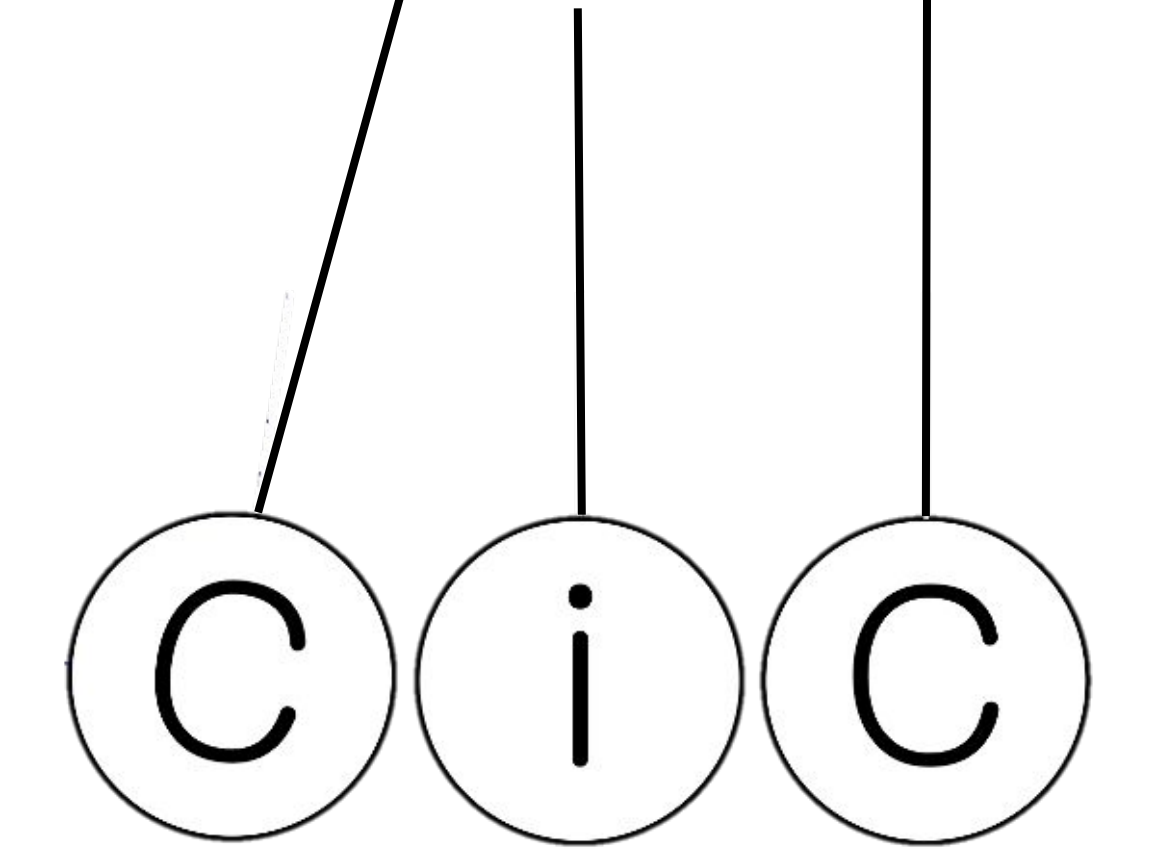




# A computational model of responsibility judgments from counterfactual simulations and intention inferences



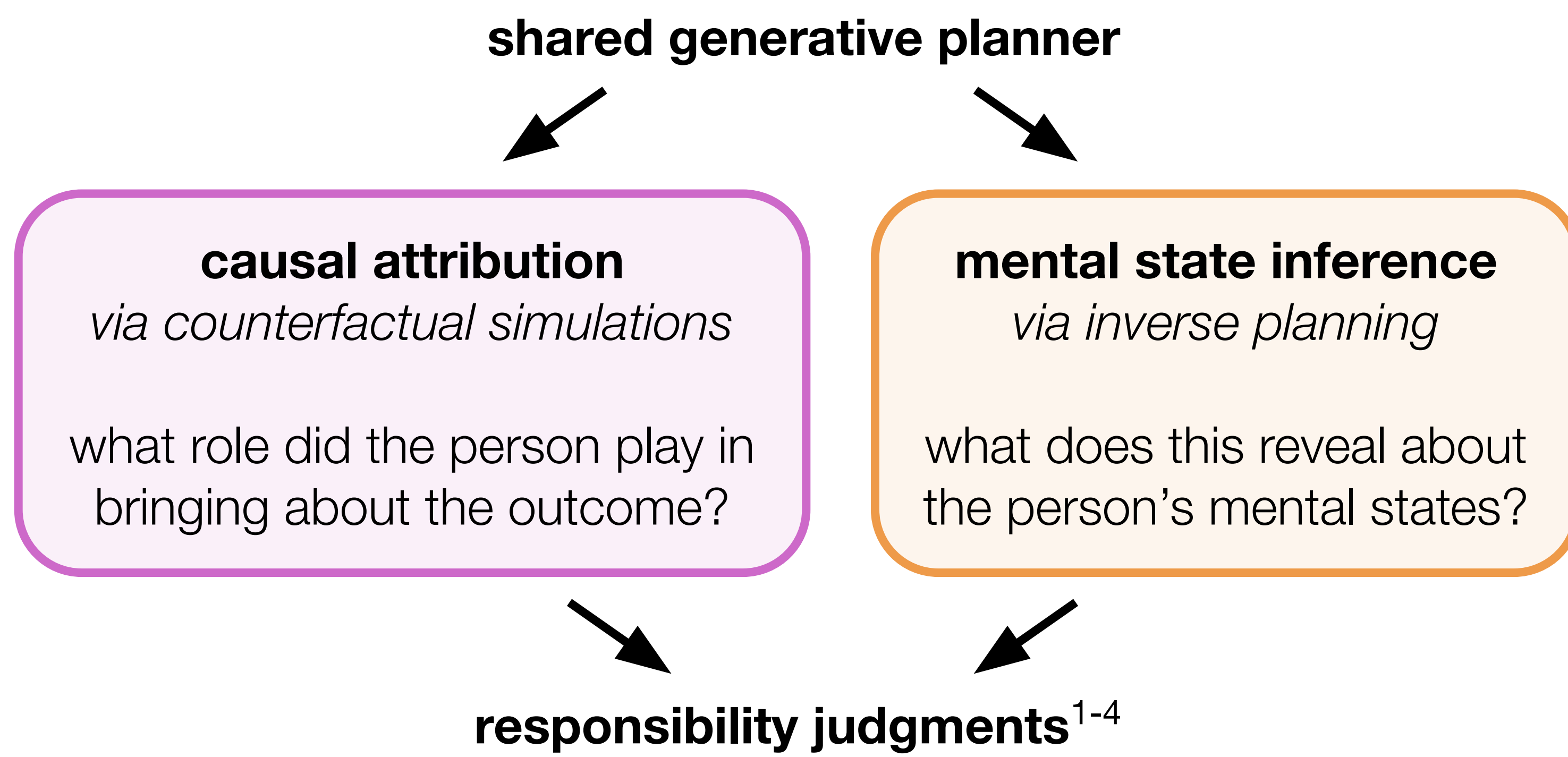
Sarah A. Wu (sarahawu@stanford.edu)<sup>1</sup>, Shruti Sridhar<sup>2</sup>, & Tobias Gerstenberg<sup>1</sup>

<sup>1</sup>Department of Psychology, Stanford University

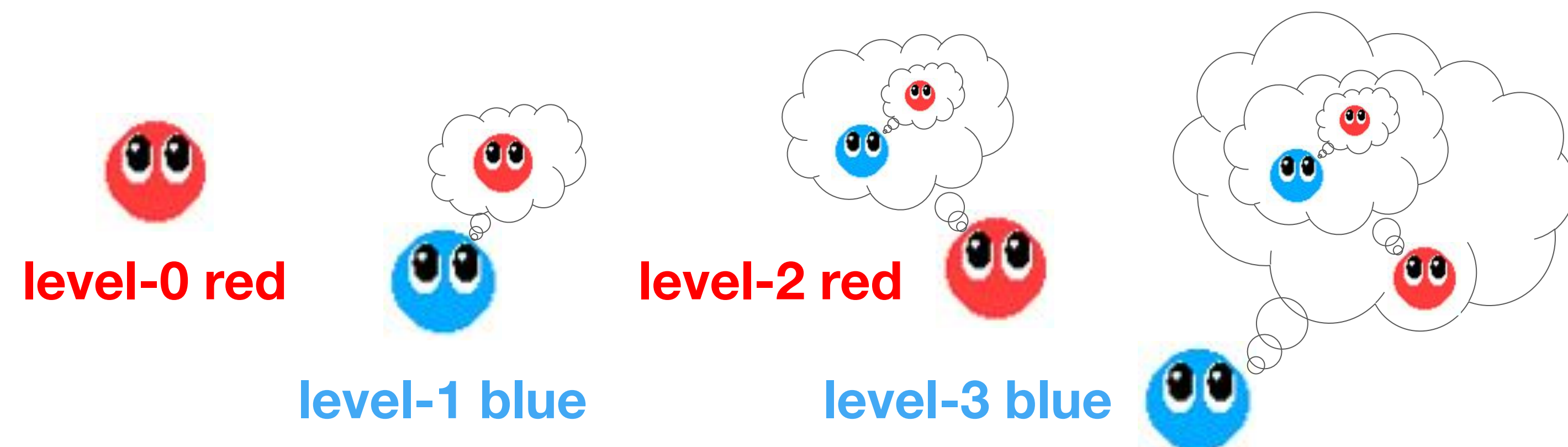
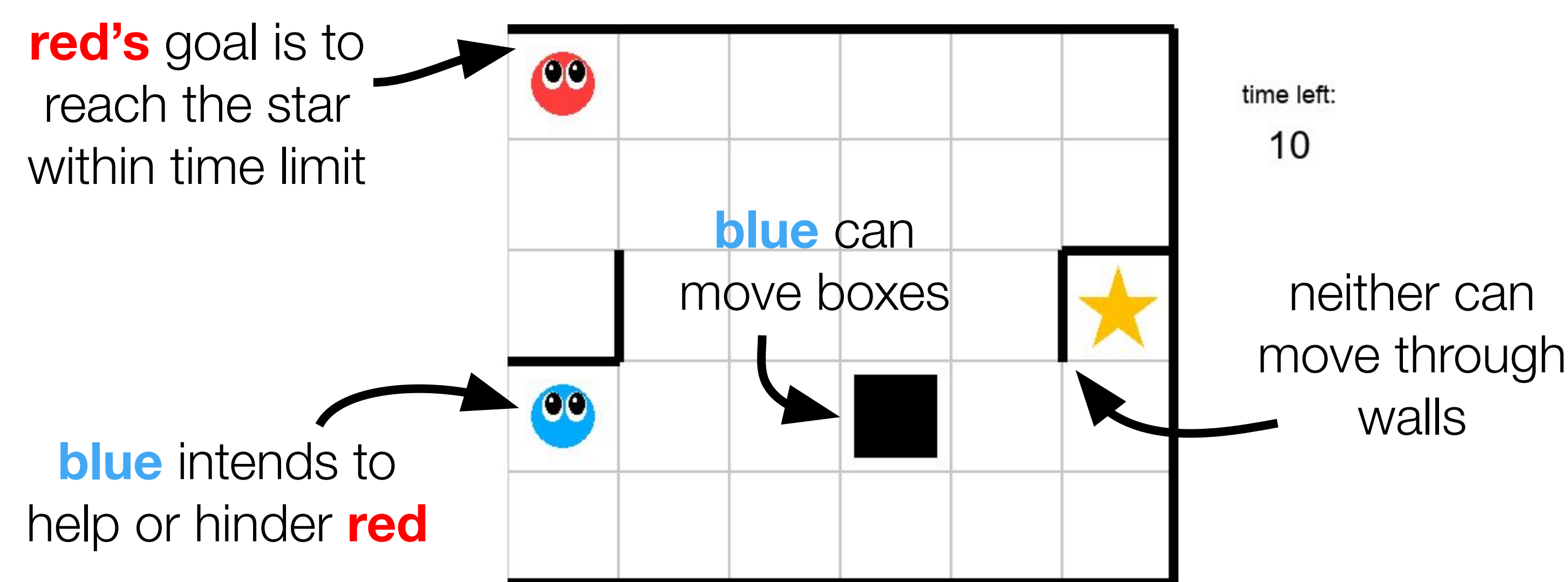
<sup>2</sup>Department of Computer Science, Stanford University

## Introduction

How do people hold others responsible in social interactions?



## Model



Environments formalized as Social MDPs<sup>5</sup>:

$$M_i^l = \langle S, \mathcal{A}, \mathcal{T}, \chi_i, g_i, R_i^l, \gamma \rangle$$

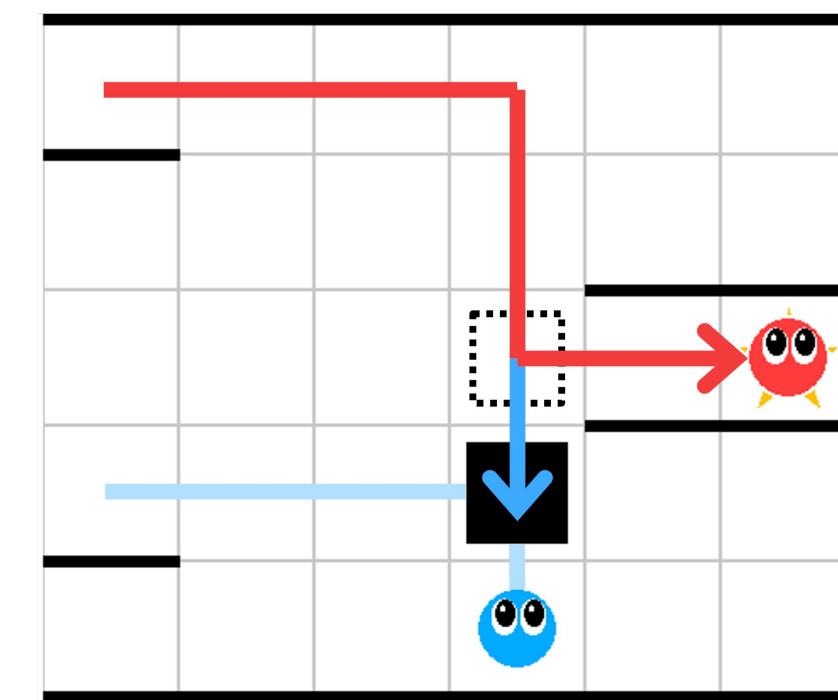
$\chi_i$  = agent  $i$ 's social goal  
 $g_i$  = agent  $i$ 's physical goal  
 $R_i^l$  =  $l$ -th level reward function for agent  $i$

**Counterfactual:** What would have happened had blue not been there?

**Mental state inference:** What was blue intending to do?

## Experiment 1

level-0 red and level-1 blue

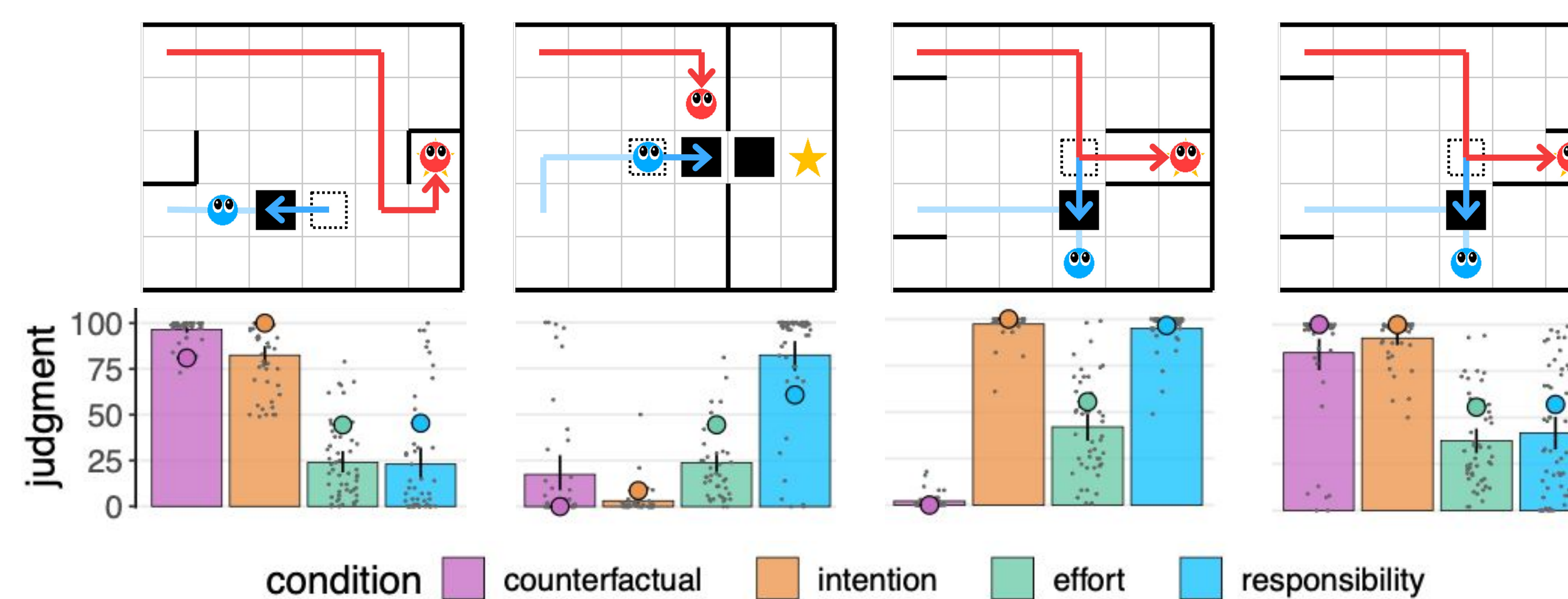


24 trials varying the actual outcome, the counterfactual outcome, and blue's intentions

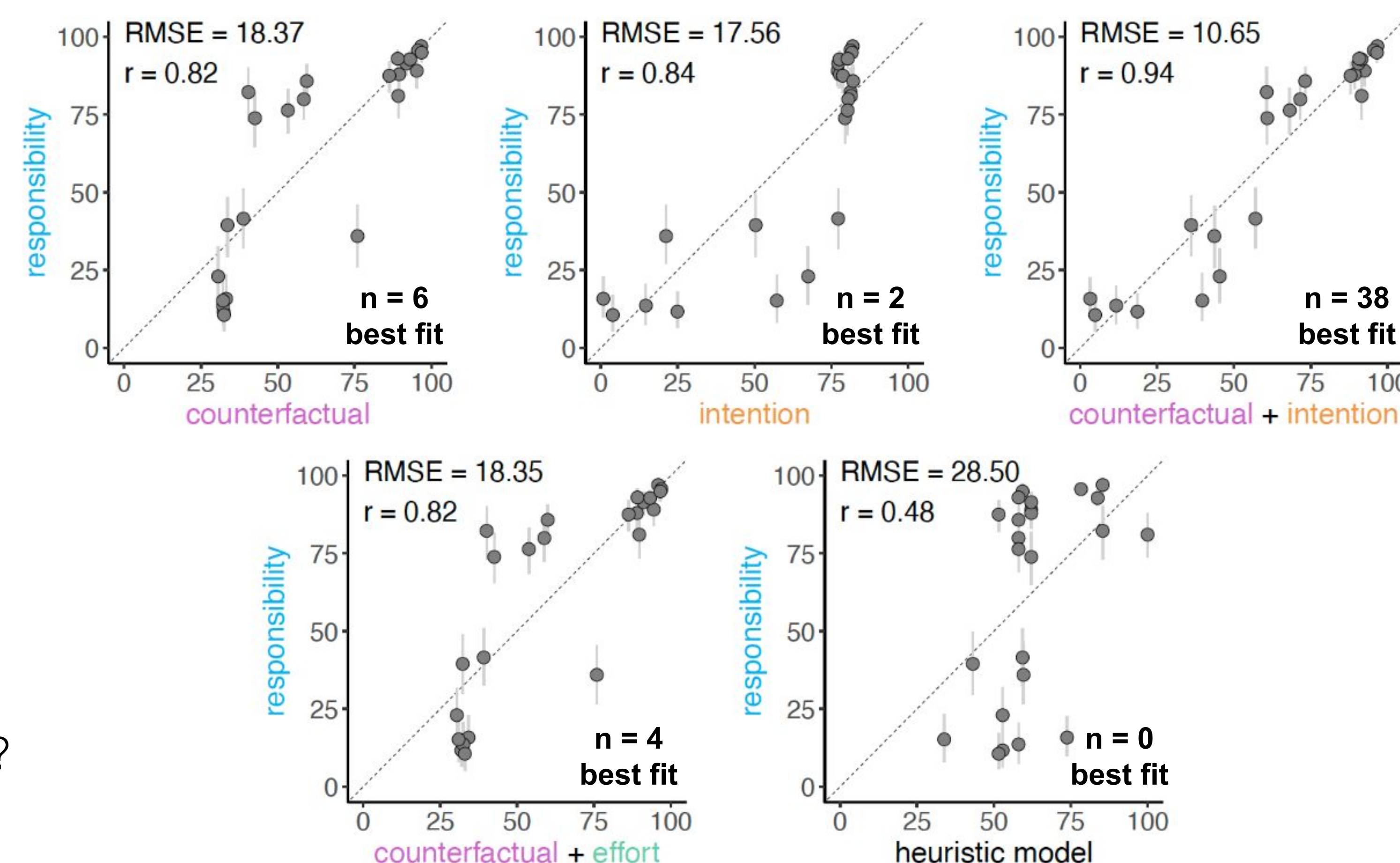
Participants in different conditions (n = 50 each) were asked:

- Counterfactual:** How much do you agree that red would (still) have succeeded if blue hadn't been there?
- Intention:** What was blue intending to do?
- Effort:** How much effort did blue exert?
- Responsibility:** How responsible was blue for red's success / failure?

Participants' judgments for select trials:

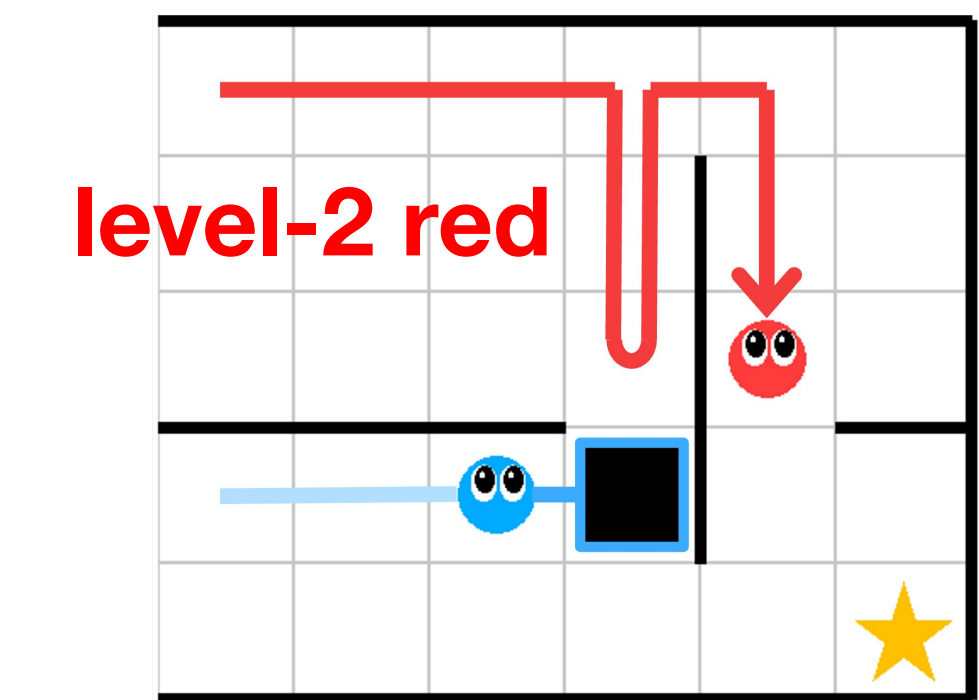
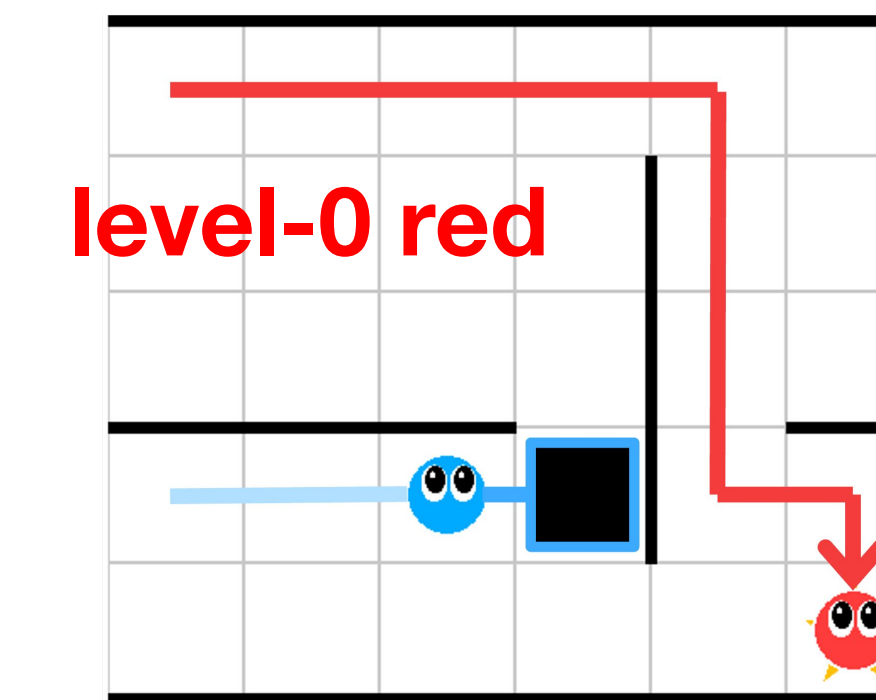


Responsibility model predictions:



## Experiment 2

includes level-2 red and level-3 blue



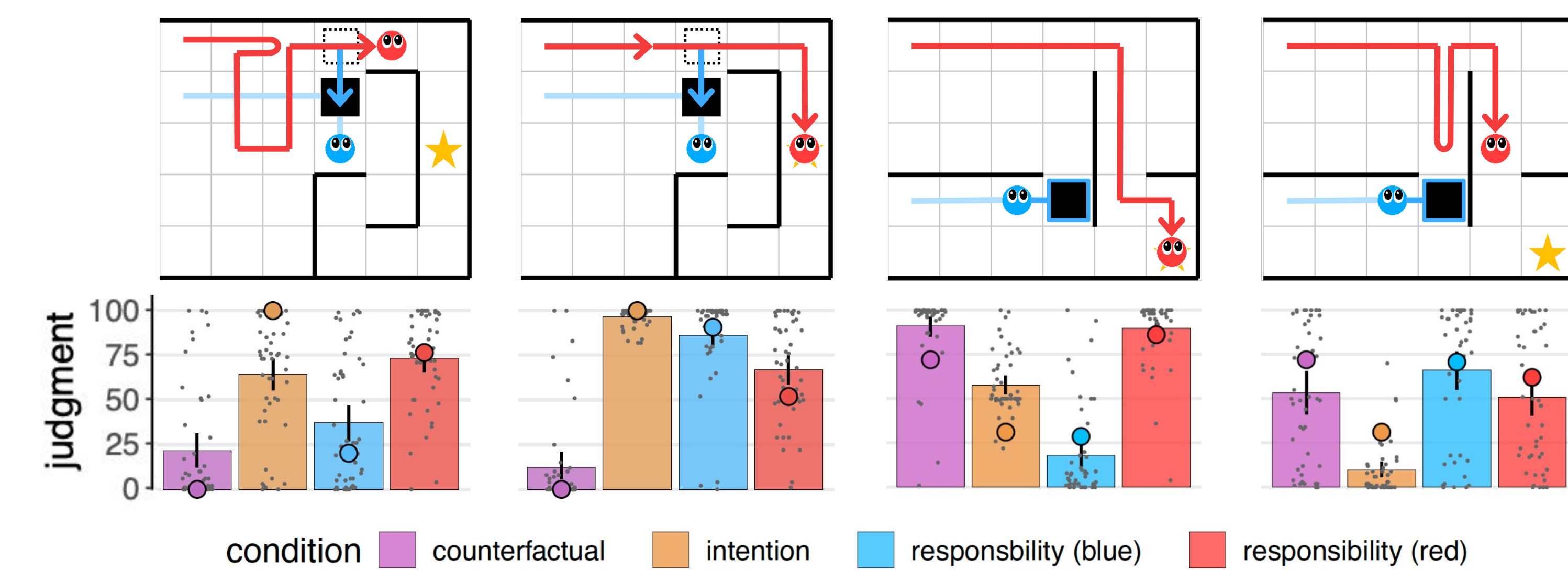
level-3 blue tricks red by appearing to be helpful, but not actually helping

12 pairs of trials differing only in whether red is level-0 or level-2

Participants in different conditions (n = 50 each) were asked:

- Counterfactual:** same as Experiment 1
- Intention:** same as Experiment 2
- Responsibility:** How responsible was blue for red's success / failure? How responsible was red for the success / failure?

Participants' judgments for select trials:



Responsibility model predictions:

- Counterfactuals + intentions model again explained responsibility judgments best (r = 0.94, lowest RMSE, n = 26/50 best fit)
- Responsibility towards blue vs. red were anti-correlated (r = -0.8)

## Discussion

Responsibility judgments are best explained by a combination of counterfactual simulations ("what would have happened otherwise?") and mental state inferences ("what was the agent intending?")

**Future work:**

- Further investigating communicative actions (signaling, deception)
- Exploring responsibility throughout repeated interactions ("fool me once, shame on you, fool me twice, shame on me!")

**References:** 1. Gerstenberg et al. (2018). *Cognition*. 2. Langenhoff et al. (2021). *Cog Psychol*. 3. Sosa et al. (2021). *Cognition*. 4. Carlson et al. (2022). *Nat Rev Psychol*. 5. Tejwani et al. (2021). *CoRL*.