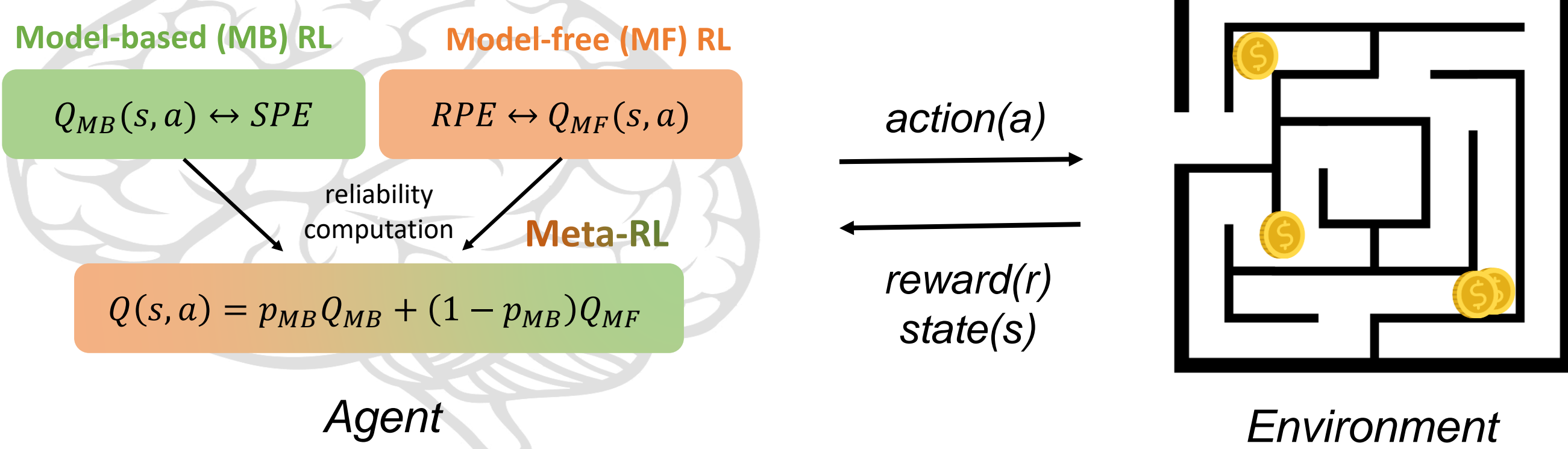


1. Introduction

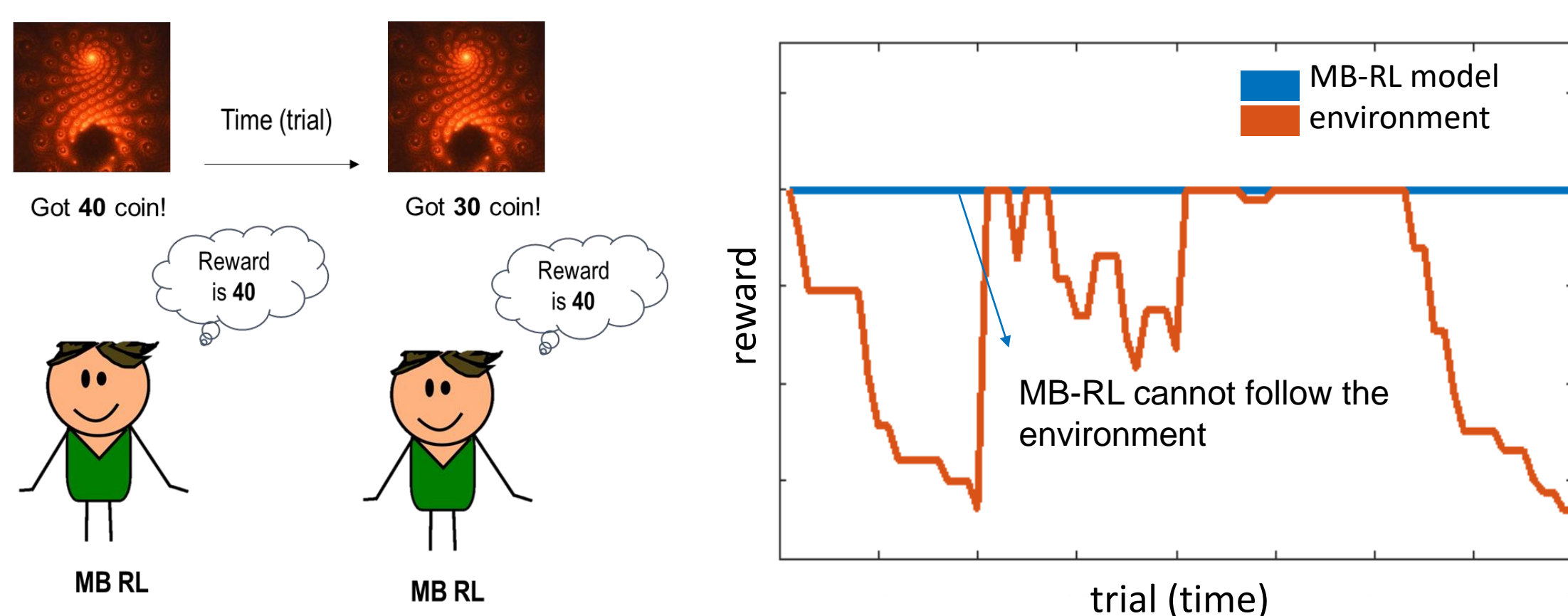
Meta-RL can explain human reinforcement learning (RL)

(Lee et al., 2014)



However, Meta-RL often fails to accommodate reward changes

Even though MB-RL possesses flexibility due to its model, it cannot follow the implicit change (e.g. rewards) in the environment

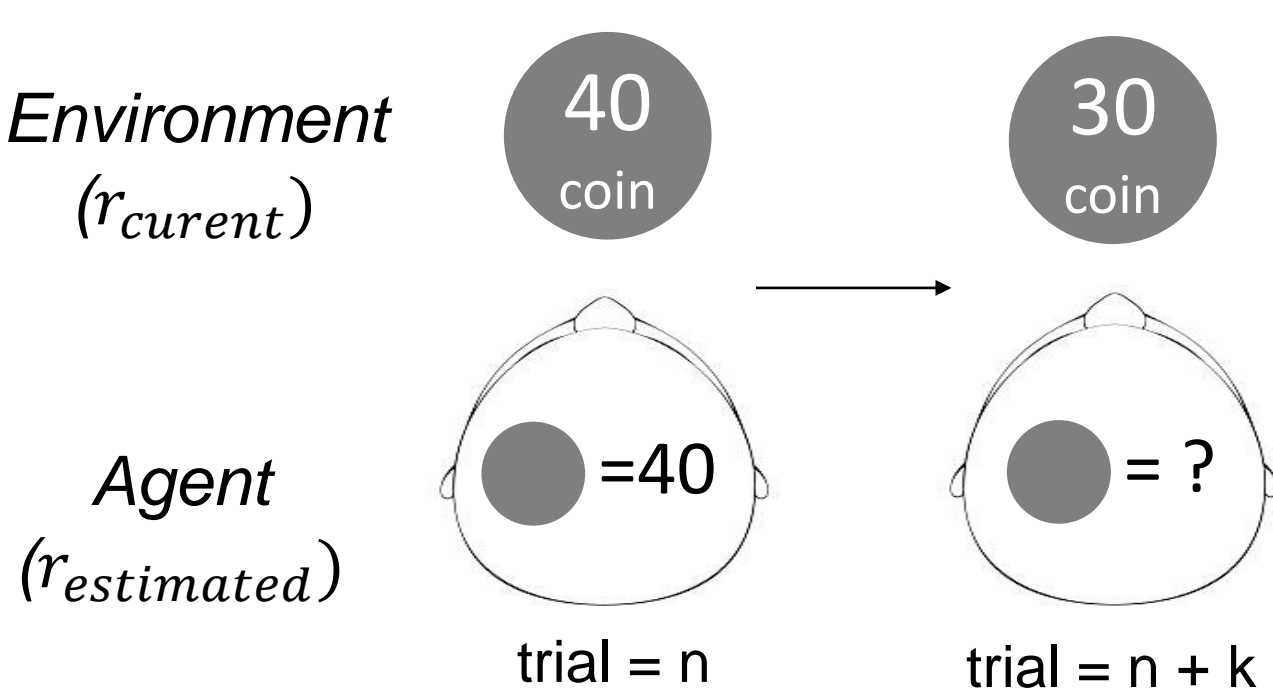


Q. Then, what strategy would a model-based system use to adapt to a dynamic environment?

2. Hypothesis Model Suggestion

Model-based system would use the temporal difference (TD) rule to update the internal reward estimation

How to update the estimation?

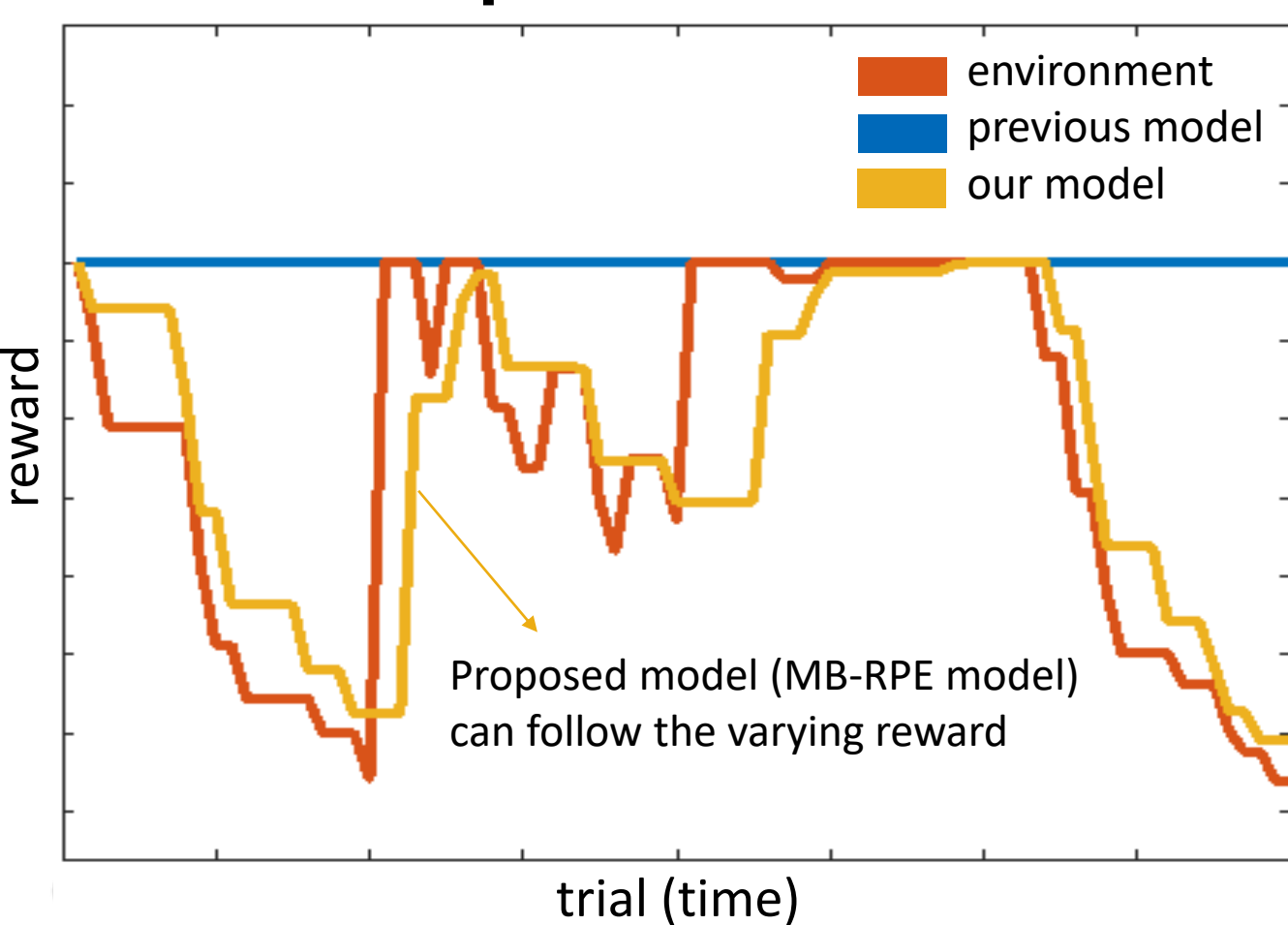


Hypothesis reward estimation

Step1. Calculate prediction error
 $\delta_{MB-RPE} = r_{estimated}(s) - r_{current}(s)$
 Step2. Update the estimated reward
 $r_{estimated}(s) = r_{estimated}(s) + \alpha \delta_{MB-RPE}$
 δ_{MB-RPE} : model-based reward prediction error (MB-RPE)
 α : learning rate / s: state

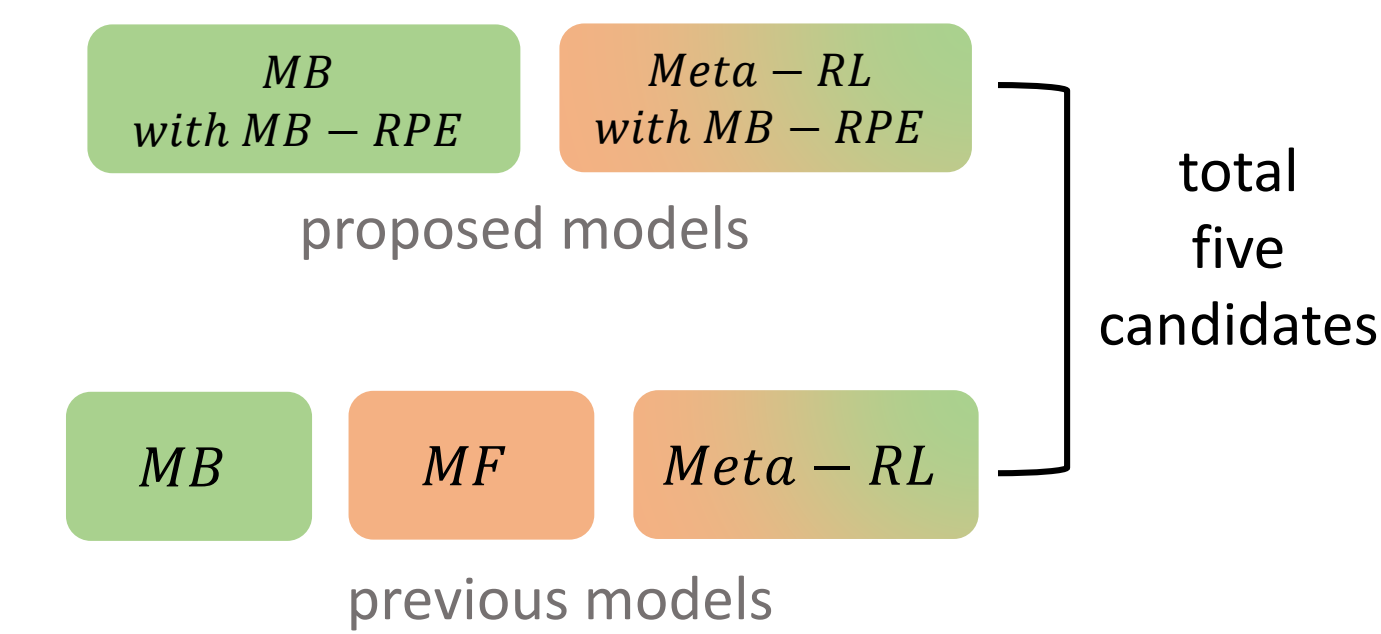
→ Using the TD rule, a model-based system would generate MB-RPE and update the estimated rewards

Our model explains subject's reward acquisition



Model comparison

Which model best explains human behavior patterns?

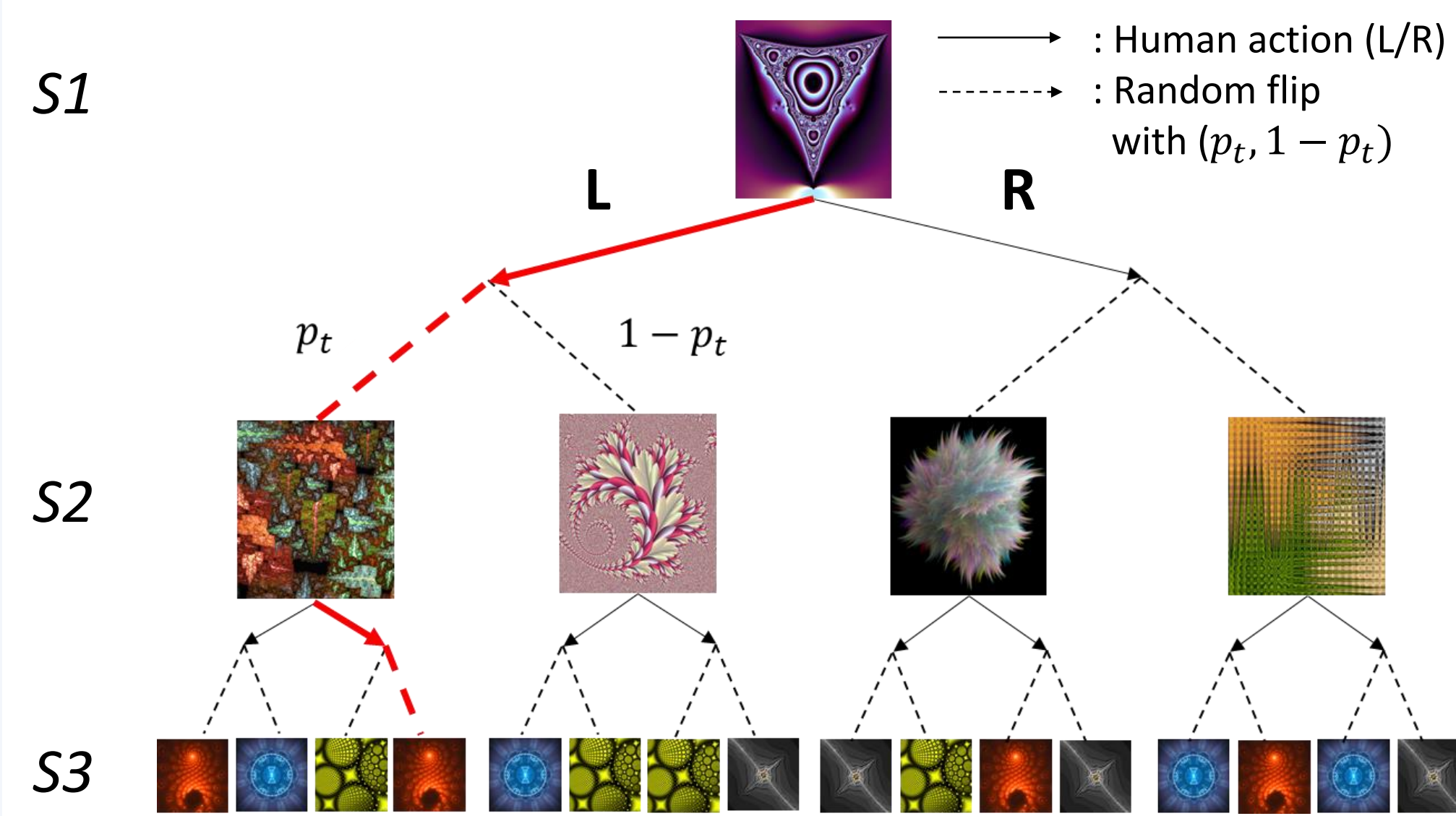


- MB = model-based / MF = model-free
- RL = reinforcement learning
- RPE = reward prediction error

Research purposes

- Test whether *Meta-RL with MB-RPE* model best explains human behavior
- Find neural evidence that brains generate MB-RPE signals

Two-step Markov decision task



- Task with context $\theta_t: \{p_t, r_t(\text{red}), r_t(\text{blue}), r_t(\text{green}), r_t(\text{grey})\}$
- Subjects make two consecutive actions (L/R) to get the reward
- Next states are also determined by an internal variable (p_t)
- State transition probability (p_t) and the final rewards (r_t) → map parameters that can change over trials
- Visited state reward decays (default setting)

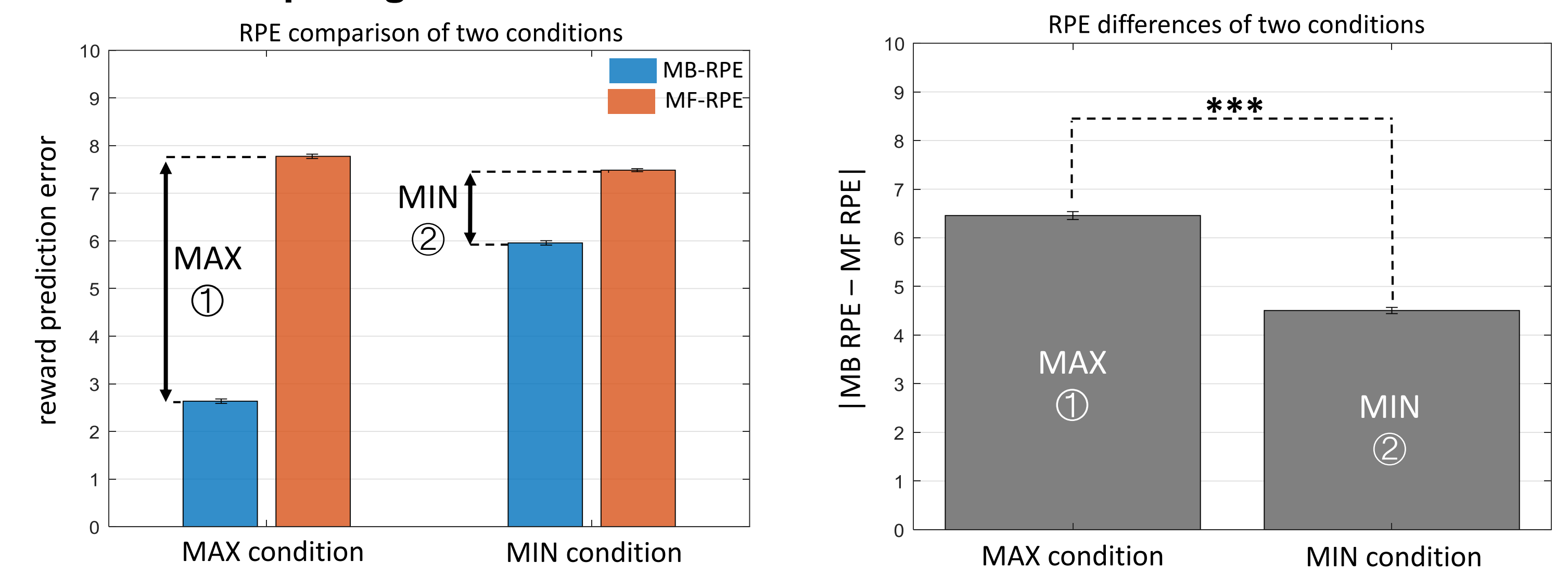
3. Task Design

Foraging rules for task design: two task conditions (MAX, MIN)

To simulate MB, MF RL either compete or cooperate, we designed two task conditions

Goals for each condition	Possible task controls per each trial	Results of task controls (20 trials)
MAX condition Maximize MB RPE - MF PRE	1: do nothing 2: shift p_t from 0.9 to 0.5 (or vice versa)	MAX condition 1 3 1 4 3 1 3 3 3 4 1 1 4 4 1 5 3 3 1
MIN condition Minimize MB RPE - MF RPE	3: reward recovery of visited state 4: reward recovery of unvisited state	MIN condition 1 1 1 1 1 1 1 1 1 1 2 1 2 1 2 2 2 1 2

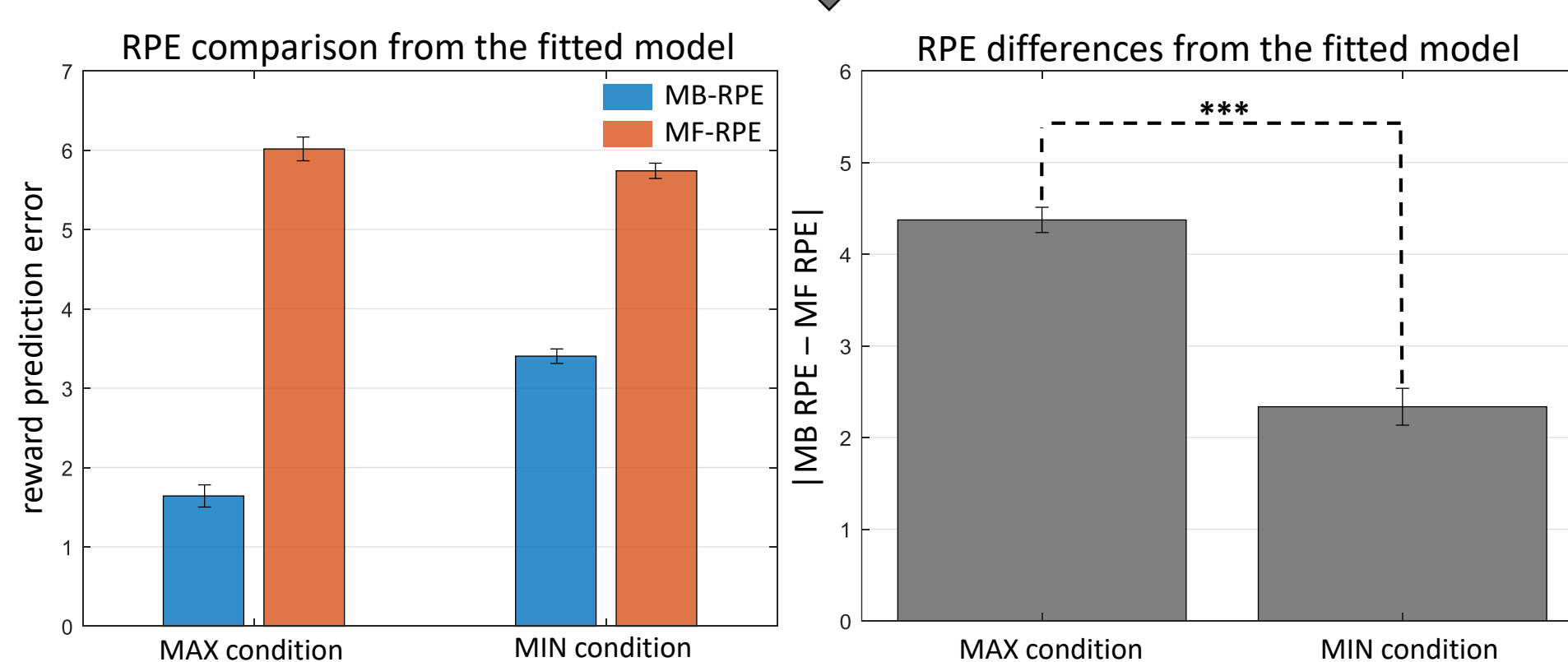
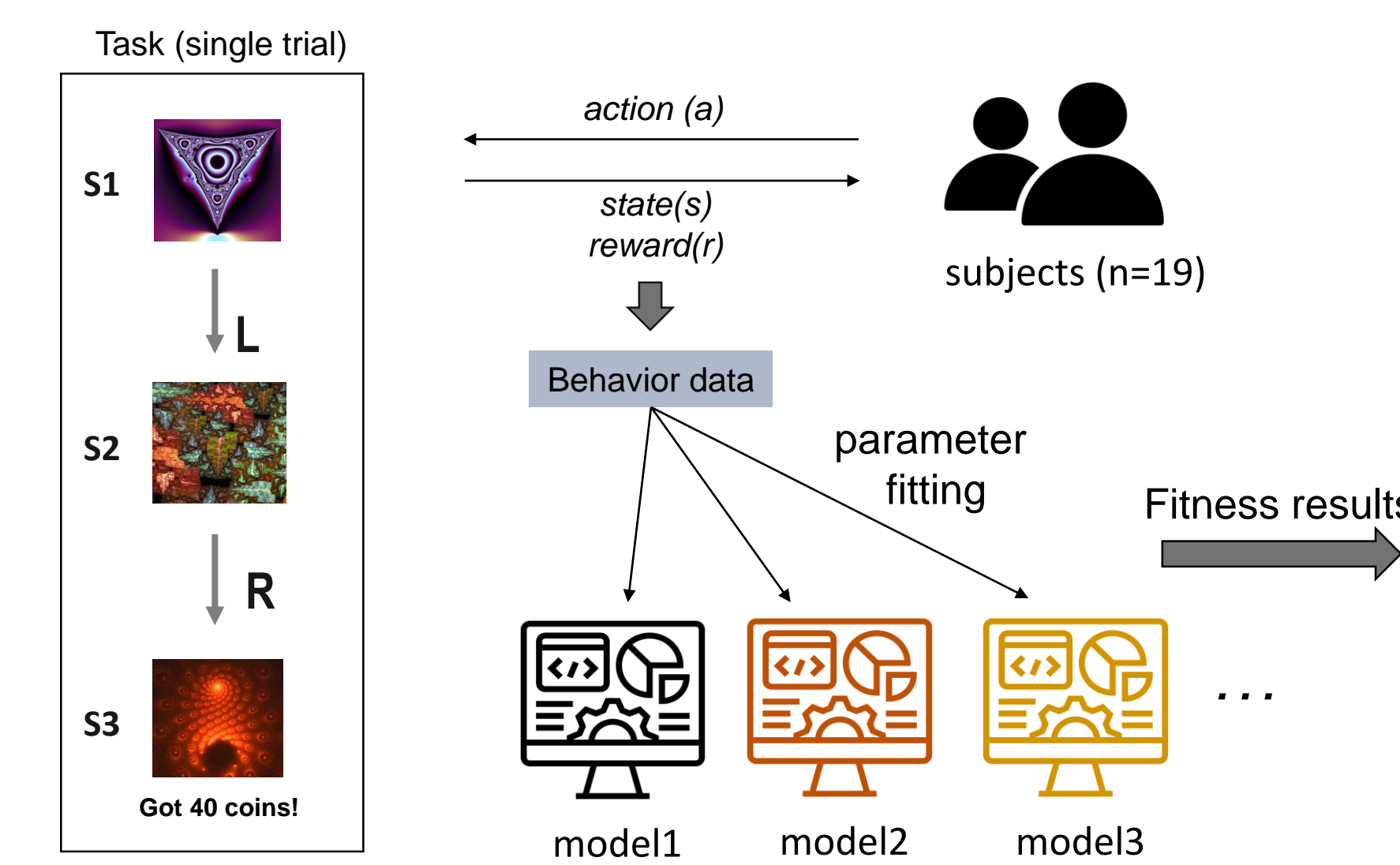
RPEs comparing in MAX and MIN conditions



4. Results

Task condition validation after model-fitting

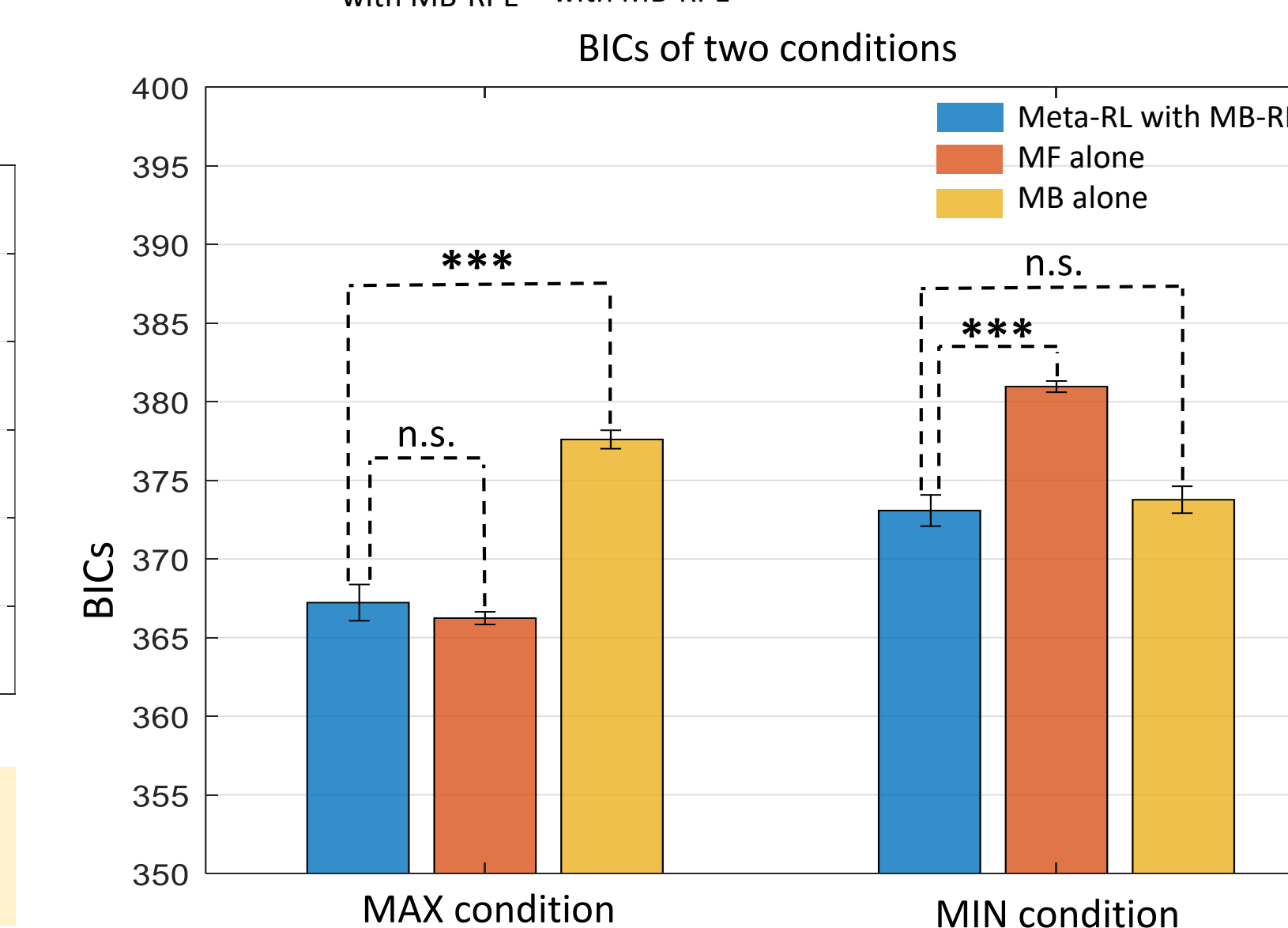
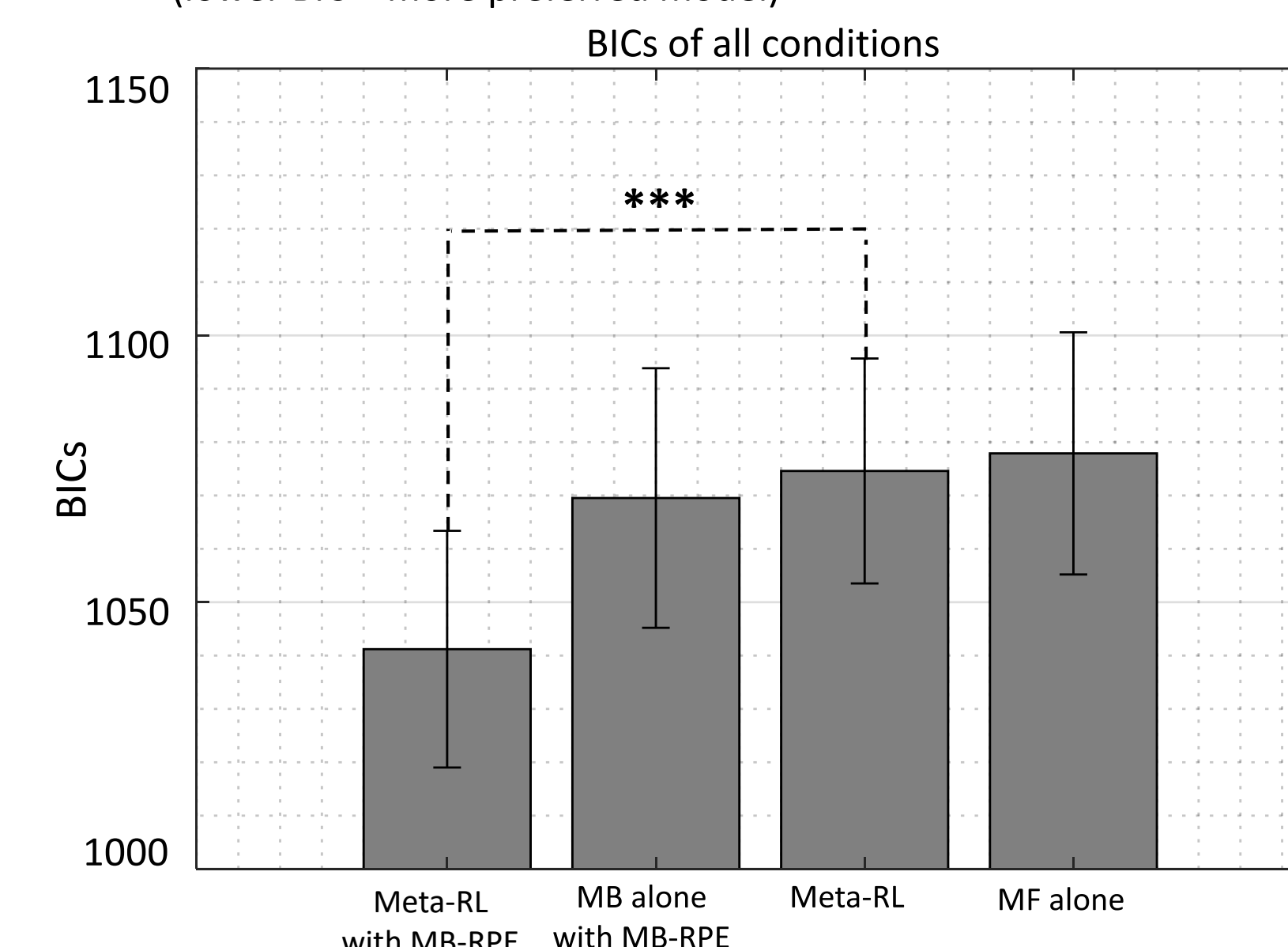
paired t-test (*: $p < 0.05$, **: $p < 1e-2$, ***: $p < 1e-3$)



Model fitting shows that two task conditions were controlled as intended initially

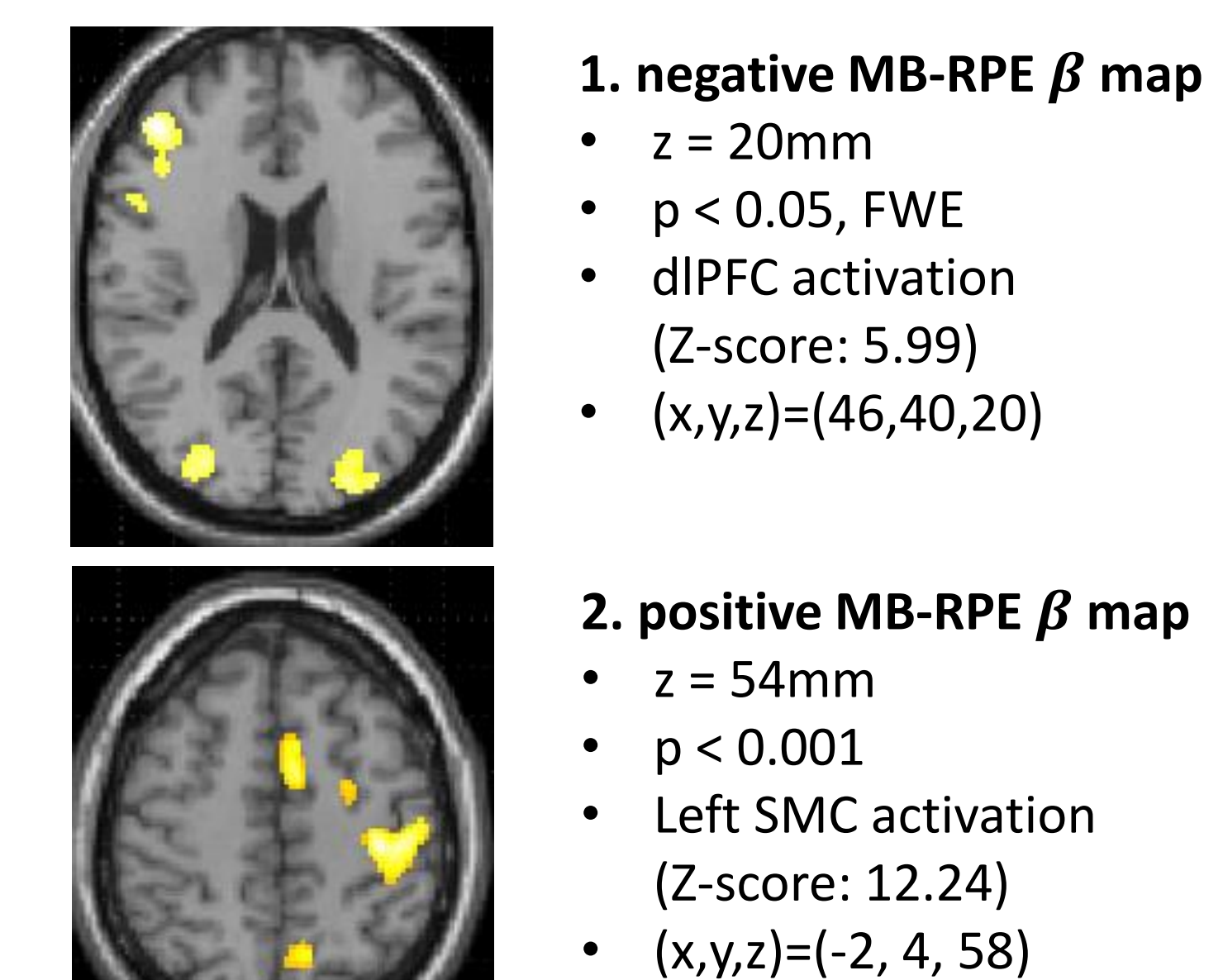
Model fitness analysis (n = 19)

BIC = Bayesian Information Criterion (lower BIC = more preferred model)



Neural data GLM analysis (n = 25)

- different subjects from behavior analysis
- 25 subjects, ten females, 23.7 ± 3.8 years



- Found neural evidence that the prefrontal cortex guides foraging
- *Meta-RL with MB-RPE* explains subjects' behavior patterns best
- *Meta-RL with MB-RPE* shows consistent accountability regardless of task conditions

5. Conclusions

- We proposed a strategy for the **model-based system to adapt to a dynamic environment** with varying rewards
- Through a simulation study, we designed **foraging tasks** with the Markov decision process with two different conditions
- From the behavior data analysis, our proposed model **best explained the human behavior data regardless of the environmental conditions**
- From fMRI data analysis, we found evidence that the **prefrontal cortex guides the foraging**