

Trust in Autonomous Human–Robot Interaction

An In-Person Pilot Study

M.C. Lau^{a,1,*}, Shauna Heron^{a,2}

^a*Laurentian University, Bharti School of Engineering,*

^b*Laurentian University, School of Social Sciences,*

Abstract

This study implements a multi-stage collaborative task system where participants collaborate with the Misty-II social robot to solve a who-dunnit type task. The system utilizes an autonomous, mixed-initiative dialogue architecture with affect-responsive capabilities.

Keywords: keyword1, keyword2

! TODO

Manually score each dialogue series.

For each interaction and stage:

- did the participant ask for help?
- how many times?
- did the robot give useful help?
- did the robot give misleading or incorrect help?
- did the robot stick to the policy?
- how many times did the robot fail to understand the participant?

For each task:

- is there evidence that the robot helped complete the task?
- is there evidence that the participant solved the problem without help?

Human–robot collaboration (HRC) has become a central topic across engi-

*Corresponding author

Email addresses: mclau@laurentian.ca (M.C. Lau), sheron@laurentian.ca (Shauna Heron)

¹This is the first author footnote.

²Another author footnote, this is a very long footnote and it should be a really long footnote. But this footnote is not yet sufficiently long enough to make two lines of footnote text.

neering, computer science, and the social sciences as robots increasingly move from controlled laboratory settings into everyday collaborative roles. In many emerging applications, collaboration depends not only on physical coordination but also on shared problem-solving through dialogue, where robots must reason, communicate, and adapt in real time. Understanding how humans perceive and respond to such systems is therefore critical for designing robots that can function as effective collaborators rather than passive tools.

A key factor shaping successful collaboration in human–robot interaction (HRI) is trust. Trust influences whether users are willing to rely on robotic systems, accept their guidance, and remain engaged during joint tasks, particularly in situations characterized by uncertainty or incomplete information. Prior work has shown that trust affects both subjective perceptions—such as perceived reliability or intent—and objective outcomes including task performance, compliance, and cooperation. As a result, a substantial body of research has focused on measuring trust in HRI, leading to the development of standardized instruments for assessing users’ evaluations of robot behaviour across industrial, medical, and social contexts.

Despite this progress, much of the existing literature on trust in HRI is based on interactions conducted under highly controlled or simulated conditions. In many studies, robot behaviour is scripted, partially simulated, or mediated through Wizard-of-Oz paradigms, where a human operator covertly controls aspects of the robot’s behaviour. While these approaches are valuable for isolating specific design factors and testing early hypotheses, they also mask many of the failures and inconsistencies that characterize autonomous systems in real-world use. Speech recognition errors, delayed or inappropriate responses, misinterpretations of user intent, and limitations of affect sensing are not peripheral issues but central features of deployed autonomous robots. These imperfections are likely to play a decisive role in shaping trust, yet they remain underexplored in empirical HRI research.

Introduce the concept of responsiveness in robot/AI systems as a moderator of trust. Prior work has suggested that robots that can adapt their behaviour based on user affect and contextual cues may foster greater trust and engagement. However, most studies examining responsiveness have relied on simulated or semi-autonomous systems, leaving open questions about how these effects manifest in fully autonomous, in-person interactions.

The present pilot study addresses this gap by examining trust and collaboration in an in-person interaction with two versions of a fully autonomous social robot operating within predefined behavioural constraints. Using a between-subjects design, participants collaborated with either a responsive or neutral robot during an immersive, dialogue-driven puzzle game in which the robot acted as a diegetic game guide and partner. The task required shared problem-solving through conversation, with participants seeking hints, advice, and support from the robot while navigating a game environment displayed on a computer screen. Crucially, all interaction management—including speech-based dialogue, task progression,

and affect-responsive behaviour was handled autonomously by the robot without human intervention.

In the experimental condition, the robot was designed to be proactive and responsive, adapting its behaviour based on participant affect—as estimated from outputs of an affect detection model—as well as conversational cues. In the other condition, the robot provided assistance only when explicitly requested, offering a more reactive interaction style. This manipulation allowed us to examine how differences in autonomy and responsiveness to human states influence trust perceptions and collaborative performance under otherwise identical task demands.

To support this interaction, we developed an autonomous spoken-language system integrated with automatic speech recognition (ASR) and affect detection on the Misty-II robot platform. The system we developed enables the robot to recognize speech, manage dialogue state, maintain conversational context, and generate coordinated verbal responses alongside LED facial expressions and head and arm movements. Rather than optimizing for flawless performance, the system was designed to reflect realistic capabilities and limitations of contemporary social robots.

By combining post-interaction trust measures with behavioural and task-level outcomes, this study aims to contribute empirical evidence on how trust is shaped in fully autonomous HRI scenarios. The focus is not on demonstrating idealized interaction under perfect conditions, but on examining trust as it emerges through realistic human–robot collaboration, where uncertainty, interactional breakdowns, and adaptive behaviour are unavoidable. In doing so, this work seeks to inform the design and evaluation of affect-responsive autonomous robots intended for real-world collaborative settings.

[also mention that this is a pilot study to inform a larger planned study with more participants and refined methods based on lessons learned here and touch on some of the lessons we learned (i.e., language issues, ASR issues, interaction design issues, etc.)]

0.1. Background

0.1.1. Trust in HRI research

Introduce the two scales here we used to measure trust in HRI: the Trust Perception Scale-HRI and the Trust in Industrial Human–Robot Collaboration scale. Discuss prior work that has used these scales and their relevance to our study.

The differences between these two scales in terms of what aspects of trust they measure (e.g., affective vs. cognitive trust, reliability vs. collaboration). In addition, discuss how these scales have been validated in prior research and their psychometric properties.

Explain why we chose to use both scales in our study to capture a comprehensive picture of trust in HRI. The reasons for selecting these scales should be linked

to our research questions and hypotheses about how robot responsiveness and affective adaptation might influence different dimensions of trust.

0.1.2. Responsiveness and affective adaptation in HRI

Discuss prior research on the role of robot responsiveness and affective adaptation in shaping trust and engagement in HRI. Summarize key findings from studies that have examined how robots that can perceive and respond to human affective states influence user perceptions and behaviours. Highlight any gaps in the literature, particularly regarding fully autonomous, in-person interactions, which our study aims to address.

Discuss theoretical frameworks that explain why responsiveness and affective adaptation might enhance trust, such as social presence theory or the CASA paradigm. Explain how these frameworks inform our hypotheses about the expected effects of robot interaction policy on trust and collaboration.

Discuss the reason for conducting the study in a fully autonomous, in-person setting rather than using Wizard-of-Oz or simulated paradigms. Emphasize the importance of examining trust in realistic conditions where interaction imperfections are present.

0.1.3. Why a pilot study?

Why did we run a pilot study? what future work will this inform? Focus on testing feasibility of the autonomous system, interaction design, and measurement approach.

Discuss pilot studies or preliminary work that informed the design of our robot interaction policies. This could include prior experiments with semi-autonomous systems, user feedback on robot behaviours, or technical evaluations of affect detection models.

0.2. Hypotheses

The primary objective of this study was to examine how differences in robot interaction policy influence trust and collaboration during fully autonomous, in-person human–robot interaction. Based on prior work linking robot responsiveness, affective behavior, and trust in HRI, we formulated the following hypotheses.

H1: Participants interacting with a responsive, affect-adaptive robot will report higher post-interaction trust than participants interacting with a neutral, reactive robot.

H2: Participants in the responsive condition will demonstrate greater engagement with the robot during the collaborative tasks, reflected in increased voluntary interaction and reliance on robot input during problem solving.

Fix H3—we didn't actually test this directly H3: Differences in trust and engagement will be most pronounced during the open-ended collaborative task, where assistance from the robot is optional rather than required.

0.3. Methods

0.3.1. Sample and recruitment

Participants ($n = 29$) were recruited from the Laurentian University community through word of mouth and the SONA participant recruitment system. Eligibility criteria required participants to be adults (18 years or older), fluent in written and spoken English, with normal or corrected-to-normal hearing and vision. Participants received a \$15.00 gift card as compensation for their time. All procedures were approved by the university's Research Ethics Board. The Misty-II robot used in this study was purchased through grant funding from the IAMGOLD President's Innovation Fund. Sample characteristics are summarized in Table 1.

0.3.2. Experimental design

The study employed a between-subjects design with robot interaction policy as the sole experimental factor. Participants interacted with the same Misty-II robot in a shared physical workspace that included both the robot and a participant-facing computer interface. The interface was used to present brief task instructions, collect participant inputs, and manage transitions between task stages. Critically, the interface did not serve as a control mechanism for the robot. Instead, the robot autonomously monitored task state and participant inputs via the interface and input from the participant. Dialogue and behavior were adapted accordingly, without any real-time human intervention (see Figure 1).

Participants collaborated with the robot during an immersive puzzle game in which the robot functioned as a diegetic game guide and collaborative partner. The interaction was fully autonomous in both conditions, and both versions of the robot were subject to the same sensory and interaction constraints inherent to real-world operation, including speech recognition variability and response timing delays. The only manipulation between conditions was the robot's interaction policy.

Participants were randomly assigned to one of two conditions:

RESPONSIVE (experimental): The robot adopted a warm, emotionally engaged, and proactive interaction style, adapting its responses based on detected participant affect, dialogue context, and task demands.

CONTROL (baseline): The robot followed a neutral, reactive interaction policy, providing information and assistance only when explicitly requested, without affect-responsive adaptation.

0.3.3. Task structure

The game consisted of five sequential stages designed to elicit interaction under differing collaboration and dependency conditions, following established approaches in HRI task design (Lin et al., 2022). Total session duration was approximately 15 minutes.

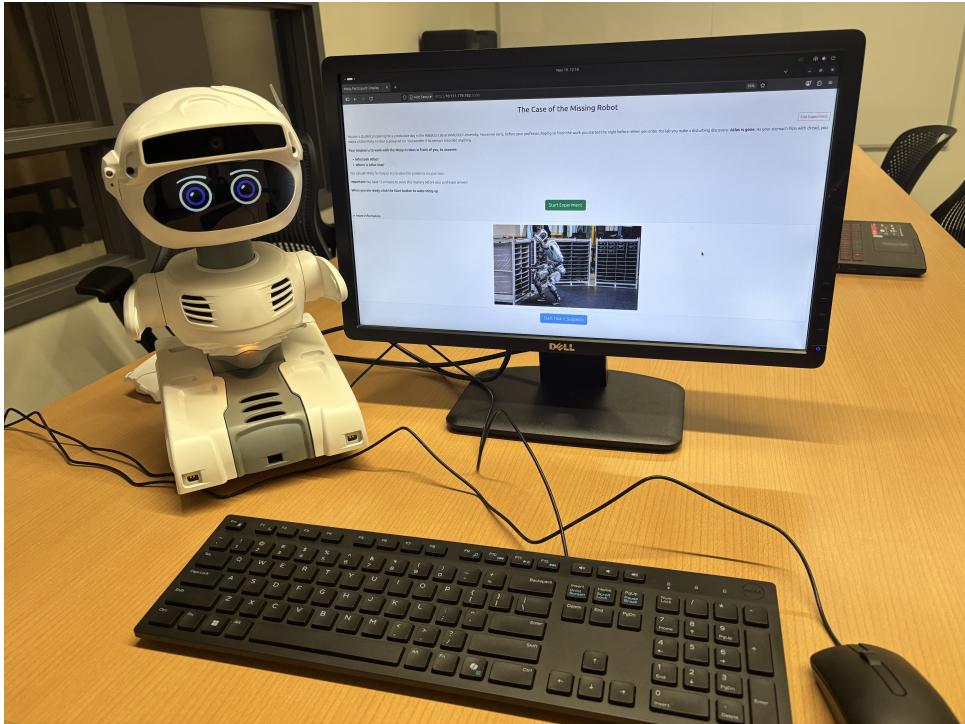


Figure 1: Experimental setup showing the autonomous robot and participant-facing task interface used during in-person sessions. Participants entered task responses and navigated between task stages using the interface, while the robot autonomously tracked task state and adapted its interaction based on participant input. No real-time human intervention occurred during the interaction.

Stage 1: Greeting. The robot introduced itself and engaged in brief rapport-building interaction.

Stage 2: Mission brief. The robot explained the narrative context and overall objectives of the task.

Stage 3: Task 1 (robot-dependent reasoning). Participants completed a constrained “who-dunnit” task.

Stage 4: Task 2 (open-ended collaborative problem solving). Participants worked to determine the location of the missing robot using technical logs.

Stage 5: Wrap-up. The robot provided feedback and concluded the interaction.

0.3.4. Task 1: Robot-dependent collaborative reasoning

In the first task, participants were required to identify a suspect from a 6×4 grid of 24 candidates by asking the robot a series of yes/no questions about the suspect’s features (e.g., hair color, accessories, clothing). The grid was displayed on the interface, while questions were posed verbally to the robot. The robot possessed the ground-truth information necessary to evaluate each question and provide correct responses.

Successful completion of this task was therefore dependent on interaction with the robot, creating a forced collaborative dynamic in which the robot served as an essential informational partner. Participants were required to coordinate questioning strategies with the robot to narrow down the correct suspect within a five-minute time limit. The structured nature of the task ensured consistent interaction demands across participants and conditions.

0.3.5. Task 2: Open-ended problem solving with advisory robot support

The second task involved a more open-ended problem-solving scenario. Participants were presented with multiple technical logs through a simulated terminal interface that were used to determine the location of the missing robot. Unlike Task 1, the robot *did not* have access to ground-truth information or the contents of the logs. The robot’s assistance was limited to general problem-solving support derived from its language model, such as explaining how to interpret logs, suggesting reasoning strategies, or prompting participants to reflect on inconsistencies.

Participants could complete this task independently or choose to solicit assistance from the robot. The robot could ask clarifying questions about what the participant observed in the logs, and participants could likewise ask the robot for guidance. This design positioned the robot as a collaborative reasoning partner rather than an authoritative source and allowed collaboration to emerge voluntarily rather than being enforced by task structure (Lin et al., 2022).

0.3.6. Wrap-up and debrief

After all responses were submitted, correct answers were displayed to participants. During the wrap-up stage, the robot engaged in a brief debriefing

interaction, acknowledging task outcomes and thanking participants for their involvement before prompting them to report back to the researcher.

0.3.7. In-person procedure

Participants completed a pre-interaction eligibility questionnaire administered via Qualtrics prior to their in-person session. This questionnaire included eligibility questions, informed consent, demographic information, the Negative Attitudes Toward Robots Scale, and a measure of Need for Cognition. Balanced random assignment was also completed in this step (need mention no-shows that threw off the balance of the group assignments). Due to variability in timing between pre-interaction questionnaires and in-person sessions, these measures were treated as baseline covariates rather than formal pre-test measures.

At the start of the in-person session, participants were seated in front of the Misty-II robot and instructed to begin the interaction by clicking a start button on the interface. They were given brief guidance on effective communication with the robot, including waiting for a visual indicator on the robot before speaking. Once participants indicated readiness, the researcher left the room and closed the door, leaving the participant and robot to complete the tasks without human presence.

Following task completion, participants exited the room and completed a post-interaction survey assessing trust using the Trust Perception–HRI scale and the Trust in Industrial Human–Robot Collaboration scale. Participants then engaged in a written and verbal debrief with the researcher. Participants were informed that they could terminate the session at any time without penalty. All participants completed the full procedure, with total session duration averaging approximately 30 minutes, and received compensation upon completion.

1. Results

1.0.1. Data exclusions and communication viability

Although participants were required to be fluent in spoken English, in-person observation on meeting participants as well as post-hoc review of interaction transcripts revealed that a small subset of participants experienced persistent communication breakdowns caused by language barriers that prevented meaningful engagement with the robot. These breakdowns were characterized by repeated speech recognition failures, fragmented responses, and task abandonment.

Because the experimental manipulation relied on language-mediated collaboration, the most extreme sessions were considered non-compliant with the intended protocol and were excluded from task-level analyses. Importantly, exclusion criteria were based on interaction viability rather than outcome measures, and results were qualitatively unchanged when these sessions were retained.

Sample eligibility and protocol adherence. Participants were required to be fluent in spoken English. During administration, the experimenter recorded

Table 1

Characteristic	N	CONTROL N = 9¹	RESPONSIVE N = 14¹	p-value²
Gender	22			>0.99
Woman		4 / 9 (44%)	7 / 13 (54%)	
Man		5 / 9 (56%)	6 / 13 (46%)	
Age Group	22			0.16
18-24		4 / 9 (44%)	7 / 13 (54%)	
25-34		2 / 9 (22%)	1 / 13 (7.7%)	
34-44		0 / 9 (0%)	4 / 13 (31%)	
45+		3 / 9 (33%)	1 / 13 (7.7%)	
Program	20			0.95
Psychology		1 / 9 (11%)	1 / 11 (9.1%)	
Engineering		2 / 9 (22%)	1 / 11 (9.1%)	
Computer Science		3 / 9 (33%)	5 / 11 (45%)	
Earth Sciences		0 / 9 (0%)	1 / 11 (9.1%)	
Other		3 / 9 (33%)	3 / 11 (27%)	
Experience w/Robots	23	5 / 9 (56%)	3 / 14 (21%)	0.18
Native English Speaker	23			>0.99
Native English		5 / 9 (56%)	8 / 14 (57%)	
Non-Native English		4 / 9 (44%)	6 / 14 (43%)	
NARS Overall	23	37 (10)	37 (7)	0.78
Need for Cognition	23	3.81 (0.76)	3.73 (0.79)	>0.99

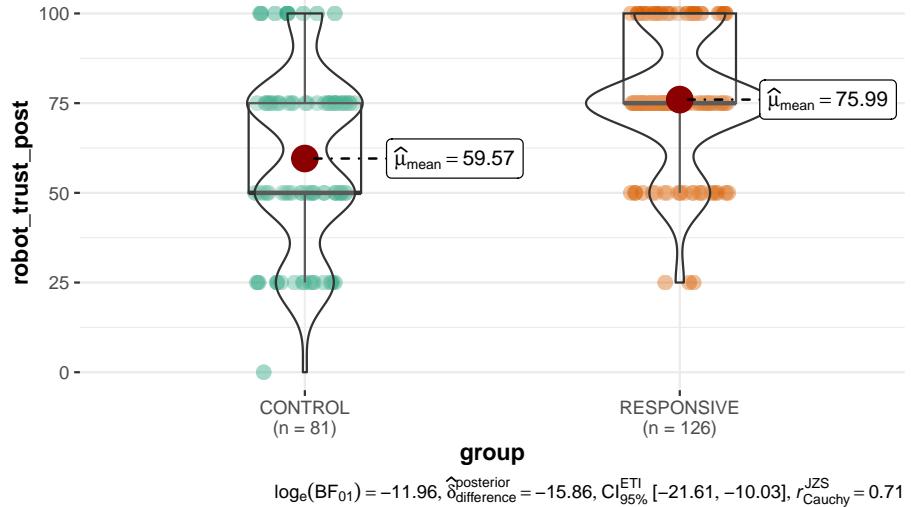
¹n / N (%); Mean (SD)²Fisher's exact test; Wilcoxon rank sum test

contemporaneous notes when conversational fluency appeared insufficient to support the language-mediated protocol. Communication viability was later assessed using objective interaction indicators extracted from system logs and manually coded transcripts (e.g., repeated speech-recognition failures, fragmented utterances, and task abandonment). Sessions showing sustained communication breakdown were treated as protocol non-adherence and were excluded from task-level analyses (n=6). All analyses are also reported with these sessions retained as a sensitivity check.

1.1. Participant characteristics and baseline measures

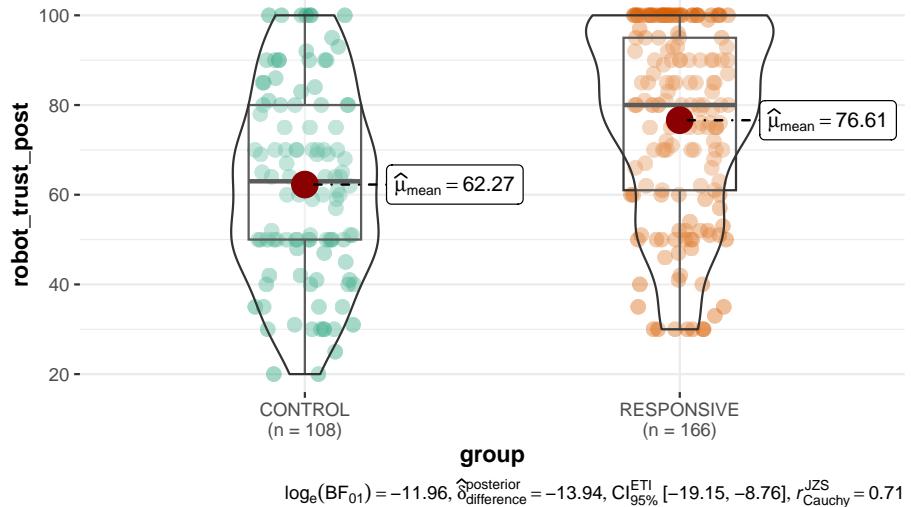
Participants in the control and responsive conditions were comparable with respect to pre-interaction demographic characteristics, academic background, prior experience with robots, and baseline attitudes toward robots. Importantly, Negative Attitudes Towards Robots (NARS) and Need for Cognition scores were similar across groups, indicating that post-interaction differences are unlikely to reflect pre-existing attitudes (see Table 1).

$$t_{\text{Welch}}(138.87) = -5.27, p = 5.04e-07, \hat{g}_{\text{Hedges}} = -0.77, \text{CI}_{95\%} [-1.07, -0.47], n_{\text{obs}} = 2$$



$$\log_e(\text{BF}_{01}) = -11.96, \hat{\delta}_{\text{difference}}^{\text{posterior}} = -15.86, \text{CI}_{95\%}^{\text{ETI}} [-21.61, -10.03], r_{\text{Cauchy}}^{\text{JZS}} = 0.71$$

$$t_{\text{Welch}}(226.24) = -5.52, p = 9.32e-08, \hat{g}_{\text{Hedges}} = -0.68, \text{CI}_{95\%} [-0.93, -0.43], n_{\text{obs}} = 2$$



$$\log_e(\text{BF}_{01}) = -11.96, \hat{\delta}_{\text{difference}}^{\text{posterior}} = -13.94, \text{CI}_{95\%}^{\text{ETI}} [-19.15, -8.76], r_{\text{Cauchy}}^{\text{JZS}} = 0.71$$

1.2. Post-Interaction Trust Differences

Descriptive comparisons of participant-level post-test scores indicated an approximately 12 point difference in post-test Trust Perception Scale-HRI scores ($M = 75$ vs $M = 63$) and a 27 point difference in the Trust in Industrial Human-robot Collaboration scale ($M = 39$ vs $M = 66$) between conditions, although in the first scale differences did not reach conventional significance under a two-sample t-test ($p=.10$), the second scale was significantly different between groups ($p=.007$).

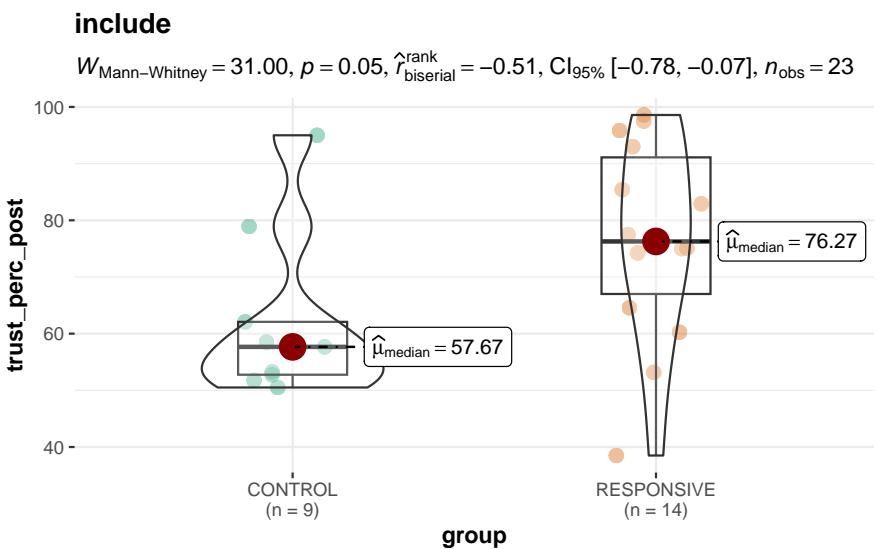
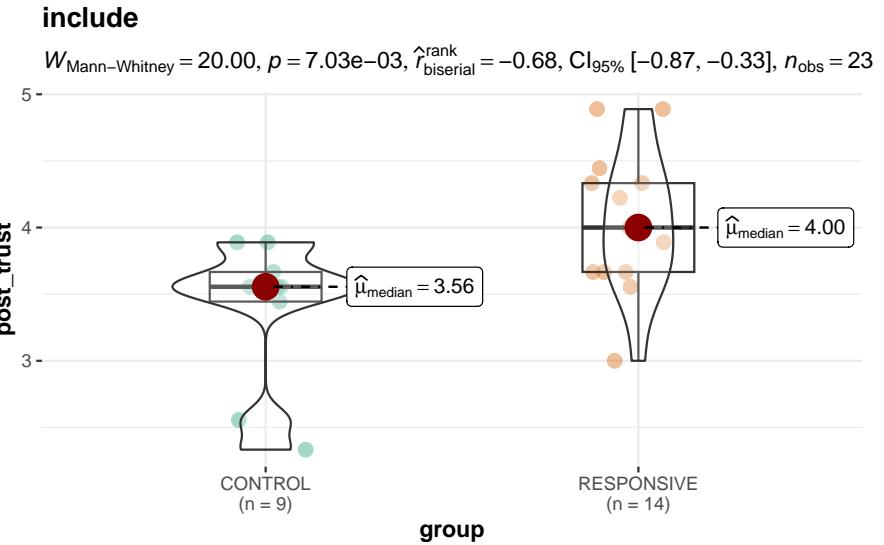
To test these findings further, we fitted several Bayesian hierarchical models were fitted (estimated using MCMC sampling with 4 chains of 4000 iterations and a warmup of 1000) to predict Robot HRI-trust and Trust in HRI Collaboration by experimental group (formula: $\text{robot_trust_post} \sim \text{group}$). The model included session_id and trust_items as random effects (formula: $\text{list}(\sim 1 \mid \text{session_id}, \sim 1 \mid \text{trust_items})$). Both models indicated higher post-interaction trust scores in the responsive robot condition across both trust-related scales (posterior median differences ~8–15 points on a 0–100 scale).

For the Trust in Industrial Robots outcome, the responsive condition showed a robust positive effect on post-task trust ratings. The estimated group difference was 14.5 points (95% CrI [5.62, 23.22]), exceeding between-session variability and remaining stable after accounting for item-level effects, with a 95% chance of being large (>6.94). In contrast, for the HRI trust perception scale, the estimated group effect was smaller at ~9 points and more uncertain 95% CrI [-1.72, 19.25], with 65% chance of being large (>6.87).

The posterior probability that the responsive condition increased trust was greater than 95% for both measures, suggesting a robust directional effect despite substantial individual variability. Sensitivity analyses using substantially wider priors yielded nearly identical posterior estimates for the group effect, indicating that results were not driven by prior specification.

In addition to directional effects, the posterior probability that the responsive condition increased trust by at least five points was 77% in HRI-trust perception and 98% in the Collaborative Trust in Industrial Robots scale, suggesting a reasonable likelihood of a practically meaningful effect. Moreover, in the latter collaboration scale, there is an 85% likelihood of an effect-size greater than 10 points.

Notably, in sessions characterized by severe communication breakdown, the responsive robot continued to provide extended assistance and meta-communication intended to repair the interaction. However, these efforts did not restore mutual understanding and may have increased participant confusion. In contrast, the control robot offered minimal, on-demand assistance, which—while less supportive overall—may have reduced cognitive overload under conditions of limited intelligibility. As a result, trust ratings in breakdown sessions did not track the intended responsiveness manipulation.



1.3. Trust subscale patterns

1.4. Interaction dynamics and task performance

1.4.1. Task performance

Objective task accuracy did not differ between conditions across any task-level measures except suspect accuracy (robot dependant task), indicating that increased trust was only attributable to improved task success when interaction was necessary to complete accurately.

Characteristic	N	CONTROL N = 9¹	RESPONSIVE N = 14¹	p-value²
post_trust	23	41 (22)	67 (21)	0.007
post_trust_reliability	23	41 (25)	65 (18)	0.022
post_trust_perception	23	39 (24)	55 (24)	0.14
post_trust_feelings	23	52 (32)	79 (22)	0.030
Post-Task Trust Perception	23	62 (15)	77 (18)	0.046
Suspect ID Accuracy	23	3 / 9 (33%)	9 / 14 (64%)	0.21
Status Accuracy	23	7 / 9 (78%)	9 / 14 (64%)	0.66
building_correct	23	6 / 9 (67%)	11 / 14 (79%)	0.64
floor_correct	23	6 / 9 (67%)	13 / 14 (93%)	0.26
zone_correct	23	5 / 9 (56%)	4 / 14 (29%)	0.38
Total Task Accuracy	23	3.00 (1.12)	3.29 (1.14)	0.49
Overall Task Accuracy	23	0.60 (0.22)	0.66 (0.23)	0.49
exclusions	23			
include		9 / 9 (100%)	14 / 14 (100%)	
Dialogue Turns	23	36 (7)	33 (5)	0.21
Avg Task Duration (mins)	23	13.82 (2.60)	15.26 (2.12)	0.16
Avg Response Time (ms)	23	13.21 (0.84)	17.24 (2.52)	<0.001
Silent Periods	23	5.67 (2.06)	4.71 (2.05)	0.29
Engaged Responses	23	2.22 (2.22)	3.50 (1.95)	0.077
Frustrated Responses	23	0.56 (0.73)	0.93 (1.21)	0.58
n_neg	23	1.00 (0.87)	0.71 (1.07)	0.28
prop_help_requests	23	0.40 (0.14)	0.36 (0.09)	0.48
prop_sentence_frag	23	0.18 (0.12)	0.20 (0.15)	0.92
prop_human_reasoning	23	0.29 (0.09)	0.35 (0.16)	0.37
prop_human_misunderstanding	23	0.02 (0.03)	0.05 (0.09)	0.85
prop_human_affective_engagement	23	0.05 (0.05)	0.10 (0.08)	0.058
prop_robot_helpful_guidance	23	0.68 (0.08)	0.84 (0.08)	<0.001
prop_robot_unhelpful	23	0.08 (0.05)	0.02 (0.03)	0.004
prop_robot_proactive_checkin	23	0.20 (0.08)	0.13 (0.08)	0.063
prop_robot_encouragement	23	0.00 (0.00)	0.36 (0.11)	<0.001
prop_robot_collaborative_lang	23	0.06 (0.04)	0.42 (0.16)	<0.001
prop_robot_reasoning	23	0.13 (0.05)	0.37 (0.14)	<0.001
prop_robot_clarification	23	0.22 (0.11)	0.48 (0.14)	<0.001
prop_robot_empathy_expression	23	0.00 (0.00)	0.13 (0.09)	<0.001
prop_comm_breakdown	23	0.21 (0.13)	0.22 (0.16)	>0.99

¹Mean (SD); n / N (%)

²Wilcoxon rank sum test; Wilcoxon rank sum exact test; Fisher's exact test; NA

Characteristic	N	include N = 23¹	p-value²
group	23		
CONTROL		9 / 23 (39%)	
RESPONSIVE		14 / 23 (61%)	
post_trust	23	57 (24)	
post_trust_reliability	23	56 (24)	
post_trust_perception	23	49 (25)	
post_trust_feelings	23	68 (29)	
Post-Task Trust Perception	23	71 (18)	
Suspect ID Accuracy	23	12 / 23 (52%)	
Status Accuracy	23	16 / 23 (70%)	
building_correct	23	17 / 23 (74%)	
floor_correct	23	19 / 23 (83%)	
zone_correct	23	9 / 23 (39%)	
Total Task Accuracy	23	3.17 (1.11)	
Overall Task Accuracy	23	0.63 (0.22)	
Dialogue Turns	23	34 (6)	
Avg Task Duration (mins)	23	14.70 (2.38)	
Avg Response Time (ms)	23	15.66 (2.84)	
Silent Periods	23	5.09 (2.07)	
Engaged Responses	23	3.00 (2.11)	
Frustrated Responses	23	0.78 (1.04)	
n_neg	23	0.83 (0.98)	
prop_help_requests	23	0.38 (0.11)	
prop_sentence_frag	23	0.19 (0.14)	
prop_human_reasoning	23	0.33 (0.13)	
prop_human_misunderstanding	23	0.04 (0.07)	
prop_human_affective_engagement	23	0.08 (0.08)	
prop_robot_helpful_guidance	23	0.78 (0.12)	
prop_robot_unhelpful	23	0.04 (0.05)	
prop_robot_proactive_checkin	23	0.16 (0.08)	
prop_robot_encouragement	23	0.22 (0.20)	
prop_robot_collaborative_lang	23	0.28 (0.22)	
prop_robot_reasoning	23	0.28 (0.16)	
prop_robot_clarification	23	0.37 (0.18)	
prop_robot_empathy_expression	23	0.08 (0.09)	
prop_comm_breakdown	23	0.22 (0.14)	

¹n / N (%); Mean (SD)

²NA

Despite similar task accuracy, interactions in the responsive condition were characterized by longer durations, slower response times, and a higher number of AI-detected engaged responses. These findings suggest that responsiveness altered the interaction dynamics and affective tone rather than task outcomes.

1.5. Individual differences and correlational patterns

As expected, we found that higher Need for Cognition (NFC) scores were negatively associated with Negative Attitudes Towards Robots (NARS), indicating that individuals who enjoy effortful thinking tend to have more positive attitudes towards robots. This relationship is consistent with prior literature suggesting that cognitive engagement is associated with openness to new technologies. In terms of NARS subscales, NFC was negatively correlated with all three subscales, but significantly so only in the domain of Situations of Interaction with Robots. This suggests that individuals with higher NFC are less likely to hold negative attitudes across various dimensions of robot interaction but especially around direct interaction with robots.

→ how to talk about post-interaction correlations w/pre-interaction measures
Several behavioural and task-level measures were correlated with post-interaction trust, consistent with the interpretation that trust judgments were shaped by interaction quality; these variables were not included as covariates in primary models to avoid conditioning on potential mediators.

Baseline negative attitudes toward robots were negatively correlated with post-interaction trust, with the strongest associations observed for affective trust subscales. In contrast, objective task performance was selectively associated with perceived reliability. Need for cognition was negatively correlated with negative robot attitudes and interaction-level negative affect, suggesting that individual differences contributed to variability in trust responses.

1.5.1. Model robustness and predictive checks

Sensitivity analyses using alternative prior specifications yielded substantively similar estimates, and leave-one-out cross-validation indicated comparable predictive performance between models with and without the group effect.

! TO DO:

- add subscale column to long format data
- run an analysis of performance by robot-dependent versus robot-independent tasks
- write up a future directions section for the planned larger study
- talk about unexpected language issues with people signing up with difficulty speaking and understanding english which caused problems with asr and interaction
- run analysis of dialogue dynamics included Bertopic or some other analysis of the actual content of the conversations/interactions

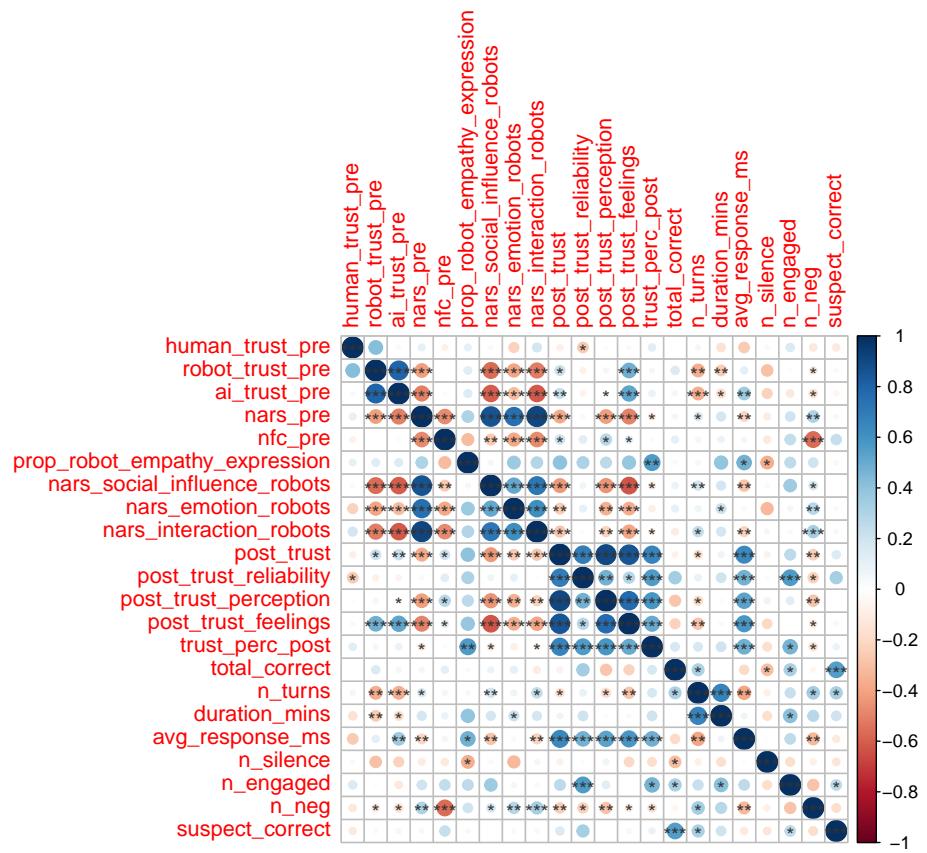


Figure 2

Need to discuss that these items were on a 0-100 scale that required sliding a bar, while the other trust scale was on a 1-5 Likert that required simply clicking. The post test was administered on a laptop with a trackpad which may have caused difficulties for some participants who found it difficult to drag the slider with the trackpad. This could have introduced additional noise into the measurement of this scale, which may explain why the effects were somewhat weaker here.

2. Discussion

Mention language confounders!!

The second task was intentionally designed to be sufficiently challenging that completing it within the allotted time was difficult without assistance. This ensured that interaction with the robot represented a meaningful opportunity for collaboration rather than a trivial or purely optional exchange. By contrasting a robot-dependent task with an open-ended advisory task, the study examined trust formation across interaction contexts that varied in both informational asymmetry and reliance on the robot.

This pilot study examined trust outcomes following in-person interaction with an autonomous social robot under two interaction policies: a responsive, affect-adaptive condition and a neutral, non-responsive control condition. By leveraging a fully autonomous dialogue system integrated with speech recognition and affect detection, the study aimed to evaluate how robot responsiveness influences trust formation in realistic human–robot collaboration scenarios.

Descriptive comparisons of post-interaction measures indicated that participants in the responsive condition reported consistently higher trust across all trust measures, with differences ranging from approximately 8 to 16 points on a 0–100 scale, although uncertainty remained high given the small sample. Notably, the responsive condition did not differ from control in objective task accuracy, suggesting that increased trust was not driven by improved task success. Instead, responsive interactions were characterized by longer durations, slower response times, and a higher number of AI-detected engaged responses, indicating a shift in interaction dynamics rather than performance.

Baseline negative attitudes toward robots were most strongly associated with affective components of trust rather than perceptions of reliability, suggesting that pre-existing attitudes primarily shape emotional responses to interaction rather than judgments of system competence. Conversely, objective task performance was selectively associated with perceived reliability, indicating that participants distinguished between affective and functional aspects of trust.

Future work with larger samples could formally test mediation pathways linking robot responsiveness, interaction fluency, affective responses, and trust judgments, as well as moderation by baseline attitudes toward robots and need for cognition.

Participants in the responsive condition also exhibited higher levels of AI-detected engagement during interaction, as indexed by a greater number of responses classified as positive affect (t-test result). This suggests that responsive behaviours altered the affective tone of the interaction itself.

2.1. Technical challenges

- Need to talk about language issues with participants who had difficulty speaking and understanding English which caused problems with ASR and interaction.
- Need to talk about issues where the AI was not able to flexibly handle when people asked a question about the suspect that was close to or another word for a ground-truth feature but not exactly the same word, causing confusion and miscommunication. E.g., “Was the suspect wearing pink?” The ground-truth feature was top: PINK, top-type: HOODIE; but the ASR and NLU did not extrapolate to understand that “wearing pink” referred to the same feature as “top: PINK”, causing confusion and miscommunication. Maybe the prompt could have included some examples of different phrasing which could improve this? To solve this issue in future work, we can expand the NLU training data to include more paraphrases and synonyms for each feature.

There was also a case where someone asked ‘is the top shirt hoodie red?’ to which the AI answered YES. It may have been confused by the multiple descriptors in the question. Future work could involve improving the NLU to handle more complex queries with multiple attributes.

Discuss future work where we will look investigate the ‘embodied’ effect of having a physical robot versus a virtual agent on trust and collaboration in HRI.

Also, prompt could include examples of what to do when dialogue appears fragmented, to remind participants to wait until the blue light is on before speaking and to switch up its phrasing if the robot seems to not understand.

Also, the control condition seemed to be somewhat neutered in terms of flexibility in responding in different ways. it would always respond with the exact same phrase when confronted with a sentence fragment or a question it could not directly answer.

Also issues with people not paying attention to the robot’s visual cues to know when to speak, leading to more fragmented dialogue. Future work could involve improving participant instructions, improved latency and ‘listening’ ... and the robot’s feedback mechanisms to better manage turn-taking.

Need to remember to flag participants who did not complete/skipped specific tasks. E.g. P56 skipped the wrapup entirely. Many skipped the brief (by advancing on their own through the dashboard).

3. Appendix

3.1. Dialogue Coding Scheme

3.1.1. Task Outcome Layer (Stage-Level)

Variable	Type	Description
task_outcome	categorical	Final task status (<code>completed</code> , <code>timeout</code> , <code>skipped</code> , <code>partial</code> , <code>abandoned</code>). Exactly one per task.
task_completed	binary	Task goal was fully completed within the allotted time.
task_timed_out	binary	Task ended due to expiration of the time limit before completion.
task_skipped	binary	Participant explicitly skipped or advanced past the task without completing it.
task_partially_completed	binary	Task progress was made, but the full solution was not reached.
task_abandoned	binary	Participant disengaged or stopped attempting the task before timeout.
task_time_remaining_sec	numeric	Time remaining (in seconds) when the task ended; 0 if timed out.
task_completed_without_help	binary	Task was completed without any help requests to the robot.
task_required_robot_help	binary	At least one robot help interaction was required for task completion.

3.1.2. Dialogue Interaction Layer (Turn-Level)

3.1.2.1. Human Turn Codes.

Variable	Type	Description
<code>human_help_request</code>	binary	Participant explicitly or implicitly asks the robot for help or guidance.
<code>human_reasoning_self</code>	binary	Participant articulates their own reasoning or problem-solving independent of the robot.
<code>human_confusion</code>	binary	Participant expresses confusion or uncertainty.
<code>human_confirmation_seeking</code>	binary	Participant seeks confirmation of a tentative belief or solution.
<code>human_ignores_robot</code>	binary	Participant proceeds without engaging with the robot's prior input.

3.1.2.2. Robot Turn Codes.

Variable	Type	Description
<code>robot_helpful_guidance</code>	binary	Robot provides accurate, task-relevant guidance.
<code>robot_misleading_guidance</code>	binary	Robot provides misleading or incorrect guidance.
<code>robot_factually_incorrect</code>	binary	Robot states information that is objectively incorrect.
<code>robot_policyViolation</code>	binary	Robot violates stated system or task constraints.
<code>robot_on_policy_unhelpful</code>	binary	Robot adheres to policy but provides vague or non-actionable assistance.
<code>robot_stt_failure</code>	binary	Robot response reflects a speech-to-text or input understanding failure.
<code>robot_clarification_request</code>	binary	Robot asks the participant to repeat or clarify their input.

3.1.3. Affective Interaction Layer (Turn-Level)

3.1.3.1. Robot Affective Behavior.

Variable	Type	Description
robot_empathy_expression	binary	Robot expresses empathy, encouragement, or reassurance.
robot_emotion_acknowledgment		Robot explicitly references an inferred participant emotional state.
robot_affect_task_aligned	binary	Robot's affective response is appropriate and supportive in context.
robot_affect_misaligned	binary	Robot's affective response is mistimed or disruptive to the task.

3.1.3.2. Human Affective Response.

Variable	Type	Description
human_affective_engagement	binary	Participant responds in a socially warm or emotionally engaged manner.
human_social_reciprocity	binary	Participant mirrors or responds to the robot's affective expression.
human_anthropomorphic_language		Participant treats the robot as a social agent.
human_emotional_disengagement		Participant responds in a curt, dismissive, or withdrawn manner.

3.1.4. Notes

- Turn-level variables are coded per dialogue turn.
- Task outcome variables are coded once per `session_id` × `stage`.
- Raw dialogue text was retained during coding and removed prior to aggregation.
- Multiple turn-level codes may co-occur unless otherwise specified.

3.2. Key Components of the System

This study implemented a multi-stage collaborative task system where participants collaborate with the Misty II social robot to solve a who-dunnit type task. The system utilizes an autonomous, mixed-initiative dialogue architecture via langchain with affect-responsive capabilities.

1. Misty-II Robot: A programmable robot platform equipped with sensors and actuators for interaction.
2. Automated Speech Recognition (ASR): A speech-to-speech pipeline that processes spoken input from users and converts it into text for LLM processing then back to speech for output on the robot.
 - STT: Deepgram API for real-time speech-to-text conversion.
 - DistilRoBERTa-base fine-tuned on emotion classification for emotion detection from user utterances
 - LLM: Gemini API for processing text input and generating contextually relevant responses in JSON format
 - TTS: Misty-II text-to-speech (TTS) engine on 820 processor.
3. Langchain Dialogue Management: A system that manages the flow of conversation, ensuring coherent and contextually appropriate dialogue within a two-part collaborative task.
4. Collaborative-Tasks
 - Task 1: Whodunnit style task where human and robot collaborate to find a missing robot via the human asking Yes/No questions (process of elimination in 6x4 suspect grid) to the robot. Robot knows ground truth but can only answer Yes/No questions about suspect features. Can not directly describe the suspect or name them. (human can choose a random suspect to solve on their own but only 1 in 24 chance of being correct without robot help)
 - Task 2: Where is Atlas? Robot collaborates with human to find Atlas by deciphering cryptic system and sensor logs. Robot does not know the answer here and can only guide the human using its expertise and knowledge of computer systems and basic logical reasoning. (human can solve on their own but very difficult without robot help depending on participants technical background).
5. Flask-gui dashboard interface: A web-based interface/dashboard that allowed participants to interact with the tasks, view task-related information and input their answers to the questions. Responses were sent to the robot to signal task progression.
 - Task 1 dashboard: Displays the suspect grid and allows the user to select suspects and view their features.
 - Task 2 dashboard: Displays system logs and allows the user to input their findings.
6. Pre and post tests:
 - PRE-TESTS: Need for Cognition Scale (short); Negative Attitudes to Robots Scale (NARS);

- POST-TESTS: Trust Perception Scale-HRI; 9 custom questions adapted from Charalambous et al. (2020) on trust in industrial human-robot collaboration;

4. Technical Specifications

4.1. System Overview

This study implements a multi-stage collaborative task system where participants collaborate with the Misty II social robot to solve a who-dunni type task. The system utilizes an autonomous, mixed-initiative dialogue architecture with affect-responsive capabilities.

4.2. Hardware Platform

Robot: Misty II Social Robot (Furhat Robotics)

- Mobile social robot platform with expressive display, arm actuators, and head movement
- RGB LED for state indication
- RTSP video streaming (1920×1080 , 30fps) for audio capture
- Custom action scripting for synchronized multimodal expressions

4.3. Software Architecture

4.3.1. Core System Components

Programming Language: Python 3.10

Primary Dependencies:

- `misty-sdk` (Python SDK for Misty Robotics API) - Robot control and sensor access
- `deepgram-sdk` (4.8.1) - Speech-to-text processing
- `ffmpeg-python` (0.2.0) - Audio stream processing
- `flask` (3.1.2) + `flask-socketio` (5.5.1) - Web interface for task presentation
- `duckdb` (1.4.0) - Experimental data logging database

4.3.2. Large Language Models

LLM Provider:

Google Gemini:

- Model: `gpt-3.5-turbo` (configurable via environment variable)
- Integration: `langchain-google-genai` with `google-generativeai` API
- Response format: JSON-only output (`response_mime_type: "application/json"`). This format is required by Misty-II for reliable parsing and for action execution.

LLM Configuration:

- Temperature: 0.7 (for balanced creativity and coherence)
- Memory: Conversation buffer memory with file-based persistence (`langchain.memory.ConversationBufferMemory`)
- Context window: Full conversation history maintained across interaction stages but reset between sessions.

4.4. LangChain Framework Integration

4.4.1. Core LangChain Components

Framework Version: langchain-core with modular provider packages

- `langchain` (meta-package)
- `langchain-community` (0.3.31)
- `langchain-google-genai` Gemini integration

4.4.2. ConversationChain Architecture

Memory Management (ConversationChain class in `conversation_chain.py`):

1. **Conversation Buffer Memory:**
 - Implementation: `langchain.memory.ConversationBufferMemory`
 - Storage: File-based persistent chat history (`FileChatMessageHistory`)
 - Format: JSON files in `.memory/` directory, one per participant session
 - Memory key: "history"
 - Return format: Message objects (full conversation context)
2. **Memory Reset Policy:**
 - Default: Reset on each new session launch
 - Archive previous session: Timestamped archive files stored in `.memory/archive/`
 - Configuration: `RESET_MEMORY` and `ARCHIVE_MEMORY` environment variables

4.4.3. Prompt Construction

Message Structure

(LangChain message types): `python [SystemMessage, *history_messages, HumanMessage]`

System Message Assembly:

- Core instructions (task framing, role definition)
- Personality instructions (mode-specific behaviour)
- Stage-specific instructions (current task context)
- Output format constraints (JSON schema specification)

```
Human Message Format:  {
  "user": "<transcribed_speech>",
  "stage": "<current_stage>",
  "detected_emotion": "<emotion_label>",
  "frustration_note": "<optional_alert>",
  "timer_expired": "<task_id>",    ...
}
```

- JSON-encoded context variables passed alongside user input
- Enables LLM to access environmental state without breaking message history

4.4.4. Memory Persistence:

- Save after each turn: `memory.save_context({“input”: user_text}, {“output”: llm_response})`
- Maintains conversational coherence across multi-stage interaction
- Enables LLM to reference previous exchanges (e.g., “As I mentioned earlier...”)

4.4.5. LangChain Design Rationale

Why LangChain for this application:

1. Memory abstraction: Automatic conversation history management without manual message list handling
2. Provider flexibility: Easy switching between Gemini and OpenAI without rewriting prompt logic
3. Message typing: Structured SystemMessage/HumanMessage/AIMessage types maintain role clarity
4. File persistence: Built-in FileChatMessageHistory enables session recovery and archiving
5. Future extensibility: Framework supports adding tools, retrieval, or multi-agent patterns if needed

Alternatives considered: Direct API calls would reduce dependencies but require reimplementing conversation history management, prompt templating, and cross-provider compatibility layers.

4.4.6. LangChain Limitations in This Context

- No chains used: Despite name ConversationChain, this is a direct LLM wrapper (no LangChain Expression Language chains)
- No tools/agents: Simple request-response pattern (could extend for future tool-use capabilities)
- Custom JSON parsing: LangChain’s built-in output parsers not used; custom extraction handles malformed responses more robustly

4.4.7. Speech Processing

Speech-to-Text (STT):

- Provider: Deepgram Nova-2 (`deepgram-sdk` 4.8.1)
- Model: `nova-2` with US English (`en-US`)
- Smart formatting enabled
- Interim results for real-time partial transcription
- Voice Activity Detection (VAD) events
- Adaptive endpointing: 200ms (conversational stages) / 500ms (log-reading task)

- Utterance end timeout: 1000ms (conversational) / 2000ms (log-reading)
- Audio processing: RTSP stream from Misty → FFmpeg MP3 encoding → Deepgram WebSocket

Text-to-Speech (TTS) - Three options:

1. **Misty Onboard TTS** (this is the one we used): Native robot voice via onboard TTS
2. **OpenAI TTS**:
 - Model: `tts-1` (low-latency variant)
 - Voice: `sage`
 - Format: MP3, served via HTTP (port 8000)
 - Ultimately chose not to use because we wanted a more robotic, non-human voice
 - Didn't want the human voice influencing trust on its own (future research could look at trust in relation to type of voice)
3. **Deepgram Aura**:
 - Model: `aura-stella-en` (conversational female voice)
 - Format: MP3, served via HTTP
 - Ultimately chose not to use because we wanted a more robotic, non-human voice

4.4.8. Emotion Detection

Model: DistilRoBERTa-base fine-tuned on emotion classification

- HuggingFace identifier: `j-hartmann/emotion-english-distilroberta-base`
- Framework: `transformers` (4.57.1) pipeline
- Hardware: CUDA GPU acceleration (automatic fallback to CPU)
- Output classes: joy, anger, sadness, fear, disgust, surprise, neutral
- Mapped to interaction states: positively engaged, irritated, disappointed, anxious, frustrated, curious, neutral

4.4.9. Multimodal Robot behaviour

Expression System: 25 custom action scripts combining:

- LLM was prompted to choose an appropriate expression from a predefined set based on context.
- Facial displays (image eye-expression files on screen)
- LED color patterns (solid, breathe, blink)
- Arm movements (bilateral position control)
- Head movements (pitch, yaw, roll control)

Nonverbal Backchannel behaviours (RESPONSIVE mode only):

- Real-time listening cues triggered by partial transcripts (disfluencies, hesitation markers)
- Emotion-matched expressions (e.g., “concern” for hesitation, “excited” for breakthroughs)

LED State Indicators:

- Blue (0, 199, 252): Actively listening (microphone open)
- Purple (100, 70, 160): Processing/speaking (microphone closed)

4.5. Data Collection

Database: DuckDB relational database (`experiment_data.duckdb`)

Logged Data:

1. **Sessions table:** participant ID (auto-incremented P01, P02...), condition assignment, timestamps, duration
2. **Dialogue turns table:** turn-by-turn user input, LLM response, expression, response latency (ms), behavioural flags
3. **Task responses table:** submitted answers with timestamps and time-on-task
4. **Events table:** stage transitions, silence check-ins, timer expirations, detected emotions

4.6. Interaction Dynamics

4.6.1. Silence Handling

Silence detection: 25-second threshold triggers check-in prompt

- RESPONSIVE: “Still working on it? No rush - I’m here if you need help!”
- CONTROL: “I am ready when you have a question.”

4.6.2. Emotion-Responsive behaviours (RESPONSIVE condition only)

Frustration tracking:

- Consecutive detection of frustrated/anxious/irritated/disappointed states
- Threshold: 2 consecutive frustrated turns triggers proactive support
- RESPONSIVE adaptation: “This part can be tough. Want me to walk you through it?”

Positive emotion matching:

- Celebratory language for curious/engaged states
- Momentum maintenance: “Yes! Great observation!”

Run Mode: Set programmatically in `mistyGPT_emotion.py` line 126:

```
RUN_MODE = "RESPONSIVE" # or "CONTROL"
```

4.7. Prompt Engineering

Modular prompt system (PromptLoader class):

- `core_system.md`: Task framing, role description, output format schema

- role_responsive.md / role_control.md: Condition-specific personality instructions
- stage1_greeting.md through stage5_wrap_up.md: Stage-specific task instructions.

Context injection: Real-time contextual variables passed to LLM:

- Current stage
- Detected emotion (if enabled)
- Task submission status
- Timer expiration notifications
- Silence check-in flags

4.8. Inter-process Communication

Flask REST API endpoints:

- GET /stage_current: Synchronize stage state with facilitator GUI
- GET /task_submission_status: Detect participant task submissions
- GET /timer_expired_status: Detect timer expirations
- POST /stage: Update stage (facilitator override)
- POST /reset_timer: Clear timer expiration flags

References

- Lin, T.H., Ng, S., Sebo, S., 2022. 2022 31st ieee international conference on robot and human interactive communication (ro-man), pp. 37–44. URL: <https://ieeexplore.ieee.org/document/9900828>, doi:[10.1109/R0-MAN53752.2022.9900828](https://doi.org/10.1109/R0-MAN53752.2022.9900828). iSSN: 1944-9437.