



Clustering Trust Dynamics in a Human-Robot Sequential Decision-Making Task

Shreyas Bhat¹, Joseph B. Lyons², Cong Shi¹, X. Jessie Yang¹

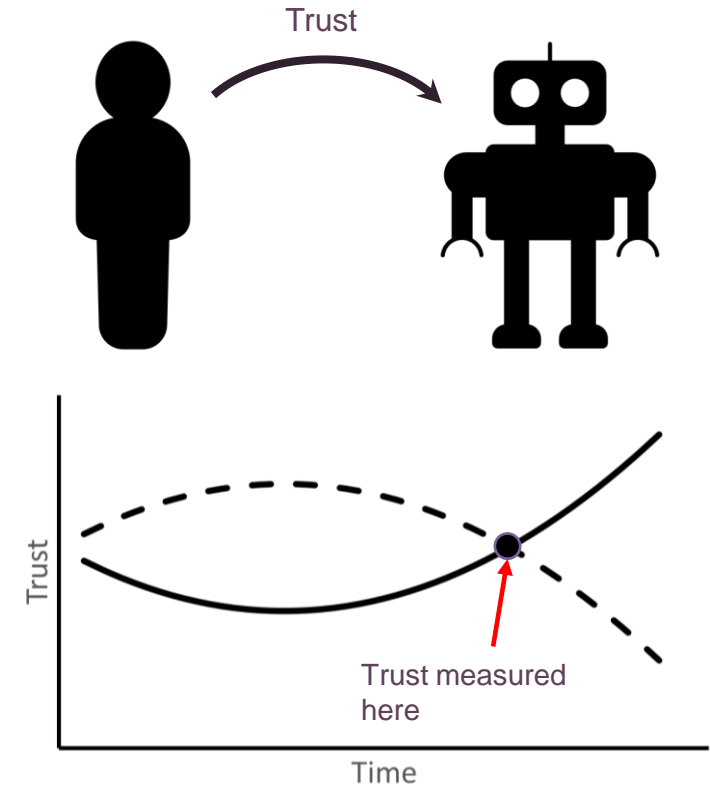
¹Industrial and Operations Engineering, University of Michigan

²Air Force Research Laboratory

1

Introduction

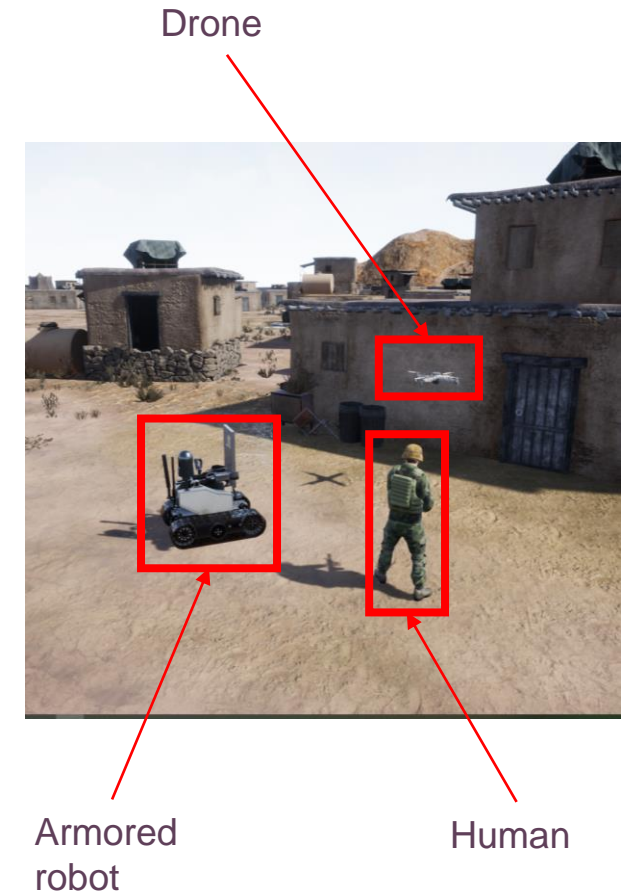
- **Trust** is a key factor for effective human-robot collaboration¹
- A snapshot view of trust is not enough. Trust can be dynamic within the interaction period²
- With a human trust-behavior model, robotic agents can be given insights into human behavior while making their decisions



1. P. A. Hancock, T. T. Kessler, A. D. Kaplan, J. C. Brill, and J. L. Szalma, "Evolving Trust in Robots: Specification Through Sequential and Comparative Meta-Analyses," *Human Factors*, vol. 63, no. 7, pp. 1196–1229, 2020.
2. X. J. Yang, C. Schemanske, and C. Searle, "Toward Quantifying Trust Dynamics: How People Adjust Their Trust After Moment-to-Moment Interaction With Automation," *Human Factors*, p. 00187208211034716, 2021.

Human-Robot Teaming Task

- Intelligence, Surveillance, and Reconnaissance (ISR) Mission
- Human-Drone team searches through N sites for potential threats
- Drone recommends whether to use an armored robot to breach the building or not
- Team receives rewards associated with health remaining and time to complete mission



Problem Formulation

- States: $t_i \sim \text{Beta}(\alpha_i, \beta_i)$

- Actions: $a_r \in \{0, 1\}$

- Human behavior model:

$$\begin{aligned}\mathbb{P}(a_h = a_r) &= t_i, \\ \mathbb{P}(a_h = 1 - a_r) &= 1 - t_i.\end{aligned}$$

- Rewards: $\mathbf{IR}_i^a = \boxed{-w_h h(a_h^i) - w_c c(a_h^i)} + \boxed{\lambda_i \cdot \mathbb{1}(A)}$

$R_i(a = a_h)$

Trust gain reward

- Transition function: $\alpha_i = \begin{cases} \alpha_{i-1} + w^s, & \text{if } P_i = 1, \\ \alpha_{i-1}, & \text{if } P_i = 0. \end{cases} \quad \beta_i = \begin{cases} \beta_{i-1}, & \text{if } P_i = 1, \\ \beta_{i-1} + w^f, & \text{if } P_i = 0. \end{cases}$

$$P_i = \begin{cases} 1 & \text{if } R_i(a_r) \geq R_i(1 - a_r), \\ 0 & \text{otherwise.} \end{cases}$$

Experiment

- 46 students from the University of Michigan participated
- Measures:
 - Big 5 Personality Traits
 - Perfect Automation Schema
 - Propensity to Trust
 - Trust after each site
 - Post-experiment Trust
 - Workload
- Participants searched through 100 sites sequentially



(a) No Threat, RARV Not Used



(b) No Threat, RARV Used



(c) Threat, RARV Not Used



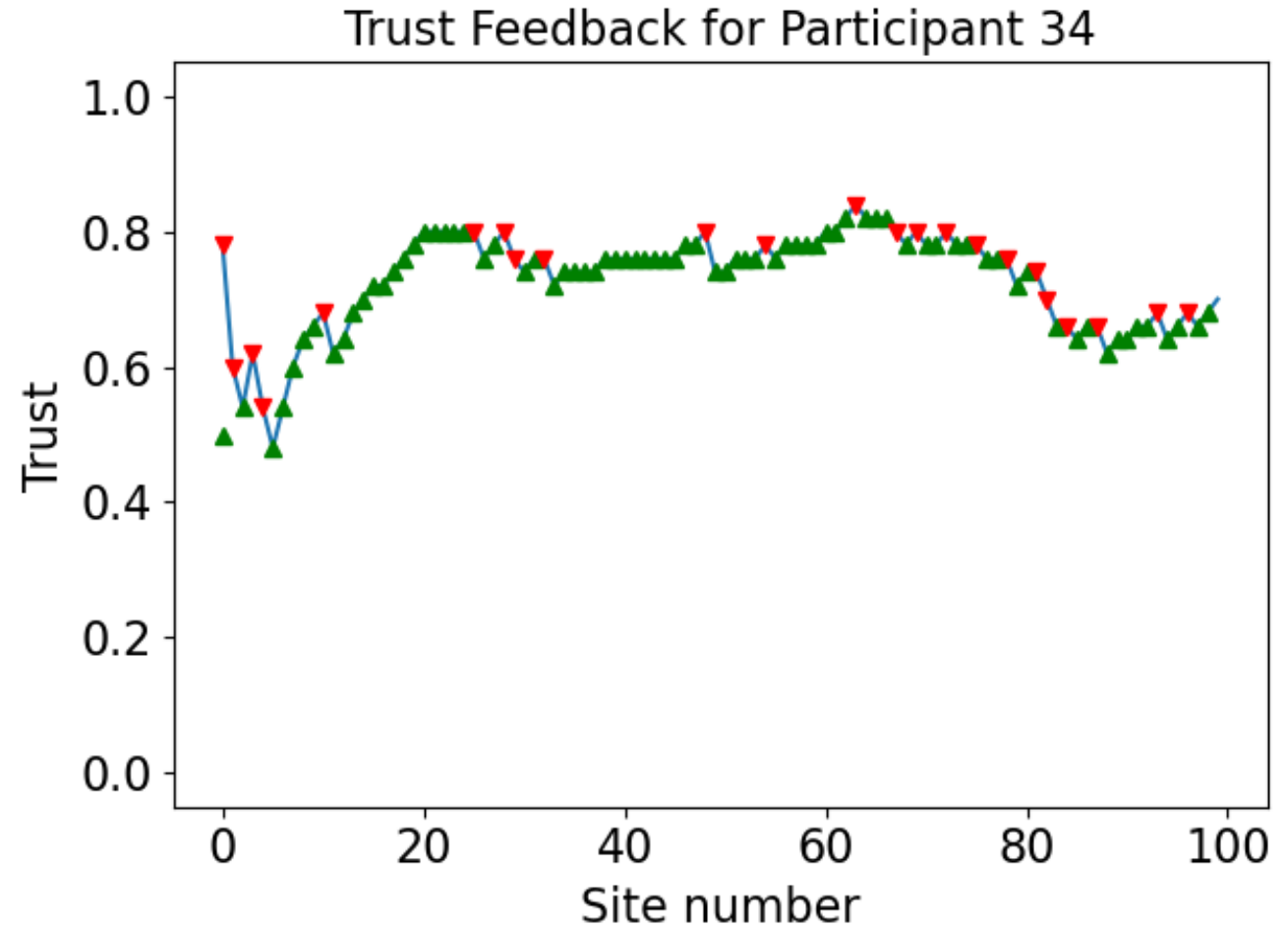
(d) Threat, RARV Used

Fig. 3. The four outcomes based on the presence of threat inside a site and the choice of action by the participant.

Results

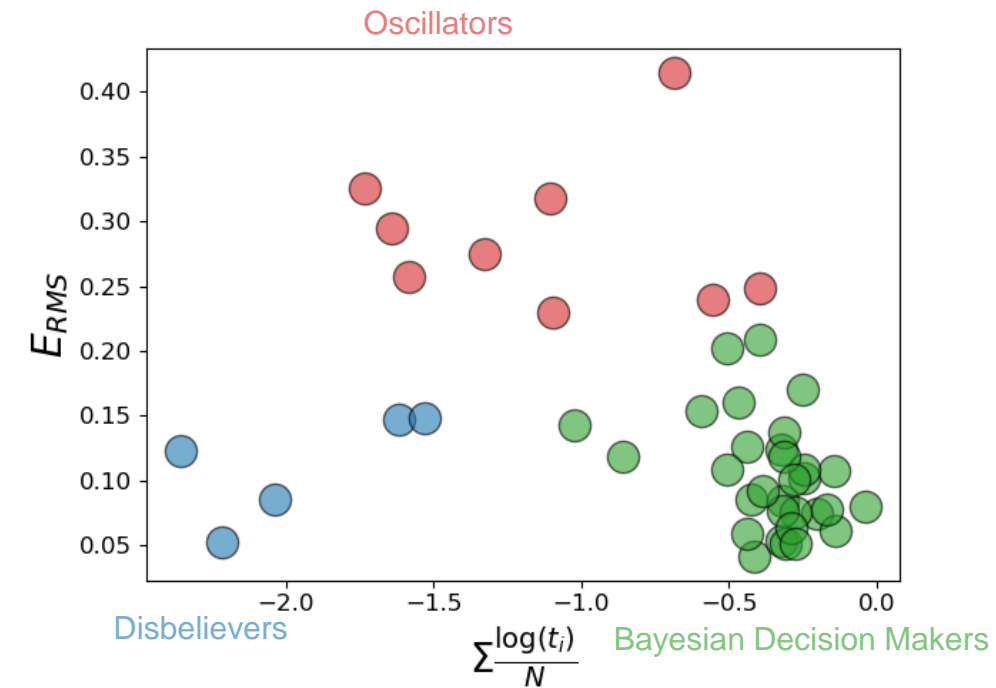
- Immediate task reward gain as a performance metric

▲ - $p_i = 1$
▼ - $p_i = 0$

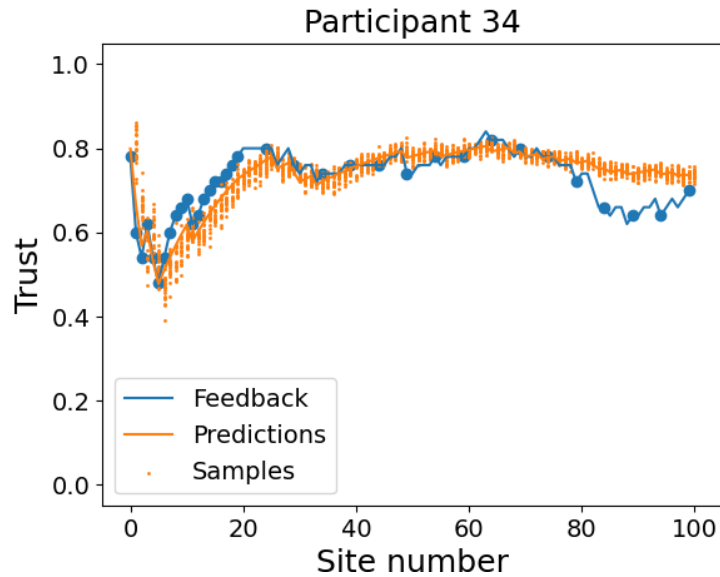


Results

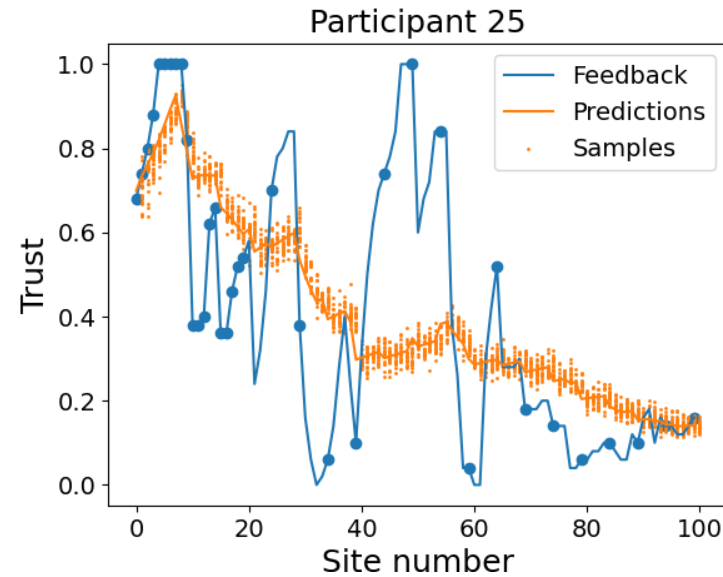
- K-means clustering analysis
- Features:
 - RMSE between feedback and predicted trust
 - Average log trust
- Elbow heuristic and silhouette scores indicate 3 significant clusters



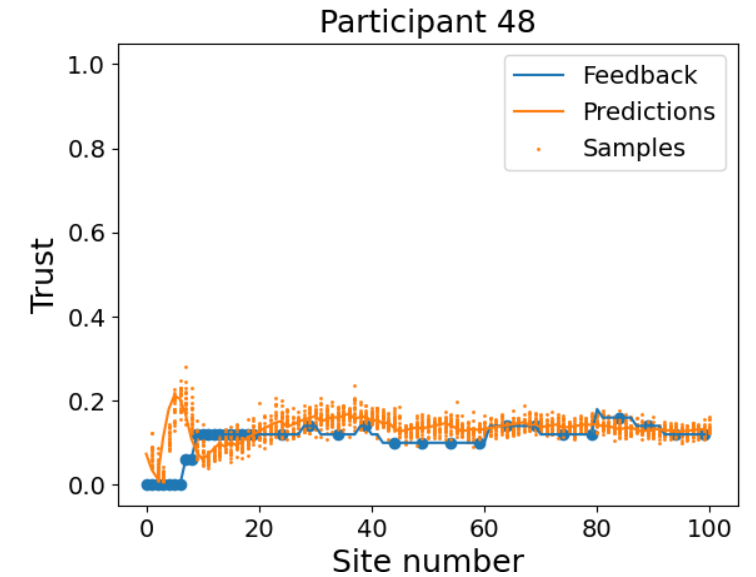
Results



Bayesian Decision Maker



Oscillator



Disbeliever

Type of trust dynamics	RMSE (SD)
Bayesian Decision Makers	0.093 (0.04)
Oscillators	0.26 (0.05)
Disbelievers	0.1 (0.04)

Results

TABLE I
MEAN AND STANDARD DEVIATION (SD) OF PERSONAL
CHARACTERISTICS BETWEEN THE THREE DIFFERENT TRUST DYNAMICS
(BDM = BAYESIAN DECISION MAKER)

Personal Characteristic	BDM	Disbeliever	Oscillator
Extraversion (/20) *	9.5 (3.3)	5.8 (2.8)	11.3 (2.9)
Agreeableness (/20) *	13.5 (2.5)	10.4 (5.0)	14.1 (1.8)
Conscientiousness (/20)	13.1 (2.7)	12.4 (3.0)	12.1 (4.5)
Neuroticism (/20)	7.9 (2.7)	6.8 (3.6)	10.2 (4.7)
Intellect/Imagination (/20) †	11.7 (2.0)	9.8 (1.8)	12.2 (1.8)
High Expectations (/28) **	12.7 (3.9)	6.4 (2.8)	12.4 (4.2)
All or None Thinking (/21)	6.6 (2.9)	6.4 (3.4)	7.1 (3.1)
Trust Propensity (/30) †	20.2 (4.4)	17.2 (4.1)	22.8 (3.2)

** – $p < 0.01$, * – $p < 0.05$, † – $p < 0.1$

Results

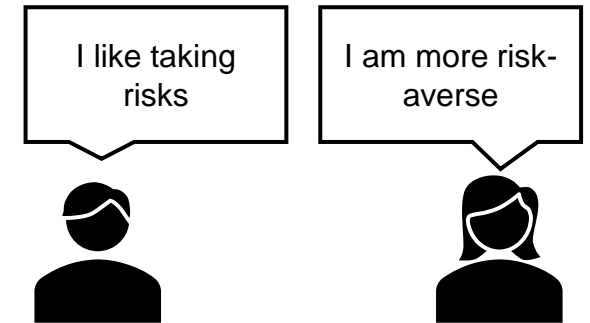
TABLE II
MEAN AND STANDARD DEVIATION (SD) OF POST EXPERIMENT METRICS
BETWEEN THE THREE DIFFERENT TRUST DYNAMICS

Personal Characteristic	BDM	Disbeliever	Oscillator
Trust (Muir) (/100) ***	65.4 (13.5)	15.8 (9.9)	44.7 (26.1)
Trust (Lyons) (/7) ***	4.5 (0.54)	3.1 (0.6)	3.6 (0.9)
Mental Demand (/100)	39.6 (25.2)	42.0 (36.6)	50.3 (28.6)
Temporal Demand (/100)	50.8 (27.4)	62.0 (24.5)	42.9 (21.3)
Performance (/100)	58.6 (19.7)	50.8 (30.0)	46.2 (31.7)
Effort (/100)	34.3 (23.0)	34.4 (17.4)	49.8 (32.2)
Frustration (/100) *	45.8 (22.2)	58.4 (25.4)	68.1 (14.3)

*** — $p < 0.001$, * — $p < 0.05$

Future Work

- Inverse Reinforcement Learning¹ to learn personalized reward functions
- Using contextual information for trust prediction
- Creating more balanced datasets



1. A. Y. Ng and S. Russell, "Algorithms for inverse reinforcement learning," in Proc. 17th International Conf. on Machine Learning. Morgan Kaufmann, 2000, pp. 663–670.



Thank you

This work was supported in part by the Air Force Office of Scientific Research
(Grant #FA9550-20-1-0406)