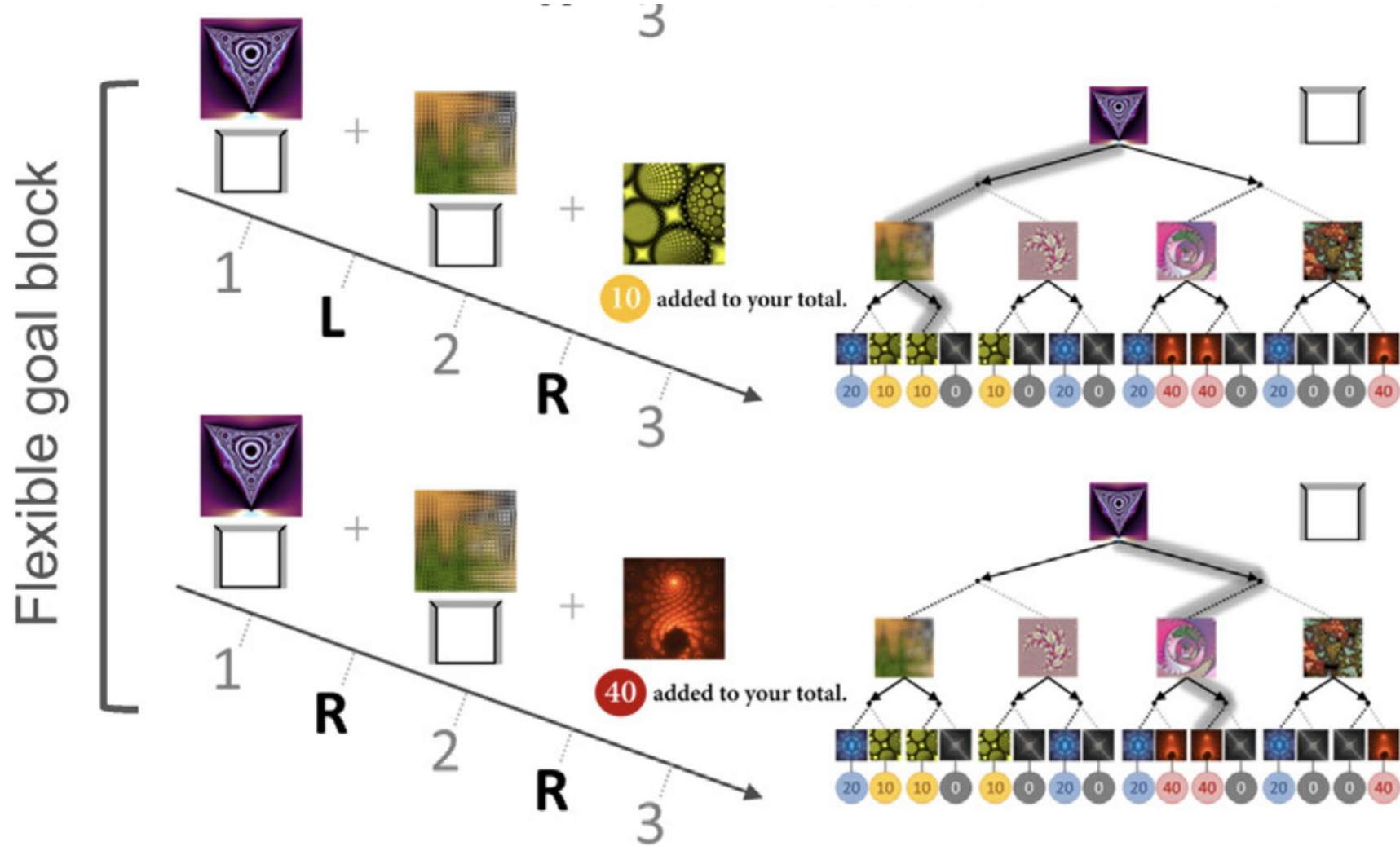
Task setting: Two Stage Task

C

Specific goal block

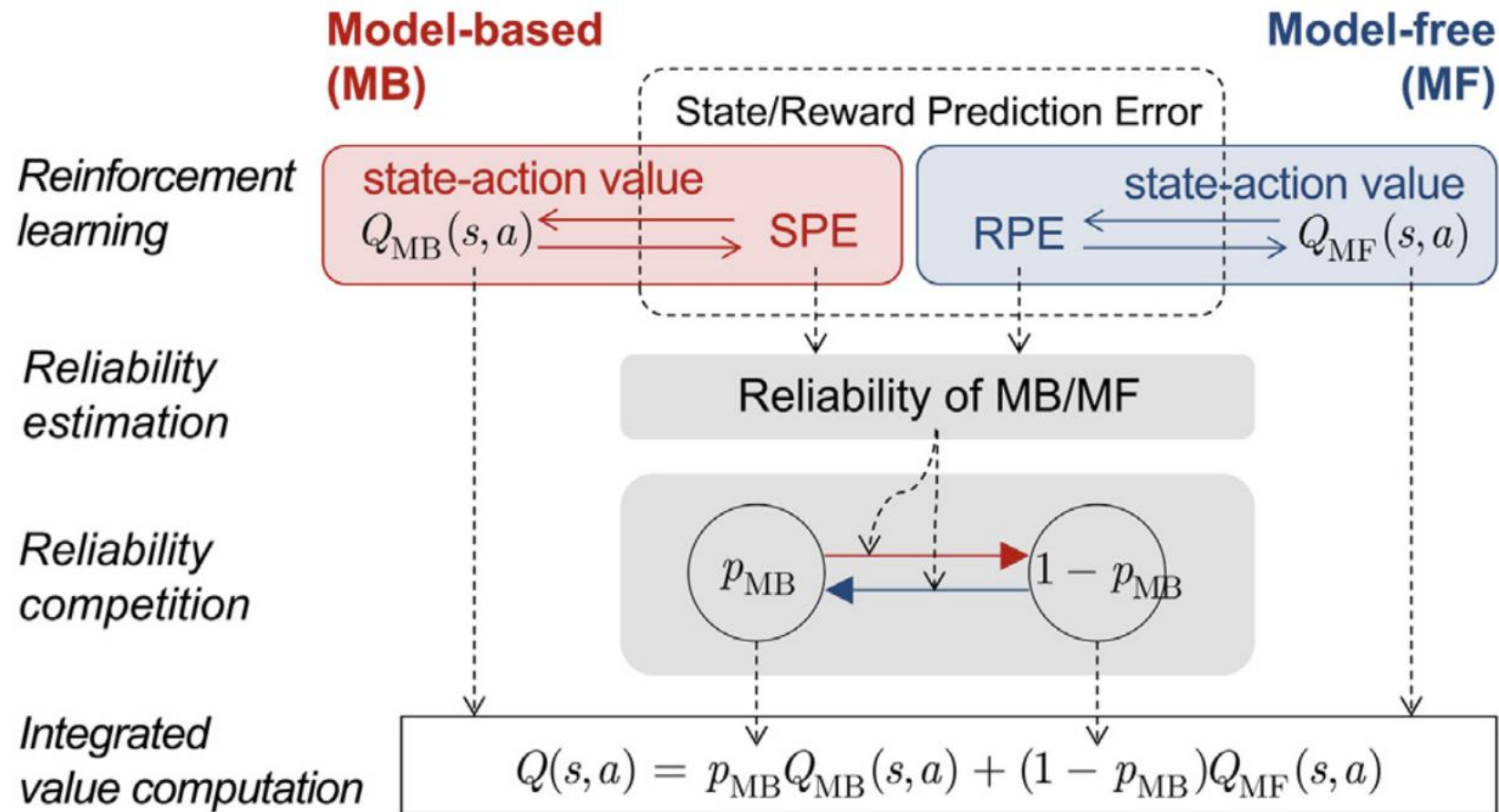


Task setting: Two Stage Task



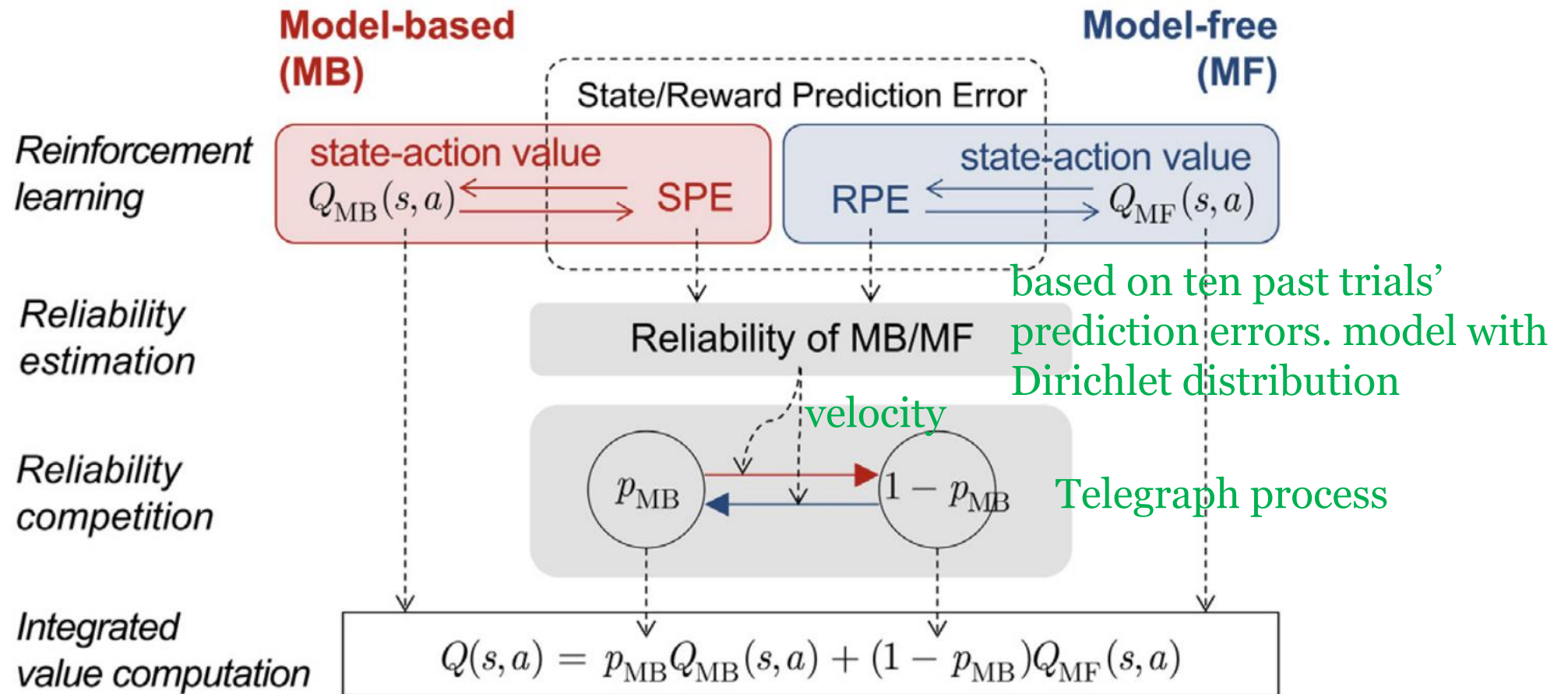
Neural Computations Underlying Arbitration between Model-Based and Model-free Learning

Sang Wan Lee,^{1,2,3,*} Shinsuke Shimojo,^{1,2,4} and John P. O'Doherty^{1,2,3}



Neural Computations Underlying Arbitration between Model-Based and Model-free Learning

Sang Wan Lee,^{1,2,3,*} Shinsuke Shimojo,^{1,2,4} and John P. O'Doherty^{1,2,3}



Loss4arbitration Model

*Integrated
value computation*

$$Q(s, a) = \overset{\downarrow}{p_{\text{MB}}} Q_{\text{MB}}(s, a) + \overset{\downarrow}{(1 - p_{\text{MB}})} Q_{\text{MF}}(s, a)$$

$w = \operatorname{argmin}_w Q_w(s, a) - Q_{\text{true}}(s, a)$ where

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

Loss4arbitration Model

*Integrated
value computation*

$$Q(s, a) = p_{\text{MB}} Q_{\text{MB}}(s, a) + (1 - p_{\text{MB}}) Q_{\text{MF}}(s, a)$$

$w = \operatorname{argmin}_w Q_w(s, a) - Q_{\text{true}}(s, a)$ where

$$Q_w(s, a) = w Q_{\text{MB}}(s, a) + (1 - w) Q_{\text{MF}}(s, a)$$

gradient descent on $\mathcal{L} = (Q_w - Q_{\text{true}})^2$

Model Fitting

minimize negative loglikelihood

subjects' data: true actions, a_t

model's prediction: $\mathcal{P} \sim \text{softmax}(Q_w)$

negative loglikelihood: $NLL = \sum \mathcal{P}(a_t)$

Model Fitting

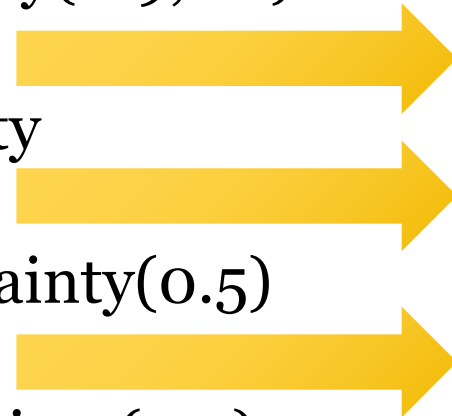
150 trials.

trial 1 ~ 37: fixed goal(20), low uncertainty(0.9,0.1)

trial 38 ~ 75: flexible goal, low uncertainty

trial 76 ~ 112: fixed goal(10), high uncertainty(0.5)

trial 113 ~ 150: flexible goal, high uncertainty(0.5)



back planning(MB)

Model Comparison

MB: pure model-based learning

MF: pure model-free learning

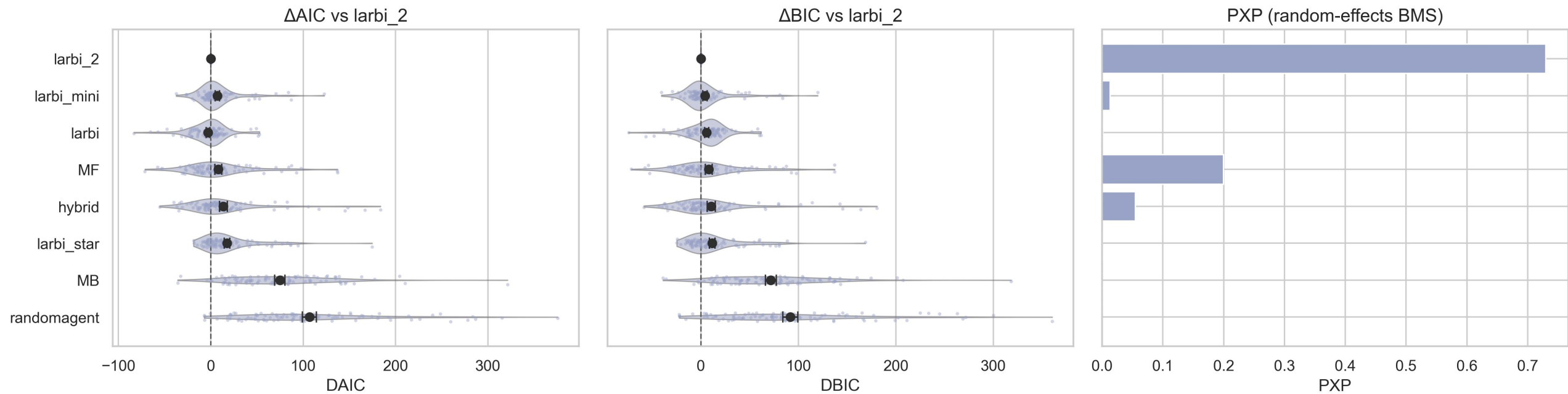
Hybrid: $Q = wQ_{MB} + (1 - w)Q_{MF}$, but fixed w

Loss4arbitration:

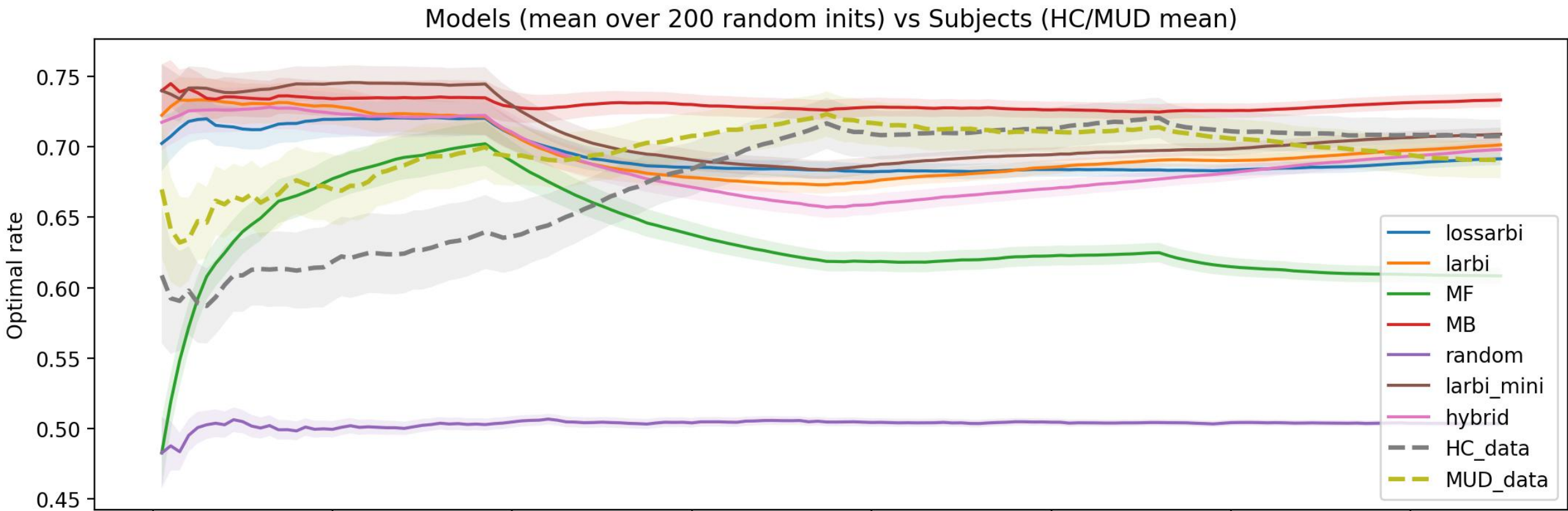
1. larbi: MF learning rates * 2, MB ... * 2, discount, w 's learning rate, exploration temperature * 2
2. larbi_mini: MF, MB, w learning rates, exploration temperature * 1
3. larbi_star: based on mini, but MB learning rate = MF's
4. larbi_2: based on mine, but a separate arbitration w at each stage

Model Comparison

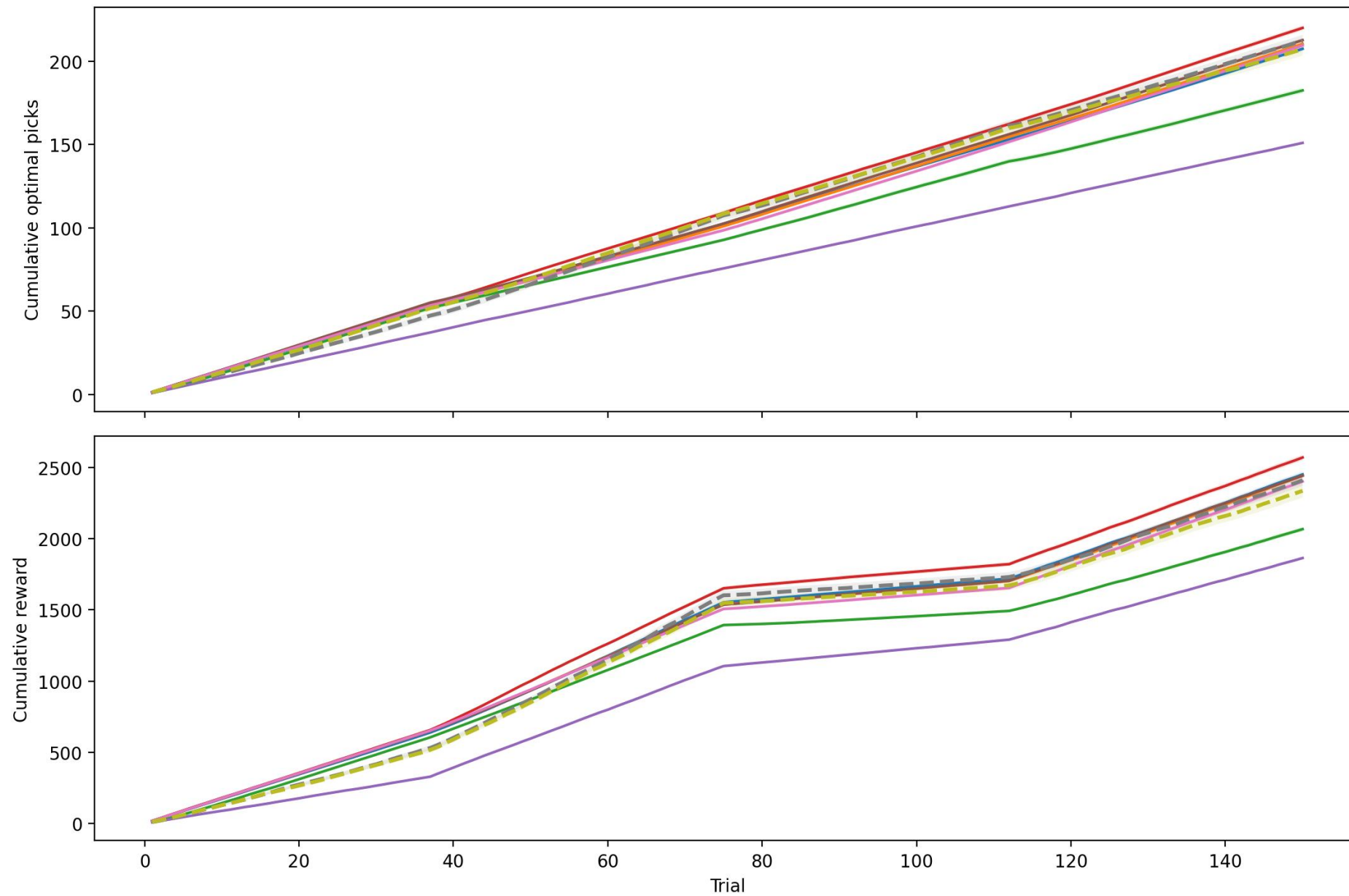
Model comparison (baseline = larbi)



Model Comparison



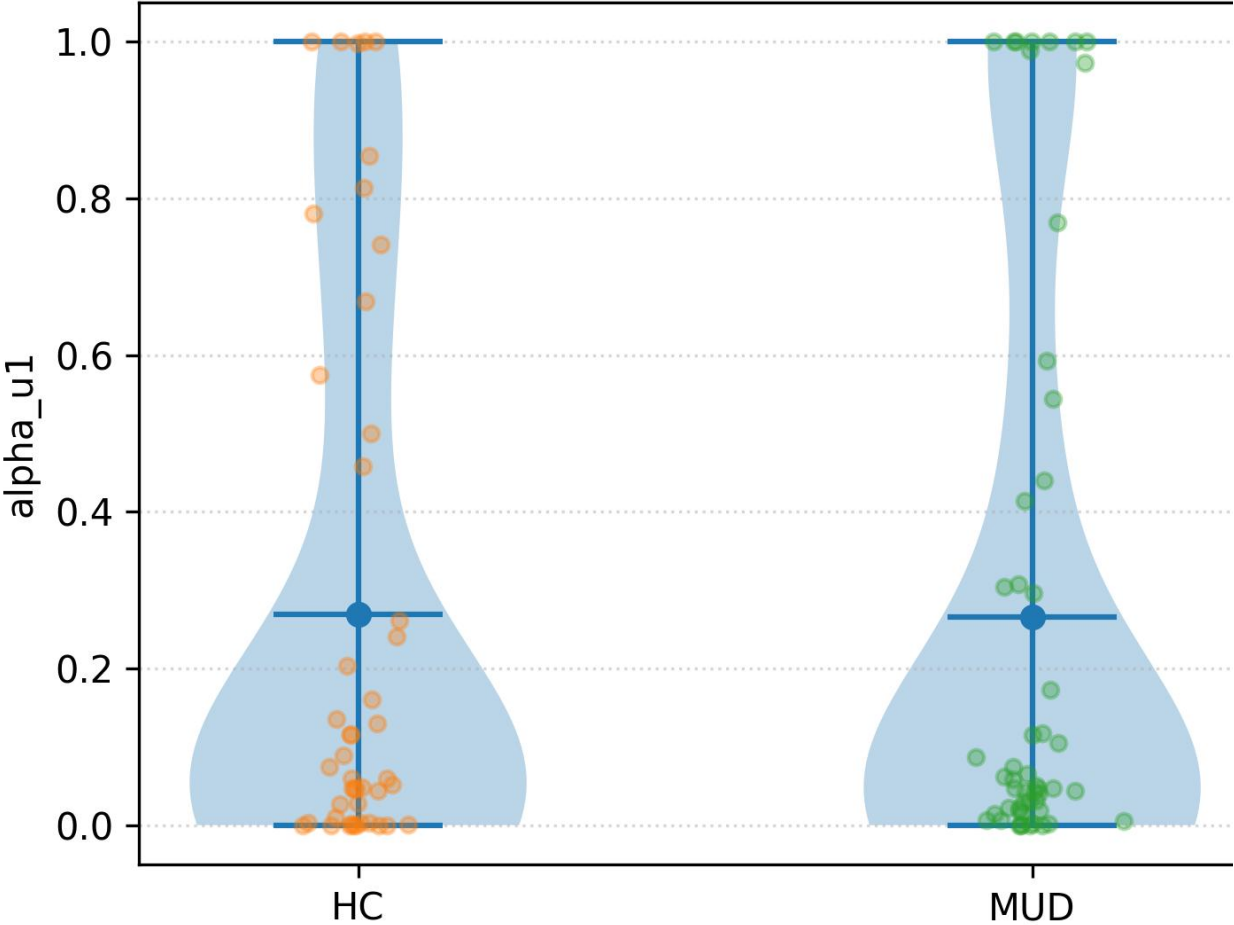
Model Comparison



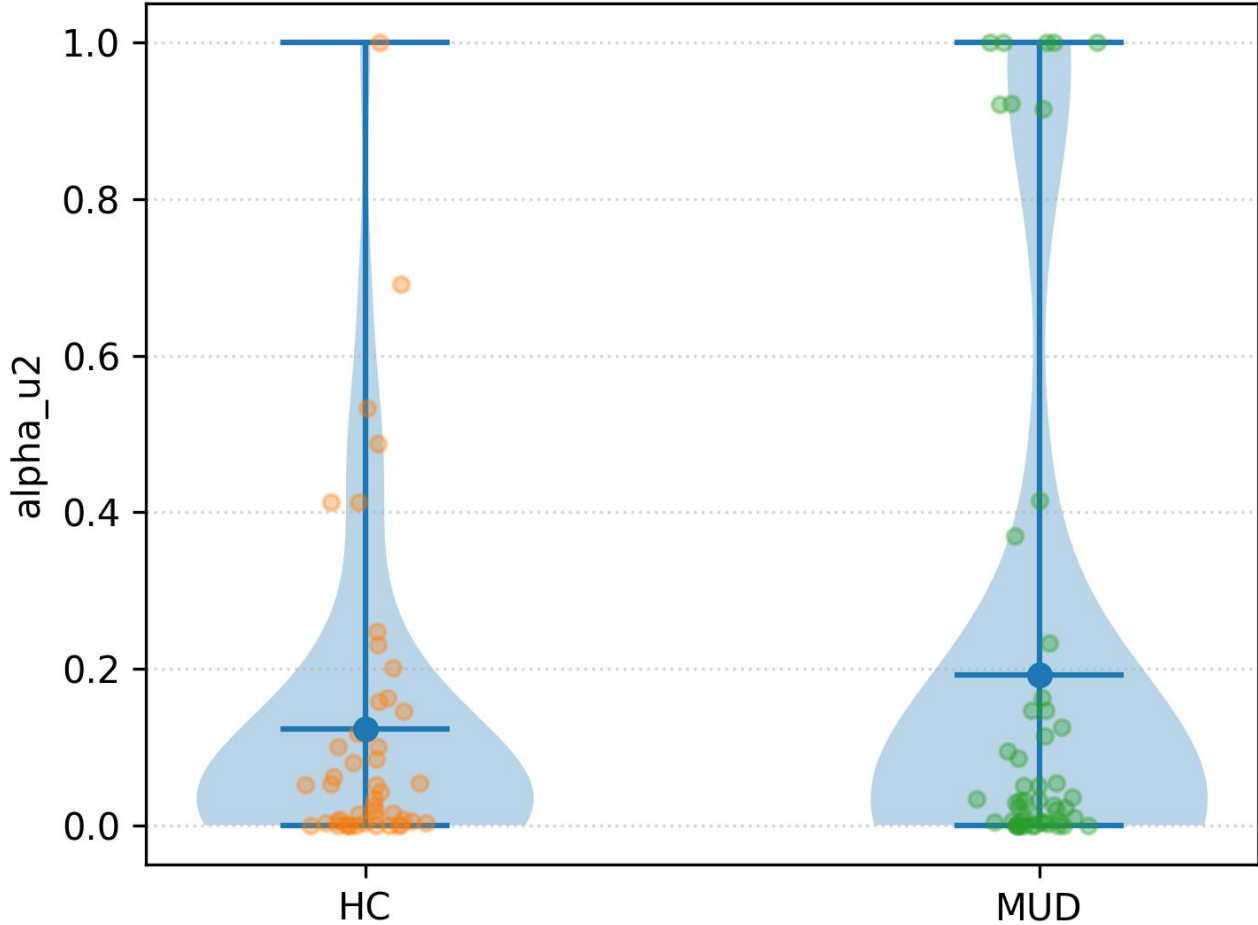
Model Comparison

Results Analysis

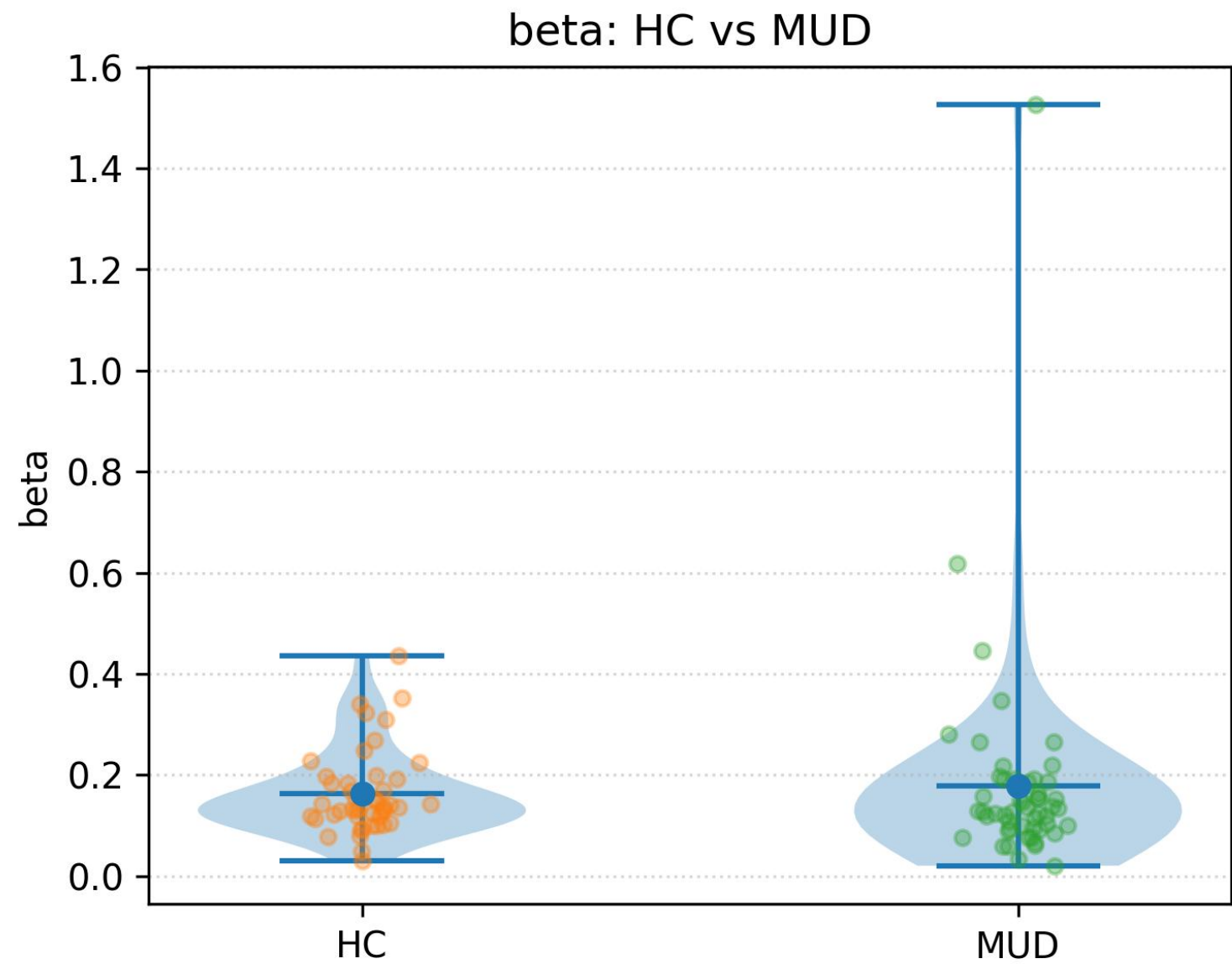
alpha_u1: HC vs MUD



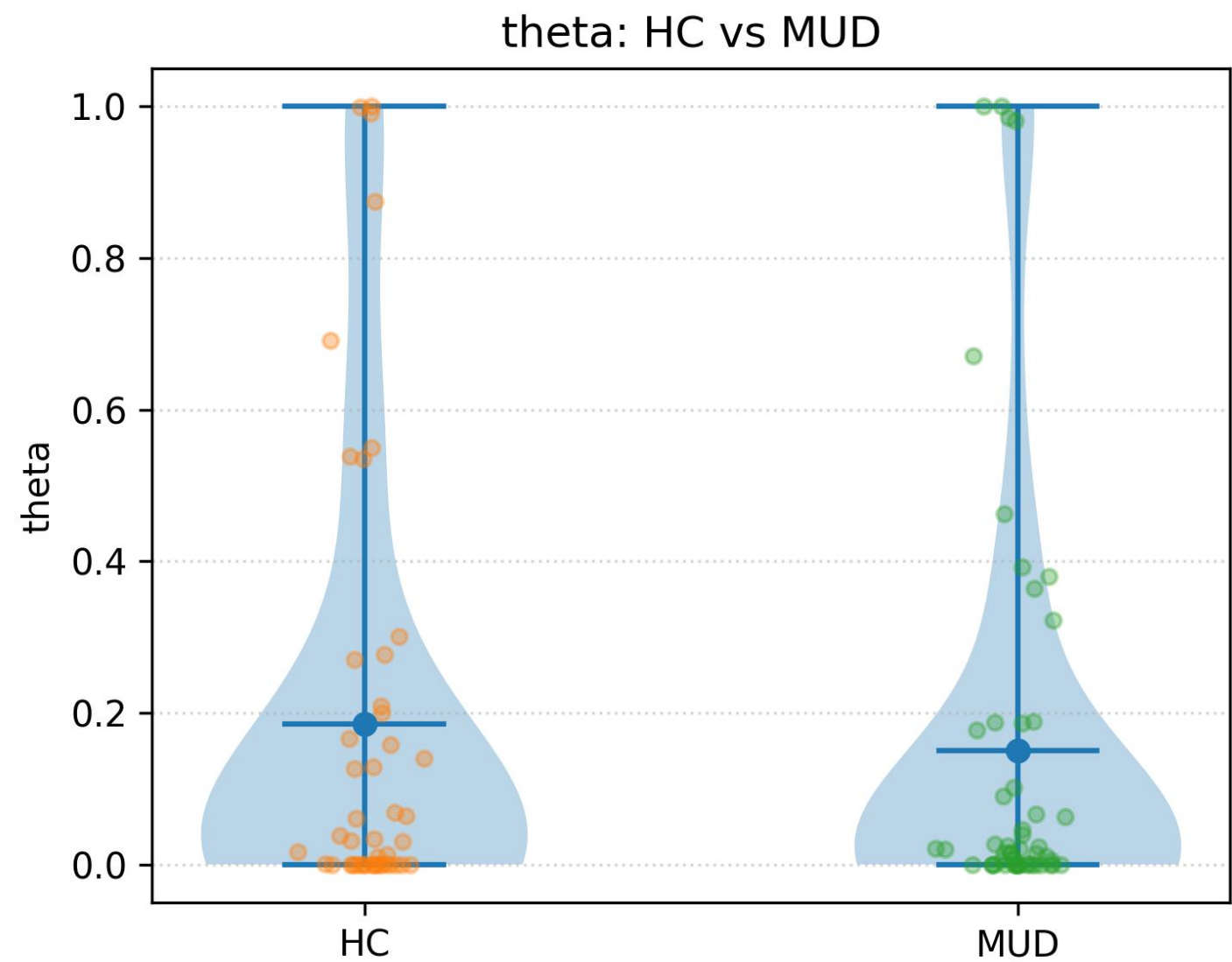
alpha_u2: HC vs MUD



Results Analysis



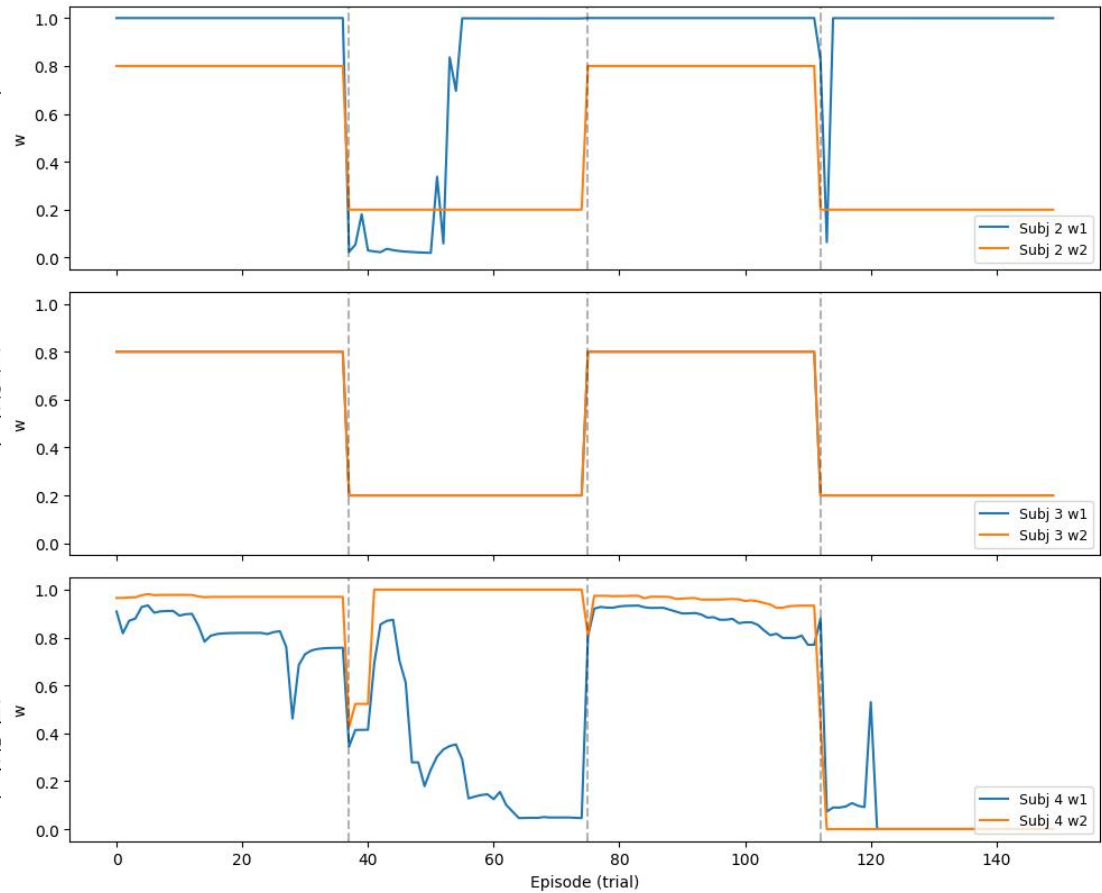
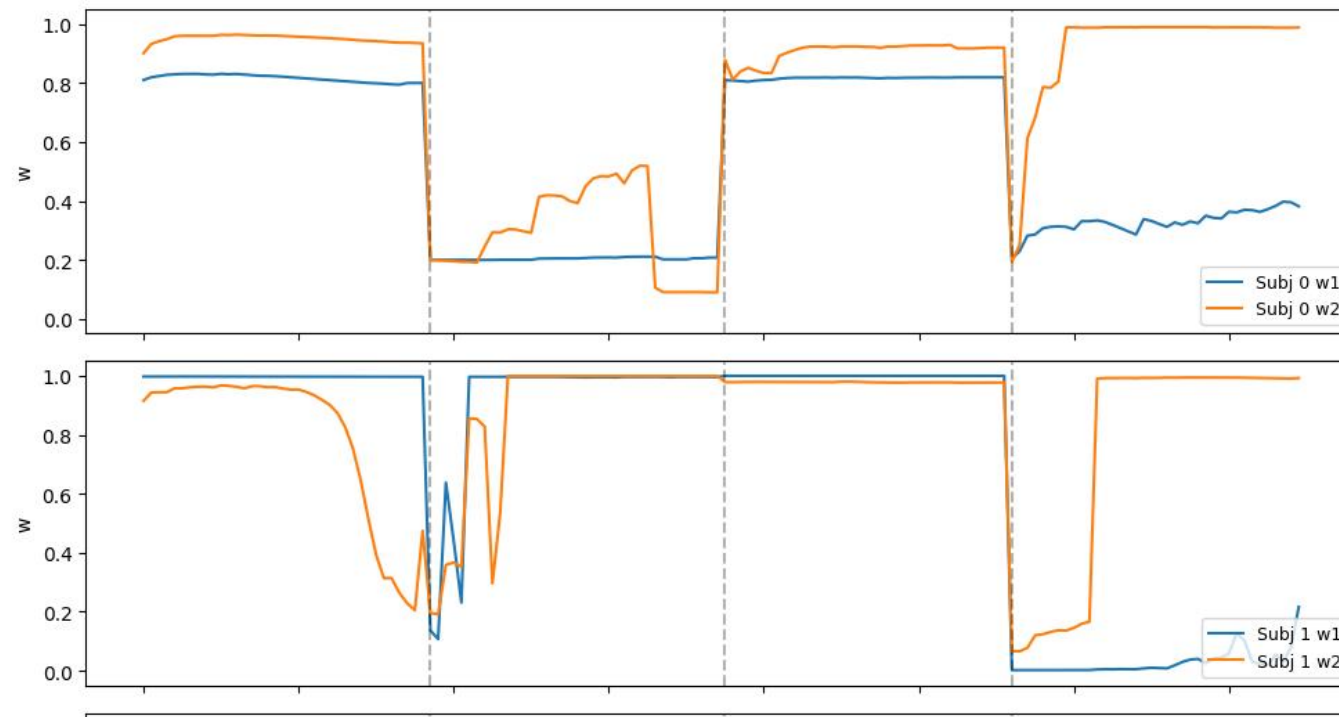
Results Analysis



around 1/3 is 0

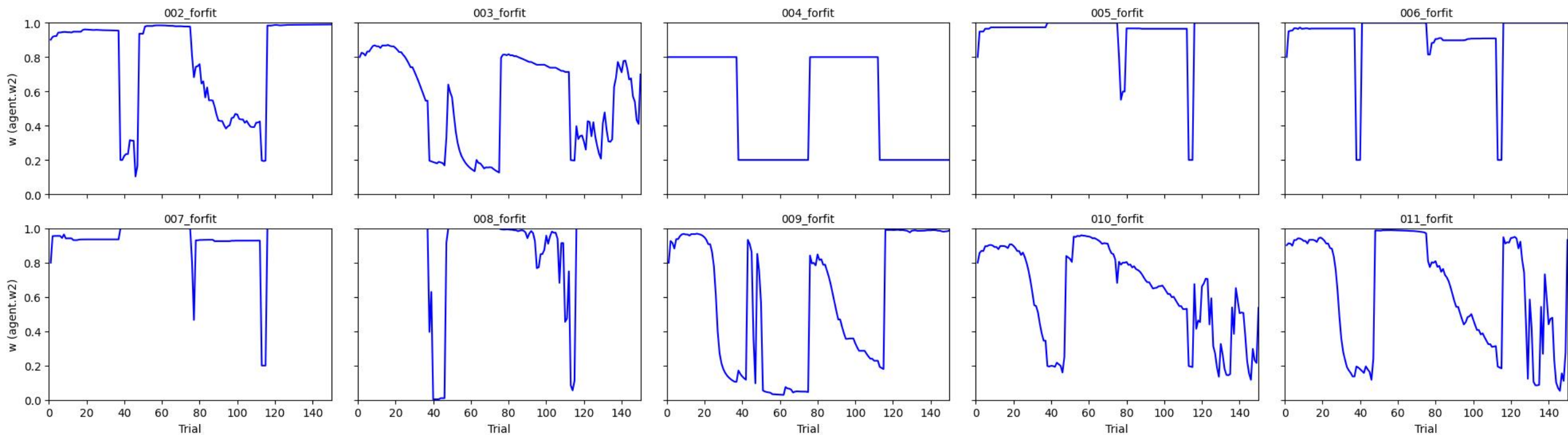
Results Analysis

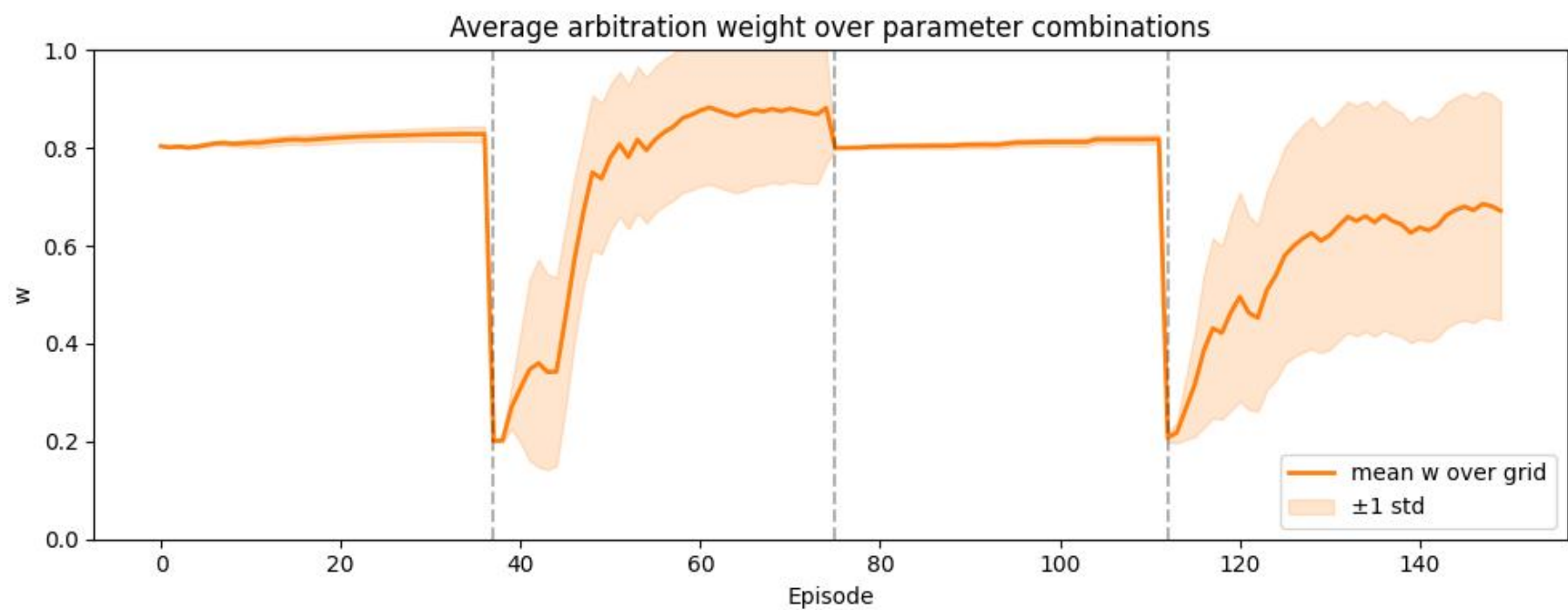
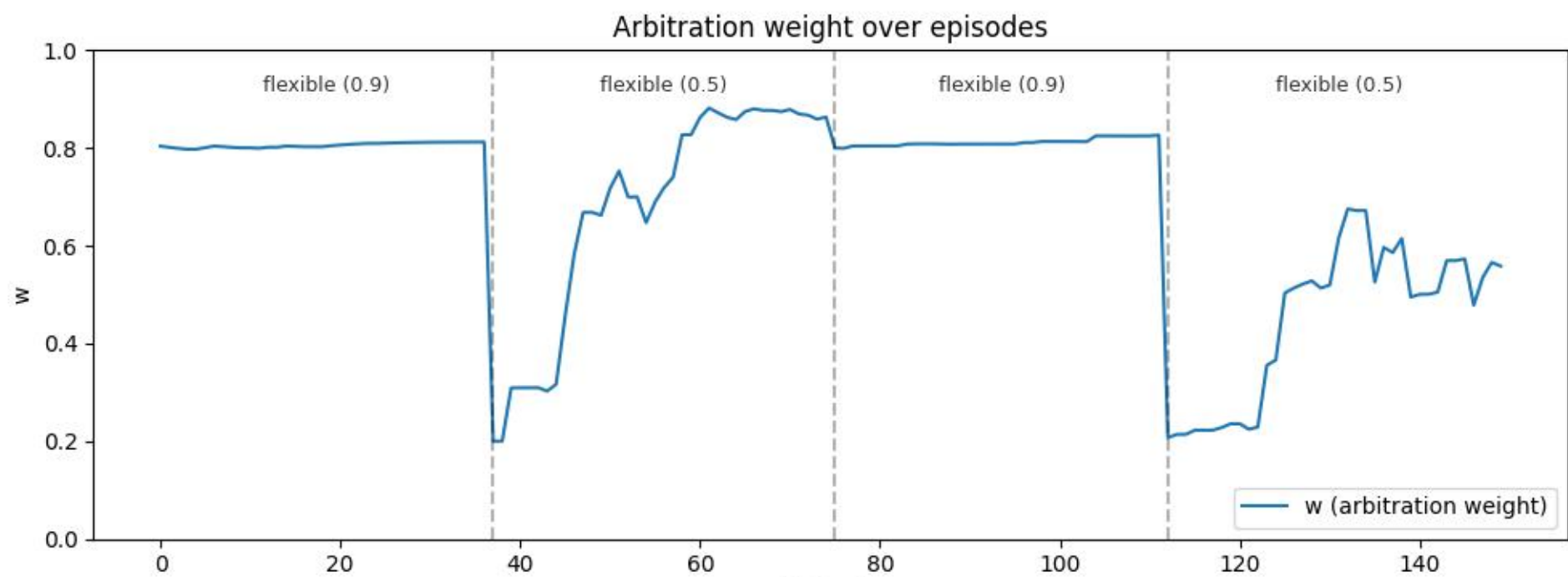
Arbitration weights w1/w2 per subject (first 5)



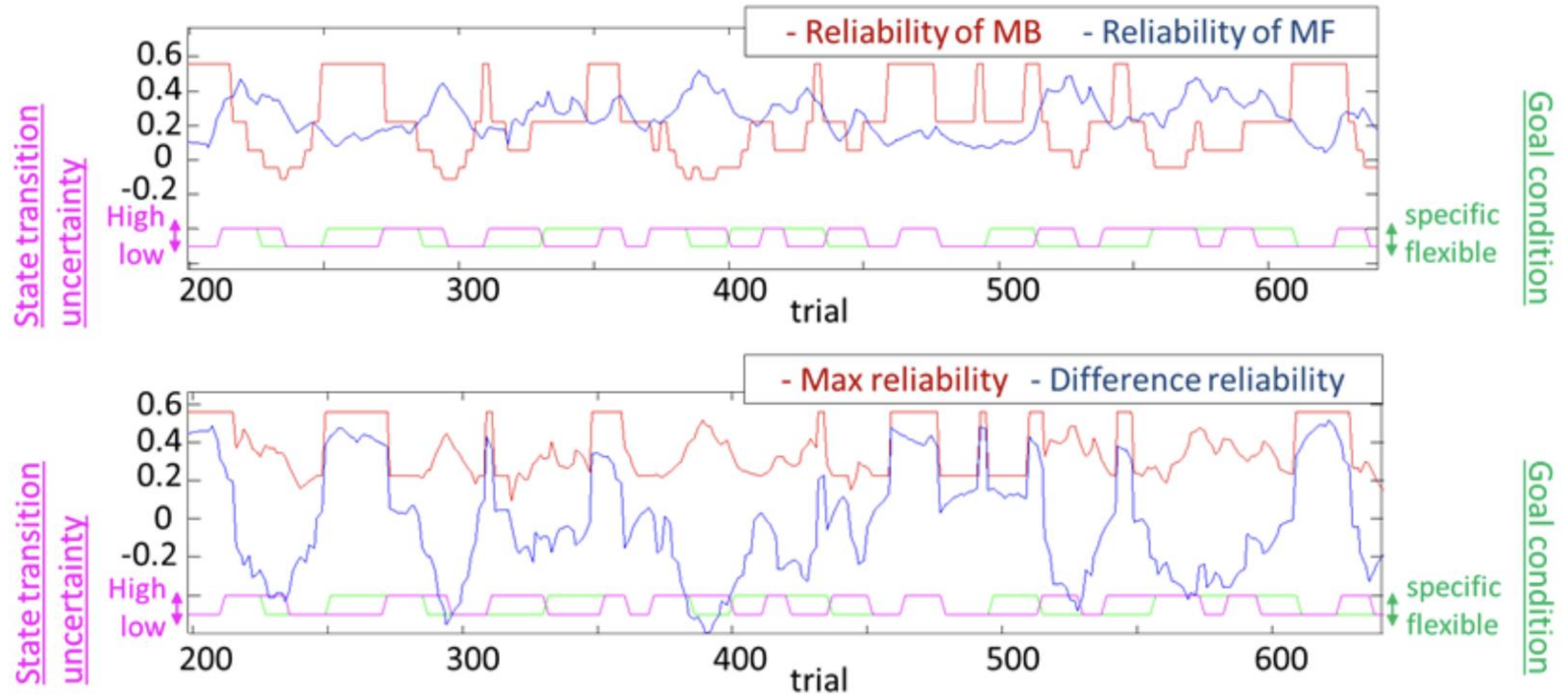
Results Analysis

w traj





Possible Explanations



Possible Explanations

However, to incentivize subjects to continue learning, throughout the task, the chances of payoff associated with the four second-stage options were changed slowly and independently, according to Gaussian random walks.

In the Lee 2014 task setting, perhaps Model-Based **is** oftentimes better? (also, from the total reward attained at each block)

The current state of research on MB-MF competition/cooperation

Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control

Nathaniel D Daw¹, Yael Niv^{1,2} & Peter Dayan¹

Speed/Accuracy Trade-Off between the Habitual and the Goal-Directed Processes

Mehdi Keramati¹*, Amir Dezfouli², Payam Piray²

1 School of Management and Economics, Sharif University of Technology, Tehran, Iran, **2** Control and Intelligent Processing Center Of Excellence, School of Electrical and Computer Engineering, University of Tehran, Tehran, Iran

When Does Model-Based Control Pay Off?

Wouter Kool^{1*}, Fiery A. Cushman¹, Samuel J. Gershman^{1,2}

1 Department of Psychology, Harvard University, Cambridge, Massachusetts, United States of America,

2 Center for Brain Science, Harvard University, Cambridge, Massachusetts, United States of America

Cost-Benefit Arbitration Between Multiple Reinforcement-Learning Systems



Wouter Kool¹, Samuel J. Gershman^{1,2}, and
Fiery A. Cushman¹

¹Department of Psychology, Harvard University, and ²Center for Brain Science, Harvard University

Beyond dichotomies in reinforcement learning

Anne G. E. Collins and Jeffrey Cockburn

MB-MF competition/cooperation

Two Stage Task

Humans are primarily model-based and not model-free learners in the two-stage task

Carolina Feher da Silva¹ and Todd A. Hare¹

¹Zurich Center for Neuroeconomics, Department of Economics, University of Zurich,
Zurich, Switzerland

(In) Relation to Memory

How much of reinforcement learning is working memory, not reinforcement learning? A behavioral, computational, and neurogenetic analysis

Anne G. E. Collins and **Michael J. Frank**

“systematically varying the size of the learning problem and delay between stimulus repetitions”

Working-memory capacity protects model-based learning from stress

A. Ross Otto^{a,1}, Candace M. Raio^b, Alice Chiang^b, Elizabeth A. Phelps^{a,b,c}, and Nathaniel D. Daw^{a,b}

(In) Relation to Memory

Hippocampal Contributions to Control: The Third Way

Máté Lengyel and Peter Dayan

`{lmate, dayan}@gatsby.ucl.ac.uk`

