

Use of Large Sample Sizes and Multiple Evaluation Methods in Human-Robot Interaction Experimentation

Cindy L. Bethel

IEEE Robotics and Automation Society Graduate Fellow

NSF Graduate Research Fellow

University of South Florida



Robotics & Automation Society



Overview

1. Introduction and Background

2. Two Issues Common in HRI Studies

3. Study Design and Planning

4. Methods of Evaluation

5. Conclusions and Recommendations

Introduction

Human-Robot Interaction is a rapidly advancing area of research. It requires strong experimental designs and methods of evaluation, in order to establish credibility and validity in HRI research.



Background

- The information presented is based on a recent HRI study conducted using:
 - 128 participants (79 females, 49 males, ages 18-62)
 - 2 different search and rescue robots
 - 2 modes of operation (standard or emotive)
 - 4 methods of evaluation
 - A high-fidelity, confined space simulated disaster site with all experiments conducted in the dark





Two Issues Common in HRI Studies

Two Issues in HRI Studies

- Lack of significant sized participant pools that accurately represent the populations being studied.
- Lack of multiple methods of evaluation to obtain convergent validity in HRI studies.



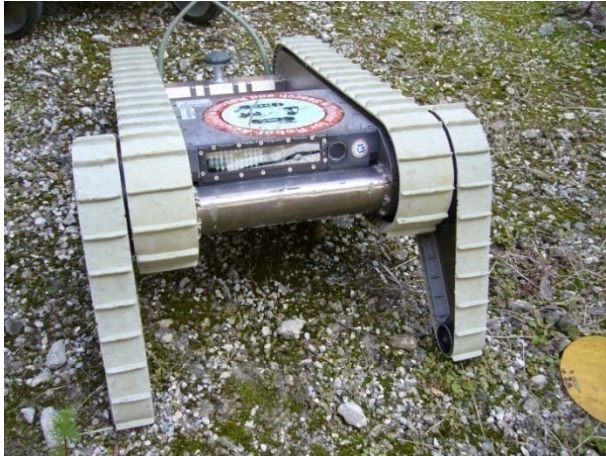
Study Design and Planning

Study Design

- ***Within Subjects*** – All participants experience all factors being manipulated. Most comprehensive, increases statistical power, reduces error variance; however participants can experience a practice and/or carryover effects.
- ***Between Subjects*** – Participants are placed into only one group and experience only one set of factors being manipulated.
- ***Mixed-Model*** - some factors or variables are set up as within-subjects and some are set up as between-subjects.

Participants interacted with:

Two Different USAR Robots



iRobot Packbot



Inuktun Extreme

... and One of Two Modes

Standard
Mode

OR

Emotive
Mode

Determining the Appropriate Sample Size Using Power Analysis

- Determine the number of groups needed for your study.
- The alpha level or significance level that you want to use (commonly $\alpha = .05$).
- The expected effect size (Cohen's levels < .1 trivial, .1 - .3 small, .3-.5 medium, > .5 large).
- The amount of statistical power (commonly .80).

Once you have this information then you can determine the appropriate number of participants needed from tables in the back of many statistics books or use software to calculate this value.

(e.g., G*Power 3 software located at <http://www.psychu.uniduesseldorf.de/abteilungen/aap/gpower3/>)

Exemplar Study – Power Analysis

- Two groups needed
- Alpha = .05
- Medium effect size (.35)
- Statistical power of .80

The results indicated a need for 64 participants per group for a total of 128 participants.

Use of Appropriate Sample Sizes

By performing an *a priori* power analysis to estimate the appropriate number of participants required for the study, it will provide a better opportunity to obtain statistically significant results that are valid, reliable, and accurate.



Methods of Evaluation

Methods of Evaluation



- Self-Assessments



- Interviews



- Observational or Behavioral Measures



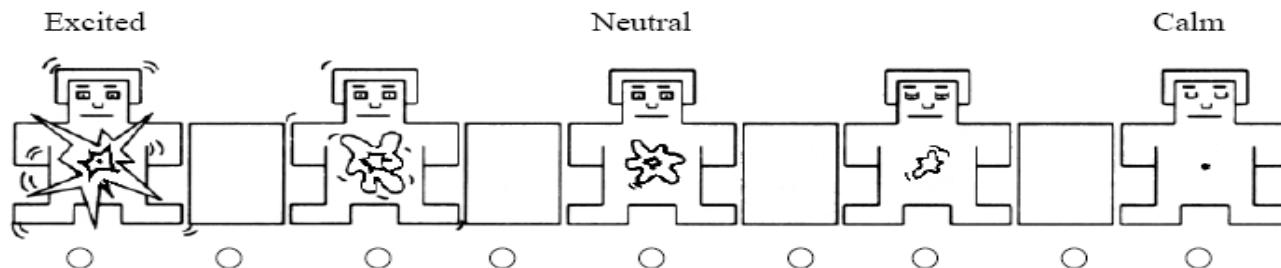
- Psychophysiology Measurements



- Performance Metrics

Self-Assessments

- Participants provide a personal assessment of how they felt or their motivations related to an object, situation, or interactions.
- **Problems:**
 - Validity and corroboration.
 - Participants may not be in touch with what they are feeling.
 - May not answer honestly but how they feel the researcher wants the questions answered.
 - Responses could be influenced by mood and emotions of participants on the day of the study.



Interviews

- This method is essentially a conversation with participants that can be a structured or freeform measurement of participants' reactions to an object, situation, or interaction.
- ***Problems:***
 - Can suffer from similar issues as the self-assessments.
 - Can be difficult to analyze with accurate and valid results.

Observational Measures

- Defined as “the watching of behavioral patterns of people in certain situations to obtain information about the phenomenon of interest.” (Johnson and Christensen 2004)
- ***Advantage:***
 - Researchers are able to record the actual behaviors of participants and do not have to rely on participants to report accurately their intended behaviors or preferences.
- ***Disadvantage:***
 - The “Hawthorne Effect” is a concern with observational studies. It is a phenomenon in which participants know that they are being observed, and it impacts their behaviors.

Psychophysiology Measurements

- ***Advantages:***

- Participants cannot consciously manipulate the activities of their autonomic nervous system.
- Psychophysiological measures offer a non-invasive method that can be used to determine the stress levels and reactions of participants interacting with technology

- ***Disadvantages:***

- Psychophysiological measurements results are not always straightforward and confounds can lead to misinterpretation of data.
- There is a tendency to attribute more meaning to results because of the tangible nature of the recordings.

- ***Recommendations:***

- Information needs to be obtained from participants prior to beginning a study to help reduce these confounds (e.g., health information, medications, food intake, state of mind, etc.).
- Multiple physiological signals should be used in order to find correlations in the results.

Task Performance Metrics

- The use of task performance metrics is evolving and becoming more common in HRI studies, especially where teams are being evaluated and/or more than one person is interacting with one or more robots.
- These metrics are designed to measure how well a person or team performs or completes tasks.

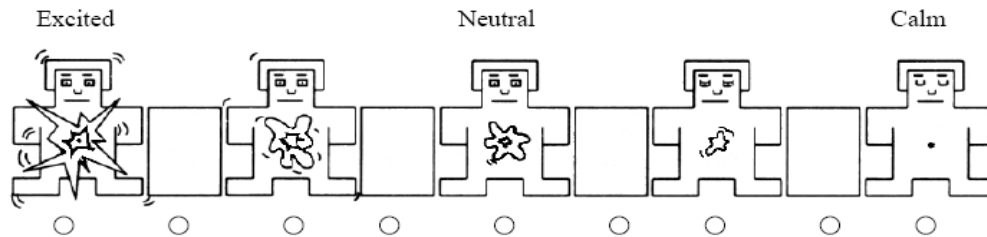
Convergent Validity

- No one method of evaluation is sufficient in and of itself to evaluate any interaction; therefore it is important to include multiple methods of evaluation in a comprehensive study to gain a better understanding of Human-Robot Interaction and obtain convergent validity.
- Within a single method of evaluation there should be multiple measures utilized.
- This will help to ensure a comprehensive study with reliable and accurate results which can be validated.

Exemplar Study

Self-Assessments and Structured Interview

How calm/aroused did you feel during your interaction with the robot?
(mark the circle under the drawing that most closely reflects your level of expected arousal)

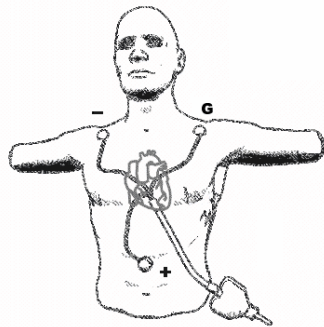


Video Observation

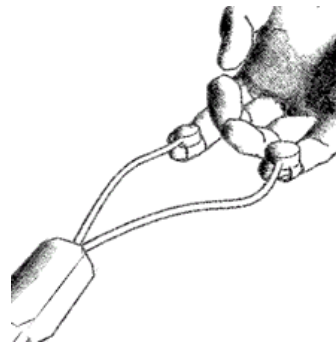


Psychophysiology Measurements

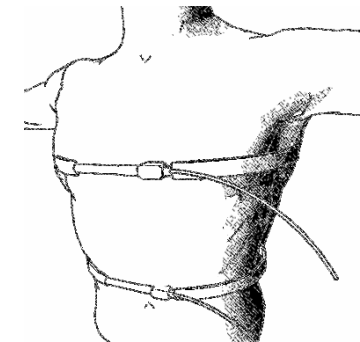
EKG: Heart Rate & IBI



Skin Conductance



**Respiration:
Thoracic and Abdominal**





Conclusions and Recommendations

Conclusions

There are two improvements that should to be made in human studies conducted in HRI and those are:

- Having larger sample sizes to appropriately represent the population being studied and so that small to medium effects can be discovered with statistically significant results.
- The use of multiple methods of evaluation to establish reliable and accurate results that will have convergent validity.

Recommendations

- Determine the most appropriate type of study for the hypotheses being investigated using either a within-subjects, between-subjects, or mixed-model design.

Recommendations

- Perform an *a priori* power analysis to estimate the appropriate number of participants required for the study in order to have a better opportunity to obtain statistically significant results that are valid, reliable, and accurate. This can be accomplished through power analysis tables or freely available software. An *a priori* power analysis is based on:
 - the number of groups in the study
 - the effect size
 - the alpha level
 - the desired statistical power
- It is recommended to add a few more participants to the estimated number to account for problems with data, participant cancellations that cannot be rescheduled, and participants that do not show up.

Recommendations

- Determine the best methods of evaluation for the hypotheses being investigated, but it is recommended that at least two or more methods should be utilized in order to obtain convergent validity in the study.
- The results from a study with multiple methods of evaluation are viewed as more reliable and accurate.
- Additionally, by incorporating multiple methods of evaluation it will overcome the inherent problems found with every method of evaluation. No one method is sufficient for accurately measuring participants' responses.

Additional Recommendations

Please refer to the paper included in the symposium handout for a more in-depth coverage of HRI study design, planning and execution recommendations developed as a result of the exemplar study conducted.

These include:

- Study protocol development
- Study location and environment
- Robots and other equipment
- Conducting the study
- Recruiting participants and volunteer research assistants
- Failures and contingencies
- Additional recommendations for conducting HRI studies

Acknowledgments

This work is supported in part under:



- A National Science Foundation Graduate Research Fellowship – Award Number DGE – 0135733



- ARL Grant Number W911NF-06-2-0041



- The Microsoft Rescue Buddy Research Project



- IEEE Robotics and Automation Society Fellowship.

Questions???



Cindy L. Bethel
cbethel@gmail.com
www.cindybethel.com