

# Examining Audio Communication Mechanisms for Supervising Fleets of Agricultural Robots

Abhi Kamboj, Tianchen Ji, and Katie Driggs-Campbell

**Abstract**—Agriculture is facing a labor crisis, leading to increased interest in fleets of small, under-canopy robots (agbots) that can perform precise, targeted actions (e.g., crop scouting, weeding, fertilization), while being supervised by human operators remotely. However, farmers are not necessarily experts in robotics technology and will not adopt technologies that add to their workload or do not provide an immediate payoff. In this work, we explore methods for communication between a remote human operator and multiple agbots and examine the impact of audio communication on the operator’s preferences and productivity. We develop a simulation platform where agbots are deployed across a field, randomly encounter failures, and call for help from the operator. As the agbots report errors, various audio communication mechanisms are tested to convey which robot failed and what type of failure occurs. The human is tasked with verbally diagnosing the failure while completing a secondary task. A user study was conducted to test three audio communication methods: earcons, single-phrase commands, and full sentence communication. Each participant completed a survey to determine their preferences and each method’s overall effectiveness. Our results suggest that the system using single phrases is the most positively perceived by participants and may allow for the human to complete the secondary task more efficiently. The code is available at: <https://github.com/akamboj2/Agbot-Sim>.

## I. INTRODUCTION

Agriculture is currently facing a human labor crisis [1], harming profitability and causing negative downstream effects. As a result, precise actions (e.g., weeding, targeted pesticides) are not feasible at the scale required for annual row crops (e.g., corn, soybeans), which dominate the Midwest and much of the US. In these settings, agriculture is only practical with heavy reliance on fertilizers, pesticides, and herbicides applied with large equipment (e.g., tractors, combines) [2], [3]. While this equipment is familiar to farmers and can be automated to alleviate labor concerns [4], such equipment is capital heavy, requires extra logistical oversight, introduces new safety risks to workers, and physically impacts the farm (e.g., soil-compaction, crop damage).

Small agricultural robots (agbots) can help alleviate the labor crisis and enable precision agriculture. These agbots are designed to be small, inexpensive, and intelligent, and have seen growing attention in recent years [5], [6]. To fully address the labor shortage, these agbots must be both easy

to use and able to be deployed at scale, where one human is supervising many robots. Despite growing interest in agbots, farmers already tend to be overwhelmed with the large number of equipment and data sources that are constantly made available to them [7]. From large suppliers automating their products to new technologies for data collection [8], farmers are being driven to manage a large set of equipment each with their own intricacies [1].

Many human-robot interaction (HRI) methodologies focus on intuitive interface designs that allow for seamless integration of robots in a society of non-experts [9]–[12]. Simulations are an effective method to test and overcome barriers preventing adoption of robots [13], [14]. We follow a design focused simulation methodology to study which type of auditory interaction most positively influences a user’s perception and productivity in a remote monitoring setting. Our simulation stands in for actual agbots navigating fields of crops. The failure cases and solutions we simulate are identifiable and previously tested on a physical robot [15].

We consider situations where a fleet of agbots are deployed in a field and a remote human operator supervises the robots. As the robots navigate through crop rows, they randomly encounter failures that require help from the operator. During each failure, the robot control center will audibly prompt the operator to provide assistance through verbal commands. Across experiments, we vary the audible prompt and measure the operator’s perception of the system and productivity of completing a secondary task.

Handling failure scenarios is of critical importance for near-term deployment, as many agriculture tasks and environments are too complex for current agbots to handle without at least occasional failure [16]–[18]. This setting follows the idea of sliding autonomy [19] and recent efforts to codify the levels of autonomy for field robots [20], which outlines how agbots with varying autonomous capabilities can interact with human operators [21]. Failure cases have a strong effect on the user’s perception and trust in a system, meaning handling failures largely impacts the overall success of an HRI system [22].

This research provides insight into the effectiveness and acceptance of audio communication interfaces when managing multiple autonomous robot failures. We studied the effect of earcons, single phrase commands, and full sentence speech on the user’s perception of the system and their efficacy in completing a secondary task. We present three contributions:

- 1) We develop a simulated control center to explore how humans interact and monitor agricultural robots deployed across a field, while potentially encountering

A. Kamboj, T. Ji, and K. Driggs-Campbell are with the department of Electrical and Computer Engineering at the University of Illinois, Urbana-Champaign {akamboj2, tj12, krdc}@illinois.edu

This work was supported by the USDA National Institute of Food and Agriculture (USDA/NIFA), through the National Robotics Initiative 2.0 (NIFA#2021-67021-33449), the AI Institute AIFARMS through the Agriculture and Food Research Initiative (AFRI) (USDA/NIFA Award no. 2020-67021-32799), as well as the Illinois Center for Digital Agriculture.

failures that require human assistance.

- 2) We demonstrate how audio signals (either tones or natural language) improve an operator's efficiency and productivity compared to a traditional visual interface.
- 3) Our user study provides insight on how well an operator perceives various auditory interaction systems in a remote robot monitoring setting, indicating which system will most effectively be adopted.

This paper is organized as follows. We review relevant literature in Section II. In Section III, we present an overview of our exploratory study, our hypotheses, and our measures. Section IV and V discuss the quantitative and qualitative findings from our user study. Finally, we conclude in Section VI.

## II. RELATED WORK

### A. Audio in Design

Auditory interfaces can be defined as bidirectional, communicative connections between two systems using audio, where audio refers to the production of sound [23]. There are various methods in which sound can be used for auditory interaction in HRI. In HRI sound either deliberately communicates an intention, such as speech, notifications, or semantic-free utterances [24], or comes without intention such as consequential sound or movement sonification [11]. Using sound to intentionally convey information has many benefits including reducing visual overload, reinforcing visual messages, and providing additional information such as direction or emotion [23]. Our study uses sound as a primary interface to reduce visual load for the human operator to complete a visually intensive secondary task. To the same degree that human-to-human correspondence involves various modes interaction such as audible, visual, and tactile, HRI also requires multi-modal interaction for successful integration of robots into society [25], [26].

Most studies about sound in HRI involve improving a user's perception of human-robot dialogue or nonverbal noises [11], [27], ignoring the wide range of possibilities that auditory interaction can offer [28]. There are four main ways data can be encoded into audio: auditory icons, earcons, sonification and speech [29]. Auditory icons are intuitive associations between a recognizable sound and a piece of information, earcons are unintuitive intentionally designed associations, sonification maps information to variations in sound, and speech conveys information verbally [29].

We derive three auditory interfaces from the literature:

- 1) **Earcons:** An association from noises to a piece of information. Earcons have often been studied in the HRI community in multimodal systems [30] such as autonomous driving [31], adaptive automation of telerobotic control [26], and healthcare [32].
- 2) **Phrases:** A truncated version of complete speech. We create this interface to balance the benefits of condensed information while maintaining some psychological or social aspects of speech. Other studies have also attempted to use some sort of hybrid interface

such as spearcons, sped up speech, [33] or audification, using data as sounds [23].

- 3) **Sentences:** A complete verbal expression most similar to conversational speech. Most of the literature studying audio in HRI uses this type of verbal dialogue.

### B. Audio in Agriculture

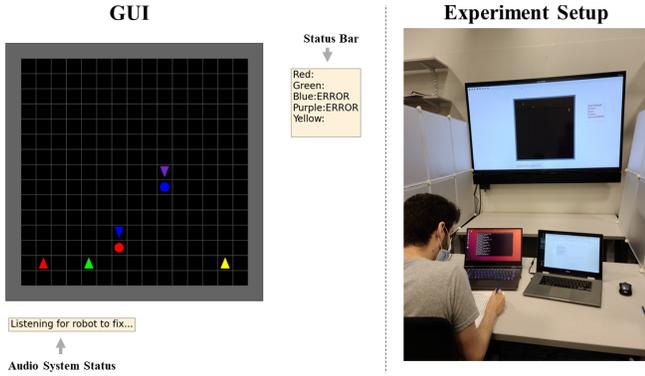
A survey on nontraditional human robot interactions in agriculture highlights some of the benefits of using speech to improve the usability of technology [17]. Moreover, recent research efforts also determined that robots have not reached a level of design that allows for effective communication of faults by untrained users [9]. Literacy has been a major barrier preventing farmers who cannot read written instructions from using robots [17]. Many studies attempt to address the issue by surveying farmers and developing audio interfaces to assist under-educated farmers in using technology and communicating with robots [34]. These studies discovered that the complexity of a robot system is one of the main challenges farmers face in adapting technology, as employees lack the necessary skills to operate complicated robots. This conclusion underscores the importance of an intuitive and effective human-centered design when it comes to robot management on an autonomous farm.

### C. Communicating Robot Failures Using Audio

Previous user studies on the failure rate in HRI tasks indicate that direct communication is more important than conversational dialogue to overcome a misperception error and complete a task collectively [35], [36]. Thus, task-oriented robots should focus on concise utterances rather than lengthy dialogue. This conclusion follows the well-researched principle of least collaborative effort in grounding dialogue, i.e. establishing a mutual understanding [37], [38].

However, more detailed speech improves perceived capability [39]. Restricting the robot's vocabulary to simple commands or sounds can negatively impact the robot's usability. Because of the robot's limited speech, the human may underestimate the capability of the robot, thus not utilizing it to its full extent which is called underperception. In overperception, the human overestimates the capability of the robot leading to the human's expectations not being met, and therefore the human is less willing to work with the robot. When a human misperceives a robot's capabilities, they misuse it, and their acceptance of the robot decreases which can decrease the success of the collective HRI task [39]. In general the habitability gap, a mismatch between human expectations and technologies' capabilities, limits speech based human-machine interaction [40].

Most literature in communicating HRI failures, including all the sources discussed thus far, deal with embodied robots, either in a video or in a physical setting. In our remote operation settings, detailed speech may not improve perceived capability as the robot is not as anthropomorphic. On the other hand, a remote failure communication system may require more collaborative effort to establish grounding. In general, strategies that work for humanoid systems do not



**Fig. 1: Left: System GUI.** The status bar indicates which robots have failed (stopped behind circles). Here the purple robot is behind a blue circle (Unrecoverable Failure) and the blue robot is behind a red circle (Row Collision). The red, green and yellow robots are unobstructed and continue to traverse their sections of the grid. The audio system status shows the system is waiting to hear a color indicating which robot to fix. **Right: Experiment setup.** The setup mimics the control center setting on an autonomous farm. The user is on an isolated desk in front of a TV screen with the GUI and is given wordsearch puzzles to work on.

transfer over to virtual robotic agents [41]. Therefore, we further investigate earcon, phrase and sentence communication systems to better understand how a user perceives such a remote operation system and what collaborative effort is needed in the dialogue.

Recent research efforts have made noteworthy progress in studying HRI in multirobot systems including coordination [42], user task-switching [43], and scalability [21], [44]. Scalability is especially a concern in failure scenarios of multi-robot systems due to the exponential growth of state and action spaces [45]. However, the effects of different forms of audio on the user’s perception of multi-robot remote supervision failure scenarios has not been well investigated.

### III. METHODS

#### A. Experimental Design

In order to study the design of an audio interface for agbot fleet monitoring, we develop an autonomous farm environment [46] as shown in Figure 1. Five robots represented as colored triangles navigate up and down columns of the grid sequentially. Each of the five robots traverses 20% of the grid and then halts to indicate the completion of the task. When a failure, represented as a circle, is reached the robot stops and prompts the user to diagnose the error. Based on the error case, the human verbally communicates to the robot on how to resolve the error to reactivate the robot.

We conduct four experiments in a random order on 13 participants. Each experiment prompts the participant differently when a robot fails: earcons, phrases, sentences, and no sound at all. Tables I-IV, are provided to the participant before each experiment. After each experiment, we measure the user’s perception of success and usability of the system through a survey and their level of productivity through a secondary task score. The order in which the participant went through the four different conditions was randomized to reduce the

**TABLE I: Scripts for each audio communication modality (refer to Table IV for earcon mappings and note how the robot’s sentence structure changes in “Sentence” to make it similar to conversational speech)**

Earcon	Phrase
Robot: “[red]”	Robot: “Error at red.”
Human: “Fix the red robot.”	Human: “Fix the red robot.”
Robot: “[robot_fail]” x 2	Robot: “Untraversable obstacle.”
Human: “Navigate around.”	Human: “Navigate around.”
Robot: “[robot_fixed]”	Robot: “Error fixed.”
Robot: “[blue], [green]”	Robot: “Errors at blue, green.”
Human: “Fix the blue robot.”	Human: “Fix the blue robot.”
Robot: “[robot_fail]” x 1	Robot: “Row collision.”
Human: “Reverse and retry.”	Human: “Reverse and retry.”
Robot: “[robot_fixed]”	Robot: “Error fixed.”
Robot: “[green]”	Robot: “Errors at green.”
...	...
Sentence	
Robot: “There is an error at the red robot.”	
Human: “Fix the red robot.”	
Robot: “The red robot is facing an untraversable obstacle.”	
Human: “Navigate around.”	
Robot: “The error has been fixed.”	
Robot: “There are errors at the following robot blue, green.”	
Human: “Fix the blue robot.”	
Robot: “Row collision has occurred at the blue robot.”	
Human: “Reverse and retry.”	
Robot: “The failure has been fixed.”	
Robot: “There are still errors at the green robots.”	
...	

bias in the results. To prevent unfamiliarity with the system from influencing the results, the participant is provided with a tutorial on how to fix the robots and can practice until they feel comfortable with the system. The audio system randomly chosen as the fourth experiment is used in the tutorial before the experiments, to discourage participants from considering the tutorial in their survey responses.

#### B. Failure Cases

In an attempt to recreate realistic failures from the agricultural domain, we consider three common failure cases: row collision, obstacle, and unrecoverable failure, detailed in Table II. The above failures can be reliably detected on the field and are assumed to be solvable from the operator’s commands [15], [18].

The participant addresses each of the failure cases using certain verbal commands shown in Table III. Although the one-to-one mapping from failures to solutions can easily be automated without a human operator, such a system can be extended to a scenario where the human operator must make an informed decision. Even in the current system, the operator chooses the order in which to fix multiple robot failures, which is not a trivial task for humans or planners [21]. Nonetheless, knowing how or which robot to address is unrelated to the type of auditory prompt the participant hears, thus the triviality of the system does not discount the merit from a user-centered design perspective.

As this study focuses on the impact of audio communication on perception and productivity, the number and type of

**TABLE II:** Failure types with description and recovery solution.

Failure Type	
Row Collision	<i>Description:</i> The robot deviates from the center line and crashes into crops due to navigation failures. <i>Solution:</i> Reverse and replan the path that tracks the center line.
Obstacle	<i>Description:</i> The robot encounters obstacles, which obstruct the center line, but still has room around to plan a collision free path. <i>Solution:</i> Navigate around the obstacle and continue the robot’s original trajectory.
Unrecoverable Failure	<i>Description:</i> The robot is in some failure scenarios where it cannot continue without human intervention, e.g, a fully blocked path. <i>Solution:</i> Send a human to the field to assist the robot to recover.

**TABLE III:** User commands for the different failure modes

Failure #	GUI Icon	Failure Type	Solution
1	●	Row Collision	“reverse and retry”
2	●	Obstacle	“navigate around”
3	●	Unrecoverable Failure	“sending human”

failures is kept constant. There are exactly 15 failures in each simulation, five of each failure type. However, the location of each failure is randomly sampled from a uniform distribution across the grid to simulate more realistic failure scenarios. A failure is not visible on the grid until a robot reaches it, inhibiting the participant from anticipating a failure.

### C. Audio Signals

The verbal interaction from the human to the robot remains constant throughout all the simulations, but the auditory interaction from the robot to the human will change as described in Tables I. When one or more robots fail, the system will prompt the participant with audio indicating the colors of the failed robots. The participant will have to verbally say a color to indicate which robot they wish to fix, and the system will indicate which type of failure the robot is facing. The participant then must say the correct command to fix the system. Once fixed, the system will notify the participant that the robot has been fixed and continue with the simulation until another robot fails.

To keep the amount of information conveyed by each audio system constant, we developed a mapping from earcons to robot colors, shown in Table IV. To convey what type of failure the robot is at, the system plays a coin noise to indicate the failure number shown in Table III: once for row collision, twice for obstacle, and thrice for unrecoverable failure. Before the earcon experiment, an earcon tutorial program played each of the earcons and their definitions as described in Table IV.

In the no sound experiment, the participants can only refer to the visual GUI to find out if a failure has occurred, which robot has encountered the failure, and what type of failure the robot is facing. Since the participant had no audio prompt, they would have to occasionally look up from the secondary task to address the failures. The GUI shows all the necessary

**TABLE IV:** Earcon mappings from sound to meaning

Robot	Red	Green	Blue	Purple	Yellow
Sound	Siren	Leaves Rustle	Splash	Violin	Taxi Honk
Condition	Robot_fixed	Robot_fail			
Sound	Ready	Coin			

information conveyed by the audio prompts, as described in Figure 1.

### D. Secondary Task and Productivity

In an autonomous robot control center setting, the operator is unlikely to fully focus on monitoring the robots. Instead, the operator would be completing other tasks and be prompted when a robot needs attention. As a controlled secondary task, a wordsearch puzzle is used in our experiment, which asks a participant to find given words going in various linear directions in a grid of letters. The wordsearch puzzle serves as a visual stimulus and cognitive load that the participant can engage with as it does not interfere with the auditory interaction of the system we wish to study. It is very easy for beginners to learn and does not give an advantage to those with more cultural knowledge or math practice (in contrast to crossword puzzles or math questions). Many psychology studies use word search puzzles to study distraction or multitasking [47], [48].

Dividing the participant’s awareness across different senses allows us to measure their productivity when switching attention between the two tasks at hand. The participant was put in a constant quiet environment such that the system was their only audio stimuli, as shown in Figure 1. User studies with physical robots often have auditory background that could affect the user’s perception of the system [49]. However, our distraction-free environment isolates the audio’s effect on the participant’s perception of the system and success in solving the wordsearch puzzles.

The productivity score of each experiment is the measure of how many words the participant found divided by the total time of the simulation (the simulation stops when all five robots reach the end of their last column of crop). This words per minute score is robust to small technical inconsistencies or pauses in the system as well as how the participant divides their attention. If the system momentarily glitches, then the participant has another second to think and find words. If the participant only focuses on finding words, the number of failed robots will increase (and eventually come to a complete standstill), thus increasing the time that the simulation takes to complete, which would give them a low productivity rate. On the other hand, if the participant does not focus on finding words at all and only focuses on the robots, they will find less words and receive a low productivity rate. Overall, the system encourages the participant to address both the wordsearch task and the robot failure task at the same time, which strengthens our findings on how audio affects efficiency when faced with a visually intensive task. At the beginning of the study, each participant

was made aware of this metric and that they should focus on both the simulation and word search.

### E. Participants

A total of 13 participants voluntarily performed the IRB approved user study, each signing a consent form beforehand. Each of the simulations took approximately eight minutes to run, but with the explanations and tutorial the entire process took around one hour. Every participant was a university affiliate; however, they had varying degrees of previous experience in robotics, with an average experience of 3.4 on a scale of 1 (unfamiliar) to 5 (expert).

### F. Survey Questions

The survey asked a few background questions before the experiment was conducted. After each experiment, the participant was asked to rate the following questions on a 7 point Likert scale:

*Q1: It was easy to diagnose and fix the errors in this system.*

*Q2: I was successful at guiding robots passed their failures.*

*Q3: This system was overwhelming to use.*

The questions were created for the purpose of this design experiment, since most standard HRI questionnaires are designed for physical experiments. However, our questions relate to the competence dimension of RoSAS [50] and the likeability dimension of Godspeed [51], which are two common HRI metrics.

After all audio systems were tested, the participant was asked to rank the notification methods from 1 to 4 (best to worse) and to leave feedback on the overall system.

### G. Hypotheses

The following three hypotheses were developed and tested:

**H1: Any audio interface will facilitate better success, usability, and productivity over a purely visual method.**

In a control center setting, the human will likely be addressing robot failures while operating other systems at once, which in our study is analogous to the word search puzzle. The sound notification acts as an interrupt, grabbing the participant’s attention as necessary which allows them to maximize the time they spend on the puzzle. With the purely visual interface, the participant must look up at the GUI occasionally to check for errors, which may disrupt their focus more frequently and result in them negatively perceiving the system and performing worse on the secondary task.

During the control experiment with no sound prompts, the participant gets to choose when to look up and address robot failures. One could argue that with sound notification a participant’s train of thought is interrupted, so looking up on their own might improve their performance. However, with sound notification, the participant can fully address a robot’s issue without ever looking up at the screen. Thus even if they lose their thought process, their eyes may still be on the word search and they still may be able to recognize words while addressing the failures.

**H2: Single phrase communication provides the best user perceived success, usability, and capability of the system.**

Speech capability improves the perception of social ability in a robot [39]. However, the principle of least collaborative effort indicates that minimal effort is socially perceived the best [37] in task-oriented HRI dialogue. Neither of these findings have been tested on a remote HRI task. We hypothesize that the phrase system provides the most intuitive balance between a socially well perceived yet efficient system.

**H3: Single phrase communication will result in the most significant improvement in productivity when a user is completing a secondary task.**

We hypothesize that in the control center scenario the participant will not want to be interrupted from the word-search with a lengthy description of the problem nor hold a conversation with the system. However, they may appreciate more easily interpretable feedback than an earcon which they have to map to a robot color. Thus, the single-word communication system will likely provide the best balance between the two extremes.

## IV. RESULTS

The survey results are shown in Table V. As predicted in H1, the system with no sound performs the worst in every metric. As the observations are independent and errors can be assumed to be normally distributed, we ran ANOVA and paired T-tests to verify the results. Using the type of audio system as the treatment groups, ANOVA was performed on the results of each question, with the following findings: Q1 ( $p \ll 0.001$ ), Q2 ( $p = 0.03$ ), and Q3 ( $p \ll 0.001$ ). Each result is significant with  $\alpha = 0.05$  which supports H1.

To verify H2, a visual comparison of the results in Figure 2 show the mean scores and standard error for each question of each system. The survey results show that the user most positively perceives the system prompting them with phrases. To test if these results are significant, paired T-tests were performed on each question of the survey results between the phrase system and the two other audio systems. These results are statistically significant ( $\alpha = 0.05$ ) as shown by the p-values in Table VI. Looking at the rankings provided by the user (Table V), phrase was ranked better than each of the other methods, which indicates that the user had the most positive impression of the phrase system. H2 is supported by the data.

**TABLE V:** Average productivity score (words per minute) and survey responses

Audio Type	Prod. Score	Q1	Q2	Q3	Rank
Earcon	1.49	5.54	6.15	2.62	2.23
Phrase	<b>1.73</b>	<b>6.54</b>	<b>6.77</b>	<b>1.85</b>	<b>1.38</b>
Sentence	1.56	5.85	5.85	3.00	2.62
No Sound	1.44	3.85	5.46	4.62	3.77

**TABLE VI:** P-values of Paired T-tests

Question	Phrase and Earcon	Phrase and Sentence
Q1	0.003	0.041
Q2	0.013	0.008
Q3	0.013	0.014

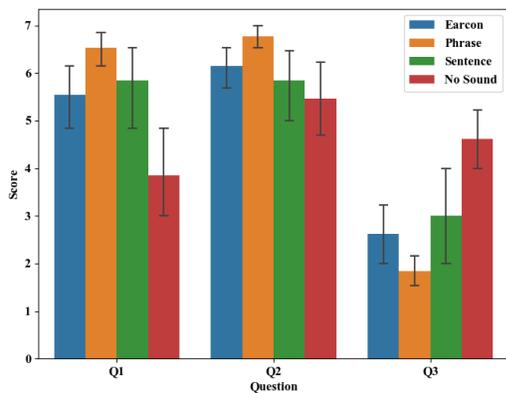


Fig. 2: Mean survey responses

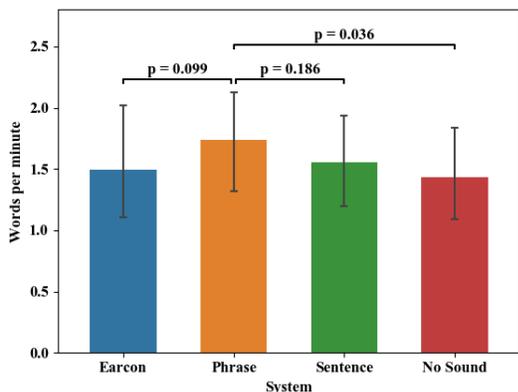


Fig. 3: Mean word search score and p-values of paired T-tests

The results of the secondary task are aggregated in Table V, and the mean of the phrase system performs slightly higher than the rest. However, ANOVA was performed across the different auditory interaction systems and indicated that the results are not statistically significant ( $p = 0.78$ ), thus H3 is not supported. We further performed paired T-tests on the productivity scores with the phrase system and every other system and display the significance values in Figure 3.

## V. DISCUSSION

### A. Explanation of Results

The data supporting H2 indicates the phrase system most positively affected people’s perception of the robots, meaning from a user-centered standpoint it is best suited for completing a human robot teleoperation task, such as fixing failures. Phrase notifications were perceived as easiest to use (Q1), most successful (Q2), and least overwhelming (Q3) and was ranked the best. One reason why sentence performed worse than phrase and earcon could be that the user overperceived the system as capable of more than it actually is since it is closer to full natural language processing. These unrealistic expectations would increase the chances of an HRI task failure as the human may try to have the system do more than it is capable of, as discovered in [39]. For example, in this study the participant often tried to interrupt the system in the middle of the sentence and was forced to repeat when it failed to register the command, making the participant frustrated. In general, imagine how many ways a human operator with

minimal robotics experience may attempt to use a system when it seems capable of full natural language processing. Intuitively, it also follows that in a task-oriented remote monitoring setting the robot does not need to have much conversational or social dialogue to be positively perceived by the participant.

Another explanation of why the phrase system scored the best is every other system may have a larger cognitive load on the participant. The earcon system required the participant to learn the mapping between earcons and robot colors during the experiment. Instead of looking at the GUI to see which robot failed, the participant usually would look at the Earcon Table IV to determine which robot failed. After the first few failures, the participant remembered the mappings and did not have to look away from the wordsearch to fix robot failures. The sentence system required the largest cognitive load out of the three audio systems. One participant reported they would “never want a system to speak in full sentences as it is highly frustrating.” Communicating one small piece of information took the sentence system a few seconds longer than necessary. Such behavior distracted the participant from the wordsearch task and often led the participant to speak before the system was done speaking. When the system did not register their response, they had to speak again. Comprehending a full sentence and being cognizant of when the system finishes the sentence increased the cognitive load on the participant. Finally, the no sound system involved the largest cognitive load overall, since the participant always had to remember to look up from the wordsearch to address the failures. During the audio experiments, the participant rarely took their eyes off the word search, which highlights how audio can aid an HRI interface and reiterates H1.

The participants preferred the earcon system over the full sentence system. One participant found “it was helpful to have an association between the sounds and the robots that needed fixing.” However, all comments about the sentence system given by the participants were negative. Furthermore, if we remove the worst performing participants (the outliers) from our analysis, the earcon system performs better than the sentence system in all measures. This indicates that earcon system may be a slightly better design than the sentence system from the user’s perspective, however this wasn’t statistically analyzed as the phrase methodology was much more positively regarded than both the other methods.

The larger cognitive load amongst the no sound, earcon, and full sentence conditions is most likely why the single phrase mechanism was the best. This cognitive load may be the effort grounding dialogue literature refers to in the principle of least collaborative effort [37], [38]. Our study extends this principle to remote robot HRI for multitasking, and shows communicating the most information in the least amount of sound is a good design choice. However, too little sound like the earcons requires the participant to still do some work to extract the information themselves so the system must take that into account.

## B. Limitations

To our surprise, the devised secondary task and its related metric did not produce significant results, leaving H3 unsupported. We hypothesized that the phrase system will improve the user's capability and increase productivity. Although H3 was not statistically supported, on average the phrase system provides a higher productivity (Figure 3), which can inform future studies investigating multitasking using audio and visual tasks. One simple explanation for the lack of significance is the small sample size that was tested. Another explanation is the imperfect speech recognition system occasionally mishearing what the human said. Such mistakes required the participant to say the command repeatedly until the robot understood them, and was often exacerbated by participants who's first language was not English. Having to repeat commands likely decreased the participants productivity enough to negate any significant increase in productivity across the sound systems. Future studies may consider having participants who are farmers for a more realistic participant group. Additionally, they may use a Wizard of Oz approach to prevent technical imperfections from influencing the results. However, this adjustment is less realistic as misunderstandings are inevitable in modern systems especially in noisy environments.

## C. Broader Impact

The support for H1 shows audio can make a significant difference in a user's perceived usability and success with the system when multitasking, and H2 proves that concise phrases is better than detailed verbal descriptions or simple sound cues. Thus, human-robot monitoring interfaces can benefit from audio over visual communication even if it cannot develop fully autonomous conversational AI systems. We study the agriculture setting; however, the design can be extended to other applications. Remote operation is also being used in the healthcare [52] and manufacturing [53]. Improving usability of robot systems allows robots to become more accessible and widespread. Enabling more people to manage multiple robots while still being able to focus on other tasks, will improve the overall productivity and quality of society.

## VI. CONCLUSION

In this paper, we studied how remote supervision of an agricultural robot through speech and audio signals affects a user's perception of the system and productivity in a secondary task. Understanding this relationship would allow robotics research to focus on developing systems that match a user's perception to improve the overall human-robot interaction. The results indicate that the average user is most likely to find single phrase notifications of a multirobot remote operation interface the easiest to interact with.

Although productivity was not improved significantly by the phrase system, it did improve participant's perception. This result indicates that the participant's perception in this remote robot management scenario had a more noticeable effect on informing the design of the system than the

participant's productivity with the system, which underscores the necessity of human-centered design. We found that participants preferred a system that communicated with simple phrases, followed by simple sound communication. Full sentence communication was ranked last by participants. These findings match our intuition that full natural language understanding capabilities may not be needed in remote robot monitoring systems as people often found them more annoying than helpful.

Understanding the optimal interface and communication mechanisms for agricultural robots will help design technology and robots that are easy to use by non-experts. What seems intuitive to an academic may not be as easy to use for a person with minimal robotics experience. Thus, the applied survey methodology is ideal for alleviating different experience levels and assuring the accessibility of technology to all types of users. By improving the design of such fleet management systems, we may increase the likelihood of adoption, paving the way for agbots to be deployed at scale.

## REFERENCES

- [1] California Farm Bureau Federation, "Survey: California farms face continuing employee shortages," <https://www.cfbf.com/news/survey-california-farms-face-continuing-employee-shortages/>, 2019, accessed: Feb. 17, 2020.
- [2] J. A. Foley, N. Ramankutty, K. A. Brauman, E. S. Cassidy, J. S. Gerber, M. Johnston, N. D. Mueller, C. O'Connell, D. K. Ray, P. C. West, C. Balzer, E. M. Bennett, S. R. Carpenter, J. Hill, C. Monfreda, S. Polasky, J. Rockstram, J. Sheehan, S. Siebert, D. Tilman, and D. P. M. Zaks, "Solutions for a cultivated planet," *Nature*, vol. 478, no. 7369, pp. 337–342, 2011.
- [3] H. C. J. Godfray, J. R. Beddington, I. R. Crute, L. Haddad, D. Lawrence, J. F. Muir, J. Pretty, S. Robinson, S. M. Thomas, and C. Toulmin, "Food security: The challenge of feeding 9 billion people," *Science*, vol. 327, no. 5967, pp. 812–818, 2010.
- [4] B. V. Ortiz, K. Balkcom, L. Duzy, E. Van Santen, and D. Hartzog, "Evaluation of agronomic and economic benefits of using rtk-gps-based auto-steer guidance systems for peanut digging operations," *Precision agriculture*, vol. 14, no. 4, pp. 357–375, 2013.
- [5] S. Yaghoubi, N. A. Akbarzadeh, S. S. Bazargani, S. S. Bazargani, M. Bamizan, and M. I. Asl, "Autonomous robots for agricultural tasks and farm assignment and future trends in agro robots," *International Journal of Mechanical and Mechatronics Engineering*, vol. 13, no. 3, pp. 1–6, 2013.
- [6] S. M. Pedersen, S. Fountas, H. Have, and B. Blackmore, "Agricultural robots: system analysis and economic feasibility," *Precision agriculture*, vol. 7, no. 4, pp. 295–308, 2006.
- [7] United Nations, Department of Economic and Social Affairs Economic Analysis, Frontier Technology Branch, "Frontier Technology Issues: Frontier technologies for smallholder farmers: addressing information asymmetries and deficiencies," 2021, accessed: March 14, 2022.
- [8] D. Vasisht, Z. Kapetanovic, J. Won, X. Jin, R. Chandra, S. Sinha, A. Kapoor, M. Sudarshan, and S. Stratman, "Farmbeats: An IoT platform for data-driven agriculture," in *14th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 17)*, 2017, pp. 515–529.
- [9] S. Honig and T. Oron-Gilad, "Understanding and resolving failures in human-robot interaction: Literature review and model development," *Frontiers in psychology*, vol. 9, p. 861, 2018.
- [10] M. Luria, J. Zimmerman, and J. Forlizzi, "Championing research through design in hri," *arXiv preprint arXiv:1908.07572*, 2019.
- [11] F. A. Robinson, M. Velonaki, and O. Bown, "Smooth operator: Tuning robot perception through artificial movement sound," in *Proceedings of the 2021 ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, 2021, pp. 53–62.
- [12] M. L. Lupetti, C. Zaga, and N. Cila, "Designerly ways of knowing in hri: broadening the scope of design-oriented hri through the concept of intermediate-level knowledge," in *Proceedings of the ACM/IEEE*

- International Conference on Human-Robot Interaction (HRI)*, 2021, pp. 389–398.
- [13] H. Choi, C. Crump, C. Duriez, A. Elmquist, G. Hager, D. Han, F. Hearl, J. Hodgins, A. Jain, F. Leve *et al.*, “On the use of simulation in robotics: Opportunities, challenges, and suggestions for moving forward,” *Proceedings of the National Academy of Sciences*, vol. 118, no. 1, 2021.
  - [14] S. Lemaignan, M. Hanheide, M. Karg, H. Khambhaita, L. Kunze, F. Lier, I. Lütkebohle, and G. Milliez, “Simulation and hri recent perspectives with the morse simulator,” in *International Conference on Simulation, Modeling, and Programming for Autonomous Robots*. Springer, 2014, pp. 13–24.
  - [15] T. Ji, S. T. Vuppala, G. Chowdhary, and K. Driggs-Campbell, “Multi-modal anomaly detection for unstructured and uncertain environments,” *arXiv preprint arXiv:2012.08637*, 2020.
  - [16] J. P. Vasconez, G. A. Kantor, and F. A. A. Cheein, “Human–robot interaction in agriculture: A survey and current challenges,” *Biosystems engineering*, vol. 179, pp. 35–48, 2019.
  - [17] A. Rodríguez, A. Fernández, and J. H. Hormazábal, “Beyond the gui in agriculture: a bibliographic review, challenges and opportunities,” in *Proceedings of the International Conference on Human Computer Interaction (HCI)*, 2018, pp. 1–8.
  - [18] T. Ji, A. N. Sivakumar, G. Chowdhary, and K. Driggs-Campbell, “Proactive anomaly detection for robot navigation with multi-sensor fusion,” *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4975–4982, 2022.
  - [19] M. B. Dias, B. Kannan, B. Browning, E. Jones, B. Argall, M. F. Dias, M. Zinck, M. M. Veloso, and A. Stentz, “Sliding autonomy for peer-to-peer human-robot teams,” in *Proceedings of the international conference on intelligent autonomous systems*, 2008, pp. 332–341.
  - [20] G. Chowdhary, C. Soman, and K. Driggs-Campbell, “Levels of autonomy for field robots,” <https://www.earthsense.co/news/2020/7/24/levels-of-autonomy-for-field-robots>, 2020.
  - [21] G. Swamy, S. Reddy, S. Levine, and A. D. Dragan, “Scaled autonomy: Enabling human operators to control robot fleets,” in *IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 5942–5948.
  - [22] S. Reig, E. J. Carter, T. Fong, J. Forlizzi, and A. Steinfeld, “Flailing, hailing, prevailing: Perceptions of multi-robot failure recovery strategies,” in *Proceedings of the ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, 2021, pp. 158–167.
  - [23] S. C. Peres, V. Best, D. Brock, C. Frauenberger, T. Hermann, J. G. Neuhoff, L. Nickerson, B. Shinn-Cunningham, and A. Stockman, “Auditory interfaces,” *HCI beyond the GUI: design for haptic, speech, olfactory, and other nontraditional interfaces*, pp. 147–195, 2008.
  - [24] S. Yilmazyildiz, R. Read, T. Belpaeme, and W. Verhelst, “Review of semantic-free utterances in social human–robot interaction,” *International Journal of Human-Computer Interaction*, vol. 32, no. 1, pp. 63–85, 2016.
  - [25] J. J. Steil, F. Röthling, R. Haschke, and H. Ritter, “Learning issues in a multi-modal robot-instruction scenario,” in *Workshop on Imitation Learning, Proc. IROS*. Citeseer, 2003.
  - [26] D. B. Kaber, M. C. Wright, and M. A. Sheik-Nainar, “Investigation of multi-modal interface features for adaptive automation of a human–robot system,” *International journal of human-computer studies*, vol. 64, no. 6, pp. 527–540, 2006.
  - [27] M. Marge and A. I. Rudnicky, “Miscommunication detection and recovery in situated human–robot dialogue,” *ACM Transactions on Interactive Intelligent Systems (TiIS)*, vol. 9, no. 1, pp. 1–40, 2019.
  - [28] C. Frauenberger, T. Stockman, and M.-L. Bourguet, “A survey on common practice in designing audio in the user interface,” in *Proceedings of HCI 2007 The 21st British HCI Group Annual Conference University of Lancaster, UK 21*, 2007, pp. 1–9.
  - [29] D. K. McGookin and S. A. Brewster, “Understanding concurrent earcons: Applying auditory scene analysis principles to concurrent earcon recognition,” *ACM Transactions on Applied Perception (TAP)*, vol. 1, no. 2, pp. 130–155, 2004.
  - [30] G. Johannsen, “Auditory displays in human-machine interfaces,” *Proceedings of the IEEE*, vol. 92, no. 4, pp. 742–758, 2004.
  - [31] N. Gang, S. Sibi, R. Michon, B. Mok, C. Chafe, and W. Ju, “Don’t be alarmed: Sonifying autonomous vehicle perception to increase situation awareness,” in *Proceedings of the 10th international conference on automotive user interfaces and interactive vehicular applications*, 2018, pp. 237–246.
  - [32] G. Rosati, A. Rodà, F. Avanzini, and S. Masiero, “On the role of auditory feedback in robot-assisted movement training after stroke: review of the literature,” *Computational intelligence and neuroscience*, vol. 2013, 2013.
  - [33] B. N. Walker, A. Nance, and J. Lindsay, “Spearcons: Speech-based earcons improve navigation performance in auditory menus,” Georgia Institute of Technology, 2006.
  - [34] S. Ghosh, A. Garg, S. Sarcar, P. S. Sridhar, O. Maleyvar, and R. Kapoor, “Krishi-bharati: an interface for indian farmer,” in *Proceedings of the IEEE Students’ Technology Symposium*, 2014, pp. 259–263.
  - [35] N. Schütte, B. Mac Namee, and J. Kelleher, “Robot perception errors and human resolution strategies in situated human–robot dialogue,” *Advanced Robotics*, vol. 31, no. 5, pp. 243–257, 2017.
  - [36] K. Fischer, B. Soto, C. Pantofaru, and L. Takayama, “Initiating interactions in order to get help: Effects of social framing on people’s responses to robots’ requests for assistance,” in *the 23rd IEEE International Symposium on Robot and Human Interactive Communication*, 2014, pp. 999–1005.
  - [37] D. Kontogiorgos and H. R. Pelikan, “Towards adaptive and least-collaborative-effort social robots,” in *Companion of the 2020 ACM/IEEE International Conference on Human-Robot Interaction*, 2020, pp. 311–313.
  - [38] S. Kiesler, “Fostering common ground in human-robot interaction,” in *ROMAN 2005. IEEE International Workshop on Robot and Human Interactive Communication, 2005*. IEEE, 2005, pp. 729–734.
  - [39] E. Cha, A. D. Dragan, and S. S. Srinivasa, “Perceived robot capability,” in *the 24th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, 2015, pp. 541–548.
  - [40] R. K. Moore, “Is spoken language all-or-nothing? implications for future speech-based human-machine interaction,” in *Dialogues with Social Robots*. Springer, 2017, pp. 281–291.
  - [41] A. Powers, S. Kiesler, S. Fussell, and C. Torrey, “Comparing a computer agent with a humanoid robot,” in *Proceedings of the ACM/IEEE international conference on Human-robot interaction*, 2007, pp. 145–152.
  - [42] Z. Yan, N. Jouandeau, and A. A. Cherif, “A survey and analysis of multi-robot coordination,” *International Journal of Advanced Robotic Systems*, vol. 10, no. 12, p. 399, 2013.
  - [43] M. A. Goodrich, M. Quigley, and K. Cosenzo, “Task switching and multi-robot teams,” in *Multi-Robot Systems. From Swarms to Intelligent Automata Volume III*. Springer, 2005, pp. 185–195.
  - [44] C. M. Humphrey, C. Henk, G. Sewell, B. W. Williams, and J. A. Adams, “Assessing the scalability of a multiple robot interface,” in *2007 2nd ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 2007, pp. 239–246.
  - [45] A. Dahiya, N. Akbarzadeh, A. Mahajan, and S. L. Smith, “Scalable operator allocation for multi-robot assistance: A restless bandit approach,” *IEEE Transactions on Control of Network Systems*, 2022.
  - [46] A. Fickinger, “Multi-agent gridworld environment for openai gym,” <https://github.com/ArnaudFickinger/gym-multigrad>, 2020.
  - [47] N. Harrell, J. McKulka, C. Kremm, and B. Mitchell, “Analysis of gaze on word search puzzles.”
  - [48] K. McMahon, B. Sparrow, L. Chatman, and T. Riddle, “Driven to distraction: The impact of distracter type on unconscious decision making,” *Social Cognition*, vol. 29, no. 6, pp. 683–698, 2011.
  - [49] E. Martinson and D. Brock, “Improving human-robot interaction through adaptation to the auditory scene,” in *Proceedings of the ACM/IEEE international conference on Human-robot interaction (HRI)*, 2007, pp. 113–120.
  - [50] C. M. Carpinella, A. B. Wyman, M. A. Perez, and S. J. Stroessner, “The robotic social attributes scale (rosas) development and validation,” in *Proceedings of the ACM/IEEE International Conference on human-robot interaction (HRI)*, 2017, pp. 254–262.
  - [51] C. Bartneck, D. Kulić, E. Croft, and S. Zoghbi, “Measurement instruments for the anthropomorphism, animacy, likeability, perceived intelligence, and perceived safety of robots,” *International journal of social robotics*, vol. 1, no. 1, pp. 71–81, 2009.
  - [52] G. Yang, H. Lv, Z. Zhang, L. Yang, J. Deng, S. You, J. Du, and H. Yang, “Keep healthcare workers safe: application of teleoperated robot in isolation ward for covid-19 prevention and control,” *Chinese Journal of Mechanical Engineering*, vol. 33, no. 1, pp. 1–4, 2020.
  - [53] L. Wang, “Collaborative robot monitoring and control for enhanced sustainability,” *The International Journal of Advanced Manufacturing Technology*, vol. 81, no. 9, pp. 1433–1445, 2015.