

Exploring Theory of Mind for Human-Robot Collaboration

M. Romeo¹, P. E. McKenna², D. A. Robb², G. Rajendran², B. Nettet², A. Cangelosi¹, H. Hastie²

UKRI Node on Trust

¹University of Manchester, ²Heriot-Watt University

Overview

The ability to impute mental states to oneself or others, or Theory of Mind (ToM), has been linked to trust between humans. Less is known about how a robot mimicking ToM affects users' trust and behaviour. We explore this through a study where we compare 3 robot personas in a cooperative maze navigation task: (i) neutral ($n=235$), (ii) explaining its reasoning in technical terms ($n=236$), (iii) mimicking ToM ($n=235$). Results show that 'robot ToM' led to more cautious navigation decisions.

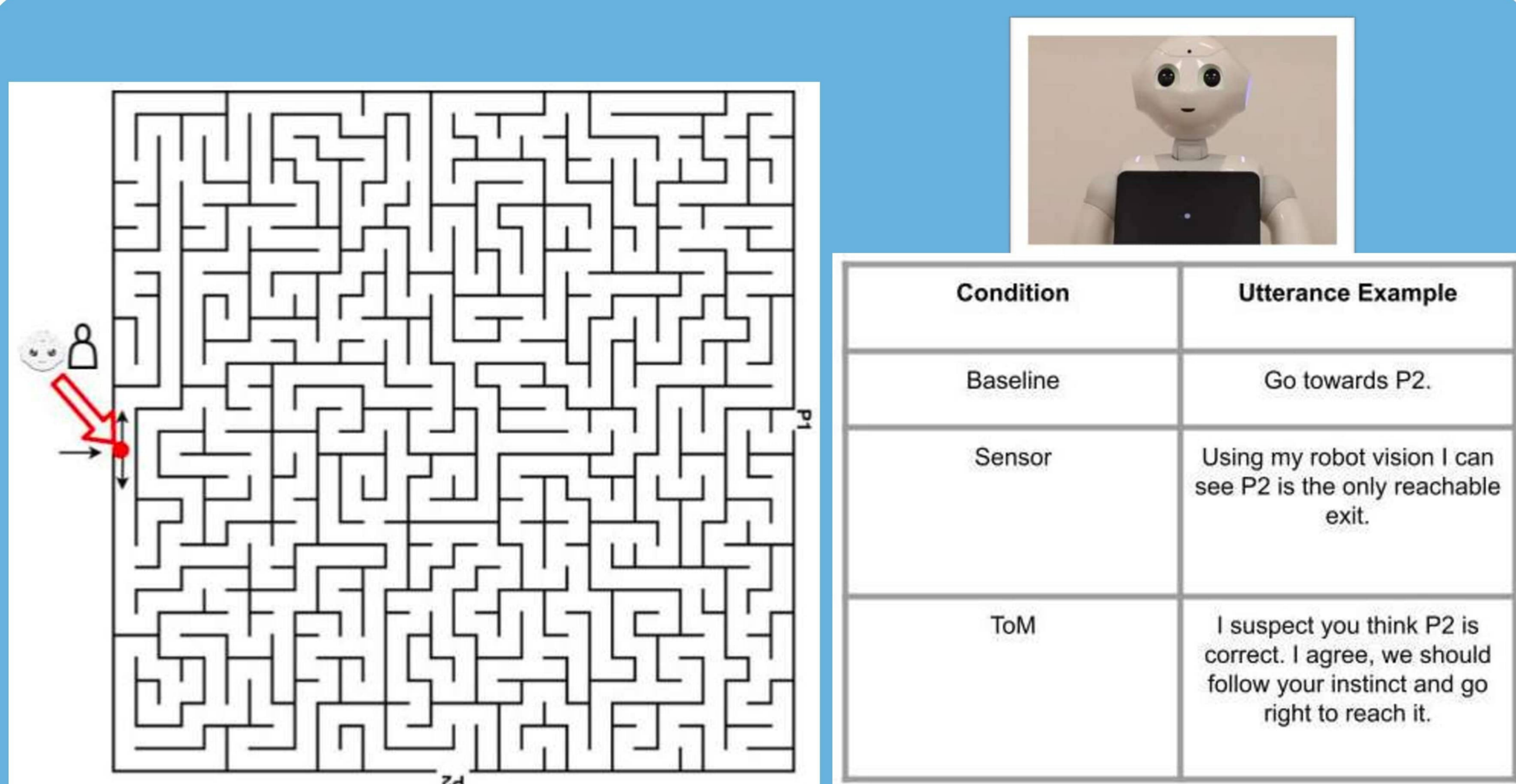


Figure 1. Maze example and Pepper's navigation suggestion utterances. ToM levels confirmed by manipulation check.

Key Features

- Design: Online, between-subjects
- IV: ToM priming & maze utterance
- DV: Subjective & behavioural measures of trust (in-task & MDMT q'nnaire)
- Task: solve 10 mazes
- Trial: follow/not follow robot's advice. T5 → Trust broken
- Pepper advice vary by persona

Results

- MDMT inconclusive
- follow/not follow did not vary per group
- Mimicked ToM: longer decision times
- Sensor: higher confidence

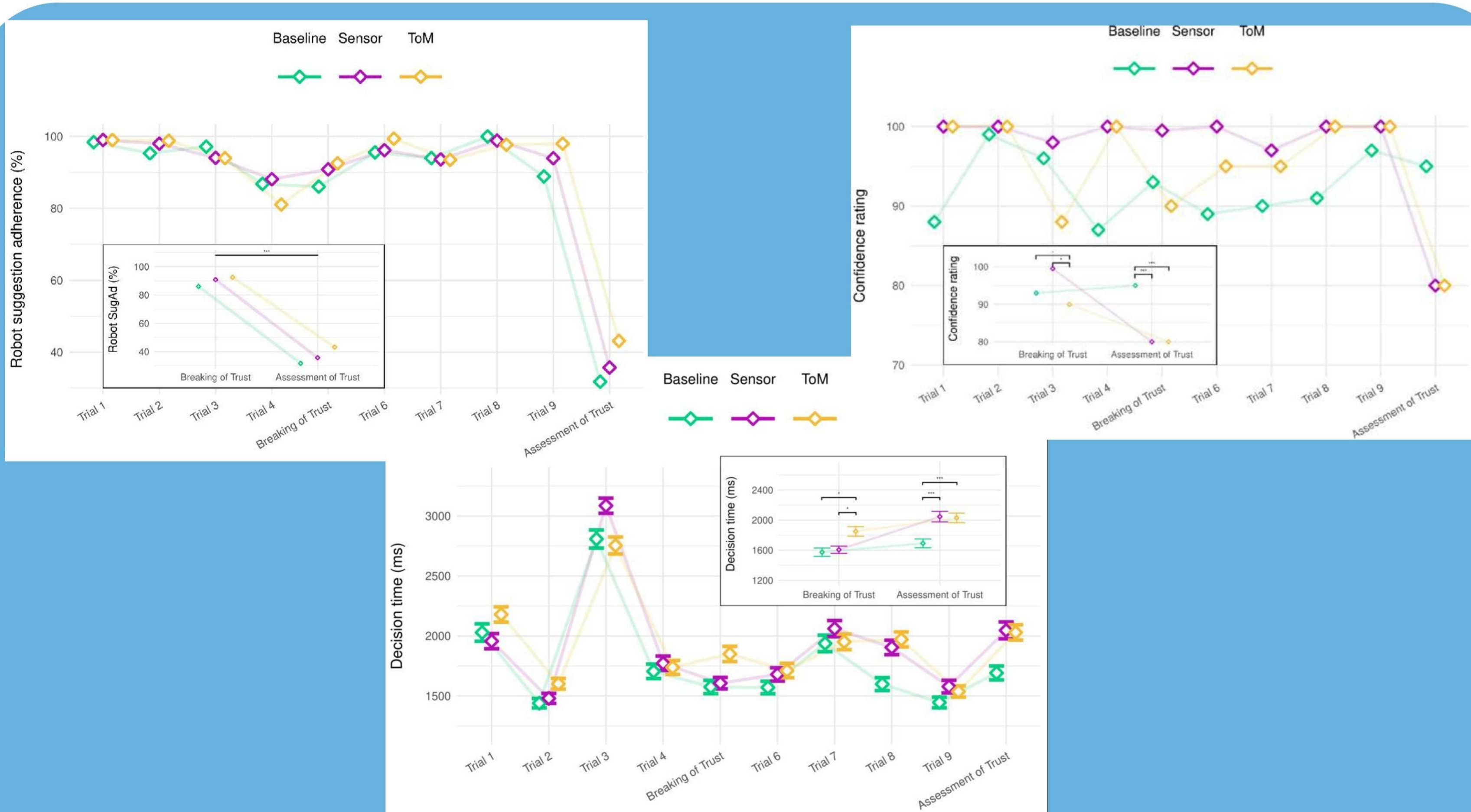


Figure 2. Results: Suggestion Adherence (left); Decision time (centre); Self-reported Confidence (right). Main plots: trend across trials. Inset plots: group comparisons & p-values.

Future Work

- Adding transparency of the interaction as a variable
- Analysis of different trust repair strategies
- Applications in: recommender systems, robot navigation, collaborative manufacturing