

Recommendations for Designing and Conducting Human-Robot Interaction Studies

LEIMIN TIAN, NICOLE ROBINSON, PAMELA CARRENO-MEDRANO, WESLEY P. CHAN, MARAM SAKR, TINA WU, ELAHE ABDI, ELIZABETH CROFT, DANA KULIĆ, School of Engineering, Monash University, Australia

Human-Robot Interaction (HRI) is a fast growing field with great potential to benefit various aspects of our personal life and our society. One of the key questions in HRI is to understand how users perceive and react to the HRI experience, which enables evaluations of the outcomes of HRI systems. Such user studies require high-quality and scientifically rigorous methodology designs to ensure the robustness of the results. In addition, standardized methodologies improve the repeatability and generalisability of the findings, as well as provide feedback for design of new systems. Thus, we are motivated to review methodology of existing user studies in HRI, as well as in the broader fields of human-computer interaction, human-centred artificial intelligence, psychology and social science. In particular, we first reviewed existing HRI methodology guidelines and review papers. We then focused on the hypothesis, participants, experiment protocols and evaluation metrics aspects of HRI studies and discussed the strengths and weakness of existing approaches. This allows us to propose a step-by-step guideline that researchers may follow when designing and conducting future HRI user studies, which leads to high-fidelity study designs and robust outcomes that advances the frontier of HRI research and applications.

Additional Key Words and Phrases: human-robot interaction, social robotics, user study, methodology, guidelines

ACM Reference Format:

Leimin Tian, Nicole Robinson, Pamela Carreno-Medrano, Wesley P. Chan, Maram Sakr, Tina Wu, Elahe Abdi, Elizabeth Croft, Dana Kulić. 2020. Recommendations for Designing and Conducting Human-Robot Interaction Studies. *ACM Trans. Hum.-Robot Interact.* 0, 0, Article 0 (2020), 41 pages. <https://doi.org/0000001.0000001>

1 INTRODUCTION

Hypotheses about HRI methodologies:

- There are shared limitations and challenges in current HRI user study methodologies.
- Shifts in methodology over time:
 - How has methodological design changed over time? Type of HRI, participants, experiment protocols, and evaluation metrics.
 - How has methodological quality changed over time?
- Introduction to methodologies from complementary fields (HCI, psychology, cog sci, etc.)

Author's address: Leimin Tian, Nicole Robinson, Pamela Carreno-Medrano, Wesley P. Chan, Maram Sakr, Tina Wu, Elahe Abdi, Elizabeth Croft, Dana Kulić, School of Engineering, Monash University, 14 Alliance Ln, Clayton, VIC, 3168, Australia, {Leimin.Tian,Nicole.Robinson,Pamela.Carreno,Wesley.Chan,Maram.Sakr, Lee.Wu,Elahe.Abdi,Dana.Kulic}@monash.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

- We can propose a set of guidelines to improve future HRI user study methodologies inspired by user studies in other fields.

1.1 A Chronology for HRI studies

Propose a set of guidelines (in the order of how they apply to each step in the HRI study chronology) that HRI researchers can follow when designing and conducting studies.

Chronology of HRI study and guidelines for each study step shown in Figure 1.¹

Interactive version of the flowchart (Google Form? Open-source / Latex / GitHub) that result in a printable summary of the study (like pre-registering of the whole study), allow people to add their own approaches in addition to the choices

Different from existing HRI methodology reviews that summarise current approaches, we focus on identifying best practices for different types of HRI study, and aim to provide specific guidelines applicable to core steps during an HRI study cycle. In particular, we will cover a variety of HRI studies, methodological approaches, and application domains in Sections 3, 4, 5, and 6 in order to gain a broad view of the field. For each aspect of HRI study methodology, our discussion will be facilitated with examples and counter-examples of existing work. To help the readers apply the proposed guidelines in their work, in Section ??, we will provide a case study of using the proposed flowchart (Figure 1). Therefore, our work provides a set of applicable and practical guidelines to improve the quality of HRI studies, especially for researchers new to the field.

1.2 Method

Borrow the Preferred Reporting Items for Systematic Reviews and Meta-Analyses guidelines for searching, screening, and evaluating existing studies [112]. All authors discuss and decide on the proposed flowchart (Figure 1), then a pair of two authors is assigned to each of the sections in Sections 2 to 6. Each pair goes on separate search of literature, then cross-validate each other's selection and summary of papers. After reaching consensus, each pair write up their section. Finally the paper is reviewed and edited by all authors.

There is a shared Zotero folder called "User Study Methodology" in the HRI-Group Shared Library. Please feel free to drop in the papers you are covering in this review. Two raters per paper, tag your other rater on Google Doc: https://docs.google.com/document/d/1qi8lj_K7tkdMmUiHLNtVc-jrHzjy6lOQZ4vcONznim0/edit?usp=sharing. Since we already have paired authors for each section, you can be the dual reviewers of the papers covered in your section.

¹Interactive version of the guideline: <https://tianleimin.github.io/HRI-Methodology-Guidelines/>

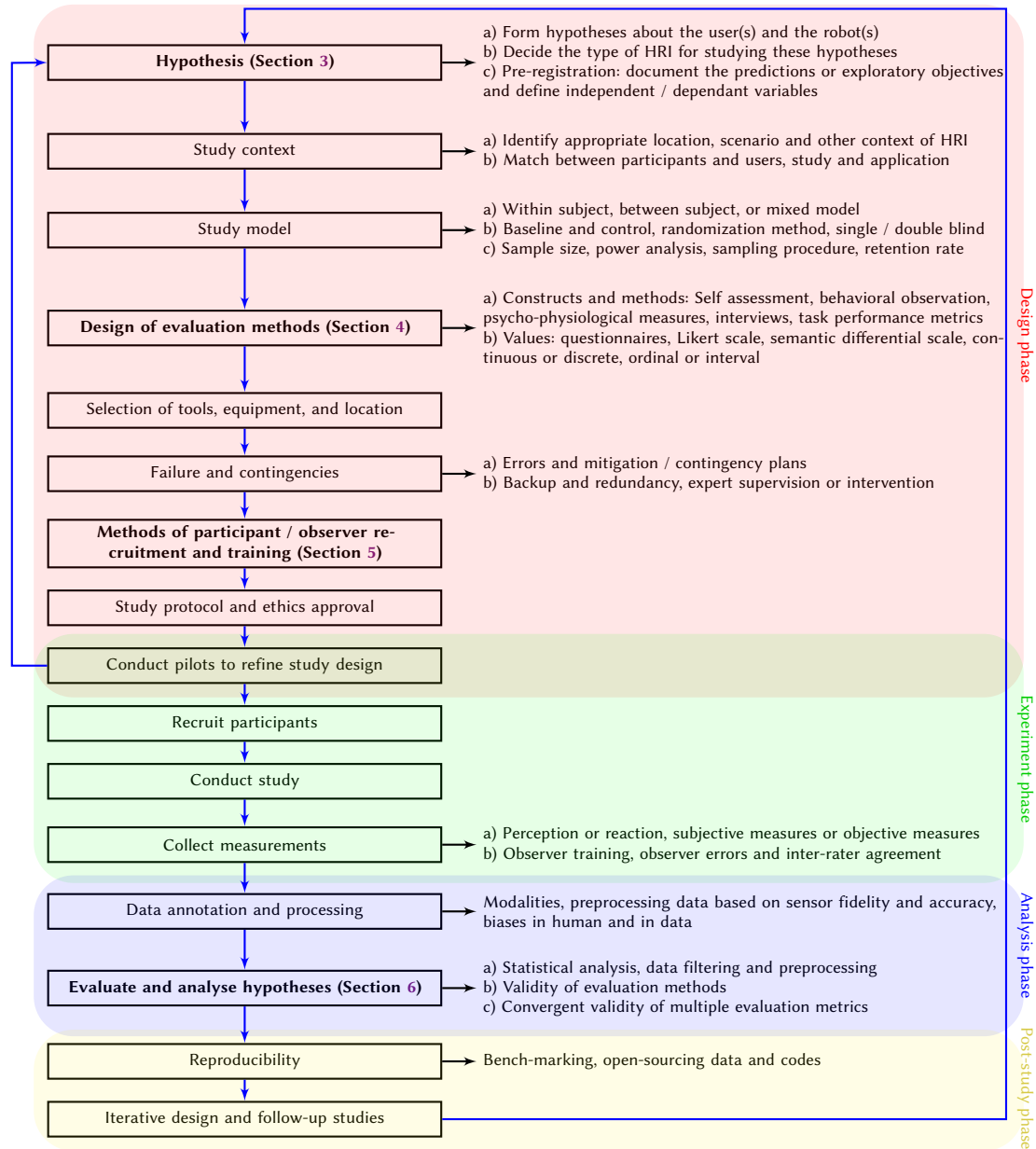


Fig. 1. Chronology and guidelines for conducting an HRI study.

2 EXISTING REVIEWS AND GUIDELINES FOR HRI STUDY METHODOLOGY

In recent years, there has been a growing effort to review the methodology in existing HRI studies. These reviews are important for the research community to improve the scientific outcomes for experimental HRI work, identify

limitations of current approaches, and recommend better practices. Existing reviews in the domain of HRI have mainly summarised current practices, rather than providing a guideline for future studies. Meanwhile, existing guidelines for conducting HRI studies either provide general advice on the study process without specifically targeting the various steps involved (e.g., Belhassein et al. [18]), or focus on a certain aspect, such as one specific type of HRI study (e.g., Wizard-of-Oz (WoZ) style studies [133]), methodological approach (e.g., psychophysiological assessment [81]), step during the HRI study cycle (e.g., evaluation [147]), application domain (e.g., healthcare [135]), or cohort of participants (e.g. elderly [184]). In this section, we will discuss representative examples of existing reviews and guidelines for HRI methodology, and provide motivation for this work. In particular, the HRI methodology guidelines we offer drew from a variety of HRI studies to provide a holistic view, while focusing on recommending suggestions applicable to core steps during an HRI study cycle.

2.1 Related work

Systematic methodological approaches have yet to permeate HRI research activity. A review of literature reveals that many HRI studies are designed and executed based on the preferences of researchers and available resources. For example, less than 10% of studies on social robots for healthcare have deployed randomised controlled trials with even fewer conducting follow-up studies [135]. In the elder care domain, two-thirds of the studies had up to 30 participants, while the majority lasted a day with only 11% lasting for 30 days or more, and most importantly, 90% of the studies lack theoretical framework [184]. That said, systematic approaches to HRI studies with a detailed study design methods are gaining more attention in recent years.

HRI studies have been categorised in multiple review papers based on a number of attributes. Themes and types of HRI studies and challenges as of year 2008 were reviewed by Goodrich and Schultz [64]. Five attributes of HRI were highlighted including autonomy, task type and structure of the human-robot team. A surge in longitudinal studies, interdisciplinary research and blended simulated/physical experiments has been observed. In a similar approach, Baxter et al. [16] categorised papers accepted at the ACM/IEEE International Conference on Human-Robot Interaction between years 2013-2015 based on level of robot autonomy, participant populations, evaluation environments, length of empirical studies, approach to statistics and replicability. Baxter et al. [16] provided a set of recommendations, highlighting the importance of clarity of the research goal, level of robot autonomy and statistical analysis, as well as justification of the subjects group size, length of interaction and replicability of the study. Establishment of standards and common metrics is reported as an emerging effort, however, years later, HRI community is yet to fully achieve this goal.

Hypotheses in HRI studies are often grounded with psychology theories. Thus, it is inevitable that the recent replication crisis in psychology has also influenced the field of HRI. For example, Irfan et al. [81] presented their failures to replicate the social facilitation theory in HRI. Based on this experience, they recommended HRI researchers who are consumers of psychology literature avoid older psychology literature with weak methods, and encouraged bench-marking efforts in HRI. Their work highlighted the challenges of observing and measuring social interaction in HRI. Similarly, Leichtmann and Nitsch [99] reviewed HRI studies on human-robot physical and social distances, and investigated the applicability crisis of informing hypotheses of HRI studies with human-human interaction (HHI) theories. They found that HRI study hypotheses are often derived from HHI under the assumption of robots being treated as social actors in a manner similar to humans. However, this assumption may not always be valid. When robots are anthropomorphized, humans would treat them as a social entity to a certain extent. Nevertheless, differences between robots and humans in displaying social cues can still result in fundamental split between HHI and HRI in many

aspects. Based on their review, the authors recommended pre-registration of hypotheses, methodology, and analysis protocol, as well as transparent reporting in HRI studies grounded with HHI theories.

Most reviews of HRI methodology focused on the evaluation aspect of HRI studies, and provided recommendations regarding the design phase of the HRI study cycle. A decade ago, Bethel and Murphy [23] identified study size and lack of multiple assessment methods as the primary issues in existing studies. They argued that a large sample size better represents the target population and exhibits a higher probability for statistically significant results. They also suggested that the deployment of three or more evaluation metrics for the same construct produces more reliable and accurate results. Currently, many metrics are just “observed” with a lack of functional/generalizable measurement mechanisms. For example, Murphy and Schreckenghost [116] proposed an HRI metric taxonomy based on a review of 29 studies. A total of 42 metrics were categorised into three groups: human (e.g., reliability, productivity and awareness), robot (e.g., plan state and time in manual/autonomous operation), and system (e.g., efficiency, safety and coactivity). Psychological assessment metrics such as cardiovascular, electro-dermal and brain activity are increasingly used to evaluate human responses in HRI [168]. When studying social interactions and human responses in HRI, one of the most commonly adopted approach is to use the Likert scale to measure subjective judgements of participants. However, a recent review on application of the Likert scale in HRI studies found that only 3 out of the 110 papers performed their analyses correctly [147]. This demonstrated the urgent need of methodological guidelines to ensure the quality and robustness of HRI studies, especially regarding the validity of evaluation methods and appropriate analyses of different metrics.

In contrast to reviews focusing on the design phase of HRI studies, Belhassein et al. [18] provided recommendations for the experiment and post-study phases of HRI user studies with a focus on participants and the replication crisis. Regarding the experiment phase, they suggested that the outcome of the HRI study can be improved by recruiting more users with diverse backgrounds, rigorous implementation of the protocol, allowing time for the user to get accustomed to the robot, ensuring physical and psychological safety, assessing suitability of the robot for the task, choosing the right and preferably standardized measurements, deployment of theoretically solid tools and measures and finally providing required details for reproducing the study. Regarding the post-study phase, they suggested standardisation of the technical components to address study reproducibility.

Many HRI systems are designed with the goal of application to education, health, and elderly care [37]. This relies on advancements in sensory technologies and algorithms in robotics [165]. Furthermore, to evaluate the outcomes of such HRI systems, it is critical to evaluate HRI in-situ. For example, Weiss et al. [174] proposed a usability, social acceptance, user experience, and societal impact (USUS) framework to evaluate human-robot collaboration in terms of usability, social acceptance, user experience, and societal impact. To envision new research questions beyond existing technical restraints, the Wizard of Oz (WoZ) approach is often applied in HRI studies, especially when studying complex social interactions. The WoZ approach also benefits experiment safety and control of variable factors in HRI. Riek [133] reviewed how WoZ has been used in HRI in order to identify valid WoZ designs. They stressed the influence of user training and instruction, as well as the formation of hypotheses on user behaviors and robot behaviors.

Existing reviews of HRI methodology often focus on one particular aspect or evaluation method. For example, Lasota et al. [96] reviewed factors influencing safety in HRI and different measurements of safety; Young et al. [183] reviewed the use of biological features and means to incorporate user variances; Prewett et al. [127] reviewed the measurement of robot teleoperator’s workload in relationship with task outcomes. It is important to have such in-depth investigation on one type of HRI or on one aspect of HRI evaluation. However, identifying a set of guidelines that inform on core

steps of an HRI study cycle is crucial to improve the overall quality of HRI studies. Therefore, in this paper we focus on providing such applicable and practical guidelines to HRI study methodology.

One of the most comprehensive overviews of HRI studies is a recent work by Bartneck et al. [14]. In this report of good practices in HRI studies, research questions (explanatory vs. confirmatory), study designs (qualitative vs. quantitative vs. mixed methods), participants (choice of population and sample size), context (location, time and social unit), robot appearance and functionality, mode of interaction (WoZ vs. physical vs. simulated) as well as direct and indirect metrics, statistical analysis and ethical considerations are covered.

The existing HRI study methodology review papers are widely focused on reporting the current common practices rather than a chronology for the HRI study cycle or a guideline regarding the best study protocols and evaluation metrics. In these studies, the discussions on (dis)advantages of each approach and analysis of the best practices are generally limited.

2.2 Summary

Because of the wide variation in research questions and experimental settings in HRI, a one-size-fits all methodological approach across all HRI studies is not a reasonable expectation. However, we have identified five shared recommendations emerging in existing reviews of HRI studies covered in this section:

- (1) **Grounding of hypotheses and the applicability crisis:** Grounding HRI hypotheses with theories in psychology and social science helps to increase the strength of the HRI hypotheses. However, researchers should also be aware of the differences between human-human interaction and HRI, which can limit the applicability of psychology or social science theories in HRI.
- (2) **The importance of context:** HRI is context dependant. Researchers are encouraged to have a holistic view when conducting HRI studies, and have a clearly defined framework for their research.
- (3) **Pay attention to evaluation:** The validity of evaluation metrics used is directly related to the validity of an HRI study. Researchers are recommended to use multiple metrics and evaluation approaches to measure a parameter, and to investigate the convergence of these metrics. In addition, researchers are encouraged to use appropriate statistical tests and conduct power analyses when studying experimental results.
- (4) **Ecological validity of a study:** Researchers are encouraged to conduct field studies and longitudinal HRI studies. In HRI experiments, it is recommended to recruit more participants from the target populations.
- (5) **Rigour and transparency are key:** To improve the reproducibility and generalizability of HRI studies, researchers are encouraged to pre-register their hypotheses, to cover both significant and non-significant results in their discussion, and to open-source the code and data used in their work.

In the following sections, we will investigate each of these aspects of HRI study and discuss the best practices for different study settings. In particular, we will cover hypotheses and study context in Section 3, evaluation methods in Section 4, participants and ecological validity in Section 5, and reporting metrics in Section 6.

3 HYPOTHESIS IN HRI STUDIES

Wesley and Elahe

Research studies can be generally classified into exploratory or confirmatory studies [13]. In exploratory studies, the researcher aims to investigate a relatively new idea or area where there is not much prior knowledge to draw expectations from. Whereas in confirmatory studies, also known as hypothesis testing, the researcher aims to confirm

certain expectations they have on how their system or users would perform or behave. Focusing on confirmatory HRI studies, a key initial step is to develop clear hypotheses about the user(s) and the robot(s).

A hypothesis is a predictive statement on how a given condition will affect a certain aspect of the outcome. It drives or guides the selection of the type of the study and the method of evaluation among others [23]. For example, a hypothesis in an HRI study comparing user perception of different robot trajectories may be:

- robot arm employing human-motion model trajectory will be perceived by users as safer to work with, compared to one employing a constant acceleration trajectory.

The above is a statement relating the given condition of employed robot trajectory, to the outcome of user perceived level of safety. It makes a prediction that with a particular trajectory type, the level of safety perceived by users will be higher when compared to the other trajectory type. As the aim of user studies is to find support for the formulated hypotheses, a good hypothesis should be verifiable in a scientific experiment. The condition stated in the hypothesis should involve a factor, often referred to as independent variable, that the researcher can vary and control independently in an experiment, and the outcome stated in the hypothesis, also known as dependent variable, should be a parameter that is measurable or observable in an experiment setup. Considering the above example, the factor of robot trajectory type can be controlled independently of other factors by the researcher. The outcome of perceived level of safety can be measured for example through questionnaires asking users to rate their perceived level of safety.

A hypothesis may be formulated based on common sense, prior knowledge, existing literature and theory [13], or parallel fields of study [33]. Various taxonomies have been proposed in the literature for classifying HRI. These include taxonomies classifying HRI by the type of study [36], the role of the robot [33], its morphology [180], level of autonomy [167] [17], behaviour [154] and characteristics [129], the role of the human [146], type of task [180] [158] and interaction [150] [121] as well as human-robot proximity [64] [3] and coordination [83] [75]. When formulating hypotheses, it is helpful to adopt a designers perspective as suggested by Goodrich and Schultz [64]. They defined HRI as "a field of study dedicated to understanding, designing, and evaluating robotic systems for use by or with humans". In their paper, they listed five attributes that the designer can change to affect HRI: (1) level and behaviour of autonomy, (2) nature of information exchange, (3) structure of team, (4) adaptation, learning, and training of people and robot, (5) shape of the task. A common strategy for formulating hypotheses is then to consider which attribute is the main concern of the study. Once the main attribute of concern has been identified, the levels for testing this attribute can be determined. These attribute settings would then naturally become the independent variables in the hypotheses. Considering the above example, the attribute of main concern is the behaviour of autonomy (robot trajectory type). The selected levels for this attribute are human-motion model trajectory and constant acceleration trajectory. The dependent variables stated in a hypothesis then often relates to what the designer hopes to achieve (i.e., improved perceived level of safety). In the following paragraphs, we take a closer look at the five attributes mentioned above.

Level and Behaviour of Autonomy. Autonomy is defined as the robot's ability to accommodate variations in its environment [167] by sensing, planning and acting to reach a goal without external control [17]. It determines the level of interaction in HRI and the type of task the robot can perform. From industrial robots to service robots, autonomy can be anywhere between manual teleoperation to full autonomy. Beer et al. have proposed a 10-point taxonomy to identify the level of robot autonomy along this spectrum [17]. This is directly related to when and how the human or robot should take the initiative in HRI. Seizing the initiative can happen in a reactive, deliberative or hybrid manner while the human-robot coordination ranges from explicit to none [83].

Nature of Information Exchange. The nature of information exchange is primarily defined by the communication channel and communication format [64]. These can be affected by human-robot interaction distance - e.g., whether they are physically attached, in the same area, or in different areas [3]. Communication channels used for human-robot interaction often mimic those in human-human interaction. Such channels, and example formats, include visual, such as gaze [2] or body gestures [105]; auditory, such as speech [182] or audio-based expressions [57]; and haptic, such as shared load [30] or vibrations [145]. A common strategy for enabling intuitive/effective HRI is by programming robots to use communication channels and formats familiar to humans [63, 67, 122].

Structure of Team. Structure of team concerns with the number (type) of humans and robots, their roles, and the configuration of communication links among them. Directly relating to these factors, Yanco and Drury's proposed taxonomy included ratio of people to robots, composition of robot teams, level of shared interaction among teams, and human role [180]. The human may take a range of roles from supervisor and operator, to mechanic, bystander, and teammate [146]. The role of the robot can be identified by comparing HRI to human interaction with other agents, including other humans, animals and objects and defining the role of the robot based on the similarity of its role to those agents [33]. This can also be affected by the robot's morphology [180]. An anthropomorphic, zoomorphic or functional robot is more likely deployed in roles matching its shape. In social robots, the ability of the robot to deceive has been suggested as an additional role capacity categorised based on the deception object, goal, and method [154].

Adaptation, Learning, and Training of People and Robot. Adaptation, learning, and training can happen for both human and robot in HRI. Most HRI designs aim to develop intuitive interfaces for interaction. Therefore, robotic systems for HRI are often designed to minimize the amount of user training required [64]. However, any system will have its own behaviours and limitations. Hence, there is always some learning and adaptation from the human taking place [65]. Short-term HRI studies commonly treat learning in humans as an undesirable carryover effect and try to mitigate it. However, in longitudinal studies, adaptation and learning warrant careful study and may be actual focuses of the studies [98]. Certain domains of application involving high risks, such as bomb disposal, can also benefit from or may require careful training of the user or operator [132]. Robot adaptation has been used to personalize interaction and promote longer term engagement with users [4, 169]. Robots can also learn tasks and skills from humans through various modes of HRI, with one common mode being Learning from Demonstration [9]. In the context of education, robots have taken the role of teachers or peer tutors for teaching students [47], while robotic pets have also been used as therapeutic animals for treating anxiety by teaching patients to adapt their breathing to the robot's [149].

Shape of the Task. Task type and task criticality (measuring importance of task outcome) [180] as well as its cognitive/physical requirements [129] are introduced as part of HRI taxonomy. The robot's level of autonomy is an influential factor in determining the type of tasks it can perform [167]. The aspects of the task the robot should perform and the extend at which the robot can perform those tasks affect the choice of the robot and the study design [17]. A sound theoretical foundation ensures that the robot is selected based on its suitability for the task rather the expected outcome [33].

For rigour in HRI studies, it is highly recommended to register the hypotheses prior to conducting the study and justify the methodology based on them rather than the expected outcome [13]. Website where researchers can register their hypotheses for this purpose include the Center for Open Science (<https://osf.io/prereg>), AsPredicted (<https://aspredicted.org>), or the U.S. National Library of Medicine (<https://clinicaltrials.gov>). Results should then be reported with reference to the initial hypotheses avoiding a selective report based on the desired outcome. Some journals (e.g., JMIR Research Protocols) may automatically generate an identifier for submitted and published research protocols

or proposals. The identifier will then be linked to all subsequent result papers. HRI literature contains good examples of this empirical approach, a few of which are highlighted in the following paragraphs.

A study on the effect of implicit and explicit communication formats in HRI teamwork efficiency provides an example of a systematic approach in defining clear hypotheses [26]. This highly cited empirical study presented three hypotheses related to transparency of robot state, task efficiency and robustness to errors and then addressed each through the experimental design, measurements and analysis steps.

Other examples are studies on personal space [163] and emotional reactions [137] in HRI that designed their experiments based on existing literature and common sense. The study in [163] examined people's behaviour (adaptation) depending on past experiences, among others factors, and their hypotheses were mainly based on a generalised model of human-human interaction, whereas [137] examined different roles a robot plays (or treatment it receives) in videos and based its hypotheses on existing literature in psychology. Both studies used factorial experimental designs with between and within subject factors² to address the hypotheses and clearly reported the results discussing the implications of their findings for theory and design.

To summarize, hypotheses are predictive statements stating how the experimental conditions are expected to affect the measured outcomes. They may be based on common sense or existing literature. When constructing experiment hypotheses for HRI studies, it is useful to consider the various design factors can be altered to affect HRI. It is highly recommended to register the hypotheses prior to conducting the study for rigour.

Clear hypotheses based on sound background is a critical initial step in any confirmatory HRI study. The hypotheses set the ground for the next steps of the study, from experimental protocol and design of evaluation methods, to recruitment of participants, evaluation metrics, statistical analysis and reporting the outcomes. Each of these steps are discussed in detail in the following sections.

4 EXPERIMENT DESIGN AND EVALUATION

Nicole and Tina This section examines the design and implementation of an experimental study in human-robot interaction. Once the hypothesis is formulated, the next step is to determine its design, parameters and select the appropriate evaluation methods that are suited to explore the proposed hypothesis. A non-exhaustive list of the study components are listed below. Each of these components will be explored in more detail, identifying advantages, disadvantages, and considerations in the following section.

- Study design
- Choice of constructs and metrics
- Choice of measurements
- Trial location (laboratory, field, and online)
- Session frequency and duration (single-session or multiple follow-up sessions)
- Type and number of robot (humanoid, zoomorphic, nonbiomimetic)
- Task behaviour of the robot (actions or behaviours)

²Factorial, between/within subject experimental designs are explained in Subsection 4.1.1

4.1 Study Model

Researchers must determine the type of study design and number of groups that is needed to explain the hypothesis. Common study design structure in a research trial can often involve one of the following: within-subjects, between-subjects, or mixed-model factorial approach [22]. This includes the consideration of a control group in the study design and will be explained in more detail below.

4.1.1 Study Design Structure.

Within-subjects Design

In a within-subjects design, there is one group of participants and each participant is exposed to all of the experimental conditions. An advantage of a within-subjects design is that a smaller sample size can be used, which is helpful in situations where participants are difficult to recruit [97]. Other advantages include experiments can be performed in a single session sequentially (i.e. more time efficient), leaving smaller statistical noise for between-group differences (i.e. the scores of the same participant is compared between each experimental conditions) [22].

On the other hand, a longer experiment duration might increase the likelihood of participants being affected by non-intended incidents, such as program glitches [22, 92]. Long experiment time and repeated task exposure can also lead to participants anticipating the experiment outcome, habituation, practice effect, and/or fatigue, resulting in careless task performance, or increased performance due to experience rather than the task at hand. Within subject design is further prone to cross-condition contamination, where exposure to one condition context affects the participant's response during the second condition, sometimes known as the Halo effect (e.g. [179]). Sequence effects can also occur based on the sequence of conditions presented to the group or individual [141]. Researchers can consider using a Latin square design and/or randomisation to control the impact of habituation and other confounding variables, and incorporate necessary breaks between experiments [97].

Recommendation Within-group design is best used when researchers want to evaluate a measure across multiple time points, or test the effects of different conditions. It is best used when you have a limited sample size. Researchers must note the possibility of sequence effects, fatigue and practice effects.

Between-subjects Design

In between-subjects designs, participants are exposed to one experimental condition only. This means that each participant only sees one condition, behavior or task whilst different participants are exposed to a different condition, behavior or task. The number of participant groups is determined based on the number of conditions. Strengths of between-subjects design are that it minimises confounding variables such as learning effect, fatigue, and frustration that occur as the byproducts of running long experimental sessions [22, 97].

Participants are different between each experimental condition and collected data will be affected by individual differences without randomisation (i.e. random allocation of participant groups into experimental conditions). As a result, it might be more difficult to detect significant differences and type II errors (i.e. false negatives) are more likely to occur [22, 85]. In a between-subjects design, the experimental conditions should also be distinct from each other for between-group differences to be perceived, e.g. [134]. Furthermore, larger sample size is required to mitigate the impact of individual differences.

Recommendation Between-group designs are best used when researchers want to test the difference between two conditions, behaviours or tasks, and researchers do not want participants to be influenced by similar conditions. The design can strongly identify effects across-groups compared to within-group. Between-group does need a larger sample size to produce statistically significant results.

Mixed-model Factorial Design

A mixed-model factorial design utilizes both between-subjects and within-subjects designs, where the between-subjects aspect is used to investigate multiple independent variables while other variables are explored using within-subjects designs [22]. This study design can simultaneously explore the effect from individual variables and the interaction effects Lazar et al. [97]. The downsides from the aforementioned discussion on each study design structure still apply and additionally, mixed-model designs often requires more subjects and randomisation is conducted properly.

Recommendation Mixed-Model designs are best used when researchers are seeking to assess both between and within-group effects. They do require a larger sample size and require more research skill to implement correctly, but produce rigorous results.

4.1.2 Baseline and Control.

Control is vital for experiments that aim to evaluate the impact of a new technology from different perspectives such as usability, functionality, and engagement level. Commonly used control conditions include waitlist/delayed (i.e. participants receiving no robot/therapy/treatment), active control, or different robot presentations. The control scenarios can take on many forms such as robot-against-human, robot-against-animal, robot-against-conventional-treatments. The purpose of a control condition is to have something to compare it to, demonstrating how the behaviour, robot or experiment performed. Interested readers can learn more about control conditions in experimental methodology books, and here we provide a brief overview.

4.2 HRI Evaluation Methodology: Constructs and Metrics

Choice of constructs (what to measure) and metrics (how to measure it) are important given that these will determine the research focus and the data to collect. Constructs and metrics will help assess the validity of the research hypothesis and the relevance of the variables controlled by the researchers, but also shed insights on the implications for the other agents involved in the interaction process or the interaction itself. We provide a summary of the common constructs and their related metrics in Table 2. Researchers can choose whichever constructs or metrics most relevant to their study, whether it is trust, acceptance, likability, response time or collaborative efficiency. However, researchers should carefully evaluate the quality and soundness of the chosen measurements as it affects the statistical analysis of the study [160].

4.2.1 Common Constructs in HRI Studies. HRI studies are often interested in validating theoretical concepts, known as constructs. Most of these constructs are not directly observable but can be linked to observable and/or measurable metrics and events [80]. Given the incredibly diverse range of applications in human-robot interaction, more than 110 different metrics have been proposed or used in the HRI literature [36]. These metrics can be grouped by the agents taking part or being measured during the interaction process, which includes the robot, the human, and the team (often referred to as the system) [116]. Given the wide scope of HRI studies, the metrics listed in this section might not apply to all types of user studies.

Robot-Related Constructs

Autonomy is a critical construct in HRI since it determines the way in which robots and humans interact with each other. Autonomy is defined as the ability of a robot to function independently [158]. Neglect tolerance [121] (i.e. how the robot effectiveness declines when the human is not attending the robot) and attention demand (percentage of time the human must control the robot) are often used as overall measures of a robot's autonomy. Beer et al. [17] suggested

to also include a subjective rating of the human intervention as a supplementary measurement of a robot's level of autonomy.

Productivity is another construct that is related to the robot's performance at a given task and/or the robot's technical capabilities (e.g., accuracy at object recognition). Time-based and error metrics such as task completion time and number of unsuccessful actions are often used to assess the former [121]. Finally, the third robot-related construct is the social attributes, which aim to capture how traits (e.g., warmth and competence) and characteristics (e.g., anthropomorphic appearance) associated with robots affect the social perception of people interacting with them. Since these attributes mostly rely on people's judgements, subjective scales such as the Robotic Social Attributes Scale (RoSAS) [29] and Godspeed questionnaires [15] are often employed as metrics.

Human-Related Constructs

Situation awareness and cognitive workload are two intertwined constructs that have been identified as particularly relevant in both automation and robotics [17, 158]. Situation awareness defines the the awareness and understanding an individual has of a situation [142]. Cognitive workload is a product of the mental resources demanded by a task and the capacity of the person performing the task. In addition, several other constructs have been used to investigate people's perception of robots, their responses to different robots under different contexts and types of interaction [14], or even predicting future use of the robot [17]. For instance, constructs such as willingness to interact with the robot ([72, 134]), trust, perceived or physiological safety, engagement, and affective states have been shown to be potential predictors of robot use [17]. The measurement or assessment of these constructs is frequently done using questionnaires (e.g., RoSAS [29], Discrete Emotions Questionnaire [68]), physiological measurements such as changes in heart rate and skin conductivity [7] and behavioral measurements.

System-Related Constructs

Most of the constructs associated with the system are defined in the context of task-oriented applications, which aim to assess how well the human(s) and robot(s) perform as a team. For instance, while productivity and efficiency [116] are associated with how well the task is completed and the time and effort required to complete the task, fluency characterizes the coordination of joint activities and actions between members of a well-synchronized team [75, 76]. Metrics such the time elapsed between the end of an agent's action and the beginning of the other agent's action (i.e., functional delay), number of unplanned human interventions or interactions, and the time required to interact with the robot (i.e., interaction effort) are often used to measure how different experimental conditions, human factors or new robot skills affect the system's performance [36].

4.2.2 Common Measurements in HRI Studies.

Multiple constructs can be of interest in an HRI study. However, existing studies often attempt to capture the interaction with a single measurement, which is insufficient Bethel et al. [22]. In addition, although many measurements are possible, specific expertise and knowledge about the phenomenon under investigation is often required for the correct measurement to be obtained and analysed. For instance, facial expression alone might not be the best indicator for evaluating emotions as expressions can be culture dependent or displayed out of context (e.g. adults learn to smile to hide disappointment or awkwardness) [12]. This example illustrates the need of finding the correct measurement and expertise to bridge the knowledge gap in an interdisciplinary HRI study. The current section highlights the five common evaluation methods listed in Bethel et al. [22] and discusses the advantages, disadvantages, and when the method should be used.

Task Performance Metrics

Task performance metrics are objective measurements of how well the robot, the participant, and the overall system

accomplish the task under investigation. The performance can be measured with respect to five different aspects of the robot's capability, which are navigation, perception, management, manipulation, and sociability [158]. Examples of these metrics include task completion time, error rate, and efficiency. More examples can be found in Table 2.

Performance measurements are useful for comparing technologies as it gives a numerical value indicating how much improvement can be achieved by adopting the technology (e.g. introducing warehouse robots to improve the sorting efficiency of packages). However, task performance metrics alone are insufficient to capture the entire HRI experience and the result might be heavily influenced by the population. For instance, people with more exposure to different technology might have a shorter task completion time in a human-robot collaboration task versus people who are less experienced with technology. The effect of population can be reduced by selecting a large, diverse group of participants and introducing a control/baseline to measure relative changes instead of the absolute values.

Behavioural Measurements

Behavioural data can be recorded under four different study contexts

- Naturalistic observation involves observing people's behaviour in a natural environment where it typically happens (quantitative or qualitative)
- Participant observation where researchers are part of the experiment they are conducting (quantitative or qualitative)
- Structured observation is when researchers only focus on gathering quantitative data of the behaviours under investigation in a particular experimental setup; and finally, case studies (detailed investigation of an individual, social unit, or events) [128]. Structured observations are typically more efficient than naturalistic/participant observation as researchers will already know what to look for during the experiment.

A metric in HRI studies might include recording the duration where participants look at the robot. In contrast with self-assessments, behavioural measurements do not rely on the participant recollection or self-interpretation of their behavior [22]. However, behavioural measurements are susceptible to the Hawthorne effect, where participants might behave differently as they normally would because they know they are being watched. Analysis of behavioural data in the form of audio/videotape requires significant time and expertise. The cost of data annotation is often too high and most researchers will not end up processing the data. With annotation comes the problem of inter-rater and intra-rater variability. Finally, behavioural measurements alone are insufficient to explain why the participant behaves the way they did [84] and should be used in conjunction with other measurements, such as interviews, quantitative or qualitative data.

Psychophysiological measures

Psychophysiological signals are physiological measurements that are affected by the participant's mental state. This includes heart rate, blood pressure, brain activity, skin conductance, muscle activity, and more [24]. These measurements are good for quantify abstract concepts such as physiological arousal using non-invasive and objective measurements. The time domain aspect of these measurements help pinpoint exactly when in the experiment a response was elicited [39]. In addition, an improvement for using psychophysiological measures compared to behaviour measures is that participants cannot easily manipulate an automatic response.

Psychophysiological measurements are informative when applied properly. Since the signal gives very specific information about the person's state (e.g. heart rate and respiration rate are indicators for arousal levels), it is important to understand what the signal is measuring and not infer meanings that may not be accurate. One challenge of obtaining the measurement is that the signal is often influenced by various confounding variables and noise (e.g. health status,

room temperature). Health conditions can prevent changes with respect to the baseline to be observed due to the "Law of Initial Values". The law demonstrates that there might be a narrower range that the measurement can possibly increase or decrease if the initial measurement is already high or low [24]. Other confounds include orienting response, defensive response, startle response, and habituation [24]. A loud noise made by the robot during the experiment can make the participant nervous, which results in a temporary increase in heart rate that is unrelated to the actual experiment. There might be scenarios where the measurement does not provide any useful information, specifically if the task might not provide a significant change from the baseline or when different tasks might elicit the same psychophysiological readings [39].

Interviews

Interviews can take on the form of informal conversational interviews, interview guide approaches, standardized open-ended interviews, and closed quantitative interviews [84, 111]. Interviews provide information that might not be collected in self-assessment and it does not force people to rate the robot on scales that are not relevant to them (e.g. how much did you trust this robot?). Qualitative responses can also be used to find new research questions and correlations between response answers and new theories. An informal interview can be tailored towards the individual and the questions would be more relevant, whereas a structured interview makes the responses more comparable. A guided interview is more comprehensive compared to an informal interview as all the topics are defined in advance. Finally, data collected in closed quantitative interviews are the easiest to analyze and compare Johnson [84].

Interviews provide a wealth of information about what were the most important factors of the interaction from the participants' perspective, and the quality of the data is influenced by the participant's response style, which can make the interaction appear overall negative, positive, or the risk of only getting socially acceptable answers. Volunteer participants might respond differently to non-volunteers, as they might inherently be more interested in robots/new technology [123]. In terms of the interview design, data collected in informal interviews are more difficult to analyze as it might require training in a particular qualitative methodology to do rigorously (i.e. interpretive phenomenological analysis or content analysis) [48, 156]. Consistency between unstructured interviews are difficult to maintain between experiments. Qualitative data analysis can be conducted on programs such as NVIVO, AtlasTI, and other natural language packages available online.

Self-assessment

Self-assessments in HRI studies can be both quantitative and qualitative, and often come in the form of paper/computer based scales, questionnaires, or surveys. A list of commonly used surveys can be found in Table 2. The assessments are easy to administer, but influenced by many factors and limited by human performance. A factor that affects the accuracy of the measurement is the timing of the survey. Since surveys are often administered at the end of the study, it means the data can only be collected on a reflective level, researchers will not be able to corroborate information from the participant directly, and participants will need to remember the feelings [24, 39]. The assessment is also influenced by societal or cultural norms where participants might answer questions based on what they think is acceptable by the research team or the society [24]. Some self-assessment measures can require substantial fees to be paid for their use, increasing the cost of research. Others have copyright that must be approved by the owner for its use, which might be difficult to obtain.

Questionnaires

Existing surveys might not always be sufficient to capture the variable at hand or the research direction that is of interest. Researchers can consider designing their own questionnaire, but it must be noted that proper validation requires extensive testing, statistical measurement and sampling. Adopting existing questionnaires is recommended if

they have been validated. Questionnaires should not be altered once validated, as they can become invalid when the order of the items or wording is changed [139].

To create a new questionnaire, the first step is to select the appropriate measurement scale (i.e. nominal data, ordinal data, interval data, or ratio; see Table 1 for more information). Once the scale is selected, types of options and the number of options need to be considered.

Table 1. Choice of measurement scales

Type of value	Definition and example	Statistical average
Nominal data	Categorized data where male participants are assigned the value of 1 and females are assigned the value of 2	Mode
Ordinal data	Ranked data where a question asks a participant to rank their engagement level during a HRI interaction between a discrete scale from 1 to 10	Median
Interval data	Data measured along a scale with an example includes participants rating their enjoyment level on a continuous scale from 0 to 100	Mean

Here we present general guideline for the design for question options as stated in Rust and Golombok [139] includes:

- A personality/mood questionnaire often includes options such as “not at all”, “somewhat” and “very much”.
- A attitude questionnaire might have “strongly disagree”, “disagree”, “neutral”, “agree”, and “strongly agree” as the options.
- A clinical symptom evaluation scale might include “always”, “sometimes”, “occasionally”, “hardly ever”, and “never” as the options.

The number of options for each questionnaire item depends on the nature of the questionnaire. The important factor is to give sufficient options to the participants so that they can express their opinions. A general rule is to include at least 4 options for rating scales and use consistent number and type of options across the survey [139]. When possible, test different survey wording and choice of the variables using a pilot study to improve the quality of the survey. Johnson [84] listed a detailed procedure to construct a questionnaire in educational research and the same principles can be adopted to HRI studies. In general, researchers should note that the customised questions have not yet been tested outside the experiment, and when feasible, provide evaluation metrics for validating the survey (see 4.2.4, 4.2.5 below). For more detailed information on how to design a questionnaire, researchers should refer to Litwin and Fink [104], Rust and Golombok [139]. Examples where customised questionnaires are designed and validated for HRI experiments can be found in [46, 76].

Table 2. Summary of Constructs and related Metrics

Agent	Construct	Measurement Type	Metrics
Robot	Autonomy [158]	Task-performance metrics	Neglect tolerance and attention demand [121]
	Social Attributes [29]	Self-reporting questionnaires	RoSAS [29]
	Productivity [158]	Task-performance metrics	Task completion time, percentage of successful actions [158]
Human	Situation awareness [146]	Task-performance metrics, self-report questionnaires, physiological measurements	SAGAT, SART [49], gaze analysis [42], secondary task performance [142]
	Cognitive workload [158]	Self-reporting questionnaires, task-performance metrics, physiological measurements	NASA-TLX; task errors and reaction time [126]; respiratory rate, heart rate and skin conductance and temperature [120]
	Acceptance [72]	Self-reporting questionnaires	Godspeed questionnaires [15], TAM-based questionnaires [72]
	Trust [158]	Self-report questionnaires, task-performance metrics	Human-Robot trust scale [144]; trust in automation scale [82]; task allocation [136]
	Affective state [94]	Self-report questionnaires, physiological measurements	Discrete Emotions Questionnaire [68]; Pleasure-Arousal-Dominance ratings [21]; cardiovascular and electromyogram (EMG) activities [94]
	Stress [96]	Self-report questionnaires, physiological measurements	Skin potential response, semantic differential (SD) questionnaire [8]
	Perceived safety [96]	Self-report questionnaires, physiological and behavioral measurement	Godspeed (safety questionnaire) [15], NARS [119], distance between human and robot [115], cardio-vascular and electrodermal activity [96]
	Engagement [158]	Physiological and behavioural measurements	Changes in heart rate and skin conductivity changes, non-verbal gestures [7]
System	Productivity [158]	Task-performance metrics, self-report questionnaires	Fan out and interaction effort [121], productivity time [36]
	Fluidity[158]	Task-performance metrics	Agent idle time, concurrent activity, functional delay [75]
	Efficiency[116]	Task-performance metrics, self-reporting questionnaires	Human-robot action or time ratio, time to complete the task [158]; perceived quality of interaction [11]
	Reliability[158]	Task-performance metrics	False alarms, number of interventions [116]

4.2.3 Other Considerations.

It is important to consider the quality of a metric, which is the extent to which it can accurately (i.e., validity) and consistently (i.e. reliability) capture the construct of interest. Donmez et al. [45] identified several other factors that should be evaluated when selecting the type of measurements and the specific metrics to be included in a HRI user study. These factors are:

- a) Experimental constraints such as the temporal and monetary resources associated to collecting and analyzing a specific metric, and the characteristics of the testing environment (e.g. gaze-based metrics make sense in controlled, in-the-lab settings but are impossible to obtain in a field context)
- b) Comprehensive understanding, or how much each selected metric explains the hypothesis of interest and the amount of additional understanding that can be obtained from the casual relationship between metrics (e.g. a decrease in task performance can be explained by an increase in the human-partner mental workload)
- c) Statistical efficiency, the selected metric should provide the researcher with a good number of measurements such that requirements on the statistical power needed to detect potential effects are met; and
- d) Measurement efficiency, unless otherwise required, the type of measurement used to collect a specific metric should not be intrusive or distracting to the participants and to the nature of the study, task and interaction.

4.2.4 Reliability.

Reliability is concerned with the degree to which metrics are free from measurement error. Formally, it is defined as the proportion of variance of a metric that is attributed to the true score variance (i.e., the metric that would be obtained if there were no measurement errors) [160]. Since true score variance cannot be calculated directly [71], in practice, reliability is often described in terms of consistency over four potential sources of error (i.e., time, forms, items and raters) and involves some type of correlation computation.

Consistency over time is often referred to as *test-retest reliability* and can be measured as the correlation coefficient between the metric values obtained under similar circumstances at two different moments in time (stability). Consistency between forms is known as *parallel form reliability* and can be obtained by measuring the correlation between the metrics collected using a different form of the original instrument in an immediate retest procedure. Consistency across items, also known as *internal consistency reliability*, is the most commonly reported type of reliability and corresponds to the correlation between people's responses across the items on a multiple-item instrument. Cronbach's alpha is the most commonly used test to determine the internal consistency of an instrument [71].

When internal raters are used to collect information that serves as the metric of interest, *inter-rater reliability* can be applied to determine the extent to which different observers are consistent in their judgements. Inter-rater reliability coefficients are sensitive to consistency in terms of relative agreement, but they do not determine absolute agreement (e.g., two observers agree on the relative rank ordering of their scores/metrics but attribute different values) [160]. Kappa is frequently used to test inter-rater reliability and detailed instructions around its use can be found in [110]. In addition to the inter-rater reliability, intra-rater reliability measures the consistency of an observer's rating over time and is also calculated using a correlation coefficient. Researchers can capture this measurement by asking a rater to complete the same rating scale at two different time points and a correlation coefficient greater than 0.7 is considered to have good agreement [104].

Independent of reliability type, it is important to note that reliability describes the quality of a set of metrics produced by a testing instrument and not the quality of the instrument itself [160]. Similarly, reliability scores are dependent on the participant sample being measured and hence researchers should always include it when reporting their studies' results.

4.2.5 Validity.

Validity is considered the most important quality of a measured dependent variable. It is defined as the degree to which empirical evidence, i.e., the scales, metrics and instruments employed in a study, actually measure and support the theoretical properties and constructs they are supposed to measure [62]. The general recommendation regarding

construct validity is to select metrics that have been the target of extended validation (e.g. NASA TLX questionnaire) [45]. In cases where non-validated metrics are chosen, construct validity and the degree to which valid inferences can be made based on the proposed metrics is assessed by considering the following four aspects [128], [130]:

- *Face validity* is the extent to which a metric or measurement method appears to be intuitively reasonable and represent that which the researcher is attempting to measure. Since face validity is based on people's intuitions about human behavior, it is considered to be a very weak evidence of construct validity [62].
- *Content validity* relates to the extent to which a measurement reflects and covers all the content that it presumably samples. In other words, does the selected measurement encompass all aspects of the construct it was designed to measure? [71]. Content validity is often assessed using expert judgements.
- *Criterion validity* relates to the extent to which different measurements (also referred to as criteria) reflect the same construct and is measured in three ways [71]: *convergent validity*, which shows that the metric of interest is highly correlated with other criteria measuring the same variable or construct; *concurrent validity*, which indicates the extent to which the measurement or metric being evaluated is related to other criteria or some other construct measured at the same moment in time; and *predictive validity*, which assess the degree to which the metric of interest predicts things they theoretically ought to predict.
- *Discriminant validity* is evaluated by the degree to which the metric being evaluated is not correlated with measures of variables that are conceptually distinct.

4.3 Trial location

Researchers will need to select an appropriate a trial location. There are a variety of different trial locations that can occur when conducting an experiment in human-robot interaction.

Laboratory Testing

Laboratory testing involves conducting the trial in a research laboratory or research team office space. Most studies are conducted in a laboratory setting due to factors such as ease of testing, convenient access to robotic equipment, and rigorous control and observation capabilities. Laboratory testing involves requesting participants to attend a testing session and to use a set-up that has been arranged for the participant to complete. Laboratory studies can generally make it easier to control for unpredictable variables, but observed interactions might be less natural and lack real-world relevance.

Field Testing

Field testing involves deploying the robot in a setting outside of a laboratory, such as a public setting. There is an increased trend to deploy robots into more realistic settings that better represent their final deployment use case, otherwise known as field trials. Field testing sites can include public spaces such as museums [125], shopping malls [86], and schools [100], or involve non-public environments such, as industrial spaces [cite]. Field trials can involve many uncontrolled variables and while experiments conducted in the field better replicate real engagement patterns, subtle evaluations can be lost in noisy data.

Online Environment

Online testing involves conducting the experiment or interaction through a digital medium, such as using tele-presence, online or computer-based mediums. There is growing interest to use robots within virtual environments given the ease of use in terms of deployment and capacity to engage large samples [22]. This method is described further in section 5.3.

4.4 Session Frequency and Duration

To determine duration and frequency of testing session, it is important to consider the following parameters:

- How many conditions?
- How many sessions?
- How long is each session?
- How long is the duration of the whole experiment?
- How frequently to run the sessions?
- What timing is best for the follow-up studies?

All parameters related to duration and frequency can be determined given the hypothesis and independent variable of the experiment as shown by three key examples. For an initial study, researchers should run a pilot proof of concept trial with a limited number of participants to best determine the frequency and duration of the session before larger sample recruitment. An increase in the number of sessions and duration of the study can help reduce the impact of the novelty effect (i.e. perception of a new technology spikes on first encounter). Alternatively, trials that run for a long time can see drop-off effects over time in relation to retention (i.e. Fernaeus et al. [56]).

4.5 Type of Robot

Researchers should consider the choice of robot for their experiment by taking into consideration different capabilities and requirements. An example of parameters could include selecting a robot based on its appearance, mobility, processing power, and level of functionality. There is a limited number of robotic systems that are commercially available and not all the systems have the features that one might be interested to explore (e.g. facial features). Researchers often select the robot model first, choosing a robot that best fits their current constraints such as budget, followed by tests to perform or data to collect. Selection of robots should take into consideration how the target sample will identify the appropriate social behaviour from the robot. To demonstrate, robot appearance can affect how people perceive the robot's sense of intelligence, sociability, likability, credibility, and submissiveness, along with other attributes Phillips et al. [124]. A humanoid robot can be more intuitive but physical resemblance to humans might set a higher expectation in robot capabilities [123]. Zoomorphic robots that look like animals can be less intuitive, but free from the user's preconceived expectations. Decisions on what robot model to use must be made based on the intended use case. Patten and Newhart [123] provides a list of design principles if researchers desire to build a customized robot for the study.

4.6 Task Behaviour of the Robot

Researchers must decide what will be the task behaviour of the robot in the experimental design. There are a variety of different options, and the role of the operator can vary depending on the experimental tasks and the capability of the robot. Robots can act autonomously, or operators can have full control over the interaction, but in some cases, they might only intervene when the robot is unable to make a decision (e.g. [86]). Researchers should specify the details of operator training and the duration of the training during experiment preparation, regardless of the task. In addition, the capability of the wizard and interaction scope need to be fully defined and stay consistent across experiments. For example, if a participant asks a question to the wizard that is outside of the scope of the interaction, the wizard should

answer by mimicking how a machine would react (i.e. I don't know), instead of inserting additional information. Any wizard error should be noted as it can affect experiment results. More information about this technique can be found in [133].

Wizard of OZ is another technique in HRI experiments that is used when current robotic systems are insufficient to handle a fully autonomous interaction in a safe or socially acceptable matter. In this design, WoZ refers to an operator (usually a member from the experiment team, known as the 'wizard') who remotely controls the robot. WoZ can replace natural language processing or other sensing requirements, generates non-verbal behaviour, navigates, localises, and performs manipulation or classification tasks as defined by the experiment [133]. A criticism of the Woz technique is that it does not reflect how a robot can behave on its own in that instance. Therefore, the experiment answers a research question about a proposed robot system that is not yet available or possible. Other ethical consideration regarding this technique both from the operator and participant's perspectives will be explored in 5.

Recommendation When possible and feasible, it is best to use robot techniques and behaviours in trial that a robot can currently use to avoid deception or assessment of effects that are not yet achievable. Wizard of Oz is best used if the researcher must evaluate a technique that cannot be delivered on its own, or the researcher is using the robot as an avatar for experiments.

4.7 Failure and Contingency

Comprehensive experimental planning requires preparing for failure and developing contingency. Failures can happen in the form of equipment failure, no-show participants, or failure to find significance in hypothesis testing.

In terms of experimental failures, if instruments or robotic equipment do not perform correctly during the experiment, the data-point should be removed to avoid influencing experimental results in a negative way. In terms of participants, researchers should aim to recruit more than what is needed to achieve statistical significance and to fill in data that needs to be removed due to experimental failures. If permitted by ethical guidelines and approval, researchers could keep a contact list to find additional participants in case of trial absence. Finally, many factors could affect the statistical results. Prior to participant recruitment, data analysis and data cleaning procedures should be established as the definition of outliers and analysis will affect the results (for more information, see Section 6). As noted by Bethel et al. [22], data management procedures should be in place, including backup equipment and data capture whenever possible to avoid data loss.

4.8 Examples

An example of a HRI study that focuses on a robot-related construct (i.e. social intelligence) can be found in [152]. This is a preliminary work into creating social intelligence, where assistive robots will be able to change behaviour based on the affect of the interactor, promoting a positive interaction feedback. Towards system-related constructs, Hoffman and Breazeal [76] pioneered the first fluency evaluation in HRI with a physical robot. The authors hypothesized that the human-robot collaboration fluency can be improved by replicating the anticipatory behaviour in human-human collaborations. The authors implemented a cognitive framework in a non-anthropomorphic robot and designed an experiment to evaluate the framework. The participants were asked to rate the robot on a customised survey in the between-subjects study. Cronbach's alpha was calculated to assess the reliability of each survey item as described in Section 4.2.4.

5 PARTICIPANTS

Tina and Dana

HRI research envisions robots in a variety of settings, interacting with users. A major requirement to validate proposed approaches is the participation of users during experiments. Therefore, recruiting the appropriate participant population is crucial. In this section, we first identify the possible types of participant involved in a HRI study and highlight the important steps to participant selection, including sampling, randomisation, and blinding. We then discuss the procedure for participant selection, recruitment, retention, and the steps for ethics application and other ethical considerations pertaining to the people involved in HRI studies.

5.1 Types of Participants

Participants of a HRI study can be classified into four categories: observers, WoZ operators, teleoperators, and interactors. The role of the observer might involve viewing videos of HRI interactions and interpreting the content based on the HRI evaluation metrics (refer to Section 4.2.2). Wizards, or operators, are participants who operate the robot and may require special training. Finally, the interactors are the participants who interact with the robot.

5.2 Participant Selection

5.2.1 Sampling.

Sampling is the process of identifying a representative group of participants suitable for the research study. An unbiased sample where participants are randomly drawn from a population is usually difficult to achieve; hence, biased (i.e. non-random) sampling techniques are often used, such as convenience sampling and purposive sampling [51]. Convenience sampling is often used when participant recruitment is limited by practical constraints, such as participant accessibility, proximity, and availability [51]. The majority of existing HRI studies use convenience samples involving university students [14, 16]. The use of convenience sample generally reduces the external validity of the research findings as the participant population is not representative of the user population in terms of social traits, cognitive performance, and attitudes towards technology [14, 16]. Geographic proximity is often another reason for using convenience sample, which comes with the trade-off of limiting the cultural diversity of the participants. In general, convenient samples introduces sampling bias (e.g. excluding certain groups) or sampling error (e.g. including a higher proportion of certain groups) [123]. When using convenience sampling (e.g. university students), researchers might consider balancing potential confounding variables, such as gender or educational background, during participant selection [14].

In purposive sampling, participants are deliberately chosen for certain characteristics that they possess, such as knowledge or experience that would enable them to assist with the research. There are different types of purposive sampling methods including maximum variation sampling, homogeneous sampling, typical case sampling, extreme/deviant case sampling, total population sampling, expert sampling and specific examples can be found in Etikan et al. [51]. Sampling affects the statistical power of the study, which is determined based the combination of the number of trials, the size of the study (number of participants), and the number of random or uncontrolled factors. A lack of statistical power can lead to both type I errors (false positives) or type II errors (false negatives) and cause the wrong conclusions to be made. Sampling more participants from a diverse population is the standard and most straightforward way to increase statistical power [43]. More discussion on the sample size can be found in Section 6.1.3.

Recommendation Balancing potential confounding variables can help with noises introduced by the participant selection process, especially when non-random sampling techniques are used.

5.2.2 Randomisation.

Randomisation and stratification can help assign participants to groups, conditions, and observers while minimising biases. Randomisation involves assigning participants with an equal chance to experimental groups to explore the effects of the condition and prevents researchers from predicting the results of the conditions. Simple randomisation is a basic form of condition assignment, which can involve methods such as a coin toss or simple computer-generated function [19]. Simple randomisation can be used for sample sizes above >200 participants, given their statistical likelihood to even out in terms of allocation, but should not be used for <100 participants as there might be an unequal number of participants in each group [20, 87, 88].

Stratified randomization can be used when the characteristics of the participants need to be equalized across the experimental groups. For instance, if age is a factor that can affect the HRI experiment, researchers can stratify the participants to ensure that similar age groups are present in each participant group. However, factors relevant to the phenomenon need to be identified and normalised, which might not always be possible. There are other stratification techniques such as block randomization [44] or response adaptive randomisation [73], which is a technique to adjust the condition assignment favoring conditions with better performance based on ongoing data collection. While these techniques are less commonly used in current HRI studies, they can be useful in longitudinal studies with socially assistive robots.

Recommendation Simple randomisation such as a coin toss or random number generator can be used for a large number of participants, but should not be used for a small participant group. Stratified randomization is best used for equal cross across multiple factors, such as age. Block randomisation and response adaptive are more sophisticated randomisation techniques and should be implemented on a needed basis.

5.2.3 Blinding and Allocation Concealment.

Blinding and allocation concealment are techniques that are used to further remove bias effects in experimental condition assignment, and must be decided on in the study design [140, 148, 178]. Blinding is the process of preventing the participants, researchers or both parties from knowing what intervention or experimental content participants are allocated to. Conventional blinding strategies include single-blinded or double-blinded studies. Blinding is often not reported Hróbjartsson et al. [78] and researchers should include this information as part of the study design.

Allocation concealment is when the researcher does not know the upcoming experimental condition until the moment of assignment. Blinding reduces the effect of observer bias (i.e. the observer's judgement is influenced by the knowledge of the group), whereas allocation concealment reduces allocation bias (i.e. how participants are assigned to each group) Ruxton [140]. Both blinding and concealment can be executed cheaply. For instance, researchers studying the effect of different robot appearances can ask a researcher unrelated to the project to randomly assign letters to each robot. The mapping between the robot and the coded letter is concealed from the researchers who are analysing the data until the data has been processed. Allocation concealment can be implemented using paper-based techniques Doig and Simpson [44] or using online tools such as SEPTRE.

Recommendation Blinding should be reported in the final report as it represents an important methodological step. Allocation concealment is an important step for demonstrating that the experiment is not influenced by the researcher's bias.

5.3 Participant Recruitment and Incentivisation

In participant recruitment, the required sample size, the best recruitment venue, and ways of motivating participants

to join the study are among the key considerations. The right sample size will depend on the study design structure and the desired statistical power as discussed in Section 6.1.3. Some participants are inherently more difficult to recruit, e.g., participants for long-term studies or studies that require special participant groups (e.g. patients, children) and researchers might inquire about the participant's future availability and willingness to partake in other relevant studies. On the other hand, researchers might need to consider when not to reuse participants. Scenarios where participants should not be reused include when the study requires first exposure to the system. Participants who have prior experience in other studies should be excluded to prevent cross-condition contamination [1].

Participants can be recruited through contacts (phone contacts, email, social media, other web-based techniques), community groups, mailing lists, volunteer banks, schools, external collaborators [97] and crowd-sourcing platforms such as Amazon Mechanical Turk [31, 32]. The techniques and the corresponding recommendations are summarized in the table below.

Regardless of the recruitment method, researchers should always pay special attention when recruiting vulnerable people. Vulnerability is not limited to disabilities but also includes people who might be exploited and abused due to an existing unequal or dependent relationship [157]. For instance, students participating in the supervisor's study or clinicians recruiting patients for trials. This also includes recruiting close friends or family, who may comply with the experimental manipulation or task to please the researcher who is conducting it. A simple way to minimise the effect of potential unequal relationship is through information disclosure (i.e. who are involved with the project, how privacy will be protected, and state that results of the study will not affect the medical care a patient might receive).

Motivating participants to enroll in the study might be another challenge. Different incentives might be provided, such as food for students, gifts for children, raffles, and participant honoraria. Participant compensation should be proportional to the amount of time required and the type of participants involved. However, biases might be introduced when a reward system is set up based on time, as financial incentives might cause data from a subset of participants to be over-represented in the entire sample [161], and interfere with the participants' ability to weigh the risks and benefits of the study [59]. When compensation is involved, the rule of thumb is that the payment should not be provided as a reward for the risk or harm [59]. Researchers should follow the standard procedure defined by the host research institute and maintain transparency. Ethics committees might only allow small amount of incidental expenses to be reimbursed. In addition, researchers can motivate participants by providing accommodations depending on the context of the study, such as offering to conduct the experiment at the participant's home to accommodate for the elders who might not be able to travel. Finally, incentives must be suitable and not cause additional risk for participation or unfairly encourage people to participate who are then not voluntarily providing consent, i.e. they are financially challenged and participating in the trial only to receive the money incentive.

Recommendation Each recruitment platform comes with different advantages and trade-offs. Researchers need to consider the platform in relation to the participants relevant to the HRI study. In addition, suitable incentivisation can be provided given that it does not affect the participant's consent.

5.4 Participant Retention

Longitudinal studies are particularly susceptible to the problem of participant retention. Participants dropping out over time can result in a type of sampling bias known as mortality. In clinical research, participants can drop out of the experiment if they do not feel the benefit of the therapy, causing the remaining sample to be biased towards the groups with more effective treatment [123]. In the HRI context, mortality is observed when participants stop using the technology over time, especially in the longitudinal, in-home studies [41, 56]. The reasons for the participant's

Table 3. Participant recruitment methods

Method	Advantages	Disadvantages
Contacts	Researcher contacts are more likely to agree to participate.	This technique might be subjected to snowball sampling, which is a biased sampling technique where one participant refers other potential participants [123]. Participants who are contacts of the researchers might be less likely to provide negative feedback.
Community groups, professional organisations	This method makes it easier to contact a large amount of specialized participants at once (good for purposive sampling).	Response rates might be low and lead to a biased sample. Researchers can consider running a short presentation at the corresponding organization to raise awareness and gauge interest.
Mailing list	This recruitment method is cheap, fast, and easy to execute.	Response rates might be low and lead to a biased sample. Cold-emailing might solicit negative perception towards the study and hence it is preferred to recruit through existing mailing lists.
Volunteer banks	Volunteer banks enable access to large number of participants.	Volunteers might be fundamentally different from non-volunteers leading to sample bias. For instance, volunteers might be more excited about technology, have prior experience with robots, or be more willing to adopt a new technology compared to non-volunteers [123].
Crowd-sourcing platforms	This method is also cheap, fast, and easy to conduct. Different filters on the platforms or sites targeted at different academic projects also enable researchers to include a more diverse population or find participants with specific skills [31]. The experimental tasks typically involve viewing pictures or videos online to replace in-person HRI [38].	In addition to issues with the data consistency and reliability, researchers might neglect the fact that crowd-sourcing workers can have prior experience with similar studies (i.e. they are nonnaïveté), bring their pre-assumed knowledge into the study, which in turns reduces the effect size of the study [32]. Workers might also discuss the experiment in the worker discussion boards where they might obtain foreknowledge about the experiment and be influenced by the opinions of other people [31]. Another HRI specific concern is that online surveys and passive observation may not reflect participants' perception and interaction with the robot in real life. Guidelines on crowd-sourcing platforms can be found in [164].

departure may not be random and it is important to record mortality in longitudinal studies and its occurrence in each group.

Recommendation Participant retention is important especially in longitudinal studies and participant dropping out should be noted in the study.

5.5 Participant Preparation and Training

Once the appropriate participants are identified and recruited for the study, planning for participant preparation and training might be needed especially for experiments that involve observers or operators. In HRI literature, commonly the details of the training program including its length and approach are not reported [133]. It might be worth considering the following questions when planning for participant preparation and training. A full description of participant debriefing is available in [5].

- What information does the participant need prior to the start of the experiment in order to give informed consent? Is deception involved?
- What training does the participant need to undergo?
- How will the observers be trained (e.g. what scale is being used to evaluate the HRI session and what are the evaluation standards)?
- Do the initial results present good inter-rater/intra-rater agreement? For instance, some experiments might require observers to undergo training, so that they can label affects from people's facial expressions. Researchers can administer a test prior to the start of the actual experiment to determine whether the participants had been trained sufficiently. If the initial results show poor rater agreement, researchers might consider extending the training or removing some observers.

Recommendation Participant preparation, such as training for WoZ operators/teleoperators, should be noted in the study. In addition, researchers should devise strategies to ensure that the participants are sufficiently trained.

5.6 Privacy and Ethical Considerations

This section discusses other considerations about participants from the procedural and ethical aspects of the experiment including privacy protection, informed consent, and special considerations for different experiment setups.

5.6.1 Privacy protection.

HRI studies follow guidelines similar to other research that deals with human participants. Maintaining the privacy of the participants is important as it protects them from potential harms, including embarrassment or fear of being judged by others. The basic guidelines noted by Lazar et al. [97] include:

- collect the minimum amount of information,
- limit the number of times data is used,
- aim for full information disclosure whenever possible,
- provide a secure data storage medium (especially when cloud-services is used),
- remove personal identifiers whenever possible,
- follow standard data disposal procedures recommended by the institution,
- and provide channels to address concerns and enforce accountability.

5.6.2 Informed consent.

Obtaining informed consent in HRI studies is similar to any other research involving human participants. Participants must understand the purpose, procedure and risks involved with their participation and their rights to data and information. For experiment setups that require deception (e.g. WoZ studies), informed consent cannot be obtained and researchers have to purposefully withhold information in the study. While the use of deception raises ethical concerns, it is sometimes needed especially when the phenomenon under investigation is prone to reactivity, which is when the

measurement changes the participant's behaviour Price et al. [128]. It is important to devise a post-study debriefing procedure during ethics application for this type of HRI studies to fulfill the researchers' duty in information disclosure.

5.6.3 Special considerations.

In addition to being aware of the participants who might require special care, accommodation, or the region specific standards when preparing the experiment protocol, this section is dedicated to highlight potential concerns relating to other special participant groups and research practices.

As HRI is a relatively new field, a lot of the phenomena are not well understood and researchers should be aware of the potential long term impacts on the participants across all age groups. For instance, Lemaignan et al. [100] conducted a study with a social robot in a classroom setting. They pointed out the potential impact of the robot, changing the dynamics between the teachers and the students in the classroom. In addition, de Graaf [40] noted that long term human-robot interactions can potentially foster human-robot relationships, where humans start forming emotional bonds towards the technology. Considerations should be given beyond risk assessment from the safety perspective. The role and utility of the technology in the user's life should also be taken into account, especially in settings such as home, long term care, workplaces, and schools. Researchers should be aware of the ethical implications involved when developing robots with human-like qualities or high level of autonomy. These functions might fundamentally be a form of deception, as the users might wrongly believe that the robots have human characteristics when they do not.

Researchers should also be aware of the potential for participants, such as Wizard of Oz operators, to have a negative experience during the study. For instance, in Rea et al. [131], the authors reviewed social HRI experiments that put the WoZ operators in stressful situations, such as having to cheat [155], enduring verbal/physical abuse, or simulating other behaviours that would otherwise be socially awkward (e.g. long periods of silence). Researchers should minimise the duration of the necessary social conflicts in the interaction script, allow the operators to have positive interactions with the participants in post experiment debriefing, and have means for the operators to be trained and express their concerns Rea et al. [131].

COVID-19 has highlighted the importance of hygiene in every procedure including experimental studies to protect both the participant and the researcher. Experimental procedures should take into account the appropriate institution and government guideline relating to social distance and minimize contact whenever possible. As some robots cannot be sanitized using alcohol-based methods, procedures should be established between the host institution and the lab to ensure proper hygiene (e.g. have people who need to work with the robots sanitize their hands before and after working with the robot).

Recommendation Considerations need to be given to all possible HRI participants, from the interactors to the researchers participating in the study, to mitigate potential negative effects of the study.

5.7 Ethics application

Most academic research projects are constrained under guidelines established by different governments or institutions. For instance, the most recent guidelines for EU projects are defined in the Horizon 2020 framework and the original principals are derived from the Charter of Fundamental Rights of the European Union, and the European Convention on Human Rights. Projects in the United States are fundamentally based on the Belmont report and Hippocratic Oath. In Canada, the Tri-Agency Policies and Guidelines govern all ethical conducts for human research. In Australia, the standards are defined under the National Statement on Ethical Conduct in Human Research. While most research projects in collaboration with industry often go through the university ethics application system, some may choose to go with commercial ethics review boards. This practice is not recommended as these commercial review boards

might be subjected to conflict of interests [101]. While pilot studies involving only internal researchers may often begin during development stage, any formal user studies involving recruited participants should only commence after ethics approval has been granted.

5.8 Examples

The study conducted by Moyle et al. [114] demonstrates sample randomisation and stratification techniques. The authors investigated the use PARO robots for treating the behavioural and psychological symptoms of dementia. The experiment took place in 28 different long term care facilities where the facilities were first stratified based on the organization type (private vs. nonprofit), then the experimental conditions were randomly assigned in blocks. A cluster randomized control trial was employed to minimize the effect of between-group contamination, as patients in different care facilities are inevitably exposed to different activities and treatment environments.

Lemaignan et al. [100] illustrate procedures for participant recruitment and study design in a study aiming to integrate robots in childhood education for both normal children and children with disabilities. This study demonstrated how an iterative study design can be used to address the challenges of validation with a vulnerable population group. The system was first validated with a more accessible population and then with the clinical group. Finally, the authors conducted two in-depth case studies at the lab with 2 children over the period of a month. Since the study took place under multiple experimental settings (i.e.: schools, the clinic, and the lab), the authors are able to identify differences in the children's behaviour in the lab compared to other scenarios.

[86] illustrate participant recruitment in a target field environment. The authors released a semi-autonomous robot in a shopping mall, where it interacted with visitors every weekday for four hours over five weeks. 332 participants were recruited by flyers that were distributed around the mall or approached by the experimenters who were onsite during the first three weeks of the study. At the end of the study, a questionnaire was mailed to each participant and 70% of the participants responded.

Wang et al. [173] illustrate the importance of considering cultural context during recruitment. The study explored how the cultural background of the participants (Chinese vs. US) affects the human-robot collaboration process when using two different communication styles (implicit vs. explicit). The research team was able to find representative samples by recruiting participants directly from two comparable universities in China and the US and conducting the experiments in the respective countries.

Yanco et al. [181] illustrate recruitment of specialised participant groups. The study evaluated HRI interfaces in-situ, developed for expert users for use during the DARPA trials. The authors used the unique event to collect data with expert users. They observed the teams participating during the trials and analyzed the HRI methodologies used for robot teleoperation.

6 DATA ANALYSIS AND STATISTICAL METHODS

Once all measurements have been collected, the next step is to analyze the data and determine whether the research hypothesis of interest is well grounded and supported by the data. To ensure that the conclusions drawn from this analysis are valid and generalizable, researchers must carefully identify any threat to the validity and generalizability of their results (e.g., extreme metric values or insufficient statistical power), take the appropriate actions (e.g., reduce error variance), and choose an adequate statistical analysis for their study design and data. Before diving into the data analysis and statistical methods that can be used for this purpose, we first cover general recommendations on how to clean, better understand, and post-treat the data. We then provide an overview of the main statistical methods found

in the literature as well as some recommendations on who to best use these methods. We end this section with some examples that showcase how statistical analysis are done in HRI studies.

6.1 Considerations Prior to Data Analysis

No matter how well designed and implemented a user study is, researchers frequently have to deal with errors from various sources and their effects on study results. As a first step in identifying these errors, researchers should employ appropriate visualization techniques to better understand and examine their data. Next, they should leverage a priori defined data cleaning strategy so as to make the data as free of errors as possible. Finally, researcher should take into consideration the type of measurements included in their data and the statistical power of their design and participant sampling scheme when deciding on an appropriate statistical method for their data.

6.1.1 Data Cleaning.

Data cleaning can be described as an iterative and repeated three-stage process in which data errors and abnormalities are screened, diagnosed and edited. However, before diving into this process, it is recommended to formulate a set of predefined rules for dealing with data errors, and missing or extreme values. Researchers must keep in mind that it is not always immediately clear whether a suspected data point is erroneous and why it is so. Similarly, missing values could be due to interruptions during data collection or unavailable information [170]. Thus, researchers should utilise their knowledge about potential technical errors (e.g., measurement errors) and expected ranges of normal values to define the rules to follow during data cleaning.

Screening: A good starting point for detecting invalid data is to do a visual scan of each participant's responses and measurements. Patterned responses, or responses completed in less time than what was observed in the across all participants are often employed as indicators of invalid data [102]. If instructional manipulation checks or attention check questions, that is, questions with obvious answers, were included among the measurements collected during the study, they can also be employed to identify potentially erroneous data points. This screening method is frequently used in studies in which participants were recruited through crowd sourcing services [70]. During screening, researchers should also look for lack or excess of data (e.g., missing values), outliers, and strange patterns in the data distributions.

Most of the the methods used to evaluate possible outliers are based on the idea that a particular proportion of valid data points will exist within a given k number of standard deviations from the population mean. By setting a threshold on the number of standard deviations considered as valid, data points above and below that threshold are considered to be outliers. Statistics such as a the first ($Q1$) and third ($Q3$) quartiles of the data as well as the interquartile range ($IQR = Q3 - Q1$) are also used to defined outlier detection bounds. Data points outside these bounds are considered to be outliers and are thus eligible to be excluded from any posterior data analysis. In the case of a sample size that is relatively small, alternative methods such as the Grubb's test (see [102] for more information) are recommended instead. When looking for outliers, it is important to keep in mind that although all outliers are characterized as extreme data points, not all extreme data points fall into the outlier category. Similarly, researchers should be aware of erroneous inliers, that is, data points generated by error but that fall well within the expected range of values.

Diagnosis: In this phase, researchers purpose is to clarify the nature of all the potentially problematic data points identified during data screening. Data points should be categorized as: erroneous, true extreme, true normal (i.e., prior expectations about normal data ranges were incorrect), or idiopathic (i.e., there is no explanation, but the data point is still suspect of error). For the suspected points for which diagnosis is less straightforward, i.e., they do not fall into the expected ranges of true extreme or outlier values, the application of a combination of diagnostic procedures is recommended [170]. Examples of these procedures are: determine whether data points were consistently the same

throughout the whole data collection procedure or collect additional information, e.g., question the experimenter in charge of data collection, and if possible, repeat the measurement.

Editing: After the identification of errors, missing values and true values, the next step is to decide what to do with these problematic data points. Overall, researchers should decide whether problematic data points should be corrected, excluded or left unchanged. In the case of impossible values, the general recommendation is to correct them if a correct value can be found, or exclude them otherwise. In the case of missing data, the researcher should decide on the amount of missing data that is acceptable before leaving a participant's data out of the analysis as well as on what acceptable substitute values should be used instead (e.g., samplewide median values are often employed to substitute a missing variable or metric) [102].

Whatever the data cleaning strategy researchers employ prior to data analysis, they should provide detailed documentation of the data-cleaning methods, error types and rates, error deletion and correction rates in their study report. Similarly, researchers should report on the differences in the outcomes of their studies with and without outliers [170].

6.1.2 Data Visualization.

Simple visualizations can be made for both individual participants and the aggregated responses with the goal of seeing the trend (the direction of data progression over time), level (changes in relative value of the data over the dependent variable), and stability/variability of the data between experiment conditions and participants [95]. Four common techniques are highlighted in this section: cumulative records, semi-logarithmic charts, bar graphs and line graphs.

Cumulative records are generated by summing the participants' responses across sessions. The method is useful when there is a progression to the experimental conditions (e.g. in each experiment, the robot learns to do a new skill in addition to what it can do previously) and a survey is administered at the end of each condition. This type of graph is a simple way to illustrate trends (increasing, decreasing, no change) in the total response within a study.

Semi-logarithmic charts are commonly used to display the rate of change or proportional changes in performance. These graphs are common in machine learning to track the accuracy of the algorithm over the training iterations. Both cumulative records and log charts are for continuous data, whereas bar graphs are common for discrete data or continuous data for studies with small sample sizes [95, 176]. Bar graphs are typically used for presenting summary statistics between experimental conditions (e.g. user ratings between the control condition versus the robot condition). However, summary statistics can be problematic as many distributions can lead to the same graph while other important features of the dataset are hidden [61, 176]. It is therefore recommended to plot the data distributions or use scatter plots instead of bar graphs. Weissgeber et al. Weissgerber et al. [176] implemented a free, online tool for data visualization, which can serve as a starting point for researchers when deciding between different visualizations. Finally, line graphs are often used when there is a temporal aspect to the data (e.g. motion tracking during an HRI experiment) [95].

6.1.3 Sample Size and Statistical Power.

A study's statistical power, that is, the probability of detecting a significant effect of the factor being manipulated on the dependent variables being measured if it exists, is directly tied to the number of samples (or measurements) being analyzed. Thus, the likelihood of detecting small or even moderate effects vanishes as the sample size decreases [113]. Similarly, the likelihood of detecting an effect when there is none (Type I error) or failing to detect an effect (Type II error) can be reduced with greater statistical power. As a result, statistical analysis are more reliable when the sample size is large and the subsequent statistical power is greater [108]. The general recommendation is to perform a pre-study power analysis to estimate the appropriate number of participants and measurements required in order to achieve

adequate statistical power or accurate parameter estimates [22, 123]. Tools like G*Power 3 can be used [54] for this purpose.

For the cases in which data from small sample sized studies still needs to be analyzed, [113] provided the following recommendations. First, in the case of categorical data, alternative exact tests such as the Fisher’s test [109] should be preferred since more common tests (e.g., the χ^2 -test) are known to lack accuracy with a small sample size. Second, in the case of continuous, interval data, it is often advised to verify key assumptions such as normality and sphericity before employing classic parametric methods such as t -tests or Analysis of Variance (ANOVA). However, in small sample sized data, the tests employed to verify these assumptions are likely to be under-powered and may lead to incorrect conclusions about the validity of these assumptions. In these cases, Welch’s extensions to the t -test and ANOVA or non-parametric options such as the Mann-Whitney or Kruskal-Wallis tests are recommended. It is important to notice however that non-parametric tests are known to have less statistical power than their parametric counterparts. Third and last, when visualizing small sample size data, the use of common visualization methods such as histograms and box plots can be misleading and hard to interpret. Scatterplots are suggested as the best choice for showing the distribution and trends of the data in this case [175].

6.1.4 On the Type of Measurement Scales.

Putting interviews aside, most of the measurements obtained in HRI studies can be categorized into continuous, interval values (e.g., behavioral measurements such as the distance between a human and a robot or physiological data such as heart rate) or ordinal data (e.g., Likert items and scales often used on self-reporting questionnaires). While the choice of statistical methods to analyze the former is a more or less straightforward process³, there is an ongoing debate on what is the correct way of analyzing the latter.

On the one hand, some researchers argue that ordinal data obtained using Likert scales can be treated as interval data and thus standard parametric statistical methods can be used when analyzing this type of measurements. On the other hand, a more conservative researchers choose to employ non-parametric tests instead even though it is known that these methods are less powerful and lack sensitivity to detect smaller effects [159].

Recently, Schrum et al. [147] provided a set of recommendations on the use and analysis of Likert scales in HRI studies. These recommendations are: i) employ summary statistics such as mode, median, range and skewness when reporting individual Likert items; ii) although there is compelling evidence suggesting that parametric tests such as ANOVA can be used on Likert scale data [159] as long as the appropriate assumptions have been tested and validated, a conservative approach in which non-parametric tests are used instead, is recommended; iii) analysis should be performed on a multi-item scale instead of on single Likert items.

6.2 Statistical Analysis

This section provides an overview of the the most common statistical analysis found in the literature: hypothesis testing methods, confidence intervals and Bayesian inference methods.

6.2.1 Hypothesis Testing.

Hypothesis testing is the most common use of statistics. We usually see in the papers or reports that a *null hypothesis* is rejected or retained with a *p-value* < 5% or *p-value* < 1%. This hypothesis test is a formal approach for deciding between two interpretations of a statistical relationship in a sample. One interpretation called the *null hypothesis* (often symbolized H_0) which refers to no relationship in the target population. The other interpretation is called the *alternative*

³The literature suggests that classic parametric tests (e.g., ANOVA) are vastly preferred if assumptions of normality and sphericity are valid, otherwise non-parametric tests (e.g., Kruskal-Wallis) are used instead [172].

hypothesis (often symbolized H_1) which refers to the existence of relationship in the population which is reflected in the population sample as well.

Significance level is the threshold at which it is decided whether the null hypothesis should be rejected or retained. This significance value is also known as *p-value* and it relates to the probability statement made about the observed sample in the context of a hypothesis, not about the hypotheses being tested [6]. The smallest significance level that is normally considered as a reasonable evidence is 5%. In practice, if the probability of observing the *null hypothesis* statement H_0 is less than 5%, the *null hypothesis* is rejected in favor of the *alternative hypothesis* H_1 . The *p-value* represents the probability of observing a similar statistical relationship if the data being analyzed was generated from random samples Kaptein and Robertson [90]. P-values do not give the probability of a hypothesis being true or false for this particular experiment, they only provide a description of the long term Type I error rate for a class of hypothetical experiments. Similarly, p-values do not indicate whether the means of the samples being analyzed are either equal or not equal [6].

Hypothesis testing methods can be classified into a) parametric and b) non-parametric tests with respect to the satisfaction of some assumptions. These two classes are discussed in detail in the following sections.

Standard Parametric Methods

Parametric tests are widely used in analyzing user study's data. These tests include t-test, analysis of variance (ANOVA), and ordinary least squares regression. Here, we focus on the ANOVA test as it is the most known and (mis)used test. ANOVA test assesses the potential difference between two or more groups on a *continuous* measurement. Before using ANOVA, there are some assumptions that need to be met in the tested data [28]. These assumptions are as follows:

- (1) Normality of data samples
- (2) Homogeneity of variance (sphericity)
- (3) Independence of samples

There are different tests for normality such as Chi-square, Kolmogorov-Smirnov, Shapiro-Wilk, Jarque-Barre, and D'Agostino-Pearson. Some researchers argue that if the sample size of all groups are equal (balanced model) and sufficiently large, the normality assumption can be relaxed provided the samples are symmetrical as noted by Blanca et al. [25] and Feir-Walsh and Toothaker [55]. For the homogeneity assumption, it can be tested by simply comparing variance values or through specific statistical tests such as Levene's, Fligner Killeen, and Bartlett's tests. ANOVA is not robust against unequal variance [66], thus if this assumption is violated, it may alter both the type I error, i.e., detecting an effect when there is none, rate [69] and the statistical power of the test [118]. For the independence assumption, there is no specific test that proves its validity. However, it must be taken into account during the design of the study itself. For instance, the observation data may be dependent if repeated measurements are collected from the same subject.

From the study design perspective (Section 4), ANOVA can be categorized based on the study model as well as the number of independent variables considered in the study [162]. In the following, we list the different types of ANOVA:

- (1) **Between-subjects ANOVA** is used when examining the differences between two or more independent groups. Within this type of ANOVA there are one-way ANOVA and factorial-ANOVA. The main difference between these two branches is the number of independent variables to be tested.
 - (a) **One-way ANOVA** is used when assessing the difference between groups with respect to one independent variable. In practice, one-way ANOVA used when studying the difference between at least three groups as in the case of two groups, t-test is usually used.
 - (b) **Factorial-ANOVA** is used when examining multiple independent variables.

- (2) **Within-subjects ANOVA**, also called repeated measures ANOVA, is used when the same subjects were tested for each experiment's condition. It is frequently used with pre- and post-test experiment design. However, it is not limited to only two time periods, it can be also used when examining the differences over two or more measurements in time. When using this type of ANOVA, it is also important to test the sphericity assumption, that is, the variance in the differences between all pairs of groups are equal, is valid [107]. This assumption can be tested using Mauchly's test.
- (3) **Mixed-model ANOVA** is also called within-between ANOVA, is used when studying the differences by group and time. That is, when the study follows both a between-subjects and within-subjects design.
- (4) **Multivariate analysis of variance (MANOVA)** is used when studying the differences between groups on multiple dependent variables. We can not replace it with simply performing multiple ANOVA's for each dependent variable as ANOVA would not account for the correlations between the dependent variables.

While ANOVA test determines whether the differences observed among groups are due to chance, it does not identify which particular differences between pairs of groups are significant/not-significant. This is what *post-hoc* tests do, they are used after the ANOVA test to assess the differences between all possible group pairs (multiple comparisons) [79]. The more groups you have in ANOVA test, the higher Type I error rate you get. Type I error rate (i.e., false positive) is defined by the significant level in the case of one comparison (two groups). The error rate inflates with the increase in the number of groups which in turn increase the number of comparison [166]. Post-hoc test constrains the experiment-wise error rate to the significant level with the adjusted p-value. A variety of methods exist in the literature for conducting Post-hoc test. Here, we focus on the most common tests available in statistical software.

- (1) **Bonferroni test** is simply a series of t-test performed on all possible pairs combination of the tested groups. In order to limit the Type I error rate to the significant level, Bonferroni test sets the significance cutoff at the significant level divided by the number of comparisons. This is called *Bonferroni correction*. Thus, Bonferroni test tends to be more conservative.
- (2) **Tukey's Honest Significant Difference (HSD)** is used when we really care about the pairwise differences in groups. It gives an estimate of the difference between the groups and a confidence interval for the estimate. This test is used with balanced data; the sample size is equal between groups.
- (3) **Tukey-Kramer HSD** the same as Tukey HSD except it can be used with unbalanced groups.
- (4) **Fisher's Least Significant Difference (LSD)** is the most 'liberal' method as it has a high probability of Type I error rate. This test is computationally identical to the Bonferroni test except for the fact that it does not employ any adjustment in Type I error [153].
- (5) **Student-Newman Keuls (SNK)** orders the groups' means difference from the largest to the smallest in a step-wise procedure. Then, it starts with the largest difference to test if it is significant, if so, it continues in order till it reaches a non-significant mean difference at which the test is terminated. It is more liable to Type II error (i.e., False negative) than a Type I error [74].
- (6) **Duncan's Multiple Range** is similar to the SNK test, but the probability of making Type I error increases with the increase in the number of comparisons.
- (7) **Scheffe's test** is very conservative method which gives the highest protection against Type I error. Also, it is robust the violations of the ANOVA's assumption [10]. However, all of these come at the cost of less power for detecting effects.

- (8) **Dunnett's test** is used when we care about the difference between control group and the other groups; not considering all the pairwise comparison. Also, it can be used with unbalanced groups as well as if the homogeneity assumption is violated.
- (9) **Binomial Sign test** uses only the signs of the differences between the pairwise comparisons.

For more details about the post-hoc test and types, we refer the reader to [77], [74], and [28].

Alternatives to Classic Parametric Methods

When deviation from parametric assumptions is a concern, sample sizes are small or there is uncertainty about the choice of the type of distribution that is the most appropriate for the data collected, non-parametric methods and modern robust statistical tests offer viable alternatives. [103].

Non-parametric methods require minimal assumptions about the underlying distribution generating the observed data. This is done through the use of the ranks of the original data instead of the data itself, and the assessment of test statistics significance through randomization tests. Ranks have the advantage of not being affected by the presence of outliers or skewed distributions and allow for test statistics distributions that do not depend on assumptions such as normality. Similarly, randomization tests, in which all ways to permute the data are considered, allow to assess a test statistic significance without making explicit assumptions regarding distributions [58]. Some of the most commonly used non-parametric tests are:

- **Wilcoxon Mann Whitney Test**, also known as the non-parametric analog of the t -test. This method is a rank-based test for comparing two populations on a continuous outcome using independent samples, e.g., compare the weight of males to that of females.
- **Wilcoxon Signed Rank Test** is the non-parametric equivalent of the paired samples t -test. This method examines whether two samples were drawn from the same population and can be used when comparing two related samples, matched samples, or repeated measurements on a single sample.
- **Kruskal-Wallis Test** is the non-parametric analog to a one-way between-subjects ANOVA. The method analyzes the population medians of the ranked form data
- **Friedman's Test** is the non-parametric equivalent of the repeated-measure ANOVA. This test examines the data based on its rank properties. Hypothesis testing for the Friedman's test may be expressed by the medians or the average ranks.

Robust methods are particularly advantageous since they provide adequate control of Type I error rates and increase the likelihood on discovering relevant differences between groups [50]. Robust methods replace traditional regression methods (i.e., ordinary least square), measures of location (i.e., the mean) and measures of associates (i.e., Pearson correlation coefficients) with robust alternatives. These alternative measures can be later used to perform hypothesis testing.

Bootstrapping methods provide precise estimates of population distributions by iteratively resampling cases from a set of observed data. They can be used to obtain accurate confidence interval estimates in the presence of outliers or strongly skewed data [138]. Rank-based methods such as rank transform ANOVA-type statistic and Wilcoxon analysis offer extensions to classic non-parametric methods and are known to produce valid results when analyzing data that is non-normally distributed and/or with different variance between groups [50].

There are still some limitations associated to use of these robust alternatives. For instance, there is a reduction in the number of degrees of freedom available for statistical tests and the estimation of sample size required for sufficient

statistical power is more complex [91]. We refer the reader to [50], [177], and [91] for a more detailed introduction to these robust alternative methods.

6.2.2 Confidence Interval.

Confidence interval (CI) use the same underlying mathematical methods as the hypothesis testing, but instead of giving a probability of a single value, it gives a range of values [43]. The concept of the CI was introduced by Jerzy Neyman in a paper published in 1937 [117]. It is powerful for giving a sense of the level of uncertainty around an estimate or prediction. It can be defined as a range of values, calculated by statistical methods, that includes the desired true parameter with a probability defined in advance; called the *confidence level*. The size of the confidence interval depends on the sample size and the standard deviation [60]. It also depends on the level of confidence we chose. For instance, if the sample size is large, the confidence interval will be narrow. On the other side, the larger confidence level is, the wider the confidence interval. The most common used value for the confidence level is 95% - similar to a 5% level of statistical significance in hypothesis test.

CIs can be one or two-sided. A two-sided CI defines the population parameter from both lower and upper bounds. A one-sided CI provides either an upper or a lower limit to the population parameter. Calculation of the CI of a sample parameter takes the general form of $CI = \text{Point estimate} \pm \text{Margin of error}$, where the margin of error is given by the product of a critical value, selected as per the required confidence limits, and the standard error of point estimate. In descriptive statistics, the CI is reported with the point estimate of the concerned parameter, indicate the reliability of the estimate.

6.2.3 Bayesian Inference Methods.

In recent years, cautionary recommendations against the use of classical statistical approach based on null-hypothesis testing have been frequently issued [171], [90], [35]. Alongside these warnings, there has been also recurrent calls to shift from frequentist methods (null-hypothesis testing included) to Bayesian inference methods. While hypothesis testing provides us with a p -value that indicates the probability of a given observation (i.e., an estimate or the difference between groups) is due to chance, Bayesian approaches provide a relative comparison of how well a null hypothesis and an alternative hypothesis account for the actual data [93]. That is, Bayesian methods allow researchers to make strong observations about the probability of a given phenomenon based on the collected evidence.

Bayesian methods offer multiple advantages such as robustness in low-power situations (e.g., small sample size) and quantification of uncertainty. Moreover, Bayesian methods allow for the inclusion of relevant context or domain information through the choice of informative priors, and they offer a power framework to build and test complex models [93]. It is important to notice however, that Bayesian methods require expertise to be used correctly since they are highly dependent in a good model specification [143]. A general recommendation is to consult with an expert especially when working with continuous data [43]. Furthermore, researchers should be aware that, compared to classical frequentist approaches, Bayesian methods also require larger sample sizes [27].

We refer the reader to [93], [52], [106], and [53] for detailed introduction to Bayesian statistics. A concrete and detailed example on how to use such methods can be found in [143].

6.3 Examples

Sena et al. [151] compare between three conditions using a within-subjects study. They conduct power analysis prior to the data collections that indicated that for a medium effect size (Cohen's $f = 0.3$) and a type I error rate = 0.05, the required sample size is 30. Then, prior to performing a repeated measures ANOVA analysis, they checked and verified that the collected data are normally distributed dependent variable, continuous dependent variable, absence of

outliers and sphericity of the data (Homogeneity of variance). They also have another study that compare between four conditions and they used mixed-factor ANOVA analysis. The power analysis indicated the required sample size is 36 or 9 per condition. However, when they check the ANOVA assumptions, they found that three out of four conditions data are not normally distributed using Anderson-Darling test. The excess Kurtosis for the four groups was 1.37, 1.60, -0.64 and 0.98, where a normal distribution would have an excess kurtosis value of zero. They chose to go with the argument about the robustness of the ANOVA for the violation of the normality assumption especially with a *minor* violation. They also tested for the homogeneity of the variance using Mauchly's test and they got $\chi^2(2) = 1.24, p = 0.54$ which indicate the validity of this assumption and no need for data correction.

On the other hand, Correia et al. [34] used non-parametric methods for analyzing their data after they found that both the normality and the homogeneity assumptions are violated. They used the Shapiro-Wilk test for checking the normality assumption and the Levene's test for checking the homogeneity assumption. They compare between two condition using the Mann-Whitney-U test. This test is similar to t-test with which a comparison between two groups can be conducted. For pairwise comparison for three conditions and more, the Kruskal-Wallis-H test can be used which is analogous to ANOVA test.

Kaplan et al. [89] used both correlation coefficient and Bayesian regression analysis for comparison between eight groups. They first conduct a pairwise correlation analysis to find how much each pair of the eight groups are correlated. Then, a Bayesian regression analysis was conducted to determine the model which best predicts the data. They used the Jeffrey's Awesome Statistical Program (JASP) for Bayesian analysis. The regression model only includes the significant predictors of the dependent variable. Any predictor variable that was not significant at the $p < 0.05$ level was excluded from the final regression model.

7 LIMITATIONS AND FUTURE DIRECTIONS

Discuss limitations and future directions.

ACKNOWLEDGMENTS

TBA

REFERENCES

- [1] 2018. When and how you should reuse research participants. <https://www.simpleusability.com/inspiration/2018/07/when-and-how-you-should-reuse-research-participants/> Library Catalog: www.simpleusability.com.
- [2] Henny Admoni and Brian Scassellati. 2017. Social Eye Gaze in Human-Robot Interaction: A Review. *Journal of Human-Robot Interaction* 6, 1 (Mar 2017), 25. <https://doi.org/10.5898/JHRI.6.1.Admoni>
- [3] Arvin Agah. [n.d.]. Human interactions with intelligent systems: research taxonomy. 27, 1 ([n. d.]), 71–107. [https://doi.org/10.1016/S0045-7906\(00\)00009-4](https://doi.org/10.1016/S0045-7906(00)00009-4)
- [4] Muneeb Imtiaz Ahmad, Omar Mubin, and Joanne Orlando. 2017. Adaptive Social Robot for Sustaining Social Engagement during Long-Term Children-Robot Interaction. *International Journal of Human-Computer Interaction* 33, 12 (Dec 2017), 943–962. <https://doi.org/10.1080/10447318.2017.1300750>
- [5] Mike Allen. 2017. *The SAGE encyclopedia of communication research methods*. Sage Publications.
- [6] Naomi Altman and Martin Krzywinski. 2017. Points of significance: interpreting P values.
- [7] Salvatore M Anzalone, Sofiane Boucenna, Serena Ivaldi, and Mohamed Chetouani. 2015. Evaluating the engagement with social robots. *International Journal of Social Robotics* 7, 4 (2015), 465–478.
- [8] Tamio Arai, Ryu Kato, and Marina Fujita. 2010. Assessment of operator stress induced by robot collaboration in assembly. *CIRP annals* 59, 1 (2010), 5–8.
- [9] Brenna D. Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. 2009. A survey of robot learning from demonstration. *Robotics and Autonomous Systems* 57, 5 (May 2009), 469–483. <https://doi.org/10.1016/j.robot.2008.10.024>

- [10] RA Armstrong and Anthony Hilton. 2004. The use of analysis of variance (ANOVA) in applied microbiology. *Microbiologist* 5, 4 (2004), 18–21.
- [11] Jimmy Baraglia, Maya Cakmak, Yukie Nagai, Rajesh PN Rao, and Minoru Asada. 2017. Efficient human-robot collaboration: When should a robot take initiative? *The International Journal of Robotics Research* 36, 5-7 (2017), 563–579. <https://doi.org/10.1177/0278364916688253> arXiv:<https://doi.org/10.1177/0278364916688253>
- [12] Lisa Feldman Barrett, Ralph Adolphs, Stacy Marsella, Aleix M Martinez, and Seth D Pollak. 2019. Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements. *Psychological science in the public interest* 20, 1 (2019), 1–68.
- [13] Belpaeme T. Eyssel F. Kanda T. Keijsers M. Sabanovic S. Bartneck, C. 2020. *Human-Robot Interaction – An Introduction*. Cambridge: Cambridge University Press.
- [14] Christoph Bartneck, Tony Belpaeme, Friederike Eyssel, Takayuki Kanda, Merel Keijsers, and Selma Šabanović. 2020. *Human-Robot Interaction: An Introduction*. Cambridge University Press.
- [15] Christoph Bartneck, Dana Kulić, Elizabeth Croft, and Susana Zoghbi. 2009. Measurement instruments for the anthropomorphism, animacy, likeability, perceived intelligence, and perceived safety of robots. *International journal of social robotics* 1, 1 (2009), 71–81.
- [16] Paul Baxter, James Kennedy, Emmanuel Senft, Severin Lemaignan, and Tony Belpaeme. 2016. From characterising three years of HRI to methodology and reporting recommendations. In *2016 11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 391–398.
- [17] Jenay M Beer, Arthur D Fisk, and Wendy A Rogers. 2014. Toward a Framework for Levels of Robot Autonomy in Human-Robot Interaction. *Journal of Human-Robot Interaction* 3, 2 (Jun 2014), 74. <https://doi.org/10.5898/JHRI.3.2.Beer>
- [18] Kathleen Belhassein, Guilhem Buisan, Aurélie Clodic, and Rachid Alami. 2019. Towards methodological principles for user studies in Human-Robot Interaction. In *Test Methods and Metrics for Effective HRI in Collaborative Human-Robot Teams Workshop, ACM/IEEE International Conference on Human-Robot Interaction*.
- [19] Elaine M Beller, Val Gebski, and Anthony C Keech. 2002. Randomisation in clinical trials. *Medical Journal of Australia* 177, 10 (2002), 565–567.
- [20] Elaine M Beller, Val Gebski, and Anthony C Keech. 2002. Randomisation in clinical trials. *Medical Journal of Australia* 177, 10 (2002), 565–567.
- [21] Alberto Betella and Paul FMJ Verschure. 2016. The affective slider: A digital self-assessment scale for the measurement of human emotions. *PloS one* 11, 2 (2016).
- [22] Cindy L Bethel, Zachary Henkel, and Kenna Baugus. 2020. Conducting Studies in Human-Robot Interaction. In *Human-Robot Interaction*. Springer, 91–124.
- [23] Cindy L Bethel and Robin R Murphy. 2010. Review of human studies methods in HRI and recommendations. *International Journal of Social Robotics* 2, 4 (2010), 347–359.
- [24] Cindy L Bethel, Kristen Salomon, Robin R Murphy, and Jennifer L Burke. 2007. Survey of psychophysiology measurements applied to human-robot interaction. In *RO-MAN 2007-The 16th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 732–737.
- [25] María J Blanca, Rafael Alarcón, Jaume Arnau, Roser Bono, and Rebecca Bendayan. 2017. Non-normal data: Is ANOVA still a valid option? *Psicothema* 29, 4 (2017), 552–557.
- [26] C. Breazeal, C.D. Kidd, A.L. Thomaz, G. Hoffman, and M. Berlin. 2005. Effects of nonverbal communication on efficiency and robustness in human-robot teamwork. In *2005 IEEE/RSJ International Conference on Intelligent Robots and Systems*. 708–713. <https://doi.org/10.1109/IROS.2005.1545011>
- [27] Marc Brysbaert. 2019. How many participants do we have to include in properly powered experiments? A tutorial of power analysis with reference tables. *Journal of Cognition* 2, 1 (2019).
- [28] Rudolf N Cardinal and Michael RF Aitken. 2013. *ANOVA for the behavioral sciences researcher*. Psychology Press.
- [29] Colleen M Carpinella, Alisa B Wyman, Michael A Perez, and Steven J Stroessner. 2017. The Robotic Social Attributes Scale (RoSAS) Development and Validation. In *Proceedings of the 2017 ACM/IEEE International Conference on human-robot interaction*. 254–262.
- [30] Wesley P Chan, Chris AC Parker, HF Machiel Van der Loos, and Elizabeth A Croft. 2013. A human-inspired object handover controller. *The International Journal of Robotics Research* 32, 8 (Jul 2013), 971–983. <https://doi.org/10.1177/0278364913488806>
- [31] Jesse Chandler, Pam Mueller, and Gabriele Paolacci. 2014. Nonnaïveté among Amazon Mechanical Turk workers: Consequences and solutions for behavioral researchers. *Behavior research methods* 46, 1 (2014), 112–130.
- [32] Jesse Chandler, Gabriele Paolacci, Eyal Peer, Pam Mueller, and Kate Ratliff. 2015. Non-Naïve Participants Can Reduce Effect Sizes. *ACR North American Advances* NA-43 (2015). <https://www.acrwebsite.org/volumes/1020052/volumes/v43/NA-43>
- [33] Emily C. Collins. 2019. Drawing parallels in human-other interactions: a trans-disciplinary approach to developing human-robot interaction methodologies. 374, 1771 (2019), 20180433. <https://doi.org/10.1098/rstb.2018.0433>
- [34] Filipa Correia, Carla Guerra, Samuel Mascarenhas, Francisco S Melo, and Ana Paiva. 2018. Exploring the impact of fault justification in human-robot trust. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*. 507–513.
- [35] Geoff Cumming. 2014. The new statistics: Why and how. *Psychological science* 25, 1 (2014), 7–29.
- [36] Praveen Damacharla, Ahmad Y Javaid, Jennie J Gallimore, and Vijay K Devabhaktuni. 2018. Common metrics to benchmark human-machine teams (HMT): A review. *IEEE Access* 6 (2018), 38637–38655.
- [37] Kerstin Dautenhahn. 2007. Methodology & themes of human-robot interaction: A growing research field. *International Journal of Advanced Robotic Systems* 4, 1 (2007), 15.
- [38] Kerstin Dautenhahn. 2018. Some brief thoughts on the past and future of human-robot interaction.
- [39] Kerstin Dautenhahn and Joe Saunders. 2011. *New frontiers in human robot interaction*. Vol. 2. John Benjamins Publishing.
- [40] Maartje MA de Graaf. 2016. An ethical evaluation of human–robot relationships. *International journal of social robotics* 8, 4 (2016), 589–598.

- [41] Maartje MA de Graaf, Somaya Ben Allouch, and Jan AGM van Dijk. 2016. Long-term evaluation of a social robot in real homes. *Interaction studies* 17, 3 (2016), 462–491.
- [42] A. Dini, C. Murko, S. Yahyanejad, U. Augsdörfer, M. Hofbaur, and L. Paletta. 2017. Measurement and prediction of situation awareness in human-robot interaction based on a framework of probabilistic attention. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. 4354–4361.
- [43] A. Dix. 2020. *Statistics for HCI: Making Sense of Quantitative Data*. Morgan Claypool. <https://ieeexplore-ieee-org.ezproxy.lib.monash.edu.au/document/9070780>
- [44] Gordon S Doig and Fiona Simpson. 2005. Randomization and allocation concealment: a practical guide for researchers. *Journal of Critical Care* 20, 2 (2005), 187–191.
- [45] Birsén Donmez, Patricia E Pina, and Mary L Cummings. 2009. Evaluation criteria for human-automation performance metrics. In *Performance evaluation and benchmarking of intelligent systems*. Springer, 21–40.
- [46] Anca D Dragan, Shira Bauman, Jodi Forlizzi, and Siddhartha S Srinivasa. 2015. Effects of robot motion on human-robot collaboration. In *2015 10th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 51–58.
- [47] Bosede I. Edwards, Idris O. Muniru, Nosiba Khougali, Adrian D. Cheok, and Rui Prada. 2018. A Physically Embodied Robot Teacher (PERT) as a Facilitator for Peer Learning. In *2018 IEEE Frontiers in Education Conference (FIE)*. 1–9. <https://doi.org/10.1109/FIE.2018.8658445>
- [48] Satu Elo and Helvi Kyngäs. 2008. The qualitative content analysis process. *Journal of advanced nursing* 62, 1 (2008), 107–115.
- [49] Mica R Endsley, Stephen J Selcon, Thomas D Hardiman, and Darryl G Croft. 1998. A comparative analysis of SAGAT and SART for evaluations of situation awareness. In *Proceedings of the human factors and ergonomics society annual meeting*, Vol. 42. SAGE Publications Sage CA: Los Angeles, CA, 82–86.
- [50] David M Erceg-Hurn and Vikki M Mirosevich. 2008. Modern robust statistical methods: an easy way to maximize the accuracy and power of your research. *American Psychologist* 63, 7 (2008), 591.
- [51] Ilker Etikan, Sulaiman Abubakar Musa, and Rukayya Sunusi Alkassim. 2016. Comparison of convenience sampling and purposive sampling. *American journal of theoretical and applied statistics* 5, 1 (2016), 1–4.
- [52] Alexander Etz, Quentin F Gronau, Fabian Dablander, Peter A Edelsbrunner, and Beth Baribault. 2018. How to become a Bayesian in eight easy steps: An annotated reading list. *Psychonomic Bulletin & Review* 25, 1 (2018), 219–234.
- [53] Alexander Etz and Joachim Vandekerckhove. 2018. Introduction to Bayesian inference for psychology. *Psychonomic Bulletin & Review* 25, 1 (2018), 5–34.
- [54] Franz Faul, Edgar Erdfelder, Albert-Georg Lang, and Axel Buchner. 2007. G* Power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior research methods* 39, 2 (2007), 175–191.
- [55] Betty J Feir-Walsh and Larry E Toothaker. 1974. An empirical comparison of the ANOVA F-test, normal scores test and Kruskal-Wallis test under violation of assumptions. *Educational and Psychological Measurement* 34, 4 (1974), 789–799.
- [56] Ylva Fernaeus, Maria Håkansson, Mattias Jacobsson, and Sara Ljungblad. 2010. How do you play with a robotic toy animal? A long-term study of Pleo. In *Proceedings of the 9th international Conference on interaction Design and Children*. 39–48.
- [57] Morten Roed Frederiksen and Kasper Stoey. 2019. Augmenting the audio-based expression modality of a non-affective robot. In *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*. 144–149. <https://doi.org/10.1109/ACII.2019.8925510>
- [58] Rudolf Freund, Donna Mohr, and William Wilson. 2010. *Statistical Methods*. Elsevier.
- [59] Craig L Fry, Wayne Hall, Alison Ritter, and Rebecca Jenkinson. 2006. The ethics of paying drug users who participate in research: A review and practical recommendations. *Journal of Empirical Research on Human Research Ethics* 1, 4 (2006), 21–35.
- [60] Martin J Gardner and Douglas G Altman. 1986. Confidence intervals rather than P values: estimation rather than hypothesis testing. *Br Med J (Clin Res Ed)* 292, 6522 (1986), 746–750.
- [61] Christopher H George, S Clare Stanford, Steve Alexander, Giuseppe Cirino, James R Docherty, Mark A Giembycz, Daniel Hoyer, Paul A Insel, Angelo A Izzo, Yong Ji, et al. 2017. Updating the guidelines for data transparency in the British Journal of Pharmacology—data sharing and the use of scatter plots instead of bar charts. *British journal of pharmacology* 174, 17 (2017), 2801.
- [62] Darren Gergle and Desney S Tan. 2014. Experimental research in HCI. In *Ways of Knowing in HCI*. Springer, 191–227.
- [63] Michael J. Gielniak and Andrea L. Thomaz. 2012. Enhancing interaction through exaggerated motion synthesis. In *2012 7th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. 375–382. <https://doi.org/10.1145/2157689.2157813>
- [64] Michael Goodrich and Alan Schultz. 2007. Human-Robot Interaction: A Survey. *Foundations and Trends in Human-Computer Interaction* 1 (01 2007), 203–275. <https://doi.org/10.1561/1100000005>
- [65] A. Green and K.S. Eklundh. 2003. Designing for learnability in human-robot communication. *IEEE Transactions on Industrial Electronics* 50, 4 (Aug 2003), 644–650. <https://doi.org/10.1109/TIE.2003.814763>
- [66] Robert J Grissom. 2000. Heterogeneity of variance in clinical data. *Journal of consulting and clinical psychology* 68, 1 (2000), 155.
- [67] Amir Haddadi, Elizabeth A. Croft, Brian T. Gleeson, Karon MacLean, and Javier Alcazar. 2013. Analysis of task-based gestures in human-robot interaction. In *2013 IEEE International Conference on Robotics and Automation*. 2146–2152. <https://doi.org/10.1109/ICRA.2013.6630865>
- [68] Cindy Harmon-Jones, Brock Bastian, and Eddie Harmon-Jones. 2016. The discrete emotions questionnaire: A new tool for measuring state self-reported emotions. *PLoS one* 11, 8 (2016).

- [69] Michael R Harwell, Elaine N Rubinstein, William S Hayes, and Corley C Olds. 1992. Summarizing Monte Carlo results in methodological research: The one-and two-factor fixed effects ANOVA cases. *Journal of educational statistics* 17, 4 (1992), 315–339.
- [70] David J Hauser and Norbert Schwarz. 2016. Attentive Turkers: MTurk participants perform better on online attention checks than do subject pool participants. *Behavior research methods* 48, 1 (2016), 400–407.
- [71] Roberta Heale and Alison Twycross. 2015. Validity and reliability in quantitative studies. *Evidence-based nursing* 18, 3 (2015), 66–67.
- [72] Marcel Heerink, Ben Krose, Vanessa Evers, and Bob Wielinga. 2009. Measuring acceptance of an assistive social robot: a suggested toolkit. In *RO-MAN 2009-The 18th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 528–533.
- [73] Trevor Higgins. 2003. Randomization in Clinical Trials: Theory and Practice. *Current Medical Research and Opinion* 19, 1 (2003), 67. <http://search.proquest.com/docview/207967010/>
- [74] Anthony Hilton and Richard A Armstrong. 2006. Statnote 6: post-hoc ANOVA tests. *Microbiologist* 2006 (2006), 34–36.
- [75] Guy Hoffman. 2019. Evaluating fluency in human–robot collaboration. *IEEE Transactions on Human-Machine Systems* 49, 3 (2019), 209–218.
- [76] Guy Hoffman and Cynthia Breazeal. 2007. Effects of anticipatory action on human-robot teamwork efficiency, fluency, and perception of team. In *Proceedings of the ACM/IEEE international conference on Human-robot interaction*. 1–8.
- [77] Susan R Homack. 2001. Understanding What ANOVA Post Hoc Tests Are, Really. (2001).
- [78] Asbjørn Hróbjartsson, Julie Pildal, An-Wen Chan, Mette T Haahr, Douglas G Altman, and Peter C Gøtzsche. 2009. Reporting on blinding in trial protocols and corresponding publications was often inadequate but rarely contradictory. *Journal of clinical epidemiology* 62, 9 (2009), 967–973.
- [79] Jason Hsu. 1996. *Multiple comparisons: theory and methods*. CRC Press.
- [80] Michael Hyland. 1981. *The Nature of Hypothetical Constructs*. Macmillan Education UK, London, 42–58.
- [81] Bahar Irfan, James Kennedy, Séverin Lemaignan, Fotios Papadopoulos, Emmanuel Senft, and Tony Belpaeme. 2018. Social psychology and human-robot interaction: An uneasy marriage. In *Companion of the 2018 ACM/IEEE International Conference on Human-Robot Interaction*. 13–20.
- [82] Jiun-Yin Jian, Ann M Bisantz, and Colin G Drury. 2000. Foundations for an empirically determined scale of trust in automated systems. *International journal of cognitive ergonomics* 4, 1 (2000), 53–71.
- [83] Shu Jiang and Ronald C. Arkin. [n.d.]. Mixed-Initiative Human-Robot Interaction: Definition, Taxonomy, and Survey. In *2015 IEEE International Conference on Systems, Man, and Cybernetics* (2015-10). 954–961. <https://doi.org/10.1109/SMC.2015.174>
- [84] Robert Burke Johnson; Larry B. Christensen; Burke Johnson. 2016. *Educational Research: Quantitative, Qualitative, and Mixed Approaches* (hardcover ed.). Sage Publications, Inc.
- [85] James M Johnston and Henry S Pennypacker. 2010. *Strategies and tactics of behavioral research*. Routledge.
- [86] Takayuki Kanda, Masahiro Shiomi, Zenta Miyashita, Hiroshi Ishiguro, and Norihiro Hagita. 2009. An affective guide robot in a shopping mall. In *Proceedings of the 4th ACM/IEEE international conference on Human robot interaction*. 173–180.
- [87] Minsoo Kang, Brian G Ragan, and Jae-Hyeon Park. 2008. Issues in outcomes research: an overview of randomization techniques for clinical trials. *Journal of athletic training* 43, 2 (2008), 215–221.
- [88] PhD. Kang, Minsoo, PhD. A.T.C. Ragan, Brian G, and PhD. Park, Jae-Hyeon. 2008. Issues in Outcomes Research: An Overview of Randomization Techniques for Clinical Trials. *Journal of Athletic Training* 43, 2 (Mar 2008), 215–21. <https://search.proquest.com/docview/206648600?accountid=12528> Copyright - Copyright National Athletic Trainers Association Mar/Apr 2008; Document feature - Illustrations; Tables; Charts; ; Last updated - 2017-11-09.
- [89] Alexandra D Kaplan, Tracy Sanders, and Peter A Hancock. 2019. The Relationship Between Extroversion and the Tendency to Anthropomorphize Robots: A Bayesian Analysis. *Frontiers in Robotics and AI* 5 (2019), 135.
- [90] Maurits Kaptein and Judy Robertson. 2012. Rethinking statistical analysis methods for CHI. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1105–1114.
- [91] Barbara Kitchenham, Lech Madeyski, David Budgen, Jacky Keung, Pearl Brereton, Stuart Charters, Shirley Gibbs, and Amnart Pohthong. 2017. Robust statistical methods for empirical software engineering. *Empirical Software Engineering* 22, 2 (2017), 579–630.
- [92] Joachim Krauth. 2000. *Experimental design: a handbook and dictionary for medical and behavioral research*. Vol. 14. Elsevier.
- [93] John K Kruschke and Torrin M Liddell. 2018. The Bayesian New Statistics: Hypothesis testing, estimation, meta-analysis, and power analysis from a Bayesian perspective. *Psychonomic Bulletin & Review* 25, 1 (2018), 178–206.
- [94] D. Kulic and E. A. Croft. 2007. Affective State Estimation for Human–Robot Interaction. *IEEE Transactions on Robotics* 23, 5 (2007), 991–1000.
- [95] Justin D Lane and David L Gast. 2014. Visual analysis in single case experimental design studies: Brief review and guidelines. *Neuropsychological rehabilitation* 24, 3–4 (2014), 445–463.
- [96] Przemyslaw A. Lasota, Terrence Fong, and Julie A. Shah. [n.d.]. A Survey of Methods for Safe Human-Robot Interaction. 5, 4 ([n. d.]), 261–349. <https://doi.org/10.1561/23000000052> Publisher: Now Publishers, Inc.
- [97] Jonathan Lazar, Jinjuan Heidi Feng, and Harry Hochheiser. 2017. *Research methods in human-computer interaction*. Morgan Kaufmann.
- [98] Kheng Lee Koay, Dag Sverre Syrdal, Michael L. Walters, and Kerstin Dautenhahn. 2007. Living with Robots: Investigating the Habituation Effect in Participants’ Preferences During a Longitudinal Human-Robot Interaction Study. In *RO-MAN 2007 - The 16th IEEE International Symposium on Robot and Human Interactive Communication*. 564–569. <https://doi.org/10.1109/ROMAN.2007.4415149>
- [99] Benedikt Leichtmann and Verena Nitsch. 2020. How much distance do humans keep toward robots? Literature review, meta-analysis, and theoretical considerations on personal space in human-robot interaction. *Journal of Environmental Psychology* 68 (2020), 101386. <https://doi.org/10.1016/j.jenvp.2019.101386>

- [100] Séverin Lemaignan, Alexis Jacq, Deanna Hood, Fernando Garcia, Ana Paiva, and Pierre Dillenbourg. 2016. Learning by teaching a robot: The case of handwriting. *IEEE Robotics & Automation Magazine* 23, 2 (2016), 56–66.
- [101] Trudo Lemmens and Benjamin Freedman. 2000. Ethics review for sale? Conflict of interest and commercial research review boards. *The Milbank Quarterly* 78, 4 (2000), 547–584.
- [102] Frederick TL Leong and James T Austin. 2006. *The psychology research handbook: A guide for graduate students and research assistants*. Sage.
- [103] Todd D. Little, Trent D. Buskirk, Lisa M. Willoughby, and Terry T. Tomazic. 2013. Nonparametric Statistical Techniques. In *The Oxford Handbook of Quantitative Methods in Psychology: Vol. 2: Statistical Analysis*. Oxford University Press, 1–68.
- [104] Mark S Litwin and Arlene Fink. 2003. *How to assess and interpret survey psychometrics*. Vol. 8. Sage.
- [105] Hongyi Liu and Lihui Wang. 2018. Gesture recognition for human-robot collaboration: A review. *International Journal of Industrial Ergonomics* 68 (Nov 2018), 355–367. <https://doi.org/10.1016/j.ergon.2017.02.004>
- [106] Alexander Ly, Akash Raj, Alexander Etz, Maarten Marsman, Quentin F Gronau, and Eric-Jan Wagenmakers. 2018. Bayesian reanalyses from summary statistics: A guide for academic consumers. *Advances in Methods and Practices in Psychological Science* 1, 3 (2018), 367–374.
- [107] John W Mauchly. 1940. Significance test for sphericity of a normal n-variate distribution. *The Annals of Mathematical Statistics* 11, 2 (1940), 204–209.
- [108] Scott E Maxwell, Ken Kelley, and Joseph R Rausch. 2008. Sample size planning for statistical power and accuracy in parameter estimation. *Annual review of psychology* 59 (2008).
- [109] John H McDonald. 2009. *Handbook of biological statistics*. Vol. 2.
- [110] Mary L McHugh. 2012. Interrater reliability: the kappa statistic. *Biochemia medica: Biochemia medica* 22, 3 (2012), 276–282.
- [111] Jane Mills and Melanie Birks. 2014. *Qualitative methodology: A practical guide*. Sage.
- [112] David Moher, Alessandro Liberati, Jennifer Tetzlaff, and Douglas G Altman. 2009. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *Annals of internal medicine* 151, 4 (2009), 264–269.
- [113] Charity J Morgan. 2017. Use of proper statistical techniques for research studies with small samples. *American Journal of Physiology-Lung Cellular and Molecular Physiology* 313, 5 (2017), L873–L877.
- [114] Wendy Moyle, Cindy J Jones, Jenny E Murfield, Lukman Thalib, Elizabeth RA Beattie, David KH Shum, Siobhan T O'Dwyer, M Cindy Mervin, and Brian M Draper. 2017. Use of a robotic seal as a therapeutic tool to improve dementia symptoms: a cluster-randomized controlled trial. *Journal of the American Medical Directors Association* 18, 9 (2017), 766–773.
- [115] Jonathan Mumm and Bilge Mutlu. 2011. Human-robot proxemics: physical and psychological distancing in human-robot interaction. In *Proceedings of the 6th international conference on Human-robot interaction*. 331–338.
- [116] Robin R Murphy and Debra Schreckenghost. 2013. Survey of metrics for human-robot interaction. In *2013 8th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 197–198.
- [117] Jerzy Neyman. 1937. Outline of a theory of statistical estimation based on the classical theory of probability. *Philosophical Transactions of the Royal Society of London. Series A, Mathematical and Physical Sciences* 236, 767 (1937), 333–380.
- [118] Kim F Nimon. 2012. Statistical assumptions of substantive analyses across the general linear model: a mini-review. *Frontiers in psychology* 3 (2012), 322.
- [119] Tatsuya Nomura, Tomohiro Suzuki, Takayuki Kanda, and Kensuke Kato. 2006. Measurement of negative attitudes toward robots. *Interaction Studies* 7, 3 (2006), 437–454.
- [120] Domen Novak, Matjaž Mihelj, and Marko Munih. 2011. Psychophysiological responses to different levels of cognitive and physical workload in haptic interaction. *Robotica* 29, 3 (2011), 367–374.
- [121] Dan R Olsen and Michael A Goodrich. 2003. Metrics for evaluating human-robot interactions. In *Proceedings of PERMIS*, Vol. 2003. 4.
- [122] Chris A. C. Parker and Elizabeth Croft. 2010. J-Strips: Haptic Joint Limit Warnings for Human-Robot Interaction. In *Volume 8: Dynamic Systems and Control, Parts A and B*. ASME, 789–795. <https://doi.org/10.1115/IMECE2010-40717>
- [123] Mildred L Patten and Michelle Newhart. 2017. *Understanding research methods: An overview of the essentials*. Taylor & Francis.
- [124] Elizabeth Phillips, Xuan Zhao, Daniel Ullman, and Bertram F Malle. 2018. What is Human-like? Decomposing Robots' Human-like Appearance Using the Anthropomorphic roBOT (ABOT) Database. In *Proceedings of the 2018 ACM/IEEE International Conference on Human-Robot Interaction*. 105–113.
- [125] Karola Pitsch, Hideaki Kuzuoka, Yuya Suzuki, Luise Sussenbach, Paul Luff, and Christian Heath. 2009. “The first five seconds”: Contingent stepwise entry into an interaction as a means to secure sustained engagement in HRI. In *RO-MAN 2009-The 18th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 985–991.
- [126] Matthew S Prewett, Ryan C Johnson, Kristin N Saboe, Linda R Elliott, and Michael D Coover. 2010. Managing workload in human-robot interaction: A review of empirical studies. *Computers in Human Behavior* 26, 5 (2010), 840–856.
- [127] Matthew S Prewett, Kristin N Saboe, Ryan C Johnson, Michael D Coover, and Linda R Elliott. 2009. Workload in human-robot interaction: a review of manipulations and outcomes. In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, Vol. 53. SAGE Publications Sage CA: Los Angeles, CA, 1393–1397.
- [128] Paul C Price, Rajiv Jhangiani, I-Chant A Chiang, et al. 2015. *Research methods in psychology*. BCCampus.
- [129] Mansour Rahimi and Waldemar Karwowski. [n.d.]. A research paradigm in human-robot interaction. 5, 1 ([n.d.]), 59–71. [https://doi.org/10.1016/0169-8141\(90\)90028-Z](https://doi.org/10.1016/0169-8141(90)90028-Z)

- [130] Paul Ralph and Ewan Tempero. 2018. Construct validity in software engineering research and software metrics. In *Proceedings of the 22nd International Conference on Evaluation and Assessment in Software Engineering 2018*. 13–23.
- [131] Daniel J Rea, Denise Geiskkovitch, and James E Young. 2017. Wizard of awwwws: Exploring psychological impact on the researchers in social HRI experiments. In *Proceedings of the Companion of the 2017 ACM/IEEE International Conference on Human-Robot Interaction*. 21–29.
- [132] Sherri A. Rehfeld. 2006. The Impact Of Mental Transformation Training Across Levels Of Automation On Spatial Awareness In Human-robot Interaction.
- [133] Laurel D Riek. 2012. Wizard of Oz studies in HRI: a systematic review and new reporting guidelines. *Journal of Human-Robot Interaction* 1, 1 (2012), 119–136.
- [134] Nicole L Robinson, Jennifer Connolly, Genevieve M Johnson, Yejee Kim, Leanne Hides, and David J Kavanagh. 2018. Measures of incentives and confidence in using a social robot. *Science Robotics* 3, 21 (2018), Article–number.
- [135] Nicole Lee Robinson, Timothy Vaughan Cottier, and David John Kavanagh. 2019. Psychosocial health interventions by social robots: systematic review of randomized controlled trials. *Journal of medical Internet research* 21, 5 (2019), e13203.
- [136] A. Roncone, O. Mangin, and B. Scassellati. 2017. Transparent role assignment and task allocation in human robot collaboration. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*. 1014–1021.
- [137] Astrid M. Rosenthal-von der Pütten, Nicole C. Krämer, Laura Hoffmann, Sabrina Sobieraj, and Sabrina C. Eimler. 2013. An Experimental Study on Emotional Reactions Towards a Robot. *International Journal of Social Robotics* 5, 1 (Jan 2013), 17–34. <https://doi.org/10.1007/s12369-012-0173-8>
- [138] Craig J Russell and Michelle A Dean. 2000. To log or not to log: Bootstrap as an alternative to the parametric estimation of moderation effects in the presence of skewed dependent variables. *Organizational Research Methods* 3, 2 (2000), 166–185.
- [139] John Rust and Susan Golombok. 2014. *Modern psychometrics: The science of psychological assessment*. Routledge.
- [140] G. Ruxton. 2017. Allocation concealment as a potentially useful aspect of randomised experiments. *Behavioral Ecology and Sociobiology* 71, 2 (2017), 1–3.
- [141] Neil J Salkind. 2010. *Encyclopedia of research design*. Vol. 1. Sage.
- [142] Paul M Salmon, Neville A Stanton, Guy H Walker, Daniel Jenkins, Darshna Ladva, Laura Rafferty, and Mark Young. 2009. Measuring Situation Awareness in complex systems: Comparison of measures study. *International Journal of Industrial Ergonomics* 39, 3 (2009), 490–500.
- [143] Daniel J Schad, Michael Betancourt, and Shravan Vasishth. 2019. Toward a principled Bayesian workflow in cognitive science. *arXiv preprint arXiv:1904.12765* (2019).
- [144] KRISTIN E SCHAEFER. 2013. *THE PERCEPTION AND MEASUREMENT OF HUMAN-ROBOT TRUST*. Ph.D. Dissertation. University of Central Florida Orlando, Florida.
- [145] Stefano Scheggi, Fabio Morbidi, and Domenico Prattichizzo. 2014. Human-Robot Formation Control via Visual and Vibrotactile Haptic Feedback. *IEEE Transactions on Haptics* 7, 4 (Oct 2014), 499–511. <https://doi.org/10.1109/TOH.2014.2332173>
- [146] Jean Scholtz. 2003. Theory and evaluation of human robot interactions. In *36th Annual Hawaii International Conference on System Sciences, 2003. Proceedings of the. IEEE*, 10–pp.
- [147] Mariah L Schrum, Michael Johnson, Muyleng Ghuy, and Matthew C Gombolay. 2020. Four Years in Review: Statistical Practices of Likert Scales in Human-Robot Interaction Studies. *arXiv preprint arXiv:2001.03231* (2020).
- [148] Kenneth F Schulz and David A Grimes. 2002. Allocation concealment in randomised trials: defending against deciphering. *The Lancet* 359, 9306 (2002), 614–618.
- [149] Yasaman Sadat Sefidgar. 2012. *TAMER: touch-guided anxiety management via engagement with a robotic pet efficacy evaluation and the first steps of the interaction design*. Ph.D. Dissertation. University of British Columbia. <https://doi.org/10.14288/1.0052145>
- [150] Johanna Seibt. [n.d.]. Towards an Ontology of Simulated Social Interaction: Varieties of the “As If” for Robots and Humans. In *Sociality and Normativity for Robots: Philosophical Inquiries into Human-Robot Interactions*, Raul Hakli and Johanna Seibt (Eds.). Springer International Publishing, 11–39. https://doi.org/10.1007/978-3-319-53133-5_2
- [151] Aran Sena and Matthew Howard. 2020. Quantifying teaching behavior in robot learning from demonstration. *The International Journal of Robotics Research* 39, 1 (2020), 54–72.
- [152] Mingyang Shao, Matt Snyder, Goldie Nejat, and Beno Benhabib. 2020. User Affect Elicitation with a Socially Emotional Robot. *Robotics* 9, 2 (2020), 44.
- [153] David J Sheskin. 2020. *Handbook of parametric and nonparametric statistical procedures*. crc Press.
- [154] Jaeeun Shim and Ronald C. Arkin. [n.d.]. A Taxonomy of Robot Deception and Its Benefits in HRI. In *2013 IEEE International Conference on Systems, Man, and Cybernetics* (2013-10). 2328–2335. <https://doi.org/10.1109/SMC.2013.398> ISSN: 1062-922X.
- [155] Elaine Short, Justin Hart, Michelle Vu, and Brian Scassellati. 2010. No fair!! an interaction with a cheating robot. In *2010 5th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 219–226.
- [156] Jonathan A Smith and Pnina Shinebourne. 2012. *Interpretative phenomenological analysis*. American Psychological Association.
- [157] Linda J Smith. 2008. How ethical is ethical research? Recruiting marginalized, vulnerable groups into health services research. *Journal of Advanced nursing* 62, 2 (2008), 248–257.
- [158] Aaron Steinfeld, Terrence Fong, David Kaber, Michael Lewis, Jean Scholtz, Alan Schultz, and Michael Goodrich. [n.d.]. Common metrics for human-robot interaction. In *Proceedings of the 1st ACM SIGCHI/SIGART conference on Human-robot interaction* (Salt Lake City, Utah, USA, 2006-03-02) (*HRI '06*). Association for Computing Machinery, 33–40. <https://doi.org/10.1145/1121241.1121249>

- [159] Gail M Sullivan and Anthony R Artino Jr. 2013. Analyzing and interpreting data from Likert-type scales. *Journal of graduate medical education* 5, 4 (2013), 541–542.
- [160] W Newton Suter. 2011. *Introduction to educational research: A critical thinking approach*. SAGE publications.
- [161] David M Swanson and Rebecca A Betensky. 2015. Research participant compensation: a matter of statistical inference as well as ethics. *Contemporary clinical trials* 45 (2015), 265–269.
- [162] Barbara G Tabachnick, Linda S Fidell, and Jodie B Ullman. 2007. *Using multivariate statistics*. Vol. 5. Pearson Boston, MA.
- [163] Leila Takayama and Caroline Pantofaru. 2009. Influences on proxemic behaviors in human-robot interaction. In *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*. 5495–5502. <https://doi.org/10.1109/IROS.2009.5354145>
- [164] Kyle A Thomas and Scott Clifford. 2017. Validity and Mechanical Turk: An assessment of exclusion methods and interactive experiments. *Computers in Human Behavior* 77 (2017), 184–197.
- [165] Andrea Thomaz, Guy Hoffman, and Maya Cakmak. [n.d.]. Computational Human-Robot Interaction. 4, 2 ([n. d.]), 105–223. <https://doi.org/10.1561/23000000049>
- [166] Bruce Thompson. 1994. Common Methodology Mistakes in Dissertations, Revisited. (1994).
- [167] Sebastian Thrun. 2004. Toward a framework for human-robot interaction. *Human-Computer Interaction* 19, 1-2 (2004), 9–24.
- [168] Lorenza Tiberio, Amedeo Cesta, and Marta Olivetti Belardinelli. 2013. Psychophysiological methods to evaluate user’s response in human robot interaction: a review and feasibility study. *Robotics* 2, 2 (2013), 92–121.
- [169] Daniel C. Tozadore, Joao P.H. Valentini, Victor H.S. Rodrigues, Fernando M.L. Vendrameto, Rodrigo G. Zavarizz, and Roseli A.F. Romero. 2018. Towards Adaptation and Personalization in Task Based on Human-Robot Interaction. In *2018 Latin American Robotic Symposium, 2018 Brazilian Symposium on Robotics (SBR) and 2018 Workshop on Robotics in Education (WRE)*. 383–389. <https://doi.org/10.1109/LARS/SBR/WRE.2018.00075>
- [170] Jan Van den Broeck, Solveig Argeseanu Cunningham, Roger Eeckels, and Kobus Herbst. 2005. Data cleaning: detecting, diagnosing, and editing data abnormalities. *PLoS Med* 2, 10 (2005), e267.
- [171] Joachim Vandekerckhove, Jeffrey N Rouder, and John K Kruschke. 2018. Bayesian methods for advancing psychological science.
- [172] Renato Paredes Venero and Alex Davila. 2020. Experimental Research Methodology and Statistics Insights. In *Human-Robot Interaction*. Springer, 333–353.
- [173] Lin Wang, Pei-Luen Patrick Rau, Vanessa Evers, Benjamin Krisper Robinson, and Pamela Hinds. 2010. When in Rome: the role of culture & context in adherence to robot recommendations. In *2010 5th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 359–366.
- [174] Astrid Weiss, Regina Bernhaupt, Michael Lankes, and Manfred Tscheligi. 2009. The USUS evaluation framework for human-robot interaction. In *AISB2009: proceedings of the symposium on new frontiers in human-robot interaction*, Vol. 4. 11–26.
- [175] Tracey L Weissgerber, Natasa M Milic, Stacey J Winham, and Vesna D Garovic. 2015. Beyond bar and line graphs: time for a new data presentation paradigm. *PLoS biology* 13, 4 (2015).
- [176] Tracey L Weissgerber, Marko Savic, Stacey J Winham, Dejana Stanisavljevic, Vesna D Garovic, and Natasa M Milic. 2017. Data visualization, bar naked: A free tool for creating interactive graphics. *Journal of Biological Chemistry* 292, 50 (2017), 20592–20598.
- [177] Rand R Wilcox and Guillaume A Rousselet. 2018. A guide to robust statistical methods in neuroscience. *Current protocols in neuroscience* 82, 1 (2018), 8–42.
- [178] Lesley Wood, Matthias Egger, Lise Lotte Gluud, Kenneth F Schulz, Peter Jüni, Douglas G Altman, Christian Gluud, Richard M Martin, Anthony J G Wood, and Jonathan A C Sterne. 2008. Empirical evidence of bias in treatment effect estimates in controlled trials with different interventions and outcomes: meta-epidemiological study. *BMJ* 336, 7644 (2008), 601. <http://bmj.com/content/336/7644/601.full.pdf>
- [179] Yuki Yamashita, Hisashi Ishihara, Takashi Ikeda, and Minoru Asada. 2016. Path analysis for the halo effect of touch sensations of robots on their personality impressions. In *International Conference on Social Robotics*. Springer, 502–512.
- [180] H.A. Yanco and J. Drury. [n.d.]. Classifying human-robot interaction: an updated taxonomy. In *2004 IEEE International Conference on Systems, Man and Cybernetics (IEEE Cat. No.04CH37583)* (2004-10), Vol. 3. 2841–2846 vol.3. <https://doi.org/10.1109/ICSMC.2004.1400763> ISSN: 1062-922X.
- [181] Holly A Yanco, Adam Norton, Willard Ober, David Shane, Anna Skinner, and Jack Vice. 2015. Analysis of Human-robot Interaction at the DARPA Robotics Challenge Trials. *Journal of Field Robotics* (2015).
- [182] Deng Yongda, Li Fang, and Xin Huang. 2018. Research on multimodal human-robot interaction based on speech and gesture. *Computers Electrical Engineering* 72 (Nov 2018), 443–454. <https://doi.org/10.1016/j.compeleceng.2018.09.014>
- [183] James E. Young, JaYoung Sung, Amy Volda, Ehud Sharlin, Takeo Igarashi, Henrik I. Christensen, and Rebecca E. Grinter. [n.d.]. Evaluating Human-Robot Interaction. 3, 1 ([n. d.]), 53–67. <https://doi.org/10.1007/s12369-010-0081-8>
- [184] Oded Zafrani and Galit Nimrod. 2019. Towards a Holistic Approach to Studying Human-Robot Interaction in Later Life. *The Gerontologist* 59, 1 (Jan. 2019), e26–e36. <https://doi.org/10.1093/geront/gny077> Publisher: Oxford Academic.

Received 2020; revised 2020; accepted 2020