

Robot Capability and Intention in Trust-based Decisions across Tasks

Yaqi Xie*, Indu P Bodala*, Desmond C. Ong[†], David Hsu*, Harold Soh*

*Department of Computer Science

National University of Singapore

Singapore

{yaqixie, indu, dyhsu, harold}@comp.nus.edu.sg

[†]A*STAR Artificial Intelligence Initiative and Institute of High Performance Computing

Singapore

desmond.c.ong@gmail.com

Abstract—In this paper, we present results from a human-subject study designed to explore two facets of human mental models of robots—inferred capability and intention—and their relationship to overall trust and eventual decisions. In particular, we examine delegation situations characterized by uncertainty, and explore how inferred capability and intention are applied across different tasks. We develop an online survey where human participants decide whether to delegate control to a simulated UAV agent. Our study shows that human estimations of robot capability and intent correlate strongly with overall self-reported trust. However, overall trust is not independently sufficient to determine whether a human will decide to trust (delegate) a given task to a robot. Instead, our study reveals that estimations of robot intention, capability, and overall trust are integrated when deciding to delegate. From a broader perspective, these results suggest that calibrating overall trust alone is insufficient; to make correct decisions, humans need (and use) multi-faceted mental models when collaborating with robots across multiple contexts.

Index Terms—Trust; Human Robot Collaboration; Capability; Intention

I. INTRODUCTION

Trust is a cornerstone of long-lasting collaboration in human teams, and is crucial for human-robot cooperation [1]. For example, human trust in robots influences usage [2], and willingness to accept information or suggestions [3]. Mismatched trust in robots can lead to poor task-allocation and unsatisfactory outcomes. In recognition of its importance, there has been a concerted effort in the research community to better understand the formation and dynamics of trust in robots and automation [4]–[10].

Nevertheless, there remains crucial gaps in our understanding of human-robot trust, particularly in the role of inferred robot “*intention*”, i.e., what people believe the robot is trying to achieve. Prior research has focussed largely on inferred *capability* [7], [9], [10], which has been shown to be a primary antecedent to trust [11], [12]. However, with the advancement of robot technology (e.g., artificial intelligence), robots are increasingly poised to achieve peer-like collaboration with humans. In this new role, robots may be afforded greater

autonomy, which involves independent decision-making, and trust in intention may surface as a critical factor [1], [13], [14].

In this paper, we seek to clarify the role of *both* inferred robot intention and capability in trust-based scenarios. Inspired by prior work on inter-human and socio-cognitive trust [13], [15], we posit that when deciding to delegate tasks to a robot, a user considers two complementary aspects: (i) whether the robot has the proper intention or motive, i.e., if it is optimizing the correct objectives, and (ii) whether the robot has sufficient capability, i.e., if it is able to carry out the task successfully under those objectives. For example, a passenger in an autonomous vehicle (AV) may trust in the AV’s capability to navigate in a complex environment, yet distrust the AV’s intention to hasten his arrival by circumventing road safety, or to value his life over the lives of others in an emergency.

We present results from a human-subject study ($n = 400$) where participants had to choose whether to delegate control to a robot under varying conditions. In particular, we sought to examine how estimations of robot intention and capability developed under one task affected trust-based decisions in the same task, and *subsequent novel tasks* (where arguably, trust plays a more important role). We created an online survey where participants had to decide whether to trust an Unmanned Aerial Vehicle (UAV) in three different tasks: (i) searching, (ii) mapping, and (iii) fire-fighting. After interacting with the robot in the searching task, participants had to decide whether to trust the robot in the latter two tasks. The searching and mapping tasks required similar robot capabilities, but involved potentially different objectives. Likewise, similar intentions (i.e., risk-behavior) would arise in the searching and fire-fighting tasks, but different robot capabilities are required.

Our primary finding is that decisions to delegate control in novel task contexts depend not only on overall trust in the robot, but also on estimations of robot capability and intention. In other words, humans appear to integrate several facets of their mental model (i.e., their estimations and beliefs) to arrive at a decision to trust. Furthermore, our results suggest that inferred capability and intention transfer (or generalize) separately to new situations, which extends previous results [9]

and suggests avenues for future work in assessing the degree of transfer.

These results suggest that human-robot trust may be similar to trust in humans, in that both inferred intention (motive) and capability play a role. However, trust in *non-robot* automation (e.g., automated alarms and decision aids) has also been shown to develop similarly to trust in humans, but with critical differences [16]. For example, people perceive automated systems to be more credible and objective sources of information compared to humans (to the point where erroneous decisions are agreed with). However, people are less tolerant of automation errors, leading to sharp declines in trust when a mistake is perceived. Here, we compared trust-based decisions when participants were informed they were working with another human player, and when they were paired with a robot. In both these cases, the agent was a confederate software program with the *same* behavioral policy. Interestingly, we did not find that humans always engaged in more trusting behavior when they thought they were working with another human (rather than a robot), *despite* reporting significantly higher trust in those situations. This suggests a complex relationship between trust, observed performance, and prior expectations given the perceived agent type.

To summarize, this work contributes a novel investigation into trust-based decisions in human-robot teams where robot capability and intention play a role in outcomes. We find evidence that humans utilize rich mental models of robot teammates when choosing to trust. Self-reported trust in the robot is an insufficient predictor of human decisions to trust the robot in different contexts. Rather, our results show that decisions to trust are based on human’s mental model of the robot, which includes inferred capability, intention, and potentially other components. Furthermore, people’s mental models appear to differ when working with a human partner or a robot—this points to a potential difference in prior expectations when working with humans versus machines. Taking a broader perspective, our results have implications for the design of human-centric robots that are able to reason about human trust and act accordingly [10]. For example, robots that aim to “teach” humans to make appropriate delegation decisions should calibrate not only general trust, but also their estimations of the robot’s intention and capabilities.

II. PRELIMINARIES: BACKGROUND AND RELATED WORK

To situate our work, we briefly review trust research, which is a large interdisciplinary endeavor spanning multiple fields including human-factors, social science, and human-robot interaction. Trust has been studied in many forms: trust in other humans, trust in organizations, and trust in machines. Here, we discuss the literature with a focus on *human trust in robots* and the aspects most relevant to this paper, i.e., the key definitions and concepts, and research in modeling trust in robots.

A. Characterizing Trust: Concepts and Definitions

Trust is a concept with varying definitions even within the same field. For example, trust has been defined as a belief—the

subjective probability whereby an agent (the *trustor*) assesses whether another agent (the *trustee*) will perform an action [17]. However, this definition has been criticized for lacking task context: for example, it fails to take into account the risk associated with the task at hand [13]. An alternative risk-related definition of trust is as the belief that a trustee will help the trustor’s goal in a situation characterized by uncertainty and trustor vulnerability [18], or just the willingness of the trustor to depend on another agent, even with the risk of possible negative consequences [19]. In this work, we adopt a recent definition of trust used in human-robot interaction, i.e., that trust is a latent (hidden) variable that summarizes past experience with an agent/robot [9], [10], which is useful for predicting future behavior of the trustee and making a decision to put oneself in a position of vulnerability. By summarization, we mean a mental abstraction or model of past interactions that is predictive of the trustee’s future behavior, since humans generally do not remember the entirety of past interactions with other agents.

There are two types of trust that differ in their situation specificity: on the one hand, there is *dispositional trust* or *trust propensity*, which is an individual difference for how willing one is to trust another. Some people may inherently be more trusting [20]. On the other hand, *situational* or *learned trust* results from interaction between the agents concerned. The more you use your new robot, the more you may learn to trust it. Dispositional trust is a trustor *trait*, i.e., it differs between trustors, and tends to be similar for the same trustor across different situations. By contrast, situational trust is a trustor *state*, i.e., it is specific to the task at hand, and may change frequently given new information. In this work, we will be concerned primarily with situational trust in robots, but we bear in mind individual differences in dispositional trust when examining empirical data.

B. Trust in Automation and Robots

Muir’s seminal work on trust in automation [5], [6] and its effect on operator control allocation found that people’s reported subjective trust ratings were sensitive to the automation’s properties—the more reliable the automation, the higher the trust ratings, which subsequently led to increased automation reliance. In other words, human trust was influenced by machine behavior, which in turn influenced human behavior. Lee and Moray [21] found that this relationship was moderated by the operator’s self-confidence, i.e., if trust in automation was greater than the operator’s self-confidence, the operators were more likely to rely on the machine.

As AI technology has matured, the research community has examined the trust in intelligent systems and robots, such as autonomous vehicles [9], [22], unmanned aerial vehicles (UAVs) [7] and medical diagnosis systems [23]. In general, trust in human-robot interactions can be influenced by robot-related factors (e.g., performance, physical attributes), human-related factors (e.g., workload, self-confidence), and environmental factors (e.g., group composition, culture, task type) [11]. A majority of prior work has focussed on

performance-related factors, particularly robot capabilities; for example, Soh *et al.* [9] examined the dynamics and transfer of trust in robot capabilities across tasks, where transfer is the ability to employ knowledge acquired in one task to improve performance in another [24]. Recent work has explored the role of the robot’s intention, e.g., its policy [14], [25] and decision-making process [26]. This work adds to this body of literature and considers both intent and capability across tasks.

In the following, it is important to distinguish between the robot’s *true* capability/intent versus the human’s *estimation* of the robot’s capability/intent. These two may differ, particularly when the human has had little experience with the robot. We will use the term *inferred capability* to refer to the human’s estimation of the robot’s ability to perform the task successfully. Likewise, *inferred intention* refers to the human’s estimation of the robot’s underlying motives or decision-making criteria, rather than the agent’s true utility function.

C. Trust Measurement

Trust is both latent (i.e., unobservable) and dynamic, which presents challenges for measurement. Multiple measurement scales have been proposed to quantify the degree of trust in a robot, including binary measures [27], continuous measures [21], [28], [29], ordinal scales [30]–[32] and an Area Under Trust Curve (AUTC) measure [28] that captures participant’s trust through the entire interaction with the robot by integrating binary trust measures over time. In this work, we use both self-reported trust measures (e.g., Schaefer’s trust scale [33]) and behavioral measures (decisions to delegate tasks to the robot).

III. HUMAN-SUBJECT STUDY: INTENTION, CAPABILITY, AND TRUST

The overarching goal of our experiments was to explore the relationship between the inferred intention, capability, overall trust, and the decision to trust. We designed a user study with a two {Grouped with Robot vs. Grouped with Human} by four {robot: High Capability/Low Capability + robot: High Risk-taking/Low Risk-taking intention} between-subjects factorial design.

A. Experimental Design

Data Collection Platform and Tasks. For this work, we developed an online survey (Fig. 1) with a delegation game; a human teammate decides whether to delegate control to an Unmanned Aerial Vehicle (UAV) performing three different mission tasks. Depending on their choice and the outcome of the task, they would gain a certain number of points. Choosing to take over control, i.e. tele-operate, would always cost some points, but this could be offset if the mission was a success. The three tasks were:

- 1) **Searching:** The UAV searches for targets by taking photos in various weather conditions. The UAV can choose between two locations, A and B, each with a different weather condition and number of targets. The weather

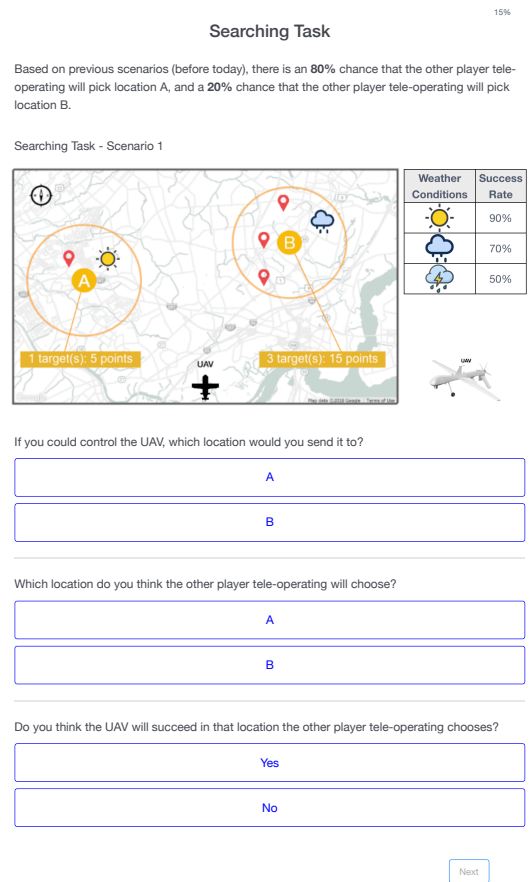


Fig. 1: Example Scenario in Task 1 (Searching).

condition affects the probability of success, which is shown to the participants (see Fig. 1).

- 2) **Mapping:** UAV creates maps by taking photos of the terrain. The UAV again chooses location A or B, but the weather in both locations are the *same*. The difference is that visiting location A results in a low-accuracy but wide coverage map. Conversely, visiting B results in a high-accuracy, low coverage map.
- 3) **Fire Fighting:** UAV puts out fires by dropping water/chemicals, again with two choices A and B. The two locations would have different fire intensities and numbers of targets.

In order to test the human trust transfer, we designed two types of capabilities and two types of intent:

- Capability 1 (Weather): Taking photos in various weather conditions. (Task 1 & 2)
- Capability 2 (Fire-fighting): Fire-fighting in various fire conditions. (Task 3)
- Intent 1 (Risk preference): Lower risk but less reward vs. Higher risk but more reward. (Task 1 & 3)
- Intent 2 (Accuracy preference): Higher coverage but low accuracy vs. Lower coverage but higher accuracy. (Task 2)

Task 1 and Task 2 required the same robot capability but different decision-making criteria (intention), while Task 1 and Task 3 involved the same intention (risk taking behavior) but different capability.

The tasks were designed to test how trust and learned mental models transfer to situations where the robot capabilities required would differ (taking photos vs. fighting fires) or where potential differences in value assignments may occur (number of targets vs. mapping). With regard to the latter, participants were incentivized to pursue certain goals via point allocations, but were informed that the other agent did not necessarily optimize the same criteria. In each task, the robot either succeeded, whereby the participant would obtain all the stated points, or failed and the participant would receive nothing (or negative points if they chose to teleoperate the robot).

Confederate Agent/Robot Types. In this study, human participants played with software agents, similar to [34]. To investigate if trust differed significantly when participants thought they were playing with another person or a machine, we randomized the participants into two groups. In (**Group: Human, GH**), participants were informed that they were paired with another player who would tele-operate the UAV, i.e. make location decisions, should participants choose to delegate control. In (**Group: Robot, GR**), participants were told they would interact with a robot. In both groups, participants interacted with one of four agent *types* which differed along two dimensions: robot capability and intention.

In Task 1, the robot's *capability* was its ability to complete a mission in different weather conditions (clear, rainy, or thunderstorm). The success probabilities for the High Capability (HC) and Low Capability (LC) robots are shown in Table I. The robot's *intention* was related to its preference for risk; high risk-taking (HR) agents would attempt to maximize the number of targets even at locations where the chance of failure was high. Conversely, low risk (LR) agents are risk-averse. The exact decision made in the different scenarios were obtained via expected utility maximization, i.e., the agents would choose actions that maximized their expected utility:

$$\mathbb{E}[U] = \sum_o U(o)p(o|a) \quad (1)$$

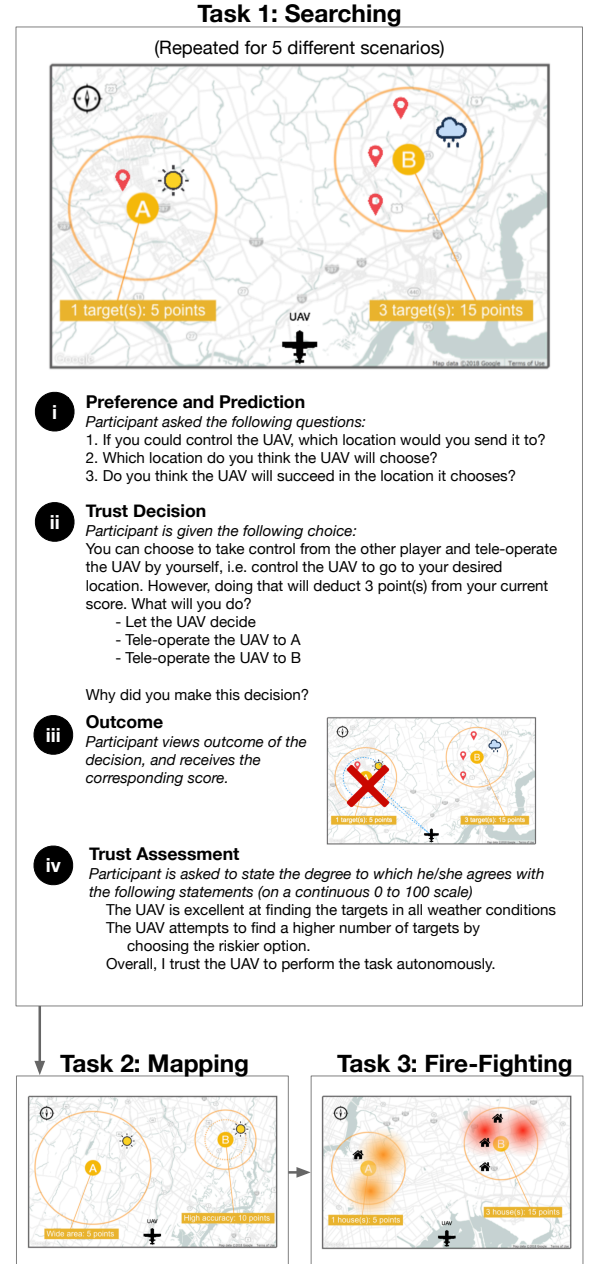
where U is a utility function and $p(o|a)$ is the probability of an outcome o given an action a . We used an exponential utility function:

$$U(o) = \begin{cases} \frac{(1-e^{-\alpha o})}{\alpha} & \alpha \neq 0 \\ \alpha & \alpha = 0 \end{cases} \quad (2)$$

which is commonly applied in economics to model risk propensity. The parameter α controls the degree of risk preference: $\alpha > 0$ for risk aversion, $\alpha = 0$ for risk neutral, and $\alpha < 0$ for risk seeking. In our work, for the risk-taking robot (HR), $\alpha_{HR} = -0.2$, and for the risk-averse robot (LR), $\alpha_{LR} = 0.8$. Note that outcomes were not shown for Tasks 2 and 3 so, neither capability specifications nor decision-making process were required.

TABLE I: Success Rates for Task 1 (Searching) in Different Weather Conditions for Agents with High or Low Capability.

Weather Type	Capability	
	High	Low
Clear	100%	66%
Raining	83%	0%
Thunderstorm	66%	0%



Tasks 2 and 3 are counterbalanced with 3 scenarios each. Participants are asked questions similar to parts (i), (ii) and (iv) in Task 1 but modified to suit the specific task.

Fig. 2: Experiment workflow comprising three tasks.

B. Methodology

Participant Recruitment. We recruited 400 participants via Amazon Mechanical Turk (AMT). Participants were required to have at least a 99% acceptance rate and were only allowed to participate in our survey once. Each survey lasted 30 minutes and participants were compensated \$3. Participants that attained a higher number of points had an opportunity to be compensated an extra \$10, which incentivized participants to pay attention and carefully consider their judgments.

Procedure. After providing consent and standard demographic data, participants were randomized to the two groups (GH or GR) and one of the four agent types. They were then presented with a description of Task 1 (Searching), and required to answer attention check questions. Participants were required to answer all questions correctly to proceed with the experiment. They were then allowed to play a trial scenario until they were ready to proceed.

The remainder of experiment workflow is shown in Fig 2. The participants were presented with five different scenarios for Task 1; each scenario consisted of a different number of targets (and hence, achievable points) and weather conditions in locations A and B. Each scenario description was followed by a sequence of four stages:

- (i) Preference and Prediction: Participants were asked to indicate their preferred location, and to predict where the robot would choose to go and the chance of success.
- (ii) Trust Decision: They then had to choose whether to delegate control to the robot, or to perform a take-over by tele-operating the UAV and overriding the UAV's choice in GR or the other player's choice in GH.
- (iii) Outcome: The outcome of their choice would then been shown, along with the points they would receive.
- (iv) Trust Assessment: Participants would then be asked to provide agreement scores on statements regarding the UAV's competence, risk-behavior, and their overall trust.

Participants then proceeded to Task 2 (Mapping) and Task 3 (Fire-fighting). The order of Task 2 and Task 3 were counterbalanced to eliminate order effects. Each task comprised 3 scenarios. Similar to Task 1, each scenario was followed by (i), (ii) and (iv); we excluded stage (iii) since providing additional observations may change participants' mental models. After completing all three tasks, participants completed a short questionnaire regarding their trust in the robot, the inferred robot capability and risk-behavior.

Dependent Variables The primary dependent variables consist of both subjective self-reported measures—e.g., of overall trust in the robot and estimations of its capability, intention—and an objective measure—whether they decide to allow the agent to perform the task autonomously. For the self-reported measures, participants indicated the agreement to statements via a continuous scale ranging from 0% to 100% where 0 indicated complete disagreement and 100 indicated complete agreement. Further details about the dependent measures are listed below:

- **Trust Decision.** This objective behavioral measure captures whether the participants trusted the agent to perform the task by itself. In our setup (see part (ii) in Fig. 2), they chose either to let the UAV decide the intended location (a decision to trust), or to tele-operate the robot to a specific location (a decision *not* to trust).
- **Self-reported Trust.** We asked patients to state their agreement with the statement, “Overall, I trust the UAV to perform the task autonomously”. We also included eight questions from Schaefer's trust scale [33].
- **Self-reported Inferred Capability and Intention Estimation.** We measured participants' perception of the robot's capability in the three different tasks using agreement statements, e.g., “The UAV is excellent at finding the targets in all weather conditions.” for Task 1. Similarly, for risk-behavior assessment, we used the statement, “The UAV attempts to find a higher number of targets by choosing the riskier option”. At the end of the survey, we asked participants to choose between one of four options about the type of agent they had interacted with in the survey: (1) “Highly capable but tends to take risk.”, (2) “Highly capable but is conservative.”, (3) “Not very capable but tends to take risk.”, or (4) “Not very capable and is conservative.”.
- **Robot Decision and Outcome Prediction** We measured the participants' predictions about the robot's choice and the outcome given the choice, via two binary answers to the questions: “Which location to you think the UAV will choose?” and “Do you think the UAV will succeed in the location it chooses?”, respectively. We also computed intent alignment by comparing the predicted robot decision to the decision the participant would have made (via their answer to the question “if you could control the UAV, which location would you send it to?”).
- **Self-reported Decision-making Similarity.** This captured the degree of perceived similarity to the agent in terms of the decision making process, measured via the agreement statement, “The [UAV's | other player's] decisions are similar to mine.”

C. Hypotheses

Our overarching hypothesis is that humans infer both robot capability and intent from observations, and use these estimations to make decisions whether to trust the robot in new, but related, tasks. We specifically hypothesized that:

- H1: Humans infer capability and intention from observations of robot performance.
- H2: Inferred capability and intention transfer separately to different tasks.
- H3: Both inferred capability and intention contribute to self-reported trust in the robot.
- H4: Both inferred capability and intention influence human decisions to trust the robot.
- H5: People are more trusting of the simulated human agents, rather than the simulated robot agents.

Hypotheses H1 and H2 capture our expectation that the different robot types and performance would engender different mental models, and that specific facets of the mental models would carry over to the other tasks. H3 and H4 are our primary hypotheses relating inferred capability and intent to trust (both self-reported and behavioral). Finally, H5 encodes our expectation that perceived agent type also has an effect on trust; recent work has found that humans tend to trust other humans more compared to intelligent decision aids [16] and software agents in economic games [34], even when the behavior of the confederate agent was identical.

IV. RESULTS

Data from 400 participants (Mean age: 38.85 years, 48% female) were included in the following analysis. Overall, we found that human estimations of robot capability and intent are important factors in determining overall trust and trust-based decisions. We analyzed the two groups GR (playing with a simulated robot) and GH (playing with a simulated human player) separately; each group contained responses from 200 participants. The statistical analyses for H1-H4 were performed using GR data only since they are defined over human-robot trust, while H5 analysis compares GR with GH.

H1: Humans infer capability and intention from observations of robot performance. Fig. 3 summarizes trust, inferred capability and intention (risk preference) scores for Task 1¹. Participants reported different inferred capability and intention for each agent ($p < 0.01$ across the four agent types). Specifically, inferred capability for HC-HR and HC-LR ($M = 0.52, SE = 0.02$) agents are significantly higher than LC-HR and LC-LR ($M = 0.29, SE = 0.02$); $t(398) = 7.62, p < 10^{-11}$. The inferred risk preference for HC-HR and LC-HR ($M = 0.66, SE = 0.03$) were significantly higher than HC-LR and LC-LR ($M = 0.42, SE = 0.03$); $t(398) = 5.99, p < 10^{-8}$. These results support H1. In addition, reported score differences between agents having the same true capability (or intent) were smaller compared to agents that differed in that dimension; the mean difference was ≈ 0.08 for both intent and capability, compared to ≈ 0.2 when the agent’s true intent/capability differed.

Participants were also largely able to discriminate between the robot types. As shown in Table II, there is a high agreement for all the cases except for LC-LR agent. Participants appeared to confuse the high capability and risk-averse agent, and the low-capability and risk-averse agent. This confusion is more severe in the simulated human group (GH). It was possible that since participants were allocated to only one agent type in our between-subjects design, they lacked a relative basis for comparison. Moreover, the low-capability and risk-averse robot would tend to choose the “safer” option, resulting in more successes.

¹Recall that the Searching task (Task 1) was used in a learning phase—outcomes are shown to participants, allowing them to update their mental representation of the robot.

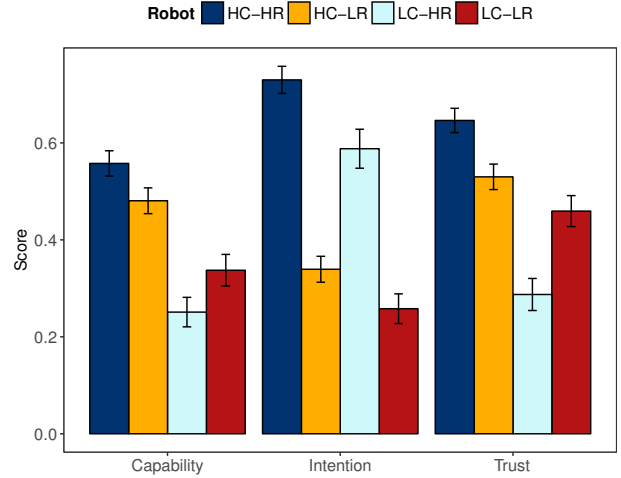


Fig. 3: Self-reported Trust, Inferred Capability and (Risk) Intention scores for the four robot types: High-Capability and High Risk preference (HC-HR), High-Capability and Low Risk preference (HC-LR), Low-Capability and High Risk preference (LC-HR), Low-Capability and Low Risk preference (LC-LR). Participants reported higher capability scores for the HC robots, and higher risk intention for the HR robots, indicating they were able to infer capability and risk preference from observed robot behavior and outcomes.

TABLE II: Confusion Matrix of Predicted Robot Types.

Simulated Robot Group (GR)				
Robot Type	HC-HR	HC-LR	LC-HR	LC-LR
HC-HR	0.65	0.2	0.15	0
HC-LR	0.25	0.575	0.075	0.1
LC-HR	0.125	0.075	0.7	0.1
LC-LR	0.1	0.225	0.15	0.525
Simulated Human Group (GH)				
Robot Type	HC-HR	HC-LR	LC-HR	LC-LR
HC-HR	0.75	0.15	0.1	0
HC-LR	0.175	0.725	0.05	0.05
LC-HR	0.275	0.025	0.65	0.05
LC-LR	0.1	0.525	0.125	0.25

H2: Inferred capability and intention transfer separately to different tasks. We expected inferred capability and intention to be separate components of a mental model that are transferred depending on the nature of the task. Recall that Task 2 (mapping) was designed such that the capability requirement of the robot was similar, i.e., operating in different weather conditions, but where value assignments were possibly different compared to Task 1 (since the risk was the same at both locations). Hence, we expected estimations of capability to transfer (i.e., be similar), but not the estimations of intention (risk-preference). Conversely, Task 3 (fire-fighting) was designed such that the capability required was different, i.e., putting out fires rather than picture-taking in different weather conditions. However, choices in both Task 3 and

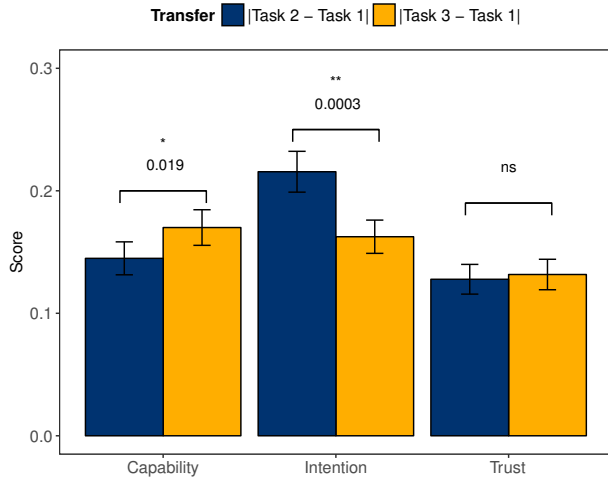


Fig. 4: Transfer of self-reported trust, inferred capability and intent as measured by the absolute difference between scores (p -values shown above the bars). Differences in inferred capability and intent were significant between tasks, indicating differences in transfer depending on the task.

Task 1 involved differences in risk. As such, we expected that estimations of intent to transfer, but estimations of capabilities not to transfer.

We measured transfer across task contexts via the absolute difference of the scores from Task 1 to Task 2 and from Task 1 to Task 3, similar to prior work [9]. Lower differences indicate higher transfer. Fig. 4 summarizes our results for all robots, which support H2. As expected, inferred capability differences were significantly smaller (greater transfer) between Tasks 1 and 2 ($M = 0.14$, $SE = 0.01$), compared to Tasks 1 and 3 ($M = 0.17$, $SE = 0.01$); $t(199) = -2.37$, $p = 0.02$. Likewise, the transfer of inferred intent was significantly greater between Task 1 to Task 2 ($M = 0.22$, $SE = 0.01$), compared to Task 1 and Task 3 ($M = 0.16$, $SE = 0.01$); $t(199) = 3.72$, $p < 10^{-2}$.

H3: Both inferred capability and intention contribute to self-reported trust in the robot. We combined the responses from all scenarios in Tasks 1, 2 and 3, yielding a total of 2000 data points (200 participants \times 10 scenarios) for analysis. Linear mixed models were used to account for repeated measures; inferred capability and intention were fixed main effects, with random intercepts for subjects, scenarios and tasks. Table III summarizes our results; significant associations were found between inferred capability ($b = 0.42$, $SE = 0.02$, $p < 2 \times 10^{-16}$) and inferred intention ($b = 0.18$, $SE = 0.02$, $p < 2 \times 10^{-16}$) with self-reported trust scores, in support of the hypothesis.

H4: Both inferred capability and intention influence human decisions to trust the robot. Similar to H3 above, we conducted mixed-effects analysis to account for dependent measures. In this case, the dependent variable was the trust decision (i.e., whether to delegate the task to the robot). Our primary model included random intercepts for subjects,

TABLE III: Mixed Effects Linear Regression Model relating Inferred Capability and Intention to Self-reported Trust.

	Coef.	SE	t value	Pr(> t)
Capability	0.42	0.02	23.39	$< 2 \times 10^{-16}$
Intention	0.18	0.02	10.99	$< 2 \times 10^{-16}$
Intercept	0.21	0.02	13.94	$< 2 \times 10^{-16}$

TABLE IV: Mixed Effects Logistic Regression Model for Trust Decisions.

	Coef.	SE	z value	Pr(> z)
Capability	0.94	0.32	2.93	3.39×10^{-3}
Intention	-1.80	0.38	-4.69	2.68×10^{-6}
Trust Residual	2.46	0.40	6.20	5.83×10^{-10}
Preference	-2.35	0.27	-8.86	$< 2 \times 10^{-16}$
Preference:Intent	3.99	0.47	8.55	$< 2 \times 10^{-16}$
Intercept	1.36	0.32	4.23	2.39×10^{-5}

scenarios, and tasks. The model’s main fixed effects were inferred robot capability and intention, as well as two additional variables: *residuals* from the self-reported trust model in H3, and participant preference. The residuals represented potential factors/components of trust that are separate from capability and intent. Participant preference was a binary indicator variable that (depending on the task) captured whether the participant preferred the higher risk option (for Tasks 1 and 3), or preferred the higher accuracy option (for Task 2).

Our results (summarized in Tbl. IV) show that decision to trust was significantly affected by both capability ($b = 0.94$, $SE = 0.32$, $p < 10^{-2}$) and intention ($b = -1.80$, $SE = 0.38$, $p < 10^{-5}$). There was a significant interaction between intention and participant preference ($b = 3.99$, $SE = 0.47$, $p < 10^{-16}$), which can be interpreted as the importance of *intention alignment*. For example, when subjects were risk-seeking, their decisions were positively associated with the degree of agreement that the robot was also risk-seeking. Conversely, when participants were risk-averse, their decisions were negatively associated with the robot’s inferred risk-preference.

The trust residuals were also significantly associated with the trust decision ($b = 2.46$, $SE = 0.40$, $p < 10^{-9}$), suggesting that other factors (e.g., robot appearance, environment or human related elements [11]) played a role in the eventual decision. However, a model with only self-reported trust as the independent variable (and otherwise the same as above), resulted in larger AIC and BIC scores (See Tbl. V where ΔAIC and ΔBIC are differences from the best model). This suggests that a multidimensional mental construct is applied to trust-based decisions. Indeed, removing different components from the initial decision model resulted in poorer quality candidate models, as shown in Tbl. V.

H5: People are more trusting of the simulated human agents, rather than the simulated robot agents. Fig. 5 shows the participants paired with simulated human agent (GH) reported significantly higher trust ($t(398) = -2.77$, $p = 0.01$),

TABLE V: ΔAIC and ΔBIC scores for four trust decision models. The “Complete Model” is the full trust decision model (described in the text) and achieves the best (lowest) AIC/BIC scores, compared to the alternative models.

Model	ΔAIC	ΔBIC
Complete Model	0	0
Without Intention	83.58	66.77
Without Capability	39.29	28.09
Without Trust Residuals	37.41	31.80
Trust Only	81.98	59.58

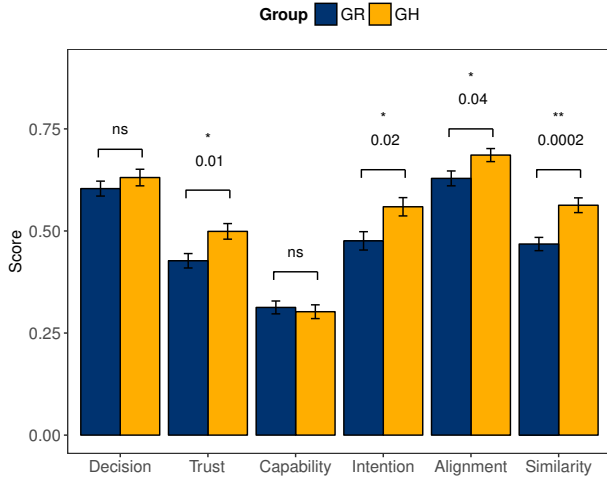


Fig. 5: Difference in dependent variables between the simulated human group (GH) and simulated robot group (GR), with p -values shown above the bars. Participants reported significantly higher trust, intention, alignment and similarity in GH compared to GR.

inferred intent scores ($t(398) = -2.63, p = 0.02$), inferred alignment ($t(398) = -2.36, p = 0.04$), and similarity in decision-making ($t(398) = -3.91, p < 10^{-2}$). However, differences in reported capability were not significant; potentially, participants considered capability as a property of the UAV since the other human player only controlled the location choice. Interestingly, differences in eventual trust decisions were also not statistically significant ($t(393.67) = -0.99, p = 0.32$); there is marginal evidence that participants tended to delegate more to the HC-HR robot ($t(98) = -1.57, p = 0.12$), but not to the other robot types ($p > 0.5$). In other words, humans tend to report higher trust in other people (rather than robots) and believe people share their decision-making objectives. But at the same time, the decision to trust appears to remain dependent on observed behavior. These mixed results suggest more investigation is needed to further elucidate the differences between inter-human and human-robot trust.

V. DISCUSSION

The results in the previous section show that both inferred capability and intent influence human decisions to trust the

robot. Humans appear to require that a robot demonstrate that it has similar intent (e.g., risk or accuracy preference in our setup), in addition to having the capability to execute the task successfully. However, inferred robot intention and capability are by themselves insufficient; other factors—captured by the overall trust—contribute towards trust-based decision-making. These findings suggest that humans make trust-based decisions using *rich mental representations* of robots, rather than relying solely on overall trust.

Trust in the simulated human agents and simulated robot agents were qualitatively similar, albeit with clear quantitative differences. Humans reported higher trust and inferred intent-alignment scores when partnered with simulated humans, which echo prior findings [16], [34]. This points to a potential difference in people’s prior expectations when working with humans versus machines. However, our failure to find statistically significant differences in decisions to trust suggests that trust decisions might vary for different agents, e.g., according to perceived similarity to human beings or specialization, and observed behavior.

Taken together, these results have implications for the development of computational trust models [7], [9] and robots that consider human trust in their decision-making process (e.g., [10]). In particular, if a robot aims to *predict* human behavior and act accordingly, tracking overall trust is inadequate when working in multiple task contexts. Instead, humans appear to internally represent agent capability and intention, allowing them to generalize appropriately to new scenarios. Our findings also advise the practical aspects of human-centric robots. For example, calibrating user inferred capability and intention in assistive robots (such as smart wheelchairs [35], [36]) may encourage adoption and proper usage.

VI. CONCLUSIONS

Trust in autonomous robots will become increasingly important as advancements in robotics are putting robots on our streets, in our factories, and in our homes. It is thus critical to study and model how people trust and delegate tasks to robots. In this paper, we presented human subjects with simulated task-delegation choices, and found that human decisions to delegate control in novel task contexts depend not only on overall trust in the robot collaborator, but also on estimations of robot capability and intention. That is, human trust in robots qualitatively mirrors human trust in other humans and is *multifaceted*, consisting of at least two important facets: capability and intention.

Our results add to a rich literature on factors that influence trust in robots (e.g., [11], [26]) and subsequent decision-making [10], [14], [37]. Future work should examine other potential facets of human-robot trust, and provide empirical evidence for transfer across a wider range of tasks contexts. We envision a more complete theory of human-robot trust would contribute towards a collaborative trust-based society comprising both human and robot agents.

REFERENCES

- [1] V. Groom and C. Nass, "Can robots be teammates?: Benchmarks in humanrobot teams," *Interaction Studies*, vol. 8, no. 3, pp. 483–500, 2007.
- [2] T. B. Sheridan and R. T. Hennessy, "Research and modeling of supervisory control behavior. report of a workshop," National Research Council Washington DC Committee on Human Factors, Tech. Rep., 1984.
- [3] A. Freedy, E. DeVisser, G. Weltman, and N. Coeyman, "Measurement of trust in human-robot collaboration," in *Collaborative Technologies and Systems, 2007. CTS 2007. International Symposium on*. IEEE, 2007, pp. 106–114.
- [4] P. A. Hancock, D. R. Billings, K. E. Schaefer, J. Y. Chen, E. J. De Visser, and R. Parasuraman, "A meta-analysis of factors affecting trust in human-robot interaction," *Human Factors*, vol. 53, no. 5, pp. 517–527, 2011.
- [5] B. Muir, "Trust in automation: Part I. Theoretical issues in the study of trust and human intervention in automated systems," *Ergonomics*, vol. 37, no. 11, pp. 1905–1922, 1994.
- [6] B. Muir and N. Moray, "Trust in automation. Part II. Experimental studies of trust and human intervention in a process control simulation," *Ergonomics*, vol. 39, no. 3, pp. 429–460, 1996.
- [7] A. Xu and G. Dudek, "OPTIMO: Online Probabilistic Trust Inference Model for Asymmetric Human-Robot Collaborations," *Proceedings of the Tenth Annual ACM/IEEE International Conference on Human-Robot Interaction - HRI '15*, pp. 221–228, 2015.
- [8] M. Desai, "Modeling trust to improve human-robot interaction," Ph.D. dissertation, University of Massachusetts Lowell, 2012.
- [9] H. Soh, P. Shu, M. Chen, and D. Hsu, "The Transfer of Human Trust in Robot Capabilities across Tasks," *RSS*, 2018.
- [10] M. Chen, S. Nikolaidis, H. Soh, D. Hsu, and S. Srinivasa, "Planning with Trust for Human-Robot Collaboration," in *Proceedings of the 2018 ACM/IEEE International Conference on Human-Robot Interaction - HRI '18*. New York, New York, USA: ACM Press, 2018, pp. 307–315.
- [11] P. A. Hancock, D. R. Billings, K. E. Schaefer, J. Y. C. Chen, E. J. de Visser, and R. Parasuraman, "A Meta-Analysis of Factors Affecting Trust in Human-Robot Interaction," *Human Factors*, vol. 53, no. 5, pp. 517–527, 2011.
- [12] K. E. Schaefer, J. Y. Chen, J. L. Szalma, and P. A. Hancock, "A meta-analysis of factors influencing the development of trust in automation: Implications for understanding autonomy in future systems," *Human factors*, vol. 58, no. 3, pp. 377–400, 2016.
- [13] C. Castelfranchi and R. Falcone, *Trust Theory*. John Wiley and Sons Ltd, 2007.
- [14] S. H. Huang, K. Bhatia, P. Abbeel, and A. D. Dragan, "Establishing Appropriate Trust via Critical States," *IROS 2018*, 2018.
- [15] R. C. Mayer, J. H. Davis, and F. D. Schoorman, "An integrative model of organizational trust," *Academy of management review*, vol. 20, no. 3, pp. 709–734, 1995.
- [16] P. Madhavan and D. a. Wiegmann, "Similarities and differences between humanhuman and humanautomation trust: an integrative review," *Theoretical Issues in Ergonomics Science*, vol. 8, no. 4, pp. 277–301, 2007.
- [17] D. Gambetta, "Can we trust trust," *Trust: Making and breaking cooperative relations*, vol. 13, pp. 213–237, 2000.
- [18] J. D. Lee and K. A. See, "Trust in Automation: Designing for Appropriate Reliance," *Human Factors: The Journal of the Human Factors and Ergonomics Society*, vol. 46, no. 1, pp. 50–80, 2004.
- [19] A. Jøsang, R. Ismail, and C. Boyd, "A survey of trust and reputation systems for online service provision," *Decision support systems*, vol. 43, no. 2, pp. 618–644, 2007.
- [20] T. Yamagishi, M. Kikuchi, and M. Kosugi, "Trust, gullibility, and social intelligence," *Asian Journal of Social Psychology*, vol. 2, no. 1, pp. 145–161, 1999.
- [21] J. LEE and N. MORAY, "Trust, control strategies and allocation of function in human-machine systems," *Ergonomics*, vol. 35, no. 10, pp. 1243–1270, oct 1992.
- [22] S. H. Huang, K. Bhatia, P. Abbeel, and A. D. Dragan, "Establishing (Appropriate) Trust via Critical States," *HRI 2018 Workshop: Explainable Robotic Systems*, 2018.
- [23] M. S. Carlson, M. Desai, J. L. Drury, H. Kwak, and H. a. Yanco, "Identifying Factors that Influence Trust in Automated Cars and Medical Diagnosis Systems," in *The Intersection of Robust Intelligence and Trust in Autonomous Systems: Papers from the AAAI Spring Symposium*, 2011, pp. 20–27.
- [24] D. G. Shapiro, T. Könik, and P. O'Rorke, "Achieving far transfer in an integrated cognitive architecture," in *AAAI*. AAAI Press, 2008, pp. 1325–1330.
- [25] S. H. Huang, D. Held, P. Abbeel, and A. D. Dragan, "Enabling robots to communicate their objectives," *Autonomous Robots*, pp. 1–18, 2018.
- [26] N. Wang, D. V. Pynadath, and S. G. Hill, "Trust calibration within a human-robot team: Comparing automatically generated explanations," in *2016 11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, mar 2016, pp. 109–116.
- [27] R. Hall, "Trusting your assistant," in *Proceedings of the 11th Knowledge-Based Software Engineering Conference*. IEEE Comput. Soc. Press, 1996, pp. 42–51.
- [28] M. Desai, "Modeling trust to improve human-robot interaction," Ph.D. dissertation, University of Massachusetts Lowell, 2012, aAI3537137.
- [29] A. Xu and G. Dudek, "Towards Modeling Real-Time Trust in Asymmetric HumanRobot Collaborations," in *Robotics Research*, 2016, pp. 113–129.
- [30] B. M. Muir, *Operators trust in and percentage of time spent using the automatic controllers in a supervisory process control task*. University of Toronto, 1989.
- [31] G. Hoffman, "Evaluating Fluency in Human-Robot Collaboration," *HRI Workshop on Human Robot Collaboration*, 2013.
- [32] J.-Y. Jian, A. M. Bisantz, and C. G. Drury, "Foundations for an Empirically Determined Scale of Trust in Automated Systems," *International Journal of Cognitive Ergonomics*, vol. 4, no. 1, pp. 53–71, mar 2000.
- [33] K. E. Schaefer, "Measuring Trust in Human Robot Interactions: Development of the Trust Perception Scale-HRI," in *Robust Intelligence and Trust in Autonomous Systems*. Boston, MA: Springer US, 2016, pp. 191–218.
- [34] M. G. Collins, I. Juvina, and K. A. Gluck, "Cognitive Model of Trust Dynamics Predicts Human Behavior within and between Two Games of Strategic Interaction with Computerized Confederate Agents," *Frontiers in Psychology*, vol. 7, no. February, p. 49, 2016.
- [35] H. Soh and Y. Demiris, "When and how to help: An iterative probabilistic model for learning assistance by demonstration," in *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*, 2013, pp. 3230–3236.
- [36] —, "Learning Assistance by Demonstration: Smart Mobility With Shared Control and Paired Haptic Controllers," *Journal of Human-Robot Interaction*, vol. 4, no. 3, pp. 76–100, 2015.
- [37] P. Robinette, A. M. Howard, and A. R. Wagner, "Effect of Robot Performance on Human-Robot Trust in Time-Critical Situations," *IEEE Transactions on Human-Machine Systems*, vol. 47, no. 4, pp. 425–436, 2017.