

Seeing through Occlusion: Uncertainty-aware Joint Physical Tracking and Prediction

Arijit Dasgupta, Andrew D. Bolton, Vikash K. Mansinghka, Joshua B. Tenenbaum, Kevin A. Smith

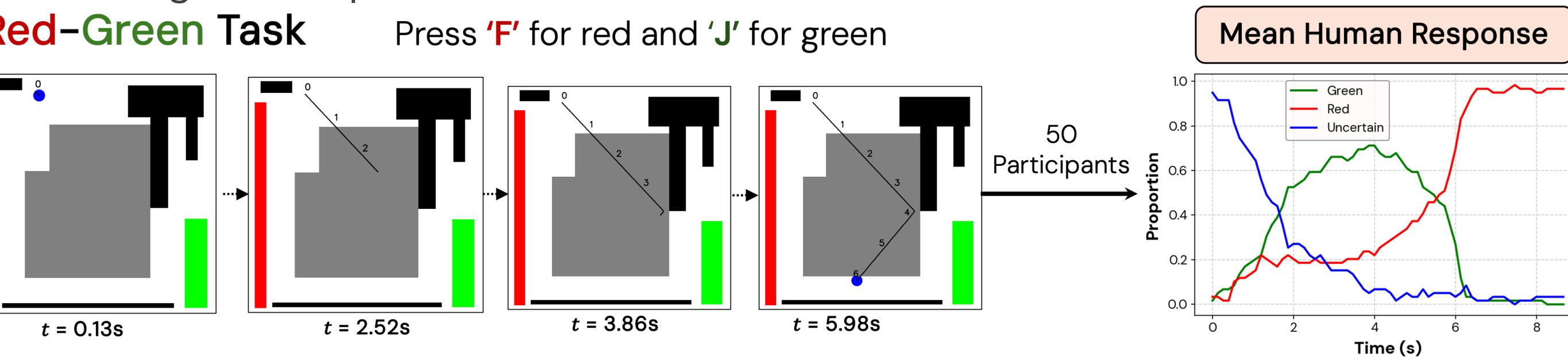
Background

- Humans can reason about hidden physical processes despite the absence of visually changing evidence.
-
- Physical tracking and future prediction have traditionally been studied **separately**
- How does **occluded motion** influence humans' future physical predictions of an object, and can this process be **computationally modeled** by integrating state estimation with prediction?

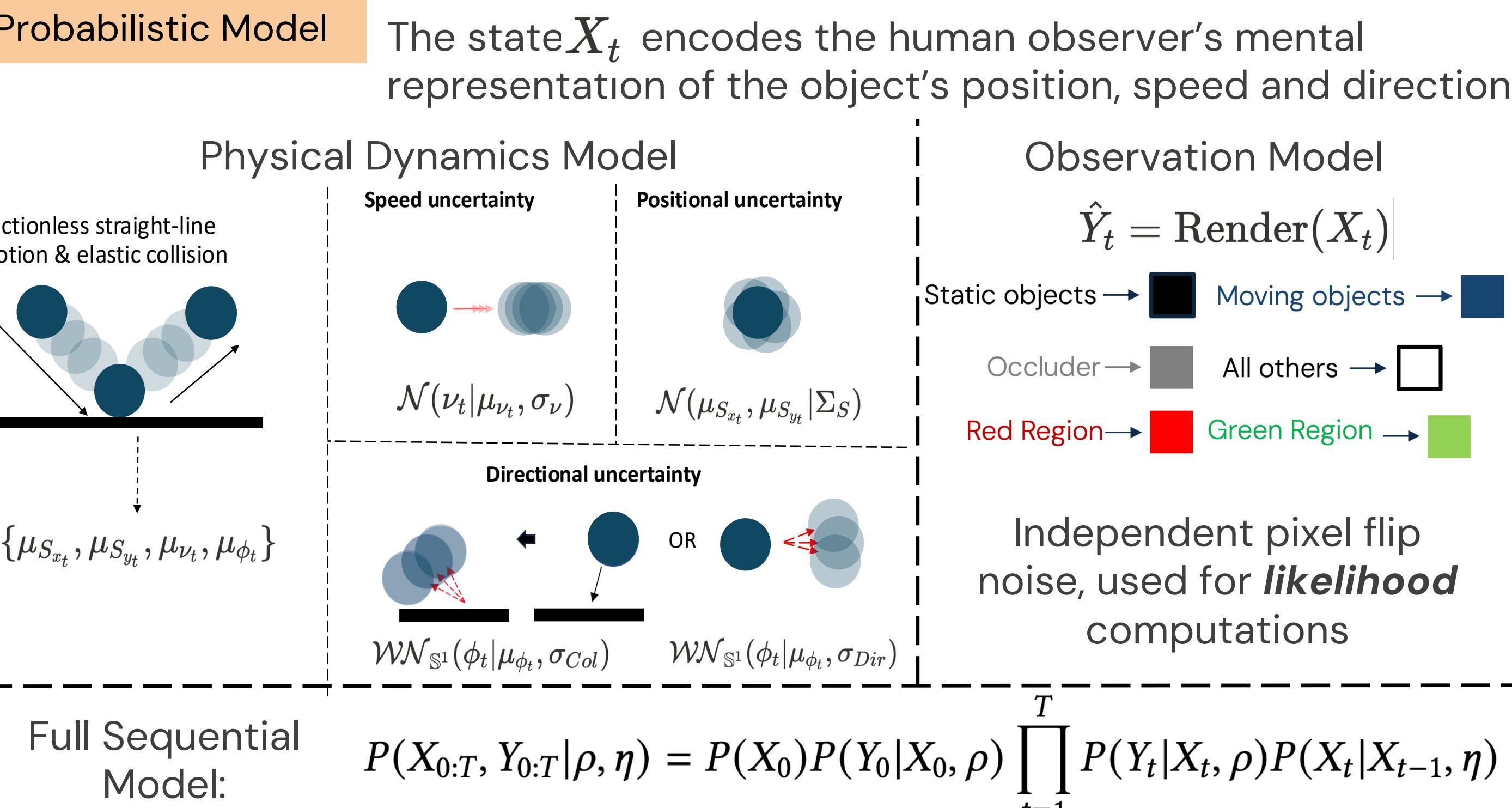
Experiment

- 59 participants recruited via Prolific (US\$ 15/hr)
- 50 video trials (34 with Occlusion, 16 without) of a ball moving in a 2.5D scene
- Objective:** continuously predict if the ball will hit **red** or **green** first
- Participants are scored between -80 to 120, penalizing mistakes and rewarding correct predictions

Red-Green Task Press 'F' for red and 'J' for green



Joint Tracking and Prediction (JTAP) Model



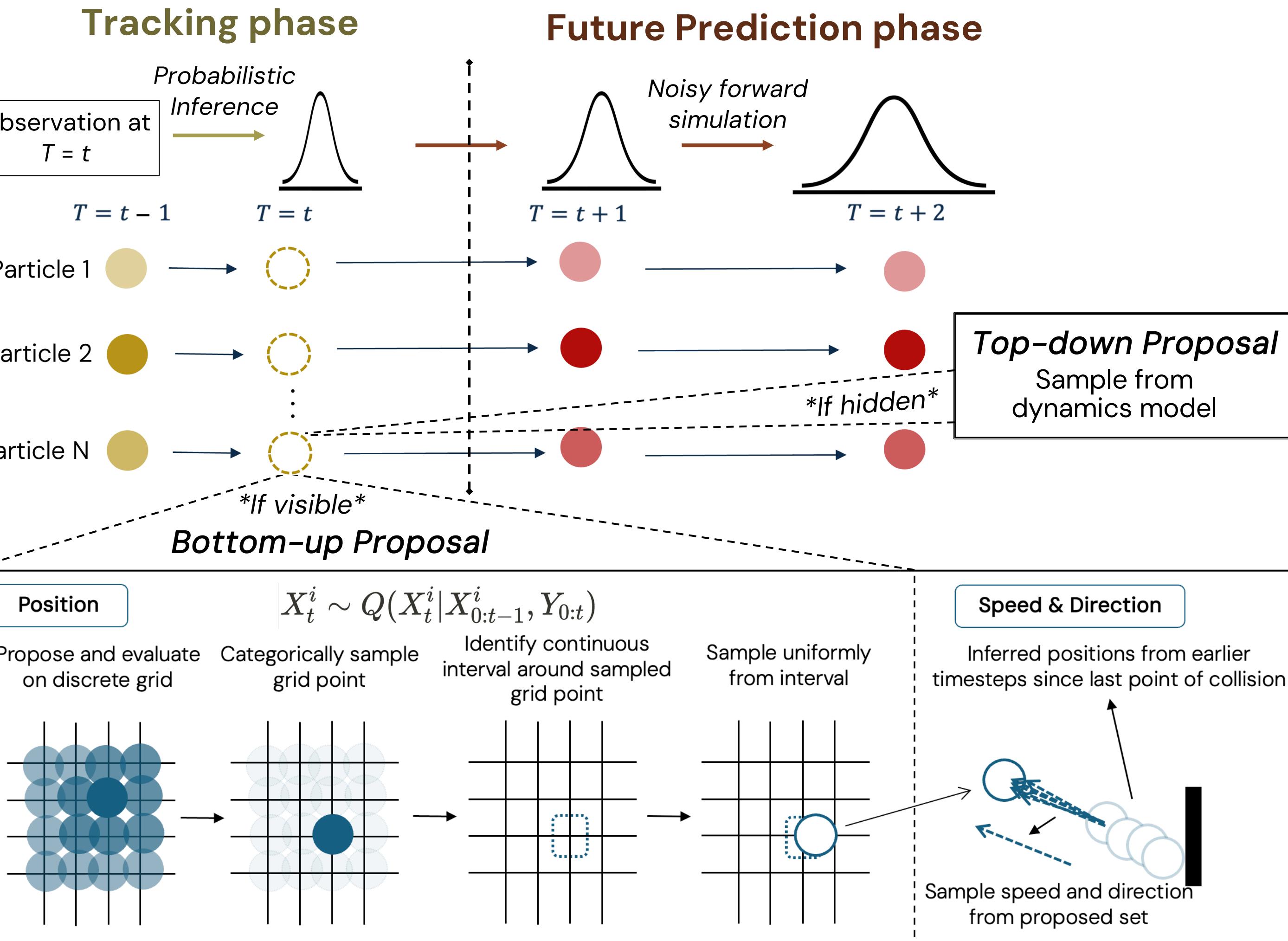
Bayesian Inference

Sequential Monte Carlo algorithm via GPU-accelerated Probabilistic Programming

Implemented in GenJAX

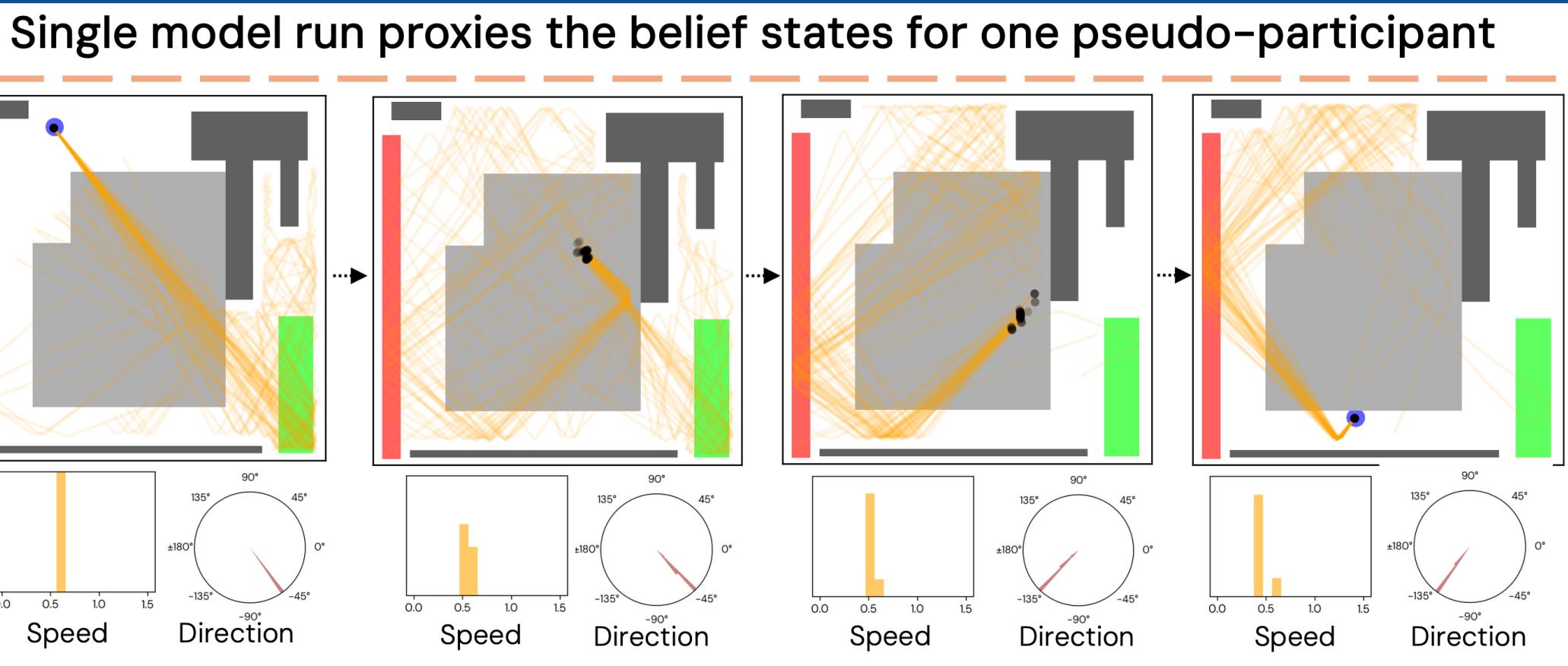
- At each timestep
1. Sample Proposal
 2. Update weights
 3. Resample **all** particles
 4. Sample dynamics model M timesteps

Tracking phase



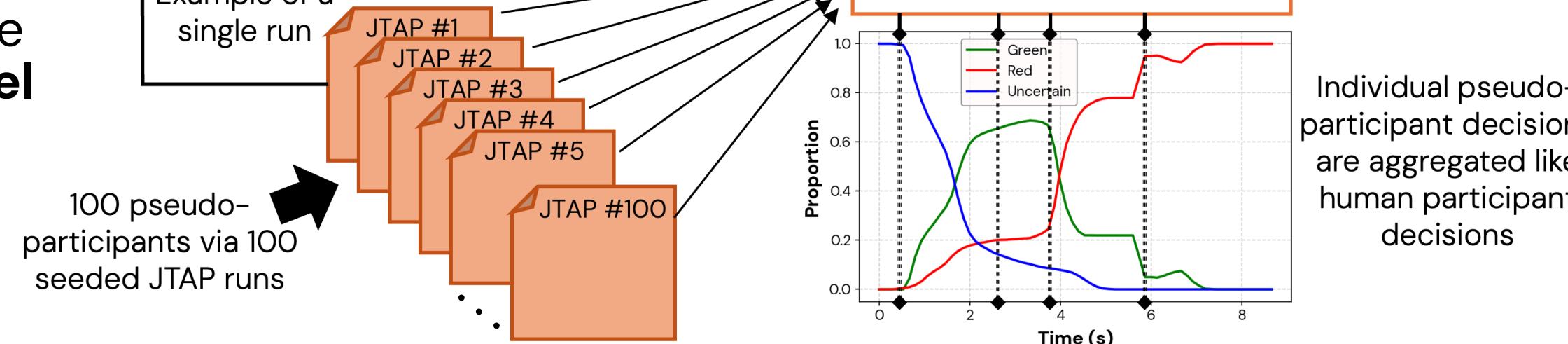
Results

Individual Model Runs

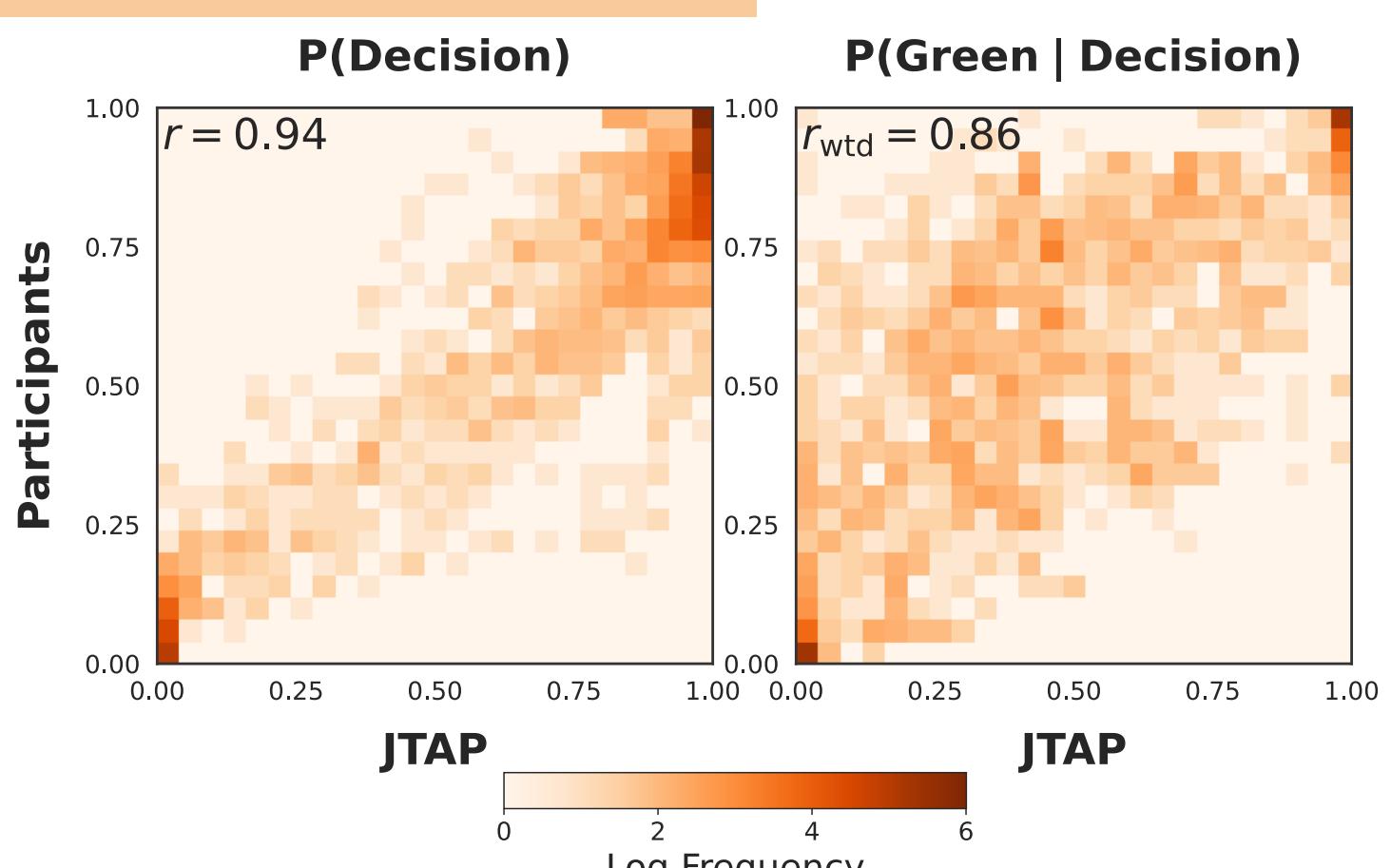


Discrete posterior approximation → View into the "mind" of inference

Modeling at the individual-level beliefs

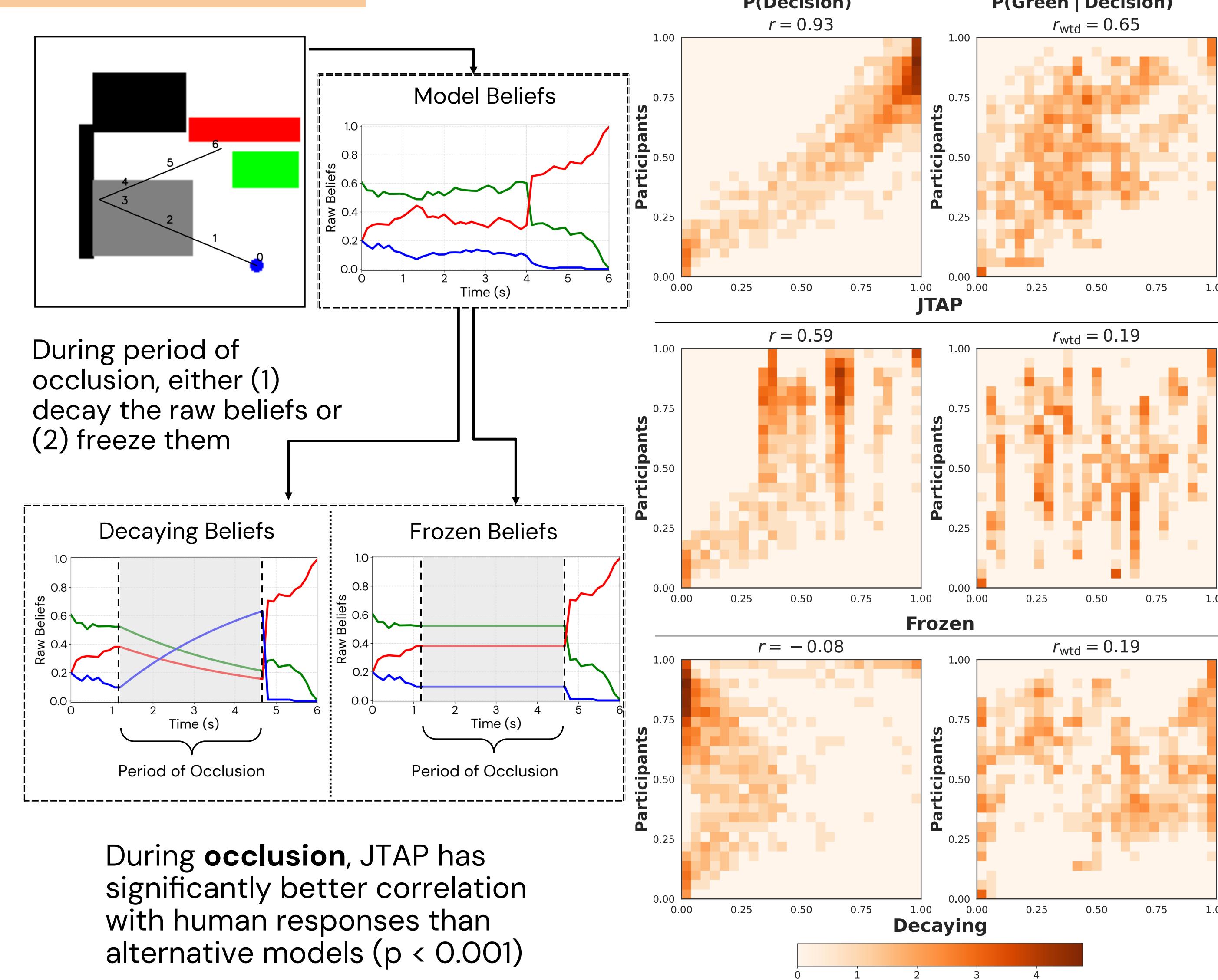


Overall Results



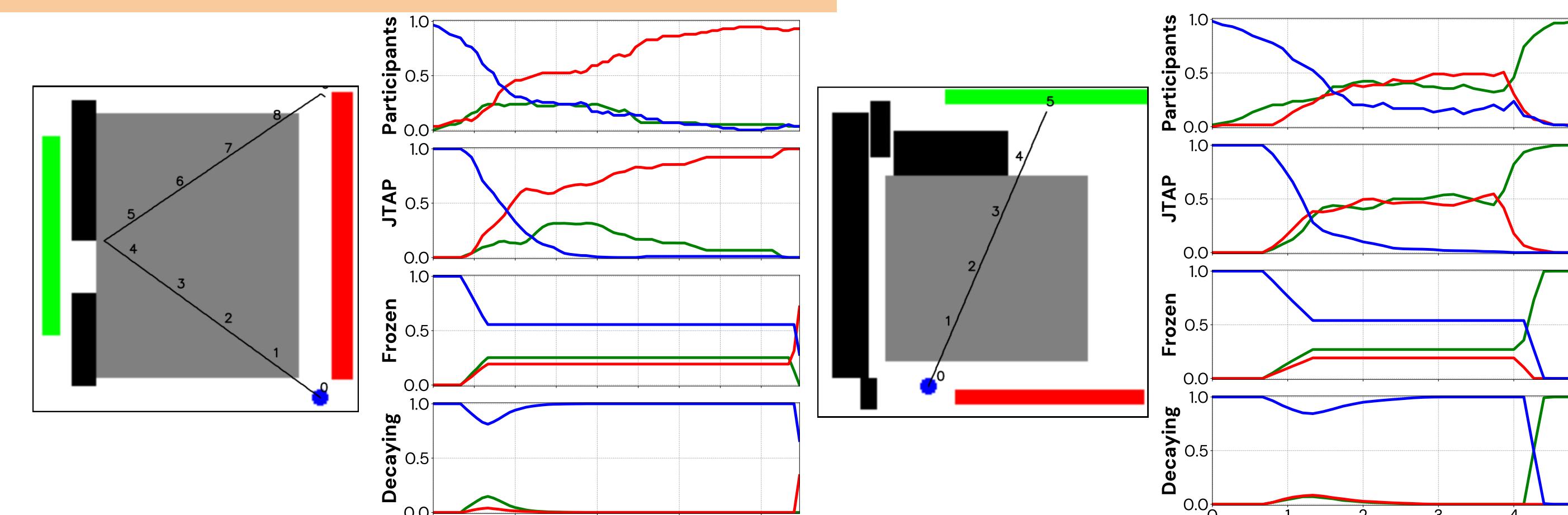
- Participants produce reliable decisions ($ICC_{1k} = 0.952$)
- JTAP explains well both **when** participants decided to push a button ($r=0.94$) and **which** button they pressed ($r=0.86$)
- Uncertainty gradations between humans are also captured by JTAP

Alternate Models



During **occlusion**, JTAP has significantly better correlation with human responses than alternative models ($p < 0.001$)

Benefits of Jointly Tracking and Predicting



The absence of visual evidence over time **is** evidence

Discussion

JTAP models human prediction of object motion during occlusion by integrating **perception**, **probabilistic** reasoning and **physical** knowledge. This approach is good at capturing human behavior overall but sometimes struggles to precisely match the timing of probability levels due to human priors.