

from trust in groups to trust in individuals

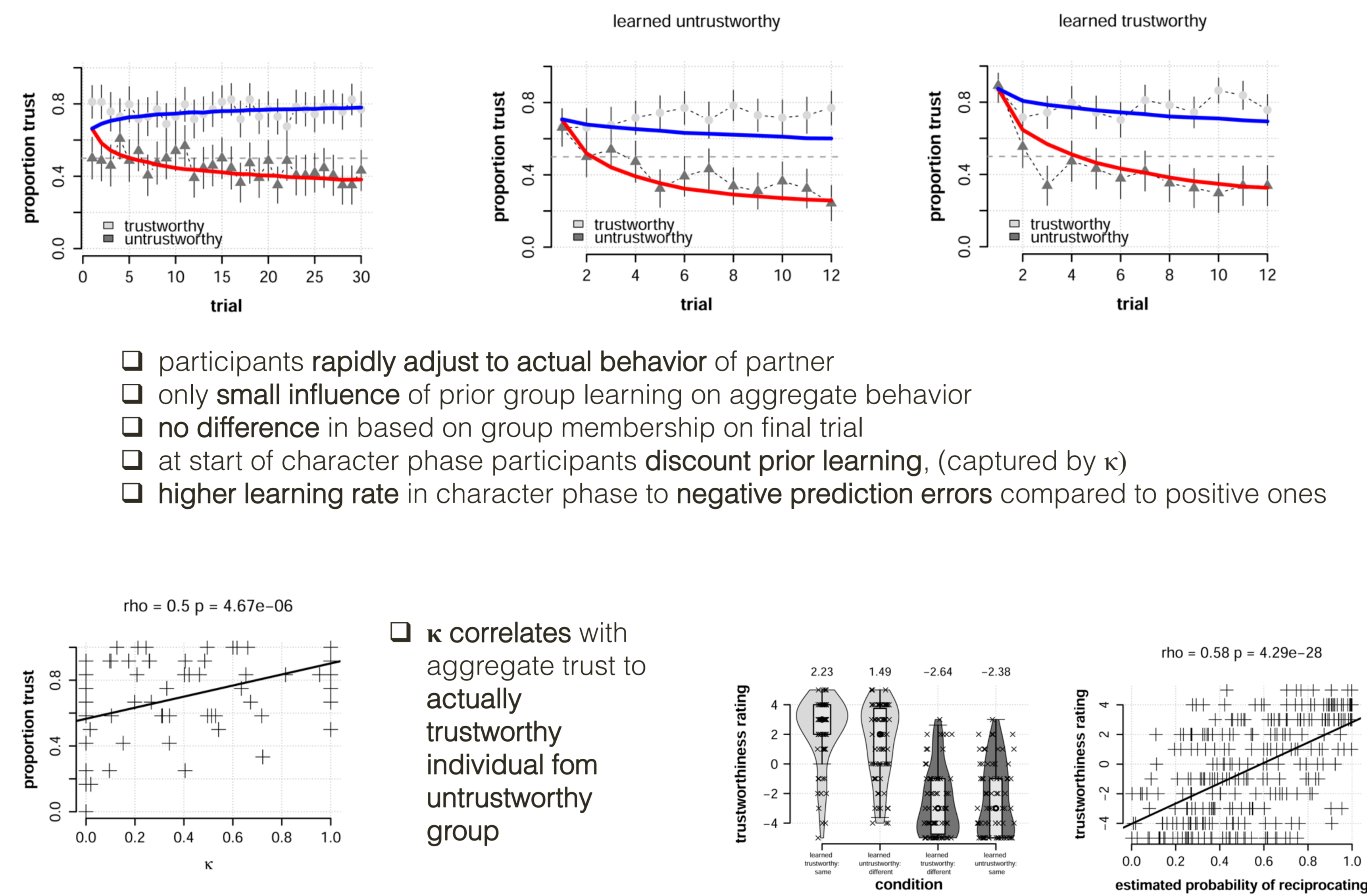
background & aim

- trust is central for social interaction.
- recent work investigated how we learn to trust other individuals by tracking their average trustworthiness using basic learning mechanisms
- however, our social world is structured by group information and little is known about how such information affects, if at all, decisions to trust

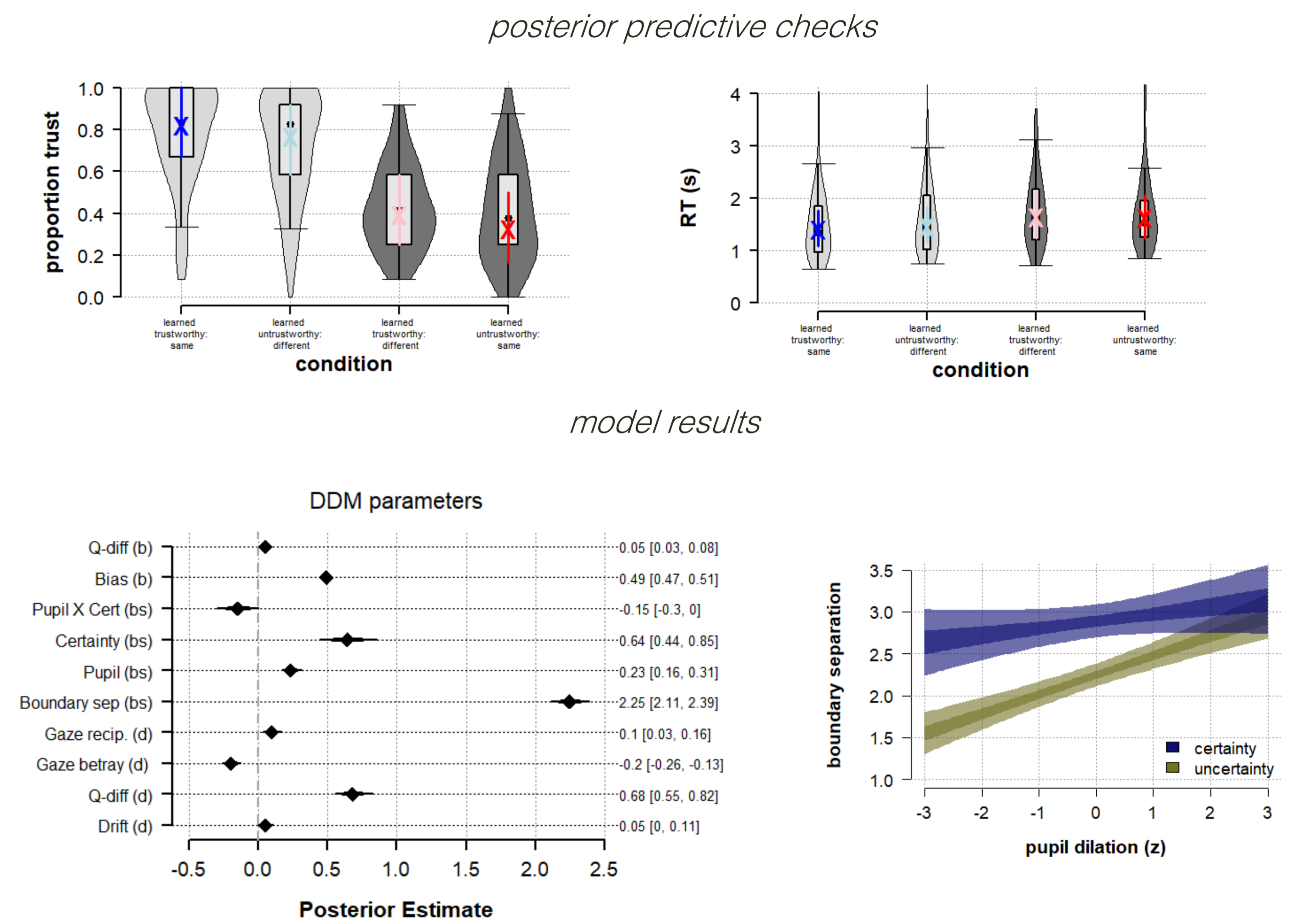
we investigated how trust in individuals is learned and unlearned in the context of minimal social groups

we aimed to characterize the dynamics both of trial-by-trial updating of trust (learning) and of within-trial variation in trust decisions

results – learning to trust



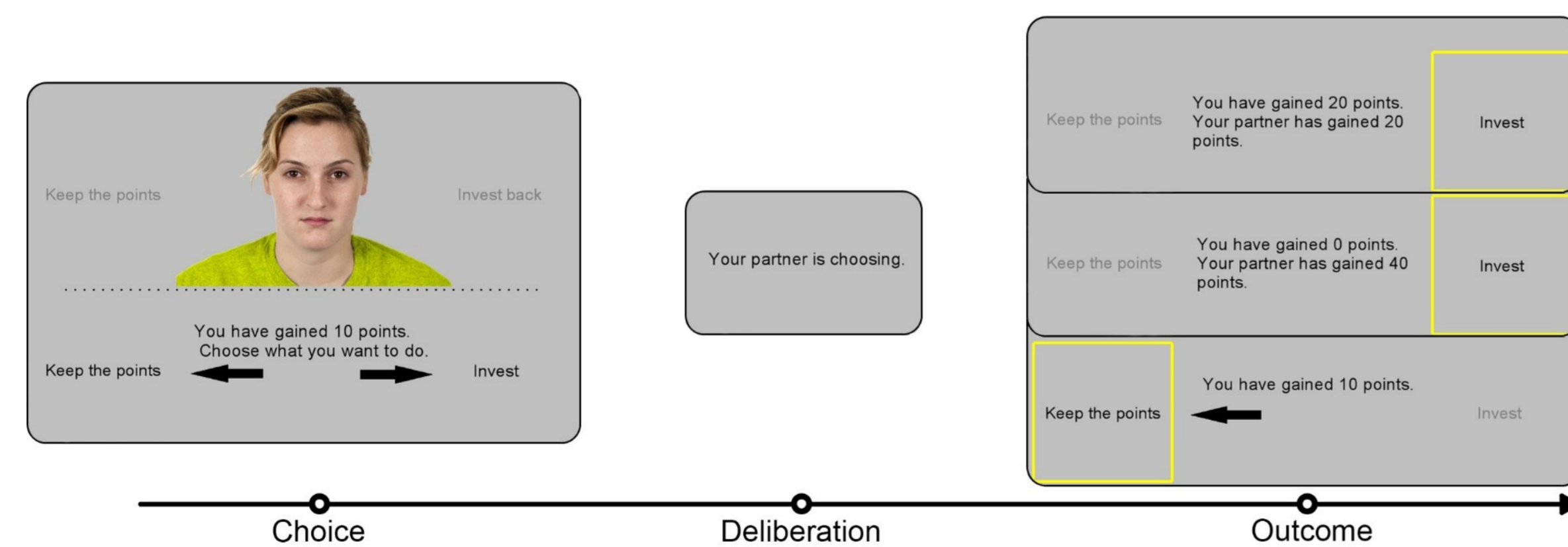
results – deciding to trust



- participants' choices & RTs can be explained by drift-diffusion model using quantities derived from RL model
- trustworthiness estimates affect both drift rates and starting points of sampling process
- looking behavior during deliberation has additive effect on drift rates
- certainty in partners' trustworthiness and pupil-linked arousal regulates boundaries of decision process

procedure

- trust games against computerized opponents
- groups identified by t-shirt color
- one group reciprocated with $P=.75$ and other group with $P=.25$
- task incentivized
- two phases of experiment (see below)



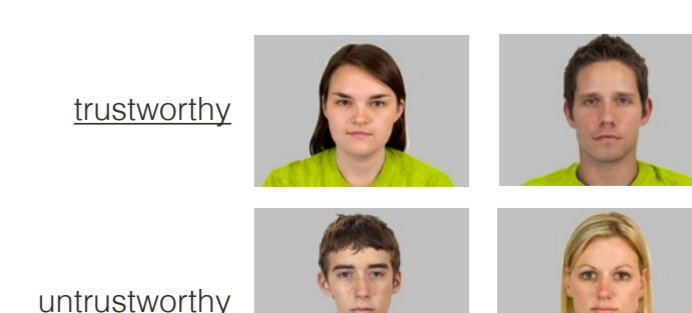
learning phase

- (learn about groups)
- multiple trust games
- trustworthiness rating of groups
- selection of group for another round



character phase

- (learn about individuals)
- repeated trust games
- trustworthiness rating of individual partners



RL model

hybrid PH-RW model

$$\text{decision} \begin{cases} EV(\text{trust}) = Pr_t * 2 + \theta \\ p(\text{trust}) = \frac{e^{EV(\text{trust})/\beta}}{e^{EV(\text{trust})/\beta} + e^{1/\beta}} \end{cases}$$

$$\text{learning} \begin{cases} d_t = R_t - Pr_t \\ Pr_{t+1} = Pr_t + \alpha A_t d_t \begin{cases} \alpha_+, & d_t > 0 \\ \alpha_-, & d_t < 0 \end{cases} \\ A_{t+1} = \lambda * |d_t| + (1 - \lambda) A_t \end{cases}$$

learning phase

$$Pr_{\text{blue}} = Pr_{\text{yellow}} = 0.5$$

$$A_{\text{blue}} = A_{\text{yellow}} = 1$$

character phase

κ : discount parameter

$$Pr_{\text{blue},1} = Pr_{\text{blue},2} = Pr_{\text{blue}} + (1 - Pr_{\text{blue}}) * \kappa$$

$$Pr_{\text{yellow},1} = Pr_{\text{yellow},2} = Pr_{\text{yellow}} + (1 - Pr_{\text{yellow}}) * \kappa$$

$$A_{\text{blue},1} = A_{\text{blue},2} = A_{\text{yellow},1} = A_{\text{yellow},2} = 1$$

- trustworthiness ratings of each partner correlate with trustworthiness estimates from RL model

conclusions

- group information (stereotypes) can be quickly overcome given individuating information in available, with minor long-term effects on aggregate behavior
- participants compute expectations of reciprocation (as predicted by RL model) and use these to guide both behavior and judgments
- combination of discounting of prior learning together with faster learning rates for negative prediction errors allow participants to avoid myopia and effectively manage learning
- participants' trust decisions depend on trial-by-trial variations in visual attention and pupil-linked arousal