

# Outcome valence and prediction error sign invariance of the reinforcement learning models

Kamil Bonna<sup>1,2</sup>, Karolina Finc<sup>1</sup>, David Meder<sup>3</sup>, Kristoffer Madsen<sup>3</sup>, Włodzisław Duch<sup>1,2</sup>, Oliver Hulme<sup>3</sup>

<sup>1</sup> Centre for Modern Interdisciplinary Technologies, Nicolaus Copernicus University, Toruń, Poland

<sup>2</sup> Institute of Physics, Faculty of Physics, Astronomy and Informatics, Nicolaus Copernicus University in Toruń, Poland

<sup>3</sup> Danish Research Centre for Magnetic Resonance, Copenhagen University, Copenhagen, Denmark



DANISH RESEARCH  
CENTRE FOR  
MAGNETIC RESONANCE

## Probabilistic reversal learning task

During the fMRI scanning session participant carried out the PRL task in two conditions: reward-seeking and punishment-avoiding. Participants were instructed to repeatedly choose between yellow and blue boxes in order to collect as many points as possible (reward condition) or loose as little points as possible (punishment condition). Participants had to learn reward probabilities from experience. Task parameters:

- box probabilities were set to 0.8 and 0.2
- reward contingency changed 4 times throughout the task
- single run consisted of 110 trials

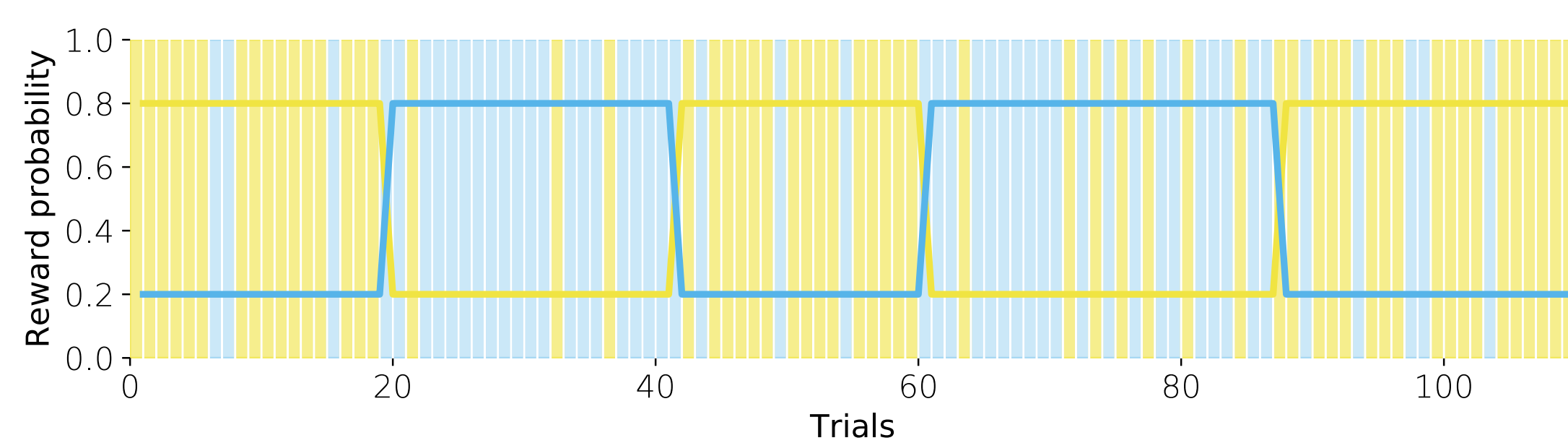
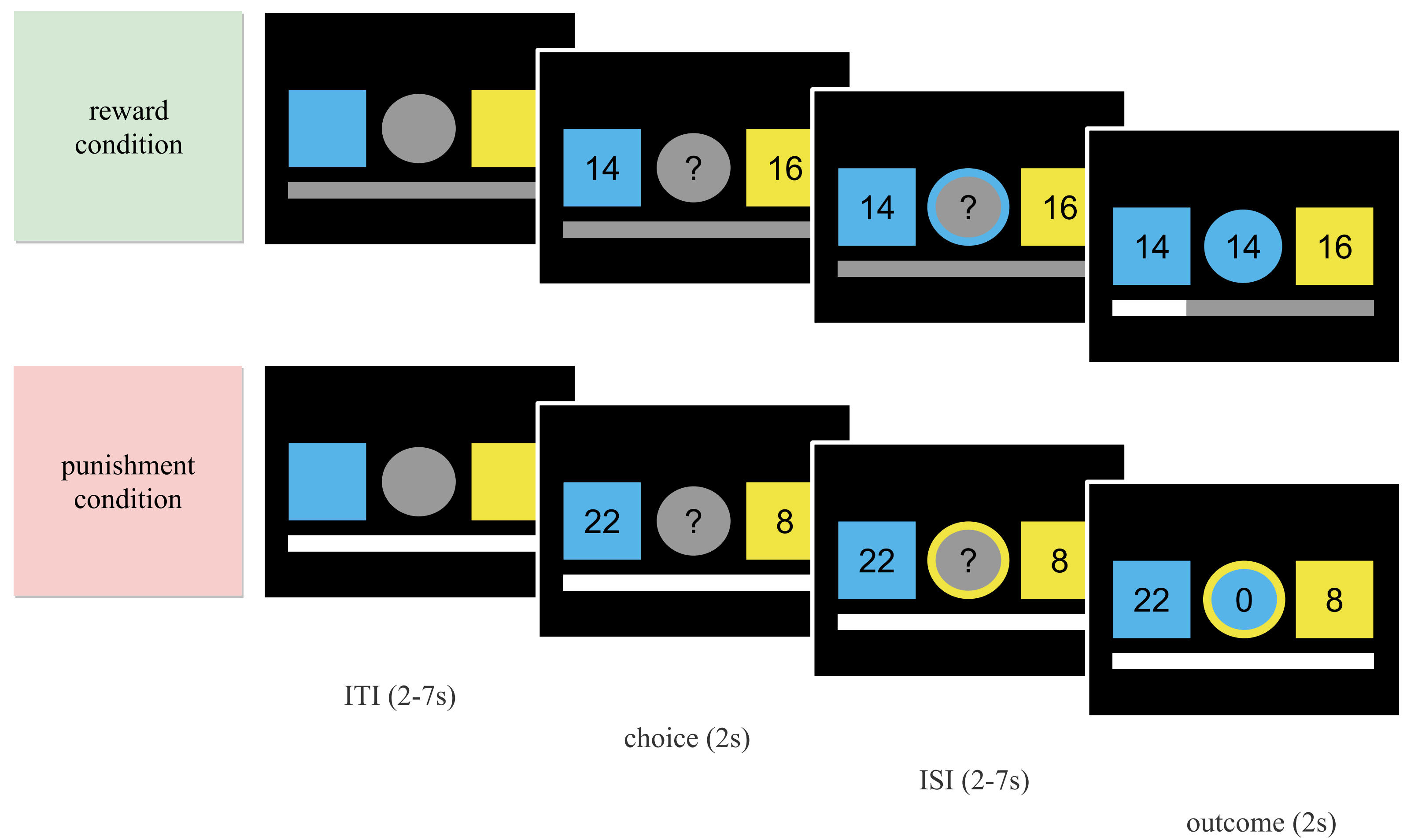
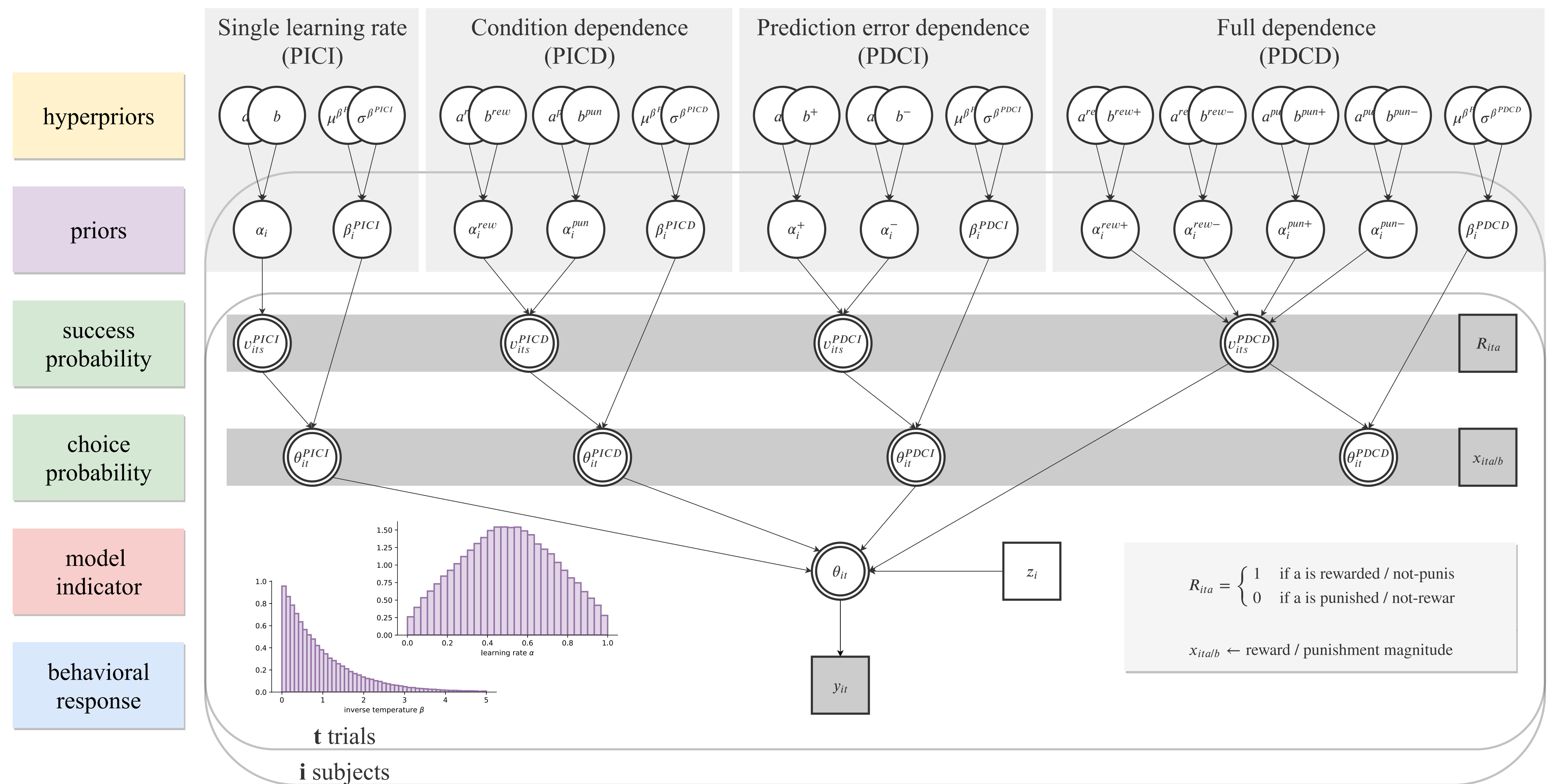
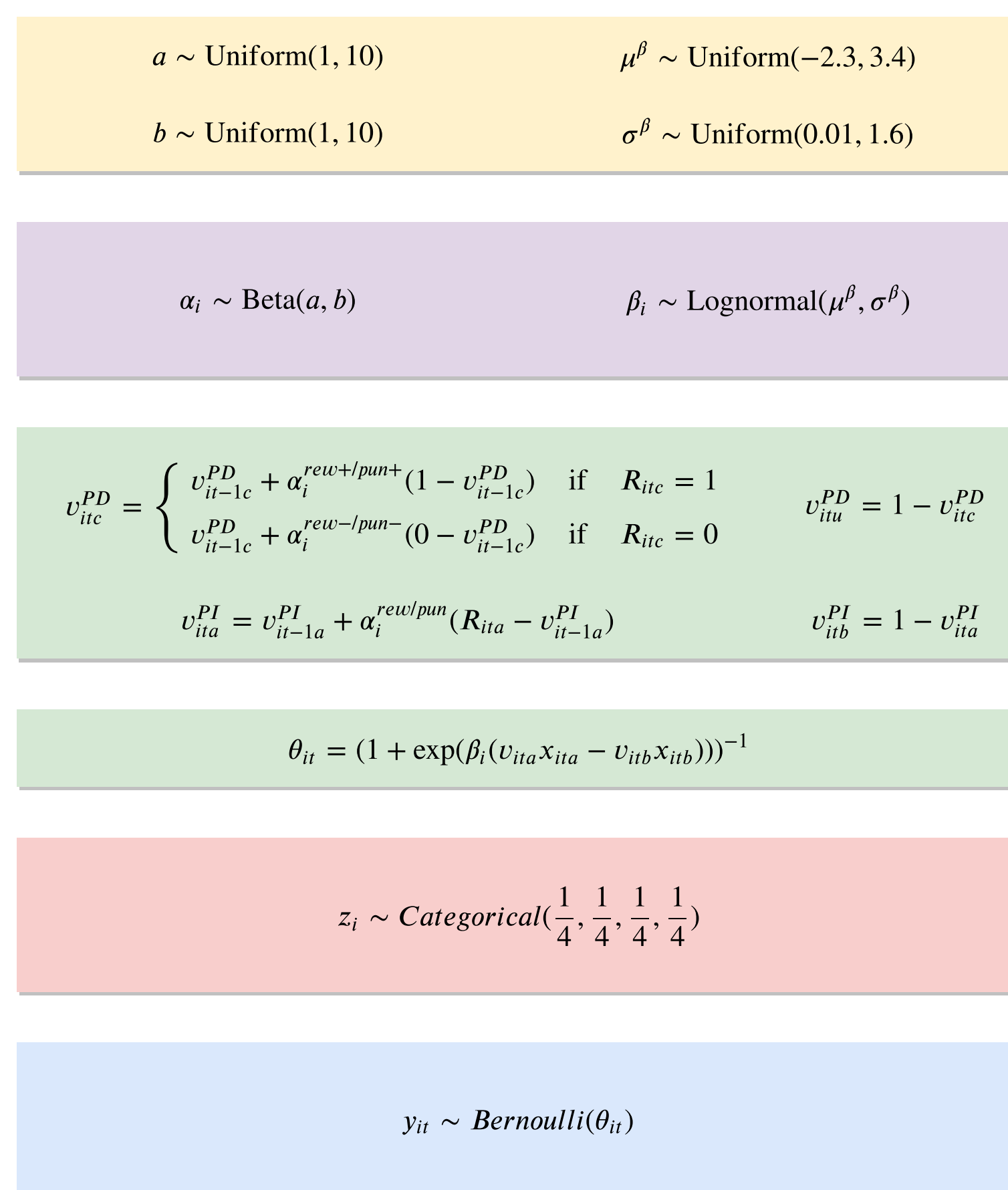


Figure 1. Example probabilistic reversal learning task structure.



## Hierarchical latent-mixture model



## Model selection

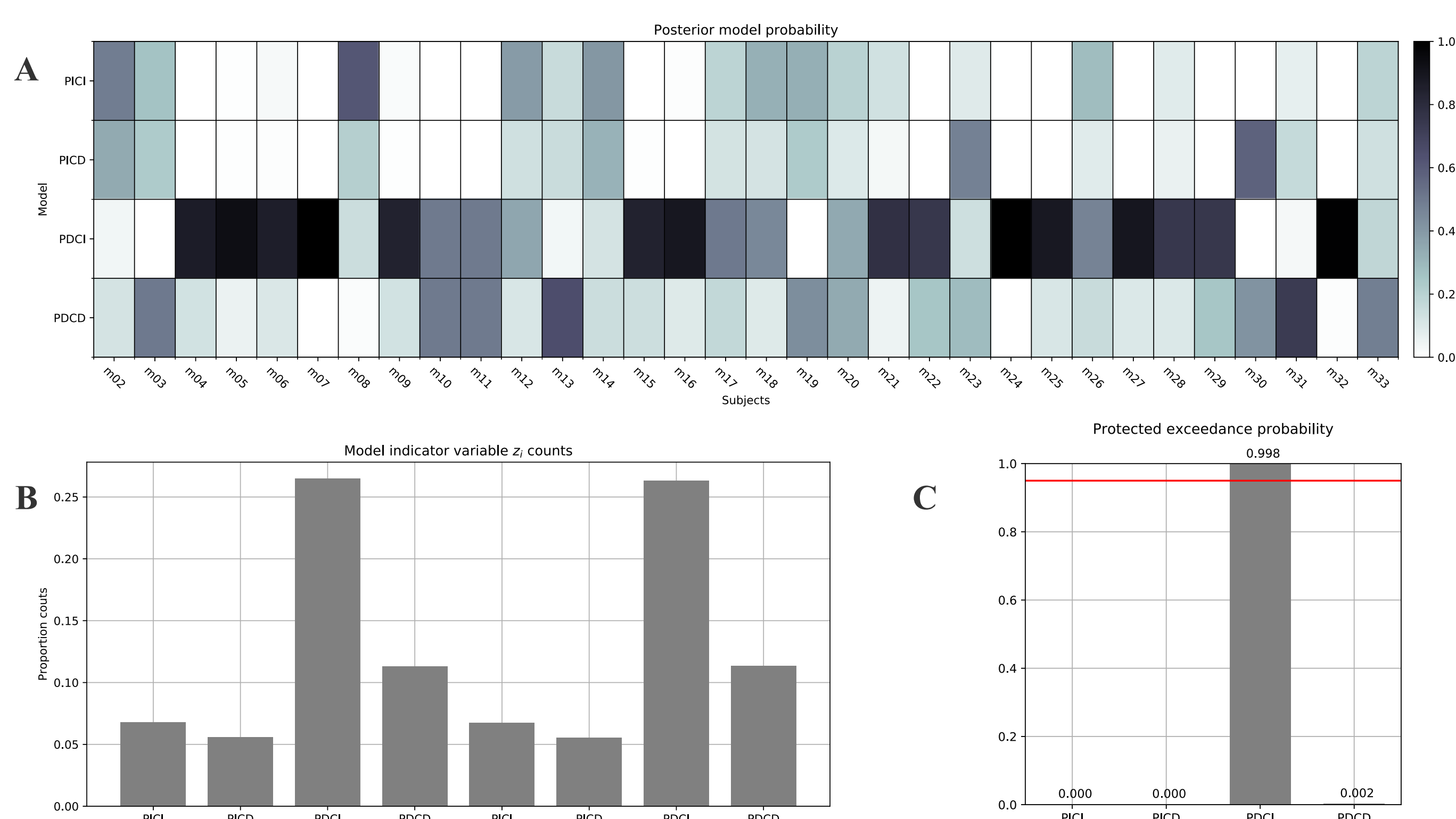


Figure 2. Model selection results. Posterior distribution of model indicator variable for (A) individual subjects and (B) marginalized over subjects. (C) Protected exceedance probabilities for all four model.

## Parameter recovery

We found that model with separate learning rates for positive and negative PEs (PDCI) was the best model (PEP > 95%) for the observed behavioral data. However, for some subjects other models had highest posterior model probability. We evaluated PDCI model separately to recover model parameters. For each subject, we calculated mean of the posterior for learning rates and inverse temperature. We found that:

- almost all subjects had extremely high values (close to 1) of learning rate for positive prediction error
- learning rates for negative prediction errors had moderate values
- difference between two learning rates was significantly correlated with mean number of reversals ( $r = -0.76$ ,  $p < 0.001$ )

Figure 3. Parameter recovery for the winning model. (A) Estimated learning rates. (B) Correlation between estimated parameters and reversal tendency.

