

Comparing Perceptions of Static and Adaptive Proactive Speech Agents

JUSTIN EDWARDS, University of Oulu, Finland

PHILIP R. DOYLE, IBM Research, Ireland

HOLLY P. BRANIGAN, University of Edinburgh, United Kingdom

BENJAMIN R. COWAN, University College Dublin, Ireland

A growing literature on speech interruptions describes how people interrupt one another with speech, but these behaviours have not yet been implemented in the design of artificial agents which interrupt. Perceptions of a prototype proactive speech agent which adapts its speech to both urgency and to the difficulty of the ongoing task it interrupts are compared against perceptions of a static proactive agent which does not. The study hypothesises that adaptive proactive speech modelled on human speech interruptions will lead to partner models which consider the proactive agent as a stronger conversational partner than a static agent, and that interruptions initiated by an adaptive agent will be judged as better timed and more appropriately asked. These hypotheses are all rejected however, as quantitative analysis reveals that participants view the adaptive agent as a poorer dialogue partner than the static agent and as less appropriate in the style it interrupts. Qualitative analysis sheds light on the source of this surprising finding, as participants see the adaptive agent as less socially appropriate and as less consistent in its interactions than the static agent.

CCS Concepts: • **Human-centered computing** → Empirical studies in HCI; Empirical studies in interaction design; **Natural language interfaces**; *HCI theory, concepts and models*.

Additional Key Words and Phrases: proactive agents, speech agent, speech interfaces, interruptions, partner model

ACM Reference Format:

Justin Edwards, Philip R. Doyle, Holly P. Branigan, and Benjamin R. Cowan. 2024. Comparing Perceptions of Static and Adaptive Proactive Speech Agents. In . ACM, New York, NY, USA, 18 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

As speech agents have become increasingly popular, users have highlighted multitasking during eyes-busy, hands-busy activities as a central motivation to trying these agents out [24]. That said, users' initial excitement for speech agents is frequently diminished to the point of disappointment and even abandonment, owing to speech agent interactions falling short of their expectations in terms of their abilities as dialogue partners [11, 24]. These expectations and internal models of speech agents as dialogue partners, termed *partner models* [5, 9] play a key role in how speech agent users understand their interactions. People's interactions with speech agents may therefore benefit from greater alignment between agent behaviour and the expectations users have for them as nearly human-like dialogue partners [6, 11]. Indeed, research participants have unfavourably compared speech agents [24] to human personal assistants, who they envision would be able to effortlessly help them multitask. In order for speech agents to meet these expectations, they will need to be able to interact proactively with users rather than waiting for the busy user to turn their attention to a speech interaction.

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

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Insofar as a benefit of proactive speech agents would come from an ability to interact with users who are already engaged in another task, proactive speech interactions must sometimes begin while users are attending to something. That is to say, speech from proactive agents will sometimes interrupt attention to an ongoing task. Recent work has investigated the characteristics of human spoken interruptions of this type, and found that people take many of the same considerations that are mentioned in proactive agent design guidelines into account when interrupting another person - seeking to limit the distraction caused by their interruptions by limiting the duration of their speech, attempting to select good moments for their interruptions, and sometimes preceding their interruptions with access rituals to make them more socially appropriate [15, 16]. But these human interrupting behaviours have not yet been implemented in the design of speech agents. By combining the well-established design principles for proactive agent interactions and the recent descriptions of proactive human speech interactions, the present study aims to investigate the effect that designing a speech agent to be proactive will have on people's partner models of that agent as compared to their partner models of proactive speech agents which, like existing speech agents, do not adapt speech behaviours to a user's context. As this work represents the first to manipulate adaptivity as an independent variable affecting perceptions of a proactive agent, adaptivity is manipulated broadly according to a variety of adaptive behaviours (explained in detail in section 3.3), seeking to establish evidence of an overall effect of adaptivity rather than isolating differential effects of particular adaptive behaviours.

2 RELATED WORK

2.1 Designing proactive agents

Some early work on agent based interaction sought to describe design principles for mixed-initiative agent-based interactions, sensitive to the principles which had guided the design of direct-manipulation user interfaces before them. Horvitz laid out 12 principles for mixed-initiative interfaces with this aim, including among others: considering uncertainty about a user's goals, considering the status of a user's attention in the timing of services, inferring ideal action in light of costs, benefits, and uncertainties, minimising the cost of poor guesses about action and timing, and employing socially appropriate behaviours for agent-user interaction [19]. Following from these early design principles, recent work on speech based proactive agents looked to test principles for the design of and implementation of proactive agents, with regards to details such as modality, timing, message content [7, 32, 37]. The present study considers proactive agents which use speech to interrupt a user who is already engaged in a task. As such, it is necessary to further consider the design details of those specific types of proactive interactions. One study on the design of a learning assistant with these characteristics proposed nine principles for proactive agent behaviour, specifying that it should be valuable, pertinent, competent, unobtrusive, transparent, controllable, deferent, anticipatory, and safe [37]. Echoing the general proactive agent design principles laid out by Horvitz [19], this set of principles again focuses on adapting interactions based on contextual information, including contexts of the agent's task, the user's environment, and the social context of a non-human agent initiating interaction with a person.

Recent studies have demonstrated the extent to which people consider the context of the task they interrupt, including the difficulty and urgency of the task (in those works operationally defined as the interrupter's perceived cost of disrupting a dialogue partner) when forming their interruption [15, 16]. This work found that interruptions of urgent tasks are delivered more quickly, owing to adaptations to speech rate and word choice. It likewise found that people try to interrupt during relatively less demanding moments of an ongoing task and that they may forego politeness norms when the ongoing task is more difficult [15]. Other research in this area has focused on tone in interruptions, finding

that assertive-voiced in-car notifications are less pleasant but more likely to elicit a response than a non-assertive voice [35]. The present study seeks to build on proactive agent design research by comparing user impressions of a proactive speech agent which adapts various characteristics of its speech according to the urgency and difficulty of an ongoing task against impressions of a proactive agent which ignores context and interacts in a static way.

2.2 Partner modelling of machine dialogue partners

Speech agent interactions are a unique form of human-computer interaction as they require users to engage in dialogue with a machine dialogue partner, making the conversational abilities of that partner central to the interaction [5]. Prior research on spoken interactions, both those with people and with machines, have established the concept of partner models, the models by which people understand the capabilities of their dialogue partners [5, 10]. Doyle and colleagues formally define partner models for machine dialogue partners as follows:

"The term partner model refers to an interlocutor's cognitive representation of beliefs about their dialogue partner's communicative ability. These perceptions are multidimensional and include judgements about cognitive, empathetic and/or functional capabilities of a dialogue partner. Initially informed by previous experience, assumptions and stereotypes, partner models are dynamically updated based on a dialogue partner's behaviour and/or events during dialogue" [13]

User studies of speech agent interactions have helped to establish that these partner models play a pivotal role in speech agent users' overall experience of these interactions, with users finding interactions particularly unsatisfying when their experience does not match their partner model [8, 24]. In a qualitative study of users of popular speech agents like Siri and Google Assistant, users remarked on the promise of human-likeness, insinuated by marketing, human-like voice synthesis, and designed personalities which mimic a human personality, which created what the researchers called the "gulf of expectations" [24], following the more general "gulfs of execution and evaluation" described by Norman [27]. Reflecting this research, it is critical for user experience that speech agents which prime human-like partner models to meet this expectation and deliver human-like capabilities.

Until recently, little research has explored the characteristics of partner models in speech agent interactions. Recent work has begun to investigate this question however, investigating the dimensions of partner models which are salient to people when engaged in dialogues with machines and with people [12–14]. This research was further developed into a validated self-report questionnaire, the Partner Modelling Questionnaire (PMQ), across three factors (perceptions of partner competence and dependability, assessment of human-likeness, and perceptions of the communicative flexibility of the system) which can be used to measure people's partner models for machine dialogue partners, indicating how strong of a dialogue partner a person views a given machine to be [12]. By designing a proactive speech agent which adapts to a user's context, this study aims to apply principles of proactive agent design to our current understanding of partner models in spoken interactions with machines. Specifically, this study aims to demonstrate that a proactive speech agent which adapts to a user's context is perceived as more competent, human-like, and flexible than existing speech agents which are not adaptive to context.

2.3 Aims and hypotheses

Prior research has highlighted the ways in which interruptions differ according to the urgency of the interruption and the complexity of task they interrupt. Holistically, these studies found that people adapt their interruptions in terms of timing, word choice, prosody, and the use of particular social markers (i.e. access rituals), taking urgency and task difficulty cues into account [15, 16]. This study aims to apply these findings to proactive non-human speech agents.

Building upon those findings as well as research on the design of proactive agents [19] and partner modelling [12], the present study hypothesises the following:

- People will rate speech interruptions from an adaptive agent as coming at better moments as compared to interruptions from a static (non-adaptive) agent (H1)
- People will rate speech interruptions from an adaptive agent as more appropriately asked as compared interruptions from a static (non-adaptive) agent. (H2),
- People will view an adaptive agent as a stronger dialogue partner than their partner models for a static (non-adaptive) agent (H3).

All hypotheses, research questions, and post-hoc analyses were pre-registered before data collection began.¹

As in recent studies of interruptions of complex tasks [15, 16], the study uses Tetris as an ongoing, complex task which a proactive agent must interrupt with speech. Insofar as those studies provided a specific description of the ways people adapt their speech to contextual information in Tetris, the same task is used here so that the interruption adaptations observed in the human participants from those studies can be directly applied to the design of the proactive agent in this study.

Rather than having participants act as Tetris players in this experiment, participants instead watched videos of interactions between the an agent and an unseen Tetris player. This video study technique is commonly used in human-robot interaction research [18, 23] owing to its benefits in being rapidly deployable to many participants including online participants, greater standardised control over the interaction, and facilitation of the use of early-stage prototypes which may lack features necessary for live interactions [36].

3 EXPERIMENTAL METHOD

3.1 Participants

80 crowdworkers (40 men, 40 women; M age = 38.4 years, SD = 11.9 years) were recruited on Prolific Academic. All participants were native speakers of English living in Ireland or the United Kingdom. 92.5% (N=74) of participants reported having used speech assistants before, with 66.3% (N=53) of participants reporting that they use a speech assistant once a week or more frequently. Participants were all familiar with Tetris, though most reported that they do not play frequently (81.3% of participants answering 3 or lower on a 7-point Likert-type question asking “If you have played Tetris before, how often do you play Tetris?”). The study took approximately 20 minutes and participants were compensated £6 through Prolific Academic for their participation. The study received ethical approval through the university’s ethics procedures for low-risk projects (Ethics code: [anonymized for review]).

3.2 Materials

Twenty-four videos of the game Tetris were created. In each video, a game of Tetris is played by an unseen player. Tetris videos were designed to match those described in prior work [15, 16], which differentiated easy and hard games of Tetris to study how urgency affected people’s interruptions in each context. Tetris videos from the present experiment were therefore either examples of easy games or hard games. Under the Tetris video, the word Urgent or Non-urgent appeared, indicating whether the video represented an urgent game of Tetris or a non-urgent game of Tetris, operationally defined as the interrupter’s perceived cost of interrupting that game. In particular, participants in those studies were told that a Tetris player would rate their appropriateness and timing of their interruption, that on urgent trials, these ratings would

¹osf.io/g8zk6/?view_only=ec53ef395bd64ff3a13dae10e94775bb

contribute much more robustly to the interrupter's total score, and that their total score determined their chance at a cash bonus prize [15, 16]. In this way, urgent trials both in prior studies and the present research could be considered something more like "safety-critical" contexts, rather than something like "time-sensitive" contexts.

A question was written at the bottom of each video, indicating the question that the proactive speech agent would be prompted to ask the Tetris player. After a fixed interval of 10 seconds, a large red dot indicator appeared in the video indicating that a proactive agent had been prompted to interrupt the player. The ten second delay was selected to give participants time to observe the Tetris game before an interruption might occur and reflects the maximum delay used in prior studies before prompting interruptions [15, 16]. After some delay (described below), a synthesised voice asked the question to the Tetris player. Videos end one second after the audio ends, with each video lasting approximately 20 seconds. 12 unique Tetris gameplay videos and prompts were sampled from [15], resulting in 12 matched trials across two within-subjects conditions.

Half of the videos selected ($n = 6$) were randomly sampled from the difficult Tetris games in [15], in which videos started with a Tetris game piece at the top of the game board, at least half of the rows of the board which already contained Tetris pieces, and the falling speed of the game piece was set to 10 rows per second. The other half of the videos selected ($n = 6$) were randomly sampled from the easy Tetris games in [15], in which videos started with a Tetris game piece at the top of the game board, at least two rows and no more than half of the rows of the board already contained Tetris pieces, and the falling speed of the game piece was set to the game minimum of 1.25 rows per second.

For each difficulty grouping per condition, half of the videos ($n = 3$) are arbitrarily marked "urgent" and the other half ($n = 3$) are arbitrarily marked "non-urgent". Tetris gameplay, interruption prompts, and urgency are all fixed throughout each block of videos and across participants - for example, Tetris gameplay video 2, which came from the easy condition from [15], was urgent and had the prompt "What was the last movie you watched?" in both blocks for all participants. This controlled against unsystematic interactions between these variables, ensuring that differences experienced by participants are only the conditional differences described below. In keeping with prior studies, the content of a prompt is unrelated to both the urgency of the trial and the content of the Tetris task [15, 16].

3.3 Experimental conditions- Adaptive vs Static Proactive Agent

The experiment followed a one-way within-subjects design. Agent condition was manipulated across two conditions with respect to the timing of their interruptions, their speech rate, and the content of their interrupting utterance, all of which were fixed for the static agent and adaptive to context (with contexts varying in terms of urgency and game difficulty) for the adaptive agent. Details of the differences between agents are explained below.

3.3.1 Adaptive agent. Previous work defined interruptible windows of Tetris games by classifying the characteristics of moments that viewers of Tetris games judged to be suitable for interruptions, in order to broadly label a variety of other Tetris games [15]. The adaptive agent varied its interruption onset and always began interruptions during one of these interruptible windows - typically moments with little active input from the player. Because prior work on human interruption adaptation manipulated task urgency and task difficulty as independent variables [15, 16], the adaptations of the agent are in relation to the way that interruptions of urgent tasks differ from interruptions of non-urgent tasks and the way interruptions of easy tasks differ from the interruption of hard tasks.

For urgent interruptions, the adaptive agent interrupted at an onset three seconds after the red dot appeared, or as close as possible to three seconds while interrupting within an interruptible window. For non-urgent interruptions, the adaptive agent interrupted at an onset five seconds after the red dot appeared, or as close as possible to five seconds

while interrupting within an interruptible window. The differences in interruption onsets was selected to reflect the difference in mean onsets observed in [15], in which urgent interruptions came after a mean onset of 3.87 seconds whereas non-urgent interruptions came at a mean onset of 4.72 seconds. This difference was slightly exaggerated in the conditions presented in this experiment with the intention of making differences more salient to an observer. The use of interruptible windows by only the adaptive agent across all trials reflected prior findings in which participants did not significantly vary their usage of these windows by urgency or by Tetris difficulty conditions [15].

The adaptive agent spoke at a 1.00 speech rate for non-urgent interruptions and a 1.10 speech rate for urgent interruptions, reflecting the difference in mean interruption durations observed in [15], in which urgent interruptions lasted for a mean of 1519ms whereas non-urgent interruptions lasted for a mean of 1596ms. For all six of the trials which featured easy Tetris games, the adaptive agent used access rituals such as “hey” and “excuse me” to lead into the interruption. For the six trials which featured hard Tetris games, access rituals were not used, reflecting the difference observed in the use of access rituals in [15], in which participants were significantly more likely to use access rituals during easy Tetris games.

Finally, the adaptive agent rephrased all of its questions, using concise language (e.g. “Got any pets?” for the prompt “Do you have any cats or dogs?”) or conversational language (e.g. “Are you right or left-handed?” for the prompt “Which hand do you write with?”) for all trials, with these styles balanced across difficulty and urgency contexts. Each rephrased question was a verbatim recreation of the way a participant phrased the corresponding interruption from [15]. Concise language and conversational language were selected to reflect the two major rephrasing strategies mentioned in qualitative data in prior works, which found neither of phrasing styles to be exclusively associated with a single urgency condition [15, 17]. All interruption audio was synthesised using Google WaveNet text to speech. Half of the participants heard voice en-GB-Wavenet-A, a feminine voice, and the other half of participants heard voice en-GB-Wavenet-B, a masculine voice, fully balanced by participant gender.

3.3.2 Static agent. For the static agent, interruptions always began 4 seconds after the red dot appeared, or as close as possible while ensuring the interruption did not begin within an interruptible window as identified in [15]. The static agent asked questions exactly as they appeared on screen, with no changes to wording, at the standard 1.00 WaveNet speech rate, and without the use of any access rituals.

Overall, the static agent condition is meant to be representative of the capabilities of current speech agents like Google Assistant or Amazon Alexa, which do not use access rituals, vary speech rates, or vary the timing of their speech based on contextual cues. The adaptive agent was designed to adapt its speech in a variety of ways representative of the ways people were observed to adapt their speech in experimental studies. While this experimental design does not allow for the analysis of any particular type of adaptation’s causal relationship with interaction outcomes, it nonetheless gives a holistic representation of the overall effect of adaptation, directly informed by the approaches that people use to adapt speech for interruptions.

3.4 Measures

3.4.1 Partner Model Questionnaire. Participants were asked to complete the 18-item Partner Model Questionnaire (PMQ). The PMQ is a validated self-report scale consisting of word pairs separated by a 7-point semantic differential scale. The scale comprises three subscales: *partner competence and dependability*, *human-likeness*, and *communicative flexibility* onto which nine, six, and three items load respectively [12]. Scores are calculated for each subscale by summing semantic differential ratings for each word pair that loads onto the respective scale, with higher numbers

corresponding to responses closer to the word more positively associated with that subscale (e.g. closer to the word "consistent" in the pair "consistent/inconsistent" which loads onto the *partner competence and dependability* subscale). Total PMQ scores are calculated by summing the three component subscale scores.

Participants were asked to complete the PMQ with the instructions "Thinking about the speech assistant you just watched, how would you rate its communicative ability on a scale between each of the following poles?". As a control, before the experiment began, participants were also asked to complete the PMQ in regards to the speech interface they are most familiar with. PMQ semantic differential item orders were randomised between participants, and 9 items were presented in reverse order (e.g. lower-scoring poles appeared on the left of the screen rather than the right) per participant, balanced across subscales. Items presented in reverse order were then reverse-scored so that presentation order did not affect scoring.

3.4.2 Single Item Questionnaires. After each trial, participants were asked to answer on 5-point Likert-type scales how much they agreed with each of two statements: "The question came at a good moment" and "The assistant asked the question in an appropriate way." These items mirrored the themes described in [15], timing and delivery, which participants identified as important features of the structure of a spoken interruption.

3.4.3 Demographic Questionnaire. Participants were asked a number of questions about themselves including their age, nationality, level of expertise with Tetris, how recently they played Tetris, their level of experience with speech agents, and which speech agents they use.

3.5 Procedure

Participants were directed to a webpage where they read an information sheet describing the study and the data rights of participants. They were then asked to indicate their consent to participating in the experiment and sharing their anonymised data. Participants were told that they would watch 12 short videos of a person playing Tetris, during which the Tetris player would be interrupted by a proactive speech agent asking them a question. Participants were told that after each video, they would be asked to answer 2 questions about the interruption that they just watched and, after all 12 videos, they would be asked to complete a 18 item questionnaire about the agent they just listened to. The informational screens explained that after completing this routine once with one agent, they would then be asked to do the same again with a different agent. Participants were told that each agent was engaged in an exercise in which it needed to ask the Tetris player a variety of questions, and its goal was to minimise disruption to the Tetris player while asking its set of questions as quickly as possible. Informational screens explained that for some trials, minimising disruption to the player is urgent as the Tetris player was rated on their play during the game shown, rated games were used to choose the winner of a cash prize, and that the Tetris player did not know which games were rated.

After this information was presented and the participants consented to take part in the study, they were asked to complete an initial PMQ questionnaire to get a baseline understanding of their views of speech agents in general. After the initial PMQ, participants saw example screenshots from a video. In the example screenshots, one image displayed a game of Tetris with a question prompt and the other screenshot displayed the same game of Tetris and question prompt with a large red dot overlaid above the Tetris game board (Figure 1). Participants were told that this visual indicator is not visible to the Tetris player, but it indicates to the observer (i.e. to the participant) that the agent has been prompted to ask a question to the Tetris player. Participants were informed that the agent could see the Tetris game and could decide when to begin its interruption any time after the interruption was prompted. After the participant viewed the

example screenshots, they were shown a practice video in which a Tetris game is played by an unseen player, the visual indicator appears after some time, and a synthesised voice asks the Tetris player a question.

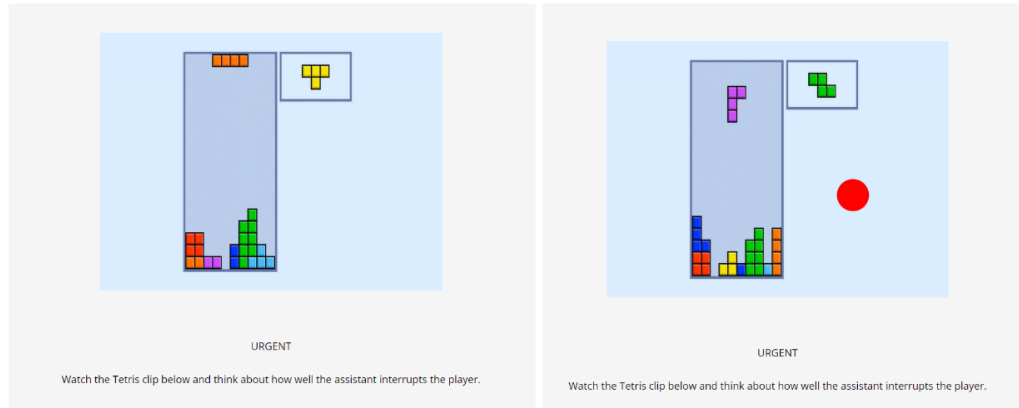


Fig. 1. Example screenshots from the experiment which participants saw as part of pre-test instructions. On the left, there is no red dot, so the agent has not yet been cued to interrupt. On the right, the red dot has appeared, signalling that the agent has been cued to interrupt.

After the participant watched the practice video, they were asked to click a button to indicate that they were ready to continue and begin their first block of trials. Blocks of trials contained 12 videos of a single agent condition, with condition order counterbalanced across participants. Within a block of 12 trials, the order of videos was randomised for each participant. Following each video, participants rated how much they agreed with each of the following statements on a 5-point Likert scale: “The question came at a good moment” and “The question was asked in a disruptive way”. After each trial (video and Likert items), a plain white screen with a black central fixation cross appeared for a short interval before the next trial began. Each trial lasted between 10 and 20 seconds.

After completing the first block of trials, participants again completed an online version of the PMQ on a single webpage. After completing the PMQ, participants were asked to confirm that they were ready for the second block of 12 trials by clicking the continue button. After completing their second block of trials, participants completed another PMQ, reflecting the agent they just saw.

After completing the final PMQ, participants were asked to complete a short demographic questionnaire. Participants were then thanked for their participation, given an opportunity to submit any other questions or comments, and debriefed on the aims of this study, including letting them know which block of trials was adaptive and which was the static agent. Finally, participants were given information for receiving payment. The full source code and materials for the experiment is provided².

4 RESULTS

4.1 Analysis approach

A total of 1920 interruption trials were viewed across the experiment by 80 participants, with participants responding to single-item questionnaires after each trial and to the Partner Model Questionnaire before the experiment and after each

²osf.io/g8zk6/?view_only=ec53ef395bd64ff3a13dae10e94775bb

Table 1. Table of means and standard deviations for single-item questionnaire responses by condition

Measure	Condition	Mean	SD
Timing	Static	2.33	1.13
	Adaptive	2.22	1.18
Appropriateness	Static	2.75	0.93
	Adaptive	2.21	1.18

of the two agent conditions. Therefore, 1920 single-item questionnaire responses and 240 PMQ responses were recorded across all participants. No data needed to be removed for technical issues or by participant request. For PMQ responses, total and subscale scores within each condition were assessed for extreme values (± 3 standard deviations from the condition means) and none were detected. For each single item questionnaire, condition means were calculated for each participant and condition means were assessed across participants for extreme values (± 3 standard deviations from the between-participant condition mean) and none were detected. This resulted in all 1920 responses for each single-item questionnaire and all 240 full PMQ responses being included in the final analysis.

Linear mixed-effects models were used to analyse the effect of agent condition on PMQ scores, single-item timing scores, and single-item appropriateness scores. Models were fit using the lme4 package version 1.1-26 [3] in R version 4.1.1 [29]. Because PMQ responses were not measured for each video stimulus, the model of PMQ responses fits the fixed effect of agent condition (pretest, static, and adaptive) with intercepts per participant. The model of each single item questionnaire score fits fixed effects of agent condition (static and adaptive) with random by-participant and by-item slopes and intercepts (by-item effects include effects of stimulus, condition order, and trial order). Each model therefore represents the maximal model for that variable. Note that the urgency and Tetris difficulty of a given trial are not modelled individually as each stimulus is fixed in terms of Tetris clip (and thus Tetris difficulty and urgency condition). For PMQ models which have three levels of agent condition, the adaptive condition was used the reference level as H3 predicts PMQ differences between the adaptive and static conditions (but not differences between PMQ scores for either condition and the pretest scores). To improve reproducibility, full model syntax and random effect outputs are included for each model [25]. Additional linear mixed-effects models were fit for each PMQ subscale as exploratory analysis to identify sources of differences between total PMQ scores. All analyses were preregistered before data collection began³.

4.2 Quantitative response data

4.2.1 Single-item questionnaires. Timing: For the first single-item questionnaire, “The assistant asked the question at a good moment”, there was no significant fixed effect of agent condition on participant ratings in a 5-point Likert-type scale [Unstandardised $\beta = -16$, SE $\beta = .08$, 95% CI $-0.26, 0.00$], $p = .053$]. H1 is rejected. Full model syntax and output are included in Table 2. Means and standard deviations of single-item questionnaire responses by condition are presented in Table 1.

Appropriateness: For the second single-item questionnaire, “The assistant asked the question in an appropriate way”, there was a significant fixed effect of agent condition on participant ratings in a 5-point Likert-type scale [Unstandardised $\beta = -0.58$, SE $\beta = 0.17$, 95% CI $-0.83, 0.-0.22$], $p = .003$]. This indicates that participants rated questions asked by the static agent as being more appropriately asked than those asked by the adaptive agent. H2 is therefore rejected as the opposite result was found. This result is visualised in Figure 2. Full model syntax and output are included in Table 3.

³osf.io/g8zk6/?view_only=ec53ef395bd64ff3a13dae10e94775bb

Table 2. Summary of fixed and random effects for timing single item questionnaire - Linear mixed effects model

Model: $Timing\ rating = Agent\ Condition + (1|subjectID) + (1 + Condition|stimulus) + (1|trialOrder)$

Fixed Effect	Std β	Unstd β	SE β	t	p
Intercept	.07	2.35	.14	16.83	.001***
Adaptive Agent	-.14	-.16	.08	-2.10	.053
Random Effects					
Group	SD	Corr			
Participant (intercept)	.53				
Participant (slope)	.27	-.08			
Stimulus (intercept)	.42				
Stimulus (slope)	.18	-.47			
Trial order	.01				

Table 3. Summary of fixed and random effects for appropriateness single item questionnaire - Linear mixed effects model

Model:

$Appropriateness\ rating = Agent\ Condition + (1 + Condition|subjectID) + (1 + Condition|stimulus) + (1|conditionOrder)$

Fixed Effect	Std β	Unstd β	SE β	t	p
Intercept	.26	2.75	.14	20.33	.002***
Adaptive Agent	-.53	-.58	.17	-3.40	.003***
Random Effects					
Group	SD	Corr			
Participant (intercept)	.72				
Participant (slope)	.66	-.65			
Stimulus (intercept)	.08				
Stimulus (slope)	.51	-.54			
Condition order	.76				

4.3 Partner model questionnaire

There was a significant fixed effect of agent condition on Partner Model Questionnaire scores, with participants having significantly stronger partner models of speech agents before the experiment as compared with after interacting with the adaptive model [Unstandardised $\beta = 6.86$, SE $\beta = 2.03$, 95% CI [2.86, 10.87], $t = 3.38$, $p = .003$] and stronger partner models of the static agent as compared to the adaptive agent [Unstandardised $\beta = 7.36$, SE $\beta = 2.03$, 95% CI [3.36, 11.37], $t = 3.63$, $p = .001$]. H3 is therefore rejected as the opposite result was found, which is visualised in Figure 3. PMQ and subscale means and standard deviations by condition are presented in Table 4. Full model syntax and output are included in Table 5. There was no difference between participants' partner models of the static agent as compared with their pretest partner model of speech agents generally. This indicates that the manipulation was successful insofar as the static agent condition matched people's prior impression of speech agents.

To better understand the source of Partner Model Questionnaire differences between the agent conditions, further models were fit to compare participants' scores across each of the three subscales of the PMQ. There was a significant fixed effect of agent condition on *partner competence and dependability* subscale scores, with participants identifying speech agents as rating higher on this factor before the experiment as compared with after interacting with the adaptive model [Unstandardised $\beta = 5.38$, SE $\beta = 1.31$, 95% CI [2.79, 7.96]] and rating the static agent as stronger on this factor as



Fig. 2. Predicted values of appropriateness questionnaire ratings by condition

Table 4. Table of means and standard deviations for PMQ total score and subscale scores by condition

Scale	Condition	Mean	SD
Total PMQ	Pretest	72.2	12.0
	Static	72.6	11.1
	Adaptive	67.0	13.9
Competence & Dependability	Pretest	42.2	7.57
	Static	42.8	6.66
	Adaptive	37.3	8.96
Human-Likeness	Pretest	20.1	5.68
	Static	19.8	6.54
	Adaptive	19.9	6.95
Cognitive Flexibility	Pretest	9.90	2.74
	Static	9.93	3.29
	Adaptive	9.85	2.88

Table 5. Summary of fixed and random effects for Partner Model Questionnaire total scores - Linear mixed effects model

Model: $PMQ = Agent\ Condition + (1|subjectID)$

Fixed Effect	Std β	Unstd β	SE β	t	p
Intercept	-.30	88.40	1.37	51.86	<.001***
Pretest	.44	6.86	2.03	3.38	.003**
Static	.47	7.36	2.03	3.63	.001**
Random Effects					
Group	SD				
Participant (intercept)	8.22				

compared to the adaptive agent [Unstandardised $\beta = 7.14$, SE $\beta = 1.31$, 95% CI [4.55, 9.73]]. This result is visualised in Figure 4. There was no difference between participants' *partner competence and dependability* subscale ratings of the static agent as compared with their pretest partner model of speech agents. Full model syntax and output are included in Table 6.



Fig. 3. Predicted values of Partner Model Questionnaire total scores by condition

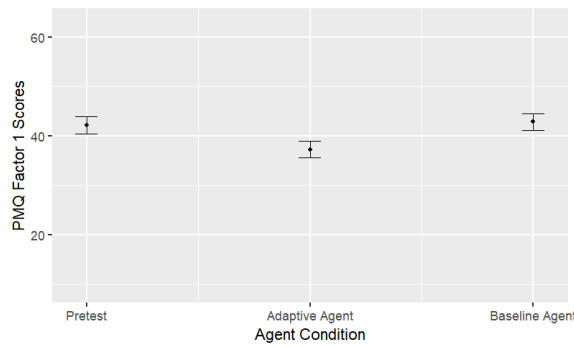


Fig. 4. Predicted values of Partner Model Questionnaire *partner competence and dependability* subscale scores by condition

Table 6. Summary of fixed and random effects for Partner Model Questionnaire *partner competence and dependability* subscale - Linear mixed effects model

Model: $PMQ F1 = Agent Condition + (1|subjectID)$

Fixed Effect	Std β	Unstd β	SE β	t	p
Intercept	-.41	51.64	1.09	47.38	<.001***
Pretest	.53	5.38	1.31	4.10	<.001***
Static	.70	5.44	1.31	5.44	<.001***

Random Effects	
Group	SD
Participant (intercept)	5.11

There were no significant fixed effects of agent conditions on either *human likeness* or *conversational flexibility* subscales, indicating that overall PMQ differences between conditions are largely explained by differences in perceived competence and dependability.

5 DISCUSSION

This study aimed to apply insights about human speech interruptions to the design of a proactive speech agent, investigating the effects of adapting speech to contexts of urgency and of ongoing task difficulty in the ways humans try to do when they interrupt. Prior work on speech agents has identified a gulf between user expectations and interaction realities [24] owing to speech agents giving cues to users that they are more capable dialogue partners than they are revealed to be through interactions [14]. The present study therefore hypothesised that participants would rate speech interruptions from an adaptive agent as coming at better moments (H1) and as more appropriately asked (H2) as compared interruptions from a static agent. It further hypothesised that participants' partner models of an adaptive agent would be rated as stronger on the Partner Model Questionnaire [12] than their partner model for the static agent (H3).

These questions were investigated through an online experiment using prerecorded interactions between agent prototypes and a Tetris player. There was no significant difference between ratings of how well the agent timed its interruptions by condition, so H1 was rejected. Interruptions from the static agent were rated as statistically significantly more appropriately asked than those from the adaptive agent, so H2 was rejected. Likewise, participants' partner models of the proactive agent were statistically significantly weaker than their partner models of the static agent or their pretest control partner model of speech agents in general, as measured by the PMQ, so H3 was rejected.

5.1 Consistency as a salient feature for adaptive agents

Contrary to expectation, the adaptive speech agent was rated lower on the PMQ by participants than was the static agent or people's pretest perception of speech agents. Post-hoc analysis revealed that differences in PMQ scores resulted from differences in perceptions of partner competence and dependability. Reflecting on the items which load onto that PMQ factor, it becomes more clear why an adaptive agent would lead to a weaker partner model across this dimension. Items such as "Dependable/Unreliable", "Consistent/Inconsistent", and "Reliable/Uncertain" [12] illustrate the importance of consistent, predictable behaviour in the formation of partner models. It may be the case that participants in this study did not have sufficient exposure to the adaptive agent to learn what contextual adaptations they could expect from the agent, leading to a poor understanding of those adaptations which cause them to seem arbitrary or inconsistent.

As commercially available speech agents are not adaptive, participants' mental models for the agents in this study would not likely lead to expectations of adaptivity. Some research on adaptive interfaces has pointed toward the benefit of explicitly describing the sorts of adaptive features that an interface has and the errors that it may cause particularly for the purpose of setting appropriate expectations [4]. Insofar as the novelty of adaptation diverges from prior experiences with speech agents, extended exposure to an adaptive agent or explicit coaching for mixed-initiative interactions may be beneficial in order to overcome perceptions of inconsistency. Returning to Horvitz's principles of mixed-initiative interface design, agent behaviour should be socially appropriate [19]. While more research is needed to determine what people consider socially appropriate interactions between speech agents and people, it may be unsurprising when interactions with novel properties like contextual adaptivity are seen as socially inappropriate when they diverge from people's prior experiences of similar interactions. Extended exposure and explicit descriptions of adaptive features should therefore be explored as ways of introducing potentially beneficial conversational features like adaptivity without introducing perceptions of inconsistency or inappropriateness.

Novel interactional elements may wane over time in the extent to which they are perceived as inconsistent. It is not clear how partner models develop over time with repeated exposure to a new partner, and the longitudinal work

required to make that determination has been identified as a challenge for over a decade [5, 10, 12]. Even if prolonged exposure to adaptive proactive agent would improve people's partner models of those agents, with better understanding of adaptive agents behaving consistently relative to particular contextual cues (rather than seeing them as inconsistent from utterance to utterance), this benefit is of little value when an early disappointment leads to abandonment of a system [11, 24]. With sensitivity toward the negative impact of novelty on partner models and the effect of poor partner models on technological abandonment, further design of adaptive proactive speech agents may need to be more incremental than this study. Introducing adaptive features piecemeal across product lines or across time spent with an agent may be less jarring to a user than interacting with a speech agent with many novel design features introduced simultaneously. In order to make progress toward adaptive speech agents, the importance of consistent interactions must be considered.

5.2 Appropriateness of adaptive proactive design

The lack of improvement on PMQ scores for the adaptive agent as compared to static agent and pretest conditions was not only surprising because of a decrease in partner competence and dependability, but also because of a lack of increase in human-likeness. Similarly, while it was hypothesised that the adaptive agent would be rated as asking questions more appropriately, the opposite was found. Each of these findings can be better understood through the lens of appropriateness in human-machine dialogue. In addressing the gulf of expectations in speech agent interactions [24], some recent work has focused on the idea of appropriateness in these interactions [1, 2, 21, 26]. This trend toward appropriateness has argued that increased human-likeness should not be a goal of itself in the design of speech agents. Instead, speech agents should be designed, in terms of voice [1, 21] and in physical appearance in the case of embodied agents [2, 26], to suit the role of the agent. Indeed qualitative work comparing human-human dialogue and human-machine dialogue has indicated that people see these two interactions as different in roles and in characteristics [28, 30], with people expressing a dislike of speech agents which try to act human-like [8, 14]. In this context, it may be clearer why participants rated the adaptive agent as no more human-like and as less appropriate in asking questions than the static agent. While adaptive interruptions may more appropriately utilise the context of an ongoing task, appropriateness also entails awareness of the social context.

5.3 Individual differences and personalisation

While neither the static nor the adaptive proactive agent were seen as significantly more human-like than participants' prior conceptions of speech agents, this may be a result of differences between participants in opposing directions rather than a lack of difference in perceptions across participants. Research on personalisation of speech agents has found a high degree of variation between people in how they would like speech agents to be designed. Some strongly prefer agents which fulfil social functions, whilst others prefer agents to only perform non-social tool-like roles [33]. Likewise, prior research comparing the roles of conversations with machines to those with humans revealed tension between some people's desires to have speech agents learn more about them to personalise interaction and others who saw speech agents building this sort of common ground with users as undesirable [8]. Tailoring speech agent design to individual users may prove especially tricky due to high variance between individuals. Recent research on how people understand the personalities of speech agents found that the popular Big Five personality types used in human personality research proved less effective for classifying machines than a more graduated model of ten personality types [34]. While more research is needed to determine differences among people's preferences for these different personality types, it is clear that the design space for machine personalities is wide and that different designs are differentiable by

the people that interact with these agents. These large individual differences in the perceptions of speech agent design decisions support the notion that personalisation of agents is both necessary and difficult.

5.4 Limitations

Individuals in this study varied not only in their preferences toward speech agents, but also in their Tetris expertise. While this study mostly involved participants with some Tetris experience (i.e. neither experts nor total novices), there is sure to be variation in skill across participants. This may impact participants' perceptions of the adaptive agent due to the differences in how expert and non-expert Tetris players perceive Tetris games [22]. Participants who have weaker understanding of Tetris gameplay and strategy were likely less sensitive to the state of the Tetris game and to the mental demand particular game states might put on the player. It may be for this reason that the adaptive agent's consistent use of interruptible windows of Tetris for initiating its interruptions went unnoticed across the study, with participants finding neither agent as significantly better at timing its interruptions. Likewise, if the adaptive agent was not perceived as better at timing its interruptions, this may further help to explain why the adaptive agent was seen as less competent and dependable than the static agent. While prior research has demonstrated that people are somewhat skilful in identifying good moments to interrupt their own discrete tasks [20], their abilities to do so for a complex task like Tetris may be much more dependent on their expertise in that task. Future research should investigate both the effect of expertise on identifying interruptible moments in complex tasks.

Another limitation of the current study is the holistic manipulation of adaptivity rather than isolating particular adaptive behaviours or contexts for adaptation. While isolating individual behavioural or contextual variables would have allowed for a more precise description of causes to changes in participant perceptions of an agent, this study was the first to look at adaptivity as an independent variable in the design of a proactive speech agent. As such, there was little theoretical basis for choosing one finding over another when considering results from prior work which demonstrated the highly varied cues and decisions people consider when producing interrupting speech [17]. This study thus presents an initial investigation into the salience and broad impact of adaptivity in this context without establishing particular causal links between particular behaviours and outcomes. Further work is needed to better understand these nuances such as the specific impact of access ritual use on perceptions of human likeness and the salience of using interruptible moments to deliver proactive speech during complex tasks. This study should be seen as an introduction to the question of how people want proactive agents to speak rather than a prescriptive set of design guidelines.

Furthermore, the PMQ has previously only been used to measure partner models formed by users of systems rather than observers who do not directly interact with the system, as in this study. As such, it may be the case that particular communicative features may be more or less salient to a person whose task was interrupted by a proactive agent as compared to an observer to the interaction. Given the level of control and more representative participant sample enabled by having participants observe an interaction rather than interacting directly, as well as the lack of adequate alternative tools for measuring partner model beliefs, we acknowledge that a limitation was introduced in the measurement of this variable. Still, insofar as hypotheses were not only rejected but directionally backward from observed results, we do not believe the magnitude of differences in the formation of partner models between an interrupted person and an observer is likely to have materially changed the findings of this work. Indeed, as partner modelling research is still in its infancy, more work is needed to understand the cost of trading the control afforded by casting participants as experimental observers for the applicability of measuring the partner models formed by people who interact with a system directly.

Finally, this study focuses solely on proactive speech interactions of a particular kind - personal questions - which interrupt a particular task - Tetris. These tasks were selected in order to maintain a high level of control over the structure of tasks and to maximally match the design of previous work in order to directly apply the findings of those studies to the design of this study. Salovaara and Oulasvirta describe the role of prototype experiments like this one as a way of evaluating possible futures [31]. In their framework, they describe design decisions as aimed at either staging - making the present have some characteristic of a possible future - or controlling - preventing particular characteristics of the present which are not expected to be part of the imagined future from becoming salient [31]. While Tetris differs in a variety of ways from the sorts of tasks that people report wanting to use speech to multitask, like cooking or driving [24], it nonetheless matches the eyes-busy, hands-busy and complex, continuous nature of each of those tasks. In this way, the Tetris task maintains the aim of experimentally staging certain features of a possible future that this study represents. That said, it is not clear that proactive agents asking questions which are irrelevant to the user's ongoing task represent a likely future use case for proactive systems. Nonetheless, the narrowing of speech tasks and stabilising of the tasks both across participants within this study serve the goal of controlling unsystematic variance, likewise helping to position this work to inform understanding of a possible future [31]. In order to explore different potential future interaction scenarios, further research will need to make carefully considered choices with regards to how to stage that future and control for unwanted present circumstances.

6 CONCLUSION

This study aimed to apply previous research insights about human spoken interruptions to the design of a proactive speech agent in order to assess people's perceptions of such an agent. Applying prior work on the design of proactive non-human agents, this study identified adaptivity as a key variable for proactive agent interactions. As prior research has identified varied contextual cues that people consider when interrupting with speech and the different modifications they make in light of those contexts, this study manipulated adaptivity by designing one agent condition around those findings with the other, static agent condition that was insensitive to context. These agent conditions were compared across three measures: the PMQ, and single-item measures of how well interruptions were timed and how appropriately they were delivered. Quantitative results revealed that participants had a stronger partner model for the static agent as compared to the adaptive agent, owing to lower ratings of the adaptive agent's competence and dependability. Participants likewise found the adaptive agent's interruptions to be less appropriate than the static agent's and detected no differences in the quality of the timing of interruptions between agents. These findings echo previous literature questioning the appropriateness of using human dialogue as a model or metaphor for nonhuman speech. While this early step toward the design of an adaptive, human-inspired, proactive speech agent revealed minimal overarching benefit to this approach, it may nonetheless serve as a guidepost for future investigation into this domain, by revealing tradeoffs encountered when using human behaviour to guide the design of speech agents.

ACKNOWLEDGMENTS

This research was conducted with the financial support of Science Foundation Ireland under Grant Agreement No. 13/RC/2106_P2 at the ADAPT SFI Research Centre at University College Dublin. ADAPT, the SFI Research Centre for AI-Driven Digital Content Technology, is funded by Science Foundation Ireland through the SFI Research Centres Programme.

REFERENCES

- [1] Matthew P. Aylett, Benjamin R. Cowan, and Leigh Clark. 2019. Siri, Echo and Performance: You have to Suffer Darling. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, Glasgow Scotland Uk, 1–10. <https://doi.org/10.1145/3290607.3310422>
- [2] Matthew P. Aylett, Selina Jeanne Sutton, and Yolanda Vazquez-Alvarez. 2019. The right kind of unnatural: designing a robot voice. In *Proceedings of the 1st International Conference on Conversational User Interfaces*. ACM, Dublin Ireland, 1–2. <https://doi.org/10.1145/3342775.3342806>
- [3] Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software* 67, 1 (2015), 1–48. <https://doi.org/10.18637/jss.v067.i01>
- [4] Matthias Beggiato and Josef F Krems. 2013. The evolution of mental model, trust and acceptance of adaptive cruise control in relation to initial information. *Transportation research part F: traffic psychology and behaviour* 18 (2013), 47–57. Publisher: Elsevier.
- [5] Holly P. Branigan, Martin J. Pickering, Jamie Pearson, Janet F. McLean, and Ash Brown. 2011. The Role of Beliefs in Lexical Alignment: Evidence from Dialogs with Humans and Computers. *Cognition* 121, 1 (Oct. 2011), 41–57. <https://doi.org/10.1016/j.cognition.2011.05.011>
- [6] Justine Cassell. 2007. Body language: Lessons from the near-human. *Genesis Redux* 346 (2007), 374.
- [7] Narae Cha, Auk Kim, Cheul Young Park, Soowon Kang, Minkyu Park, Jae-Gil Lee, Sangsu Lee, and Uichin Lee. 2020. "Hello There! Is Now a Good Time to Talk?": Opportune Moments for Proactive Interactions with Smart Speakers. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 3 (2020), 28. <https://doi.org/10.1145/3411810>
- [8] Leigh Clark, Cosmin Munteanu, Vincent Wade, Benjamin R. Cowan, Nadia Pantidi, Orla Cooney, Philip R. Doyle, Diego Garaialde, Justin Edwards, Brendan Spillane, Emer Gilmartin, and Christine Murad. 2019. What Makes a Good Conversation?: Challenges in Designing Truly Conversational Agents. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems - CHI '19*. ACM Press, Glasgow, Scotland Uk, 1–12. <https://doi.org/10.1145/3290605.3300705>
- [9] Benjamin R Cowan and Holly Branigan. 2017. They Know as Much as We Do: Knowledge Estimation and Partner Modelling of Artificial Partners. *Proceedings of CogSci '17* (2017), 6.
- [10] Benjamin R. Cowan, Holly P. Branigan, Mateo Obregón, Enas Bugis, and Russell Beale. 2015. Voice anthropomorphism, interlocutor modelling and alignment effects on syntactic choices in human-computer dialogue. *International Journal of Human-Computer Studies* 83 (Nov. 2015), 27–42. <https://doi.org/10.1016/j.ijhcs.2015.05.008>
- [11] Benjamin R Cowan, Nadia Pantidi, David Coyle, Kellie Morrissey, Peter Clarke, Sara Al-Shehri, David Earley, and Natasha Bandeira. 2017. "What can I help you with?" infrequent users' experiences of intelligent personal assistants. In *Proceedings of the 19th International Conference on Human-Computer Interaction with Mobile Devices and Services*. 1–12.
- [12] Philip R. Doyle. 2022. *The Dimensions and Adaptation of Partner Models in Human-Machine Dialogue*. PhD Thesis. University College Dublin. School of Information and Communication Studies.
- [13] Philip R. Doyle, Leigh Clark, and Benjamin R. Cowan. 2021. What Do We See in Them? Identifying Dimensions of Partner Models for Speech Interfaces Using a Psycholexical Approach. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, Yokohama Japan, 1–14. <https://doi.org/10.1145/3411764.3445206>
- [14] Philip R. Doyle, Justin Edwards, Odile Dumbleton, Leigh Clark, and Benjamin R. Cowan. 2019. Mapping Perceptions of Humanness in Intelligent Personal Assistant Interaction. In *Proceedings of the 21st International Conference on Human-Computer Interaction with Mobile Devices and Services - MobileHCI '19*. ACM Press, Taipei, Taiwan, 1–12. <https://doi.org/10.1145/3338286.3340116>
- [15] Justin Edwards. 2023. *Using Speech to Interrupt Complex Tasks Understanding Human Spoken Interruptions and Designing Interruptions for Proactive Speech Agents*. Ph. D. Dissertation. University College Dublin, Ireland.
- [16] Justin Edwards, Christian P. Janssen, Sandy J. J. Gould, and Benjamin R. Cowan. 2021. Eliciting Spoken Interruptions to Inform Proactive Speech Agent Design. *Proceedings of the 3rd Conference on Conversational User Interfaces* (2021). <https://api.semanticscholar.org/CorpusID:235352489>
- [17] Justin Edwards, He Liu, Zhou Tianyu, Sandy J. J. Gould, Leigh Clark, Philip R. Doyle, and Benjamin R Cowan. 2019. Multitasking with Alexa: How Using Intelligent Personal Assistants Impacts Language-based Primary Task Performance. In *Proceedings of the 1st International Conference on Conversational User Interfaces*. Dublin, Ireland.
- [18] Moojan Ghafurian, Gabriella Lakatos, Zhuofu Tao, and Kerstin Dautenhahn. 2020. Design and Evaluation of Affective Expressions of a Zoomorphic Robot. In *Social Robotics*, Alan R. Wagner, David Feil-Seifer, Kerstin S. Haring, Silvia Rossi, Thomas Williams, Hongsheng He, and Shuzhi Sam Ge (Eds.). Springer International Publishing, Cham, 1–12.
- [19] Eric Horvitz. 1999. Principles of mixed-initiative user interfaces. In *Proceedings of the SIGCHI conference on Human factors in computing systems the CHI is the limit - CHI '99*. ACM Press, Pittsburgh, Pennsylvania, United States, 159–166. <https://doi.org/10.1145/302979.303030>
- [20] Christian P. Janssen, Duncan P. Brumby, and Rae Garnett. 2012. Natural Break Points: The Influence of Priorities and Cognitive and Motor Cues on Dual-Task Interleaving. *Journal of Cognitive Engineering and Decision Making* 6, 1 (March 2012), 5–29. <https://doi.org/10.1177/1555343411432339>
- [21] Sébastien Le Maguer and Benjamin R. Cowan. 2021. Synthesizing a human-like voice is the easy way. In *CUI 2021 - 3rd Conference on Conversational User Interfaces*. ACM, Bilbao (online) Spain, 1–3. <https://doi.org/10.1145/3469595.3469614>
- [22] John K. Lindstedt and Wayne D. Gray. 2019. Distinguishing experts from novices by the Mind's Hand and Mind's Eye. *Cognitive Psychology* 109 (March 2019), 1–25. <https://doi.org/10.1016/j.cogpsych.2018.11.003>
- [23] Manja Lohse, Marc Hanheide, Britta Wrede, Michael L. Walters, Kheng Lee Koay, Dag Sverre Syrdal, Anders Green, Helge Huttenrauch, Kerstin Dautenhahn, Gerhard Sagerer, and Kerstin Severinson-Eklundh. 2008. Evaluating extrovert and introvert behaviour of a domestic robot - a video

- study. In *RO-MAN 2008 - The 17th IEEE International Symposium on Robot and Human Interactive Communication*. 488–493. <https://doi.org/10.1109/ROMAN.2008.4600714>
- [24] Ewa Luger and Abigail Sellen. 2016. "Like Having a Really Bad PA": The Gulf between User Expectation and Experience of Conversational Agents. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems - CHI '16*. ACM Press, Santa Clara, California, USA, 5286–5297. <https://doi.org/10.1145/2858036.2858288>
- [25] Lotte Meteyard and Robert A.I. Davies. 2020. Best practice guidance for linear mixed-effects models in psychological science. *Journal of Memory and Language* 112 (June 2020), 104092. <https://doi.org/10.1016/j.jml.2020.104092>
- [26] Roger K Moore. 2017. Appropriate Voices for Artefacts: Some Key Insights. *1st Int. Workshop on Vocal Interactivity in-and-between Humans, Animals and Robots* (2017), 5.
- [27] Donald A Norman. 1983. Some observations on mental models. In *Mental models*. Psychology Press, 15–22.
- [28] Martin Porcheron, Joel E Fischer, Stuart Reeves, and Sarah Sharples. 2018. Voice interfaces in everyday life. In *proceedings of the 2018 CHI conference on human factors in computing systems*. 1–12.
- [29] R Core Team. 2020. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>
- [30] Stuart Reeves. 2019. Conversation considered harmful?. In *Proceedings of the 1st International Conference on Conversational User Interfaces (CUI '19)*. Association for Computing Machinery, New York, NY, USA, 1–3. <https://doi.org/10.1145/3342775.3342796>
- [31] Antti Salovaara, Antti Oulasvirta, and Giulio Jacucci. 2017. Evaluation of Prototypes and the Problem of Possible Futures. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*. ACM, Denver Colorado USA, 2064–2077. <https://doi.org/10.1145/3025453.3025658>
- [32] Rob Semmens, Nikolas Martelaro, Pushyami Kaveti, Simon Stent, and Wendy Ju. 2019. Is Now A Good Time?: An Empirical Study of Vehicle-Driver Communication Timing. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems - CHI '19*. ACM Press, Glasgow, Scotland UK, 1–12. <https://doi.org/10.1145/3290605.3300867>
- [33] Sarah Theres Völkel, Penelope Kempf, and Heinrich Hussmann. 2020. Personalised Chats with Voice Assistants: The User Perspective. In *Proceedings of the 2nd Conference on Conversational User Interfaces*. ACM, Bilbao Spain, 1–4. <https://doi.org/10.1145/3405755.3406156>
- [34] Sarah Theres Völkel, Ramona Schödel, Daniel Buschek, Clemens Stachl, Verena Winterhalter, Markus Bühner, and Heinrich Hussmann. 2020. Developing a Personality Model for Speech-based Conversational Agents Using the Psycholexical Approach. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu HI USA, 1–14. <https://doi.org/10.1145/3313831.3376210>
- [35] Priscilla N. Y. Wong, Duncan P. Brumby, Harsha Vardhan Ramesh Babu, and Kota Kobayashi. 2019. Voices in Self-Driving Cars Should be Assertive to More Quickly Grab a Distracted Driver's Attention. In *Proceedings of the 11th International Conference on Automotive User Interfaces and Interactive Vehicular Applications (AutomotiveUI '19)*. Association for Computing Machinery, Utrecht, Netherlands, 165–176. <https://doi.org/10.1145/3342197.3344535>
- [36] Sarah N Woods, Michael L Walters, Kheng Lee Koay, and Kerstin Dautenhahn. 2006. Methodological Issues in HRI: A Comparison of Live and Video-Based Methods in Robot to Human Approach Direction Trials. In *ROMAN 2006 - The 15th IEEE International Symposium on Robot and Human Interactive Communication*. 51–58. <https://doi.org/10.1109/ROMAN.2006.314394>
- [37] Neil Yorke-Smith, Shahin Saadati, Karen L. Myers, and David N. Morley. 2012. The Design of a Proactive Personal Agent for Task Management. *International Journal on Artificial Intelligence Tools* 21, 01 (Feb. 2012), 1250004. <https://doi.org/10.1142/S0218213012500042>