

Robot Navigation in Risky, Crowded Environments: Understanding Human Preferences

Aamodh Suresh*, Angelique Taylor**, Laurel D. Riek**, Sonia Martínez*

Abstract—Risky and crowded environments (RCE) contain abstract sources of risk and uncertainty, which are perceived differently by humans, leading to a variety of behaviors. Thus, robots deployed in RCEs, need to exhibit diverse perception and planning capabilities in order to interpret other human agents’ behavior and act accordingly in such environments. To understand this problem domain, we conducted a study to explore human path choices in RCEs, enabling better robotic navigational explainable AI (XAI) designs. We created a novel COVID-19 pandemic grocery shopping scenario which had time-risk tradeoffs, and acquired users’ path preferences. We found that participants showcase a variety of path preferences: from risky and urgent to safe and relaxed. To model users’ decision making, we evaluated three popular risk models (Cumulative Prospect Theory (CPT), Conditional Value at Risk (CVAR), and Expected Risk (ER). We found that CPT captured people’s decision making more accurately than CVaR and ER, corroborating theoretical results that CPT is more expressive and inclusive than CVaR and ER. We also found that people’s self assessments of risk and time-urgency do not correlate with their path preferences in RCEs. Finally, we conducted thematic analysis of open-ended questions, providing crucial design insights for robots in RCE. Thus, through this study, we provide novel and critical insights about human behavior and perception to help design better navigational explainable AI (XAI) in RCEs.

I. INTRODUCTION

Robots are increasingly being deployed in everyday risky and crowded environments (RCE), including shopping malls, museums, streets, and sidewalks (i.e., autonomous cars) [1]. These environments are often crowded, contain multiple sources of risk (e.g., dynamic and chaotic human-motion trajectories) and uncertainty (e.g. noisy sensor measurements, including those from camera ego-motion [2]). As robots become more integrated into such environments, they need to appropriately deal with these challenges and navigate in a safe and socially-acceptable manner [1], [3], [4], [5].

Modeling of how humans perceive risk [6] can help us understand and close this gap. These models differ on the degree of rationality assumptions made on the human when subject to risky choices. These range from the consideration of human behavior as completely rational and possibly risk-averse (e.g., Expected Risk (ER), Conditional Value at Risk (CVaR) [7]) to non-rational and possibly risk-insensitive (e.g., Cumulative Prospect Theory (CPT) [8]).

However, little is known about the validity of these models in a risky social navigation setting, as well as how they

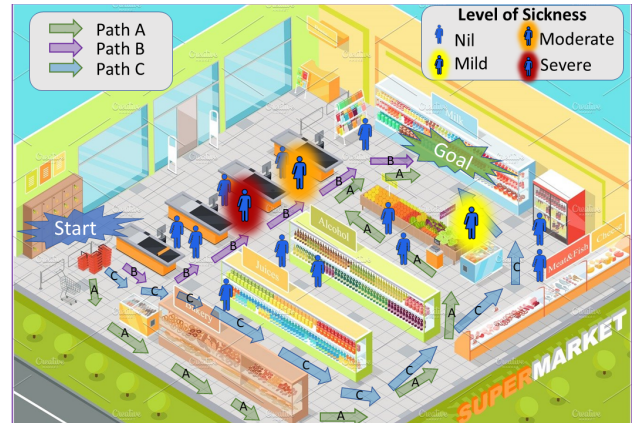


Fig. 1: Grocery store environment used in user studies. Participants select one of three paths (A, B, and C), to go from the entrance (shown as the ‘Start’) to the milk section (shown as the ‘Goal’). The supermarket is crowded with people with levels of sickness ranging from ‘Nil’ to ‘Severe’.

compare with humans’ self perception of risk. In particular, we are interested in understanding how robots can reason with humans, explaining their behaviors and actions, also known as Explainable Artificial Intelligence (XAI) [9]. XAI “explains” itself by opening up its reasoning to human scrutiny, resulting in better, faster, more accurate and more aligned human-robot decisions [10], [11].

Prior Work: Risk is a relevant notion of urgency used to design navigation algorithms in robotics [12]. Accordingly, various models have been employed to quantify and reason about risk. CVaR is one such popular model adopted from finance in robotics [13], [12], which captures risk aversion (i.e., “play it safe”) by employing linear and rational notions of decision making. While this is analytically convenient, it cannot capture non-linear and non-rational decision making that humans usually exhibit [14], [15], [16].

Recently, CPT methods [8] have been proposed [17], [18] to address this shortcoming. Theoretically, it has been shown that CPT is more “expressive” [17], “versatile”, and “inclusive” [18] than CVaR and Expected Risk (ER), thus capturing a wider range of risk profiles of humans. Preliminary evidence that CPT better captures human decision making under risk can be found in applications of traffic intersection management and routing [19], and resource management settings by operators [20]. In practice, these approach is yet to be evaluated extensively in user studies pertaining navigation in RCE.

To do so, user studies that employ natural or explainable metrics to humans need to be developed. Unfortunately, commonly used risk variables such as money [21], time [22], or collision probabilities [23], do not satisfy this criterion for all cases. In fact, recent studies have found that humans are often

This work was supported by grants XXX-XXXX-XXX

* A. Suresh and S. Martínez are with the Department of Mechanical and Aerospace Engineering, University of California San Diego, La Jolla, CA 92093, USA. E-mail: {aasuresh,soniamd}@eng.ucsd.edu.

** A. Taylor and L. Riek are with the Department of Computer Science and Engineering, University of California San Diego, La Jolla, CA 92093, USA. E-mail: {amt062,lriek}@eng.ucsd.edu.

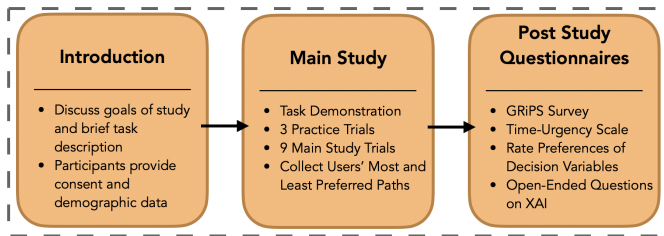


Fig. 2: An overview of our study design.

sub-optimal in planning paths in such situations [24]. These studies assume that the human is either “noisy-rational” or do not have correct environment models to choose optimally.

A few other avenues of using risk for planning paths in RCE include fall risk assessment [25], [26], [27], risk of localization and mapping systems [28], [29], and planning risk in search and rescue operations [30]. These arguments are from a robot’s perspective which acts in an expected manner and also expects the human to do so. However, from a human-centered and XAI perspectives, the robot’s “expected” behavior might lead to mistrust and confusion [31], [32]. To the best of our knowledge, general studies pertaining to everyday scenarios that employ more abstract cost interpretations are lacking, and are needed for better explainable AI design for robots in RCEs.

Contributions: In this work, through the design of a novel user study, we bridge the gap in existing literature by characterizing human perception of risk in RCEs, comparing theoretical risk models with observed human responses, and by exploring the consistency of human perception of risk and time urgency with standard survey responses. In addition, we provide new valuable insights for XAI design. Specifically, our work aims to address the following research questions:

RQ 1: What is the relationship between participants’ path preferences and those arising from standard risk models?

RQ 2: What is the relationship between participants’ self-risk and self-time-urgency perception and their actual path choices?

RQ 3: How do humans relatively weight time and risk to make navigational decisions in everyday scenarios?

RQ 4: What are the users’ preferences to interact with robots navigating in everyday scenarios?

We conducted a large scale online study ($n = 82$) and found that most participants do not make decisions in an expected manner (in accordance with expected risk metric) and that CPT as a risk model captures the observed responses better than CVaR and ER. Interestingly, through the application of standard questionnaires, we find that there is a mismatch between humans’ self-risk/self-time-urgency assessment and their actual choices. Additionally, participants generally give a higher weight to risk than time while choosing paths.

Finally, we provide valuable insights to design XAI for robots in RCEs. For example, we found that most participants want robots that can explain its rationale behind decision-making and they also suggested user interface design to have a two-way motion intention communication between users and robots. Thus, equipped with these results and insights, XAI design can be improved to enable robots to operate and adapt to human preferences in RCEs.

II. METHODOLOGY

We conducted an IRB-approved (approval code: 201638) within-subjects study on the Qualtrics¹ survey platform. The study was designed to evaluate people’s risk perception, and was inspired by the COVID-19 pandemic. This provides an easy-to-relate context for participants to think about decision making under risk. Thus, we consider a grocery-store shopping scenario, where the risk is characterized by being coughed at by potentially infected people. Participants were asked to imagine being an “Instacart² Shopper” who needs to go from the entrance of the store to the milk section. Time-urgency is characterized by the need to complete shopping quickly in order to get better ratings and tips.

This scenario is illustrated in Figure 1. Here, participants had three paths to choose from. Each path had a varying intensity of risk and time urgency (discussed in detail in Section II-B). The participants indicated their most and least preferred paths for each scenario. In the following paragraphs, we explain the scenario methodology and the list the post-study questions that we use (see Tables I and II).

Participants: We recruited 82 participants affiliated with a university campus through university list-serves and via word of mouth. The participants consisted of 27 females, 49 males, 1 binary/third gender and 5 that preferred not to answer this question. The ages ranged from 21-32 (mean = 25.6, SD = 2.5) and their educational background had a distribution of 68 in Engineering, 3 in Mathematics, 5 in Basic Sciences, 1 in Management, and 5 from other fields.

A. Measures and metrics employed

We obtained the following measures and metrics from each participant through nine trials and post-study questionnaires:

1) *Path preferences:* For each trial the participant revealed their path preference order by providing their most preferred path (MPP) and Least preferred path (LPP). The MPP is denoted as $M \triangleq \{m_1, m_2, \dots, m_9\}$, for all 9 trials. For a trial i , we encode the user’s MPP choice as $m_i \triangleq 0$ for path A (resp. $m_i \triangleq 1$ if they chose path B, and $m_i \triangleq 0.5$ if they chose C), indicating the level of relative risk and time-urgency. Similarly the Least preferred path (LPP) is denoted as $L \triangleq \{l_1, l_2, \dots, l_9\}$, for all 9 trials. We encode $l_j \triangleq 0$ for path A (resp. $l_j \triangleq 1$ if they chose path B, and $l_j \triangleq 0.5$ if they chose C), indicating the level of relative risk and time-urgency). Using these measures we describe the risk taking behavior of the participants.

2) *Observed risk taking behavior:* We characterize participants’ decision making into three different categories based on the risk taking behaviors expressed. First, we consider “expected behavior” using Expected Risk (ER), from which we get M^{exp} and L^{exp} . These are the MPP and LPP choices when risk is defined in an expected manner. Next, we observe “risk aversion” through CVaR, where we similarly obtain M^{av} and L^{av} . Finally, we consider “risk insensitive” behavior using CPT with M^{ins} and L^{ins} as the MPP and LPP choices, respectively. We note that from our previous theoretical results [17], [18], we have shown that CPT is the most

¹<https://www.qualtrics.com/>

²A grocery delivery service (www.instacart.com).

TABLE I: The 8-item the General Risk Propensity Scale (GRiPS) [33] that we administered to participants after engaging in our study.

GRiPS Survey Questions
1. Taking risks makes life more fun
2. My friends would say that I'm a risk taker
3. I enjoy taking risks in most aspects of my life
4. Taking risks is an important part of life
5. I commonly make risky decisions
6. I am a believer of taking chances
7. I would take a risk even if it meant I might get hurt
8. I am attracted, rather than scared, by risk

inclusive model and can capture all three perceptions. CVaR can capture expected and risk averse perception, whereas ER only captures the expected behavior. In the following, we list the corresponding metrics, and how they are calculated.

- Average MPP score, \bar{M} : This is the average risk and time-urgency of the participants' MPP, $\{m_1, m_2, \dots, m_9\}$.
- Average LPP score, \bar{L} : This is the average risk and time-urgency of the participants' LPP, $\{l_1, l_2, \dots, l_9\}$.
- Deviation from expected behavior, $J^{exp} = \sum_{i=1}^9 |m_i - m_i^{exp}|$: Larger values indicate greater deviations from the expected behavior.
- Deviation from risk aversion, $J^{av} = \sum_{i=1}^9 |m_i - m_i^{av}|$: Larger values indicate greater deviations from risk-averse (CVaR) behavior.
- Deviation from risk insensitivity, $J^{ins} = \sum_{i=1}^9 |m_i - m_i^{ins}|$: Larger values indicate greater deviations from risk insensitive behavior.
- Deviations from RPMs: For ER we have $J^{ER} = J^{exp}$. For CVaR we have $J^{CVaR} = \min\{J^{exp}, J^{av}\}$ and for CPT we have $J^{CPT} = \min\{J^{exp}, J^{av}, J^{ins}\}$.

Next, we look at the self-reported measures that the users provide us through post-study questionnaires.

3) *Self-reported measures*: After collecting user path choices, we then conducted these post-study questionnaires.

GRiPS: General Risk Propensity Scale (GRiPS) [33], which measures the participants' self risk-taking abilities, i.e. it evaluates how risk-averse or risk-taking they think they are in their daily lives. GRiPS is a self-report measure (see Table I) of general risk and pro-social behavior consisting of 8-items which participants answer on a Likert scale from 1 (Strongly Disagree) to 5 (Strongly Agree). We denote the responses for the GRiPS questionnaire as $R \triangleq \{r_1, r_2, \dots, r_8\}$ for the 8 questions.

Time Urgency: The second questionnaire is called the "Time Urgency Scale" [34], which measures participants' self assessment of how time-urgent and urgent they think they behave in everyday scenarios. It is a self-report measure (See Table II) of general time-related behavior consisting of 6-items (as commonly used [35]), which participants answer on a Likert scale from 1 (Strongly Disagree) to 5 (Strongly Agree). We denote the responses for the Time-Urgency questionnaire as $T \triangleq \{t_1, t_2, \dots, t_6\}$ for the 6 questions.

We employ metrics to evaluate the participants' self perception of risk and their observed risk-taking behavior, which are described next.

- Risk score* \bar{R} : The average response (normalized between 0 and 1) from the GRiPS survey responses R

TABLE II: The 6-item the Time Urgency Scale [34] that we administered to participants after engaging in our study.

Time Urgency Survey Questions
1. I find myself hurrying to get places even when there is plenty of time.
2. I often work slowly and leisurely.
3. People that know me well agree that I tend to do most things in a hurry.
4. I tend to be quick and energetic at work.
5. I often feel very pressed for time.
6. My spouse or a close friend would rate me as definitely relaxed and easy going.

(as commonly used [36]). Here, a $\bar{R} = 1$ indicates an adventurous perception, whereas $\bar{R} = 0$ indicates a risk-averse perception.

- Time urgency score* \bar{T} : The average response (normalized between 0 and 1) from the Time-urgency survey responses T . Here a $\bar{T} = 1$ indicates a hasty behavior, whereas $\bar{T} = 0$ indicates a relaxed behavior.
- Risk similarity score* $R^{sim} = \bar{M} - \bar{R}$: measures the deviation from users' self-risk perception and their observed risk perception. $R^{sim} \in [-1, 1]$, where $R^{sim} \approx 0$ indicates user's self-perception of risk and their observed risk taking characteristics are similar. A larger positive (negative) value indicates users chose riskier (safer) paths than their self perception through the questionnaire.
- Time-urgency similarity score* $T^{sim} = \bar{M} - \bar{T}$: Shows if people choose paths according to the time-urgency survey responses. As above, a $T^{sim} \in [-1, 1]$, $T^{sim} \approx 0$, indicates similarity in user's self perception and observed behavior of time urgency. Larger positive (negative) values indicate users choose more urgent (leisured) paths than their self perception.

Decision Variable Preferences: In addition, after the study trials and questionnaires are administered, we ask participants how they relatively weighed (as a %) each of the four variables in making their decisions: that is, the time taken, the number of sick people, the level of sickness, and the chance of being coughed at for a particular path. The participants' relative reliance (in %) on each decision variable is denoted by $\{v_1, v_2, \dots, v_4\}$ corresponding (Table III) to the four variables.

4) *Open ended questions*: We finally asked open-ended questions (listed below) to better understand human preferences towards designing robot navigation in RCE.

- Q1: Would you like to know how robots make decisions and plan paths? If so, how do you want a robot to explain its thought processes behind its decisions and what modality (e.g., speech, expressions) would you prefer?
- Q2: How do you want the robot to communicate its movement intentions (e.g., moving right or left)?
- Q3: Would you like a robot to know how you are making decisions and planning paths? If so, how do you want to explain your intentions and what modality (e.g., speech, touchscreen) would you choose?

To summarize, the trial data helped us answer RQ 1. Question RQ 2 can be assessed from analyzing the trial data along with the questionnaire responses. The data on decision-variable preferences helped us answer RQ 3. Whereas the responses to open-ended questions helped us answer RQ 4.

TABLE III: Description of decision variables and their ranges for each path in every scenario.

No.	Decision Variables	Range of values presented
1	Time Taken	Path A : 20 mins Path B : 5 mins Path C : 10 mins
2	Number of Sick people	Path A : 0-1 Path B : 2-3 Path C : 1-2
3	Level of Sickness	0-3 for each path
4	Chance of being coughed at	0-100 % for each path

B. Scenarios

Participants were presented with three choices of paths to choose from, including paths A, B, and C (see Figure 1). Path A was the longest, path B was the shortest, and path C had a length that was in between. We used the situation of “being coughed at by sick people” to elicit risk for each path in every scenario. This risk was described by four decision variables: “time taken”, “number of sick people”, “level of sickness”, and “chance of being coughed at” (see Table III). The time taken varied from 5 to 20 minutes, the number of sick people from 0 to 2, the level of sickness from 0 to 3, and the chance of being coughed at was expressed as a percentage for each sick person encountered. Regarding the level of sickness, we used the following terminology: 0-Nil, 1-Mild, 2-Moderate and 3-Severe (see Figure 1). We purposefully kept the consequences of being exposed to a sick person abstract, in order to extract realistic risk perceptions from participants. The variables are summarized in Table III.

We administered nine trials with different values for decision variables, aimed at capturing a wide range of scenarios. In each trial, the shortest path (Path B) had the most risk and uncertainty, while the longest path (Path A) had the least risk and uncertainty. Participants choose their most and least preferred paths, thus providing a preference order. The risk variables for each trial are designed in such a way that the best expected choice³ (A or B or C) varies across trials. We then group the nine trials into three levels of uncertainty (w.r.t. number of sick people). Comparing the chosen preferences with those based on expected risk, and across different levels of uncertainty, helps us understand how human choices are affected by uncertainty.

C. Procedure

An overview of our study design is available in Figure 2. The study began with a consent form providing a brief description of the study. After giving informed consent, participants provided their demographic information including age, gender, occupation, and area of expertise. The main study then started, with a discussion the goal of the study, which is to investigate how people choose paths in risky situations. We further explained the user interface (UI) and elements of the study through a demonstration trial.

Participants then saw a demonstration scenario of the grocery store (Fig. 1), describing the four decision variables pictorially, in sentences and through a summary table⁴ (see

³That is, best according to the expected risk metric.

⁴These decision variables were selected after a few rounds of pilot studies which indicated different people prefer different variables to describe the scenarios.

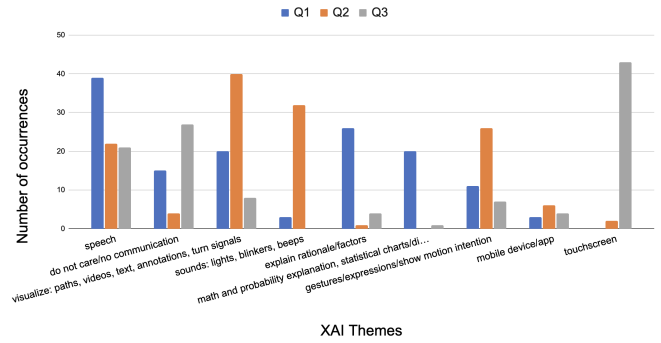


Fig. 3: This shows the number of occurrences of each explainable AI (XAI) theme identified in the data collected in responses to three open-ended questions from Section II-A to address RQ4.

Table III). Then, they selected their most and least favorite paths. Based on their most preferred path choices, we randomly selected a risk outcome and display the final results⁵ (e.g., ‘You encountered no sick people’).

Next, participants engaged in three practice rounds with three different scenarios and selected their most and least favorite paths. After the practice rounds, they then participated in the main study, which presented 9 different combinations of “risk” in each scenario. To remove ordering effects and the influence of regret, we randomized the ordering of the trials across all participants.

At the conclusion of the study, we administered the GRiPS [33] and Time-Urgency [34] questionnaires, to respectively measure self risk-taking and time-urgency perception among participants. In addition, we asked participants to relatively weigh (as a percentage of) each of the four variables to make their path choices. Finally, we asked open-ended questions (listed in Section II-A) to better understand their preferences towards designing explainable AI for path planning in everyday scenarios.

D. Analysis

RQ1 Analysis: We examined the relationship between participants’ path choices and the the corresponding choices that result from a risk model. To do this, we analyzed the descriptive statistics of the objective metrics that measure deviation from various behaviors and risk models. (As we had only 9 trials, the sample was insufficient for parametric correlation analysis.)

RQ2 Analysis: We explored the relationship between participants self-perception of risk and time-urgency, compared to the expected (baseline) risk associated with their path choices in terms of most and least preferred paths. In this regard, we analyzed the descriptive statistics of the subjective metrics (\bar{R} and \bar{T}) and the similarity scores (R^{sim} and T^{sim}). We also conducted a correlation analysis to understand the interaction between self-perception of risk and time-urgency. We also performed a pairwise comparison of means (paired t-test) to reveal additional trends between variables representing participants choices and self perception.

RQ3 Analysis: We studied the relative importance that users give to each decision variable (in Table III) in order to

⁵Through pilot studies we learned that displaying the trial results enhanced user engagement.

TABLE IV: Descriptive statistics of variables to compare users’ decisions with standard risk-model decisions to address RQ 1.

Variable	Mean	95% confidence interval of mean	Standard Deviation	Range
\bar{M}	0.60	0.56 to 0.66	0.20	0.00 to 1.00
\bar{L}	0.40	0.35 to 0.47	0.30	0.00 to 0.94
J^{exp} for MPP	0.40	0.35 to 0.40	0.10	0.11 to 0.61
J^{av} for MPP	0.30	0.31 to 0.36	0.10	0.00 to 0.61
J^{ins} for MPP	0.30	0.28 to 0.36	0.20	0.00 to 0.89
J^{CVaR} for MPP	0.30	0.31 to 0.36	0.10	0.00 to 0.61
J^{CPT} for MPP	0.20	0.21 to 0.26	0.10	0.00 to 0.44

choose their most preferred and least preferred paths. To do this, from this data, we created two new variables to measure the relative importance of time and risk used to make decisions. This was done by first averaging the three variables representing risk, and then expressing it as a percentage w.r.t. the total of time taken and average risk percentages. We provide descriptive statistics on the responses to post-study question on decision variable preferences.

RQ4 Analysis: We determined user preferences for XAI design through open-ended questions. Two members of our team performed thematic coding on open-ended question responses following grounded theory [37] as commonly done in the literature. This process involved reviewing responses (see Section II-A), generating high-level themes to capture key ideas in the data, reviewing the thematic codes with the team, negotiating them based on key ideas in the data, and repeating the process until all codes have been agreed upon. Next, we coded a total of 246 codes with our final set of codes shown in the x-axis of Figure 3. We computed the inter-rater agreement using Krippendorffs-Alpha as we used multiple codes for each quote. This is advantageous because it supports categorical and ordinal data. We found an IRR of 1.0 which is considered high agreement. We believe this is due to sparse responses with 2126 words total across 82 participants. There was an average of 8.64 words per response, a minimum of 1, a maximum of 102, and a median of 4 words per response across the dataset.

III. RESULTS

We provide descriptive statistics of the relevant variables and metrics, including the mean, median, standard deviation, and 95% confidence interval. To study the correlation between two variables, we calculate the Pearson’s correlation coefficient, along with null-hypothesis significance testing with threshold p-value = 0.05.

A. RQ1: Comparing users’ risk perception

In order to compare risk models with users’ decision-making, we provide descriptive statistics of the relevant variables, which are summarized in Table IV. Recall that MPP (similarly, LPP) for a j^{th} trial with $m_j = 0$ is path A with lowest risk and most leisured. Whereas $m_j = 1$ is path B with the highest risk and most time-urgency, and Path C is in between in both risk and time-urgency with a value of $m_j = 0.5$.

Path choice characteristics: From Table IV, the mean, median, and confidence interval are over 0.5 for the average \bar{M} , and under 0.5 for the average \bar{L} . This indicates a preference towards Path B (more risky and time urgent), and

TABLE V: Descriptive statistics of similarity scores of users’ decisions in trials compared to their questionnaire responses to address RQ 2.

Variable	Mean	95% confidence interval of mean	Standard Deviation	Range
\bar{R}	0.6	0.55 to 0.64	0.2	0.00 to 1.00
\bar{T}	0.5	0.48 to 0.56	0.2	0.08 to 0.96
R^{sim} for MPP	0.10	-0.05 to 0.09	0.30	-0.84 to 0.94
R^{sim} for LPP	0.00	-0.08 to 0.07	0.30	-0.91 to 0.94
T^{sim} for MPP	0.10	0.04 to 0.15	0.30	-0.58 to 0.67
T^{sim} for LPP	0.10	0.00 to 0.14	0.30	-0.58 to 0.67

a disinclination towards Path A (more safe and leisured). Additionally, \bar{M} exhibits a full range from 0 – 1, whereas \bar{L} has a range 0 – 0.94. Thus, although their preferences were scattered across the full spectra, no participant disliked path B across all trials.

Behavioral Characteristics and model comparisons:

From Table IV, we note that the deviation from the expected behavior J^{exp} is in general greater than J^{av} and J^{ins} , indicating that expected behavior (from expected risk) is the least aligned with participants’ preferences. Also, from row 3, the min is > 0 , indicating that not a single participant showed expected behavior across all their trials. The almost-similar statistics for risk-averse and risk-insensitive behaviors show that participants’ exhibited both of these behaviors equally frequently. Also, the deviation J^{CPT} is the least, showing that CPT is a better model to approximate the participants’ decision making. Similar deviation statistics were correspondingly obtained for LPP choices; hence, we omit its discussion here.

B. RQ2: Users’ Self-Perception of Risk vs. Expected Risk

We provide descriptive statistics of the relevant variables (Table V) and conduct correlation studies next. From Table V, the statistics of the average survey responses \bar{R} indicate that, in general, participants are more inclined towards taking risks, as the mean, median, and confidence interval are all over 3. However, a wide range and high standard deviation suggest a fairly diverse set of risk-taking behaviors. A similar trend is observed for time-urgency \bar{T} . To compare the participants’ decision making characteristics with their self perception of time-urgency and risk, we first consider the similarity scores R^{sim} and T^{sim} (Table V). The risk similarity score R^{sim} for MPP and LPP is balanced with the mean, median, and confidence interval close to 0, but have a high range and standard deviation. Hence, there are a variety of people with different perceptions, and the GRiPS responses may not fully represent their decision making in the study. A similar trend is observed for time urgency, with a slightly greater inclination of participants to act more urgently than what they indicate in the survey.

Next, we try to identify trends between the GRiPS and Time-Urgency survey responses and path choices by measuring correlation and performing linear regression. We highlight the significant results here. There was a significant interaction between the average survey responses for risk similarity score R^{sim} and time urgency similarity score T^{sim} for MPP and LPP with and $p < 0.05$ and a effect size < 0.5 . This implies that people who acted more/less riskier than they indicated in the GRiPS survey, also acted correspondingly more/less time urgent than they indicated in

the time-urgency survey. This can arise because of the study construction, where paths which are shorter (time urgent) are also riskiest and vice versa.

Interestingly, there was an insignificant interaction between the normalized GRiPS response \bar{R} with MPP and LPP scores \bar{M}/\bar{L} with $p = 0.936$ and $p = 0.655$, and effect size close to 0. This reveals that participants' self risk evaluation and their path choices are not related. The trend was similar in the case of time-urgency. Thus, relying solely on the GRiPS and/or time-urgency survey to depict participants' behavior may not be effective.

Table VI shows the paired t-test statistics of various pairwise means comparisons. Interestingly we found significant evidence that \bar{L} is greater than \bar{R} ($p < 0.001$) with a medium effect size ($d > 0.5$). On average, \bar{L} is greater than \bar{R} by 0.18 (95% CI: 0.1, 0.26) units. By this, we can infer that participants had more disinclination towards riskier paths than their indicated risk appetite through the GRiPS survey.

C. RQ3: Users reliance on decision variables

We measured the relative importance that participants give to the four decision variables (Table III) to choose paths in terms of percentages. The data is described as boxplots (first 4 from the left) in Figure 4a. The last two boxplots represent the relative time and risk consideration, respectively. We note that the relative consideration of time has a mean 40.8% and median 42.8%, as opposed to relative consideration of risk which has mean 59.2% and median 57.2. So, in general, participants seem to consider risk factors more importantly than time factor while making decisions. However, with a large standard deviation of 30.2 and a full range of 0 – 100 for both variables, the generalization may not apply to many participants. This again reflects the diversity regarding time and risk consideration for making path planning decision by humans. Thus, XAI needs flexible models in this regard.

D. RQ 4: Users' Interaction Preferences

We asked participants three open-ended questions from Section II-A to gauge their preferences w.r.t. how they would like robots to communicate their intentions. After conducting the analysis as described in Section II-D, we found many redundant responses for the three questions. Thus, we discuss descriptive statistics for each question and describe the overarching themes we found in all responses.

Descriptive statistics of open-ended questions: We are interested in learning about users' preferences for XAI-capable robots in everyday settings. However, we found it challenging to fully address our research questions without also collecting qualitative data about users' preferences. Thus, we asked the open-ended questions from Section II-D.

Figure II-A shows the descriptive statistics of the XAI themes used to code the data across all open-ended questions. In summary, there are ten XAI themes identified in the data (x-axis of Figure II-D). From here, 73/82 participants indicated interest in XAI systems, and 8/82 participants did not want robots that explain their motion intentions. On the other hand, 69 out of 82 participants wanted robots that understand how users in their environment plan paths, 10 out of 82 do not, and 3 out of 82 participants indicated

maybe. The 'touchscreen' theme achieved the highest occurrence count. Several themes achieved the lowest occurrence score of 0 which include themes 'sounds, lights, blinkers, beeps' for how the robot knows what users want, 'math, probability, explanation, statistical charts/diagrams' for how to communicate intent, or 'touchscreen' in terms of how they would like robots to plan paths.

Thematic coding of Open-Ended Questions: We identified four overarching themes that inform users' preferences for XAI robotic navigational systems. These themes emphasize the importance of data variables used to interact with XAI systems for mobile platforms, how robots should communicate their navigation plans, and users' concerns about using XAI for mobile robots in everyday environments.

What data robots should use to communicate with users: Users identified a range of preferences for data robots can use to communicate its motion intentions. The most popular modality discussed by 62 out of 82 participants was speech, as it is convenient to communicate naturally. More specifically, 60 out of 82 participants discussed the need to state rationale visually and mathematically such as with pie charts, using arrows, or visual representation of weights used in AI decision-making. Additionally, 26 out of 82 participants envisioned robots that can explain the factors behind their decision, consistent with prior work in XAI [10], [11] e.g., level of sickness during COVID-19. Lastly, participants discussed the need for mobile robots that adapt to users' movement over time. P39 stated, "I would want to know what the robot is 'thinking' so that I know how to adjust my own behavior/path/location."

Robots that communicate motion intentions with body cues: Our analysis of the open-ended questions indicate that participants envision robots that communicate non-verbally using bodily cues. For instance, 26 out of 82 users discussed the need for full-body motion. They discussed robots that preemptively gesture in the direction they plan to move in before doing so. Also, 6 out of 82 participants envisioned robots that can provide hand gestures to indicate the direction they plan to move in. P84 said, "Something like a turn indicator on a car. It should flash/get attention visually and make a noise if possible so that people with low vision/hearing would be able to notice it".

Preferred communication devices with robots: Our results show that users preferred several devices for communicating with robots. Participants envisioned intelligent user interfaces to facilitate interaction with robots inspired by Google Maps as a top-down map of the robot and people around it in real-time. Building on this, they discussed a feature that enables them to view a ranked list of paths the robot is considering and why it chose its current path in terms of factors. P9 said, "I would like to know the top options the robot considered and why it decided to go with its final choice." Another idea was for robots to have 'car-like' features like turn signaling with lights or using a human-actuated hand to point or gesture in the direction the robot intends to move in. Lastly, 48 out of 82 participants envisioned physical devices to visualize robot motion intention information including on their phone or touchscreen devices.

Concerns about XAI for mobile robots: 14 out of 82

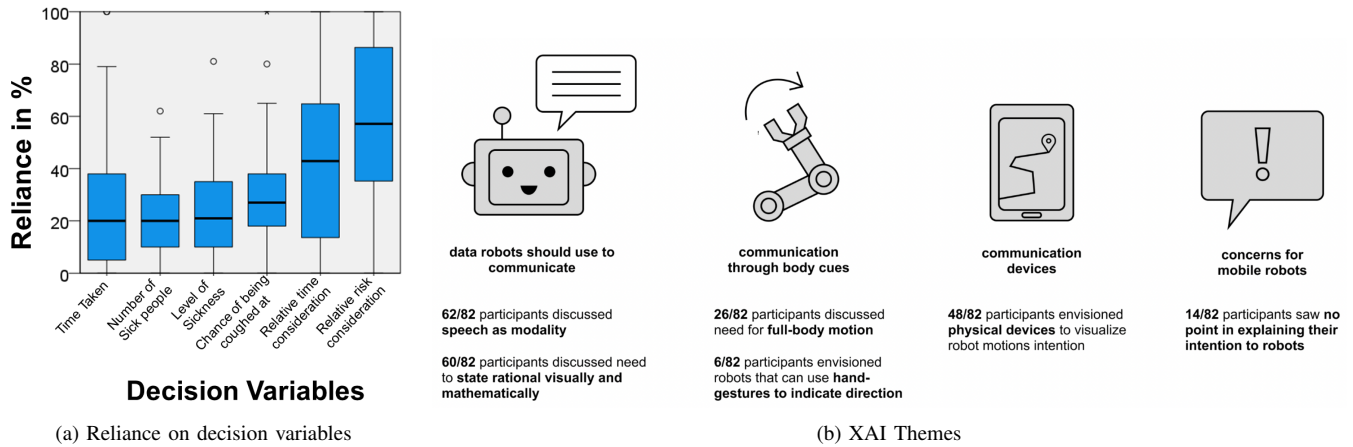


Fig. 4: (a) Boxplots showing the users’ preferences of decision variables (Table III) in percentage to choose paths in the trials to address RQ 3. (b) Themes that inform users’ preferences for XAI systems in social navigation settings from Section II-A to address RQ 4

TABLE VI: Paired t-test statistics of relevant variables to address RQ2

Variable X	Variable Y	P-Value	Effect Size (Cohen’s d)	Difference Between Means (Y - X)	95% confidence interval of difference
T	R	0.029	0.249	0.070	0.010 to 0.140
T^{sim} for MPP	R^{sim} for MPP	0.029	0.249	-0.070	-0.140 to -0.010
M	R	0.544	0.068	-0.020	-0.090 to 0.050
L	R	<0.001	0.503	0.18	0.100 to 0.260
M	T	0.006	0.318	-0.090	-0.160 to -0.030
L	T	0.008	0.302	0.110	0.030 to 0.190

participants saw no point in explaining their intentions to robots. Instead, they envisioned robots that adapt to them. Furthermore, they wanted robots to be passive actors in their environment instead of an active agent that they can interact with, indicating “[...] this should be inferred. I don’t want to change my natural behavior” -P21. One salient reflection resulting from the analysis is that some users expressed concern for robots using speech to communicate and they foresaw it as ‘creeper’ or ‘annoying’. P9 said, “Speech is a convenient way [...], but it could get annoying if the robot is constantly talking about why it’s going where it’s going.” Overall, there were a range of preferences identified in the data which highlights the need for systems with diverse XAI capabilities that adapt to users’ preferences over time.

IV. DISCUSSION

We have provided crucial insights about human behavior and perception to help design better navigational XAI in RCE. We summarize our findings from the user study. Participants tend to prefer riskier and time urgent paths, while they least prefer safer and time relaxed paths. Participants tend to show risk averse and risk insensitive behavior more often than expected behavior with respect to risk. CPT risk model captures participants’ decision making better than CVaR and ER, while ER performs the worst of the three models, in accordance with the theoretical findings from previous work [17], [18]. On an average the survey responses effectively captured participants’ risk and time urgency behavior from the trials. However, they may not be fully indicative of people’s risk and time urgency preferences as there was a large standard deviation in the similarity scores. Participants tend to prefer safer paths the least compared to their indicated risk appetite for the survey. There was no significant correlation between participants’ risk propensity and time urgency indicated in the respective surveys.

Our study reveals that humans act in a diverse manner in RCE, thus motivating the need to equip XAI with more inclusive risk perception models to foster cogent interaction and integration in such environments. We also revealed that standard questionnaires that measure risk propensity and time-urgency scale of individuals may not be very reliable to consistently predict human decision making in everyday RCE. Thus XAI may need to rely on other techniques such as online learning to understand the human perception of risk for efficient communication.

Our research demonstrates insights about users’ preferences for XAI-capable mobile robots that communicate their motion intentions. Overall, we found that most participants want robots that can explain the rationale behind decision-making, factors they considered, and how those factors were weighed against each other to make navigation decisions. Also, many participants envisioned intelligent user interfaces to communicate their motion intentions and for robots to express their motion intention on a Google maps-like interface. For those few participants that indicated no interest in explainable mobile robots, they envisioned XAI interaction could be useful in case of accidents or situations that require a written report. Overall, we found much potential for intelligent user interfaces that enable robots to communicate their risk perception to humans and hope that users find this information useful to enable safe and trustworthy mobile robot systems.

Limitations and Future Work: Here we will discuss some of the main limitations of our work and then mention possible future work for improvement. Due to the online nature of the study, participants’ path choices may not be fully representative of their actions in real world settings. Also, there might have been an inadvertent sampling bias resulting in a relatively younger age group of a student participant population. A more in-person study in a real world supermarket or similar setting may be needed to further validate our claims.

Although, we had a large sample size in terms of number of people, we only collected limited data (only 9 trials) per participant to minimize fatigue. More data points are needed to explicitly compare and characterize decision making between participants and risk models. This limitation can be

again alleviated by conducting in-person user studies where data can be collected in a natural and continuous manner (entire paths), which can then be used to perform more rigorous comparisons between various risk models.

We also performed correlation studies (point-biserial correlation) to determine the interaction effects between deviation from expected behavior in MPP and LPP using GRiPS score and time-urgency score. We found no statistical significance in these cases. We believe generic risk and time-urgency questionnaire responses may not reflect the participants' decision making in our particular environment. An in-person study might throw some more light on these interactions. Also, questionnaires focused on navigation in RCE may be needed to better capture human path choices.

V. CONCLUSIONS

In this work, we found that people tend to exhibit a variety of risk perceptions and behaviors in a crowded social navigation setting. We found that risk models like CPT, that are more expressive and inclusive, can better depict the observed human behavior, which thus support the previous theoretical findings. We also found that existing standard questionnaires to determine a users' risk-taking and time-urgency traits were not consistent with the exhibited behavior. In addition, we provided novel insights to consider for future XAI development for social navigation scenarios.

REFERENCES

- [1] A. Taylor, S. Matsumoto, W. Xiao, and L. D. Riek, "Social navigation for mobile robots in the emergency department," *In Proc. of the IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
- [2] A. Taylor and L. D. Riek, "Regroup: A robot-centric group detection and tracking system," in *Proceedings of the 2022 ACM/IEEE International Conference on Human-Robot Interaction*, 2022.
- [3] A. Taylor, S. Matsumoto, and L. D. Riek, "Situating robots in the emergency department," in *AAAI Spring Symposium on Applied AI in Healthcare: Safety, Community, and the Environment*, 2020.
- [4] A. Taylor, M. Murakami, S. Kim, R. Chu, and L. Riek, "Hospitals of the future: Designing interactive robotic systems for resilient emergency departments," in *In Proc. of the ACM Conference on Computer Supported Collaborative Work (CSCW)*, 2022.
- [5] A. Suresh and S. Martínez, "Human-swarm interactions for formation control using interpreters," *International Journal of Control, Automation and Systems*, vol. 18, pp. 2131—2144, 2020, doi: 1007/s12555-019-0497-3.
- [6] M. Kwon, E. Biyik, A. Talati, K. Bhasin, D. P. Losey, and D. Sadigh, "When humans aren't optimal: Robots that collaborate with risk-aware humans," in *Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction*, 2020, pp. 43–52.
- [7] R. T. Rockafellar and S. Uryasev, "Optimization of conditional value-at-risk," *Journal of risk*, vol. 2, no. 1, pp. 21–42, 2000.
- [8] A. Tversky and D. Kahneman, "Advances in Prospect theory: Cumulative representation of uncertainty," *Journal of Risk and Uncertainty*, vol. 5, no. 4, pp. 297–323, 1992.
- [9] F. Xu, H. Uszkoreit, Y. Du, W. Fan, D. Zhao, and J. Zhu, "Explainable ai: A brief survey on history, research areas, approaches and challenges," in *CCF international conference on natural language processing and Chinese computing*. Springer, 2019, pp. 563–574.
- [10] A. Das and P. Rad, "Opportunities and challenges in explainable artificial intelligence (xai): A survey," *arXiv preprint arXiv:2006.11371*, 2020.
- [11] G. Papagni and S. Koeszegi, "Understandable and trustworthy explainable robots: a sensemaking perspective," *Paladyn, Journal of Behavioral Robotics*, 2021.
- [12] V. D. Sharma, M. Toubeh, L. Zhou, and P. Tokekar, "Risk-aware planning and assignment for ground vehicles using uncertain perception from aerial vehicles," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 11 763–11 769.
- [13] A. Choudhry, B. Moon, J. Patrikar, C. Samaras, and S. Scherer, "Cvar-based flight energy risk assessment for multirotor uavs using a deep energy model," *arXiv preprint arXiv:2105.15189*, 2021.
- [14] S. S. Stevens, "Neural events and the psychophysical law," *Science*, vol. 170, no. 3962, pp. 1043–1050, 1970.
- [15] —, "On the psychophysical law," *Psychological Review*, vol. 64, no. 3, pp. 153–181, 1957.
- [16] S. Dhami, *The Foundations of Behavioral Economic Analysis*. Oxford University press, 2016.
- [17] A. Suresh and S. Martínez, "Planning under non-rational perception of uncertain spatial costs," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 4133–4140, 2021.
- [18] —, "Risk-Perception-Aware Control Design under Dynamic Spatial Risks," *IEEE Control Systems Letters*, vol. 6, pp. 1802 – 1807, 2022.
- [19] M. Kwon, E. Biyik, A. Talati, K. Bhasin, D. P. Losey, and D. Sadigh, "When humans aren't optimal: Robots that collaborate with risk-aware humans," in *2020 15th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 2020, pp. 43–52.
- [20] P. E. U. de Souza, C. P. C. Chanel, M. Mailliez, and F. Dehais, "Predicting Human Operator's Decisions Based on Prospect Theory," *Interacting with Computers*, vol. 32, no. 3, pp. 221–232, 2020.
- [21] A. Majumdar and M. Pavone, "How should a robot assess risk? towards an axiomatic theory of risk in robotics," *arXiv:1710.11040*.
- [22] S. Gao, E. Frejinger, and M. Ben-Akiva, "Adaptive route choices in risky traffic networks: A prospect theory approach," *Transportation Research Part C: Emerging Technologies*, vol. 18, no. 5, pp. 727–740, 2010.
- [23] A. Hakobyan and I. Yang, "Wasserstein distributionally robust motion planning and control with safety constraints using conditional value-at-risk," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 490–496.
- [24] S. Reddy, A. Dragan, and S. Levine, "Where do you think you're going?: Inferring beliefs about dynamics from behavior," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [25] R. S. Novin, A. Yazdani, A. Merryweather, and T. Hermans, "Risk-aware decision making in service robots to minimize risk of patient falls in hospitals," *arXiv preprint arXiv:2010.08124*, 2020.
- [26] K. Koide and J. Miura, "Collision risk assessment via awareness estimation toward robotic attendant," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 11 011–11 016.
- [27] J. Ballesteros, J. M. Peula, A. B. Martinez, and C. Urdiales, "Automatic fall risk assessment for challenged users obtained from a rollator equipped with force sensors and a rgb-d camera," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 7356–7361.
- [28] O. A. Hafez, G. D. Arana, and M. Spenko, "Integrity risk-based model predictive control for mobile robots," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 5793–5799.
- [29] G. D. Arana, O. A. Hafez, M. Joergler, and M. Spenko, "Localization safety validation for autonomous robots," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 6276–6281.
- [30] V. Shree, B. Asfora, R. Zheng, S. Hong, J. Banfi, and M. Campbell, "Exploiting natural language for efficient risk-aware multi-robot sar planning," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 3152–3159, 2021.
- [31] T. Kruse, A. K. Pandey, R. Alami, and A. Kirsch, "Human-aware robot navigation: A survey," *Robotics and Autonomous Systems*, vol. 61, no. 12, pp. 1726–1743, 2013.
- [32] B. Shneiderman, *Human-Centered AI*. Oxford University Press, 2022.
- [33] D. C. Zhang, S. Highhouse, and C. D. Nye, "Development and validation of the general risk propensity scale (grips)," *Journal of Behavioral Decision Making*, vol. 32, no. 2, pp. 152–167, 2019.
- [34] F. J. Landy, H. Rastegary, J. Thayer, and C. Colvin, "Time urgency: The construct and its measurement," *Journal of Applied Psychology*, vol. 76, no. 5, pp. 644–657, 1991.
- [35] S. Mohammed and S. Nadkarni, "Temporal diversity and team performance: The moderating role of team temporal leadership," *Academy of Management Journal*, vol. 54, no. 3, pp. 489–508, 2011.
- [36] M. A. Nadeem, M. A. J. Qamar, M. S. Nazir, I. Ahmad, A. Timoshin, and K. Shehzad, "How investors attitudes shape stock market participation in the presence of financial self-efficacy," *Frontiers in Psychology*, vol. 11, 2020.
- [37] V. Braun and V. Clarke, *Thematic analysis*. American Psychological Association, 2012.