

Understanding User Perceptions of Response Delays in Crowd-Powered Conversational Systems

TAHIR ABBAS, Eindhoven University of Technology, Netherlands

UJWAL GADIRAJU, Delft University of Technology, Netherlands

VASSILIS-JAVED KHAN, Digital Transformation, Sappi Europe, Belgium

PANOS MARKOPOULOS, Eindhoven University of Technology, Netherlands

Crowd-powered conversational systems (CPCS) are gaining considerable attention for their potential utility in a variety of application domains, for which automated conversational interfaces are still too limited. CPCS currently suffer from long response delays, which hampers their potential as conversational partners. The majority of prior work in this area has focused on demonstrating the feasibility of the approach and improving performance, while evaluation studies have primarily focused on response latency and ways to reduce it. Relatively little is currently known about how response delays in a CPCS can affect user experience. While the importance of reducing response latency is widely recognized in the broader field of human-computer interaction, little attention has been paid to how response quality, response delay, conversational context, and the complexity of the task affect how users experience the conversation, and how they perceive waiting for responses in particular. We conducted a between-subjects experiment ($N = 478$), to examine the influence of these four factors on the overall waiting experience of users. Results show that users 1) evaluated the waiting experience more negatively when the response delay was longer than 8 seconds, 2) underestimated the elapsed time but experienced more frustration in tasks with high complexity, 3) underestimated the elapsed time and experienced less frustration with high quality bot's utterances, 4) judged response delays to be slightly longer, and experienced more frustration in an emotion-centric CPCS compared to a task-centric CPCS. Our insights can inform the design of future CPCSs with regards to defining performance requirements and anticipating their potential impact on the user experience they can facilitate.

CCS Concepts: • **Human-centered computing** → **Natural language interfaces**; **Empirical studies in HCI**.

Additional Key Words and Phrases: Crowd-powered Conversational Agents, Response Latency, Perceived Response Latency, Time Perception, Waiting Tolerance, Cognition, Affect

ACM Reference Format:

Tahir Abbas, Ujwal Gadiraju, Vassilis-Javed Khan, and Panos Markopoulos. 2021. Understanding User Perceptions of Response Delays in Crowd-Powered Conversational Systems. In *CSCW '21: ACM conference on Computer Supported Cooperative and Social Computing, October 23–27, 2021, Virtual Conference*. ACM, New York, NY, USA, 42 pages.

1 INTRODUCTION

Conversational user interfaces have attracted considerable attention in human-computer interaction (HCI), already since the early days of this field, as they can be deployed quickly, they have low training costs and thereby can support a wide variety of applications, such as tutoring, companionship, and information retrieval, among others. Despite great progress, current artificial intelligence

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 ACM 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.

CSCW '21, October 23–27, 2021, Virtual Conference

© 2021 Association for Computing Machinery.

ACM ISBN XXX-1-XXXX-XXXX-X/XX/XX...\$15.00

(AI) techniques and natural language processing are not yet capable of dealing with the full complexity of free-form conversational interactions, often resulting in conversation breakdowns [9]. Crowd-Powered Conversational Systems (CPCS) [48, 69] have been proposed as a remedy to these shortcomings of AI. CPCS can be more robust than current AI and can handle interaction with users in a fluid, multi-turn conversation. CPCSs make use of sophisticated recruiting, rewarding and user interface techniques to reduce latency of crowd input from hours to a few seconds [11]. A pioneering example of a CPCS is Chorus [69], which is a text-based conversational agent that assists end-users with information retrieval tasks by conversing with an online group of workers in real time. Since such CPCSs rely (fully [69] or partially [48]) on operation by humans, significant response delays can be expected that might frustrate end-users and compromise the overall user experience. The typical workflow of CPCS is depicted in Figure 1, highlighting locations where latency may occur.

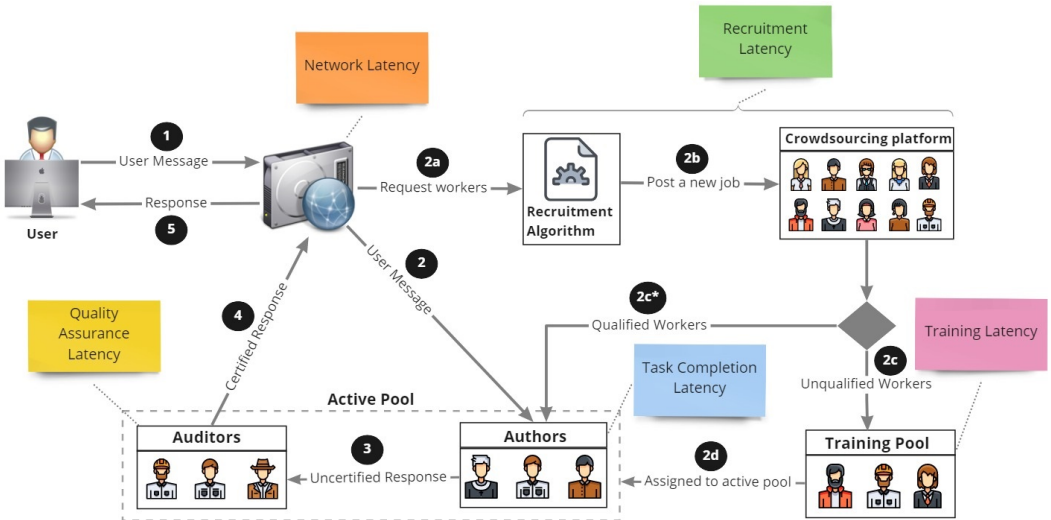


Fig. 1. This diagram depicts a typical workflow of CPCS. (1) First, user sends a message to the server. This message may take some time to reach the crowd workers depending on network bandwidth, resulting in *network latency* [67]; (2) The message is then sent to a pool of active workers; (2a-2b) if there are no workers in the active pool when user sends a message, then server requests recruitment algorithm to hire a sufficient number of workers from the crowdsourcing platform. The time it takes for an interested crowd worker to accept a newly posted task is referred to as *recruitment latency* [43]; (2c) When workers take their first task, they are often required to complete tutorials or qualification tasks before being able to perform actual tasks. This causes *training latency* [43]; (2d) After training, they are assigned to a pool of active workers; (2c*) If they are already qualified for the job, they can begin immediately, albeit this is rare. Once in the active pool, they are allocated to one of two roles: one set of workers (authors) creates responses, while another group of workers (auditors) validates those responses before they are returned to the users. The *task completion latency* is the time required for a worker to finish a task, which varies according to the workers' skill, fatigue, and a variety of other factors [52]. (3) The messages generated by one group of workers (authors) are forwarded to another group of workers (auditors) for validation, introducing *quality assurance latency*; (4) This certified message is then transmitted to the server, where it may encounter delays owing to network latency; (5) finally, the message is reached to the intended user.

Research on time perception in HCI suggests that with longer response time, users show more frustration and stress [104], perform poorly [18], demonstrate dissatisfaction with the system [27]

99 and may even be affected physiologically [98] (e.g., elevated blood pressure). Although there is
100 wide agreement concerning the negative effects of long response times on end-users, views differ
101 on what waiting times are acceptable; some researchers advocate 10 seconds as an upper limit
102 [85], while others suggest a stricter 8 seconds limit [35]. The notion of an acceptable time delay
103 refers directly to how users experience delays rather than on the objectively measurable time delay
104 between user utterances and system responses. The majority of work in CPCS has only focused on
105 reporting or reducing the “actual” or system response latency (See section 2.2 for more details), but
106 relatively little is known concerning the subjective experience of such response delays. Here we
107 examine the “perceived” latency, which is a subjective measure pertaining to the user experience
108 which pertains to the perception of the waiting time between conversational turns. Evaluating the
109 perceived waiting times in CPCS and how they influence user’s attitudes towards the system can
110 help prioritize requirements and guide the design of these systems, leading to our first research
111 question:

112 **RQ1:** *How do end users perceive the response latency of a CPCS while interacting with*
113 *it?*

114 One can expect to receive responses of mixed quality from CPCSs. This is in part due to incon-
115 gruity of workers’ skills and expertise with the task, biased interests, unfair incentives, poorly
116 designed tasks, insufficient workers’ training, among others [3, 6, 26, 61]. A high-quality response
117 is one that is appropriate to the user’s utterance, contains helpful and concrete advice and gives
118 individualized attention to the user [123]. We hypothesize that a high-quality response from CPCS
119 can positively influence the overall waiting experience. For this paper, *waiting experience* refers
120 to perceived waiting time (believing that less or more time has passed than the actual duration),
121 cognitive appraisals of the wait [70, 93, 116] (i.e., evaluations based on instrumental beliefs, such as
122 performance and usefulness), affective appraisals of the wait [70, 93, 116] (i.e., evaluations based
123 on feelings, emotions, and moods) and perception of overall interaction quality. We refer to the
124 cognitive and affective appraisals collectively as *appraisals of the wait*. Our selection of measures
125 for the appraisals of the wait is based on prior research in time perceptions that has shown an
126 inverse relationship between utilitarian evaluations (e.g. performance and lengthy waiting time
127 [18]) and emotional states (e.g. stress and lengthy waiting time) [104], and between the quality of
128 the response contents and negative affect states [125]. Thus, in the case of high-quality responses,
129 users may perceive time to pass faster, may experience less stress, and a better overall interaction
130 quality. This line of reasoning leads us to the following research question:

131 **RQ2:** *How does the quality of responses generated by a CPCS influence the waiting*
132 *experience of users?*

133 Recently, researchers have adopted machine learning to quantify task clarity [37] and task
134 complexity [124] for crowdsourcing tasks and studied their impact on the worker performance
135 [36]. In our study, we want to examine whether the complexity of a conversational task has any
136 impact on the time perceptions within the CPCS context. For instance, take an example of an
137 information retrieval based CPCS where a simple search task requires only retrieving factual
138 information from an information resource (e.g., What is the name of the deepest point in the
139 ocean?). Next to that, a more complex search task involves also organizing factual information
140 into a new pattern or generating a new plan (e.g., Help me to put together two thirty-minute
141 low impact exercise programs for my 90-year-old granny). Similarly, a stress management CPCS
142 requires answering user complaints involving complicated and multi-layered stressors and entails
143 higher cognitive complexity and time investment than a simple stressor concerning weight loss
144 or exam stress. Correspondingly, workers may need more time in generating a response for a
145 complex case and the user might expect the CPCS to take longer to respond. Arguably, this could
146
147

148 influence the overall waiting experience positively, with users underestimating the elapsed time
149 and experience comparatively less affective responses to the delay when the complexity of task
150 is high. This notion is supported by prior research that explains that increase in cognitive load
151 can decrease the perceived time delay [5, 13]. It should be emphasized that task complexity can
152 be related to both the difficulty of the task [122] and the context of the conversation [107] (see
153 next paragraph). Since one of the study's overarching aims was to gain a better understanding of
154 the influence of task complexity as a broad construct, we employed two types of tasks to reflect
155 complexity in terms of the conversational context in addition to the complexity specified at the
156 task level, as opposed to considering either of these alone. In the context of CPCSs this leads to the
157 following research question:

158
159 **RQ3:** *How do conversational tasks of varying complexity affect the waiting experience*
160 *of users in CPCS?*

161
162 Furthermore, we wanted to explore the effects of two different task types on the waiting expe-
163 rience for the following three reasons: (1) We created two distinct task types in light of a prior
164 research indicating that user expectations regarding delays may vary according to the type of
165 conversational context [89, 107]. For instance, in our recent study [1], we evaluated the usefulness
166 of several engaging interventions or waiting-time fillers in lessening the impression of time. We
167 discovered that fillers, notably those that demand users to complete a secondary activity while
168 waiting, were more helpful in Information Retrieval (IR) activities than the stress mitigation task.
169 Likewise, when prolonged delays in the service of a robot are inevitable, the research [89] found
170 that consumers in the chitchat settings were marginally less satisfied with lengthier waits than
171 users in the IR conditions. Thus, we were intrigued if this distinction remained true when we
172 enabled end users to interact with various sorts of activities in CPCS, when no fillers were offered.
173 This knowledge would aid the developer of CPCS in determining the upper and lower limits of
174 an acceptable waiting time in CPCS for different task types and in creating applicable mitigation
175 measures depending on the conversational context. (2) Furthermore, we focus on both stress miti-
176 gation and IR tasks due to the pervasiveness of these tasks in contemporary CPCSs. For instance,
177 Panoply [83] and CoZ [3] are crowd-powered conversational systems that provide on-demand
178 emotional support to people. In the context of IR, Chorus [69] and Evorus [48] are crowd-powered
179 text-based chatbots that assist users with IR tasks. The two tasks differ in breadth and depth of
180 focus, duration, and speed [42]. In stress mitigation, workers engage with users through multi-turn
181 conversation by asking focusing or evoking questions, expressing empathy, and reflecting and
182 validating the users' concerns. Therefore, users are expected to experience shorter delays, giving the
183 impression of a continuous and fluent dialogue, enhancing therapeutic engagement, and sustained
184 usage. On the other hand, IR tasks can be completed with either single or fewer dialogue exchanges
185 with considerable delays because workers sometimes have to find the relevant information on the
186 Web. In Table 1, we have provided some examples of task-centric and emotion-centric dialogues
187 to help readers distinguish between the two tasks. (3) Finally, we choose to explicitly add the
188 stress mitigation task because of its ubiquitous relevance. While technology-assisted behavioral
189 interventions can be implemented more easily than in-person interventions, they are limited in
190 their ability to comprehend or express emotions and provide contextualized assistance [62, 87]. As
191 a result, researchers have lately begun to examine crowd-sourced emotional support solutions [83].
192 However, the bulk of existing tools are oriented on the asynchronous approach, which places a low
193 premium on response time. However, research indicates that when creating conversational bots
194 for mental health applications, response time is a crucial characteristic to consider [23]. Thus, we
195 sought to ascertain the acceptable waiting time for real-time crowd-powered affective applications,
196

which are gaining traction in response to global mental health challenges. This potential influence of task type motivates the following research question:

RQ4: *How does the nature of the conversational task affect the waiting experience of user in CPCS?*

Table 1. Examples of task-centered and emotion-centered dialogues

| IR based bots | Examples | Topics of Motivational Interviewing (MI) / Active listening | Examples [88, 120] |
|-----------------------|-------------------|--|--|
| Task-centric Dialogue | Chorus [69] | Hi, I want to know how to make lasagna? | Focusing Question How would you feel about that? |
| | Evorus [48] | Can you find me some good restaurants in Pittsburgh? | Evoking Question What worries you about your current situation? |
| | Guardian [47] | hello, I like to know some information about the movie Titanic | Reflective listening If I heard you correctly, this is what I think you are saying: *. |
| | Humorous Bot [19] | I want to understand how rocks are formed | Expressing Empathy / MI-adherent statement Anyone would feel this way if they were in your situation. |
| | | Emotion-centric Dialogue | |

To answer **RQ4**, we developed two chatbots (Figure 2 & 3) to simulate the CPCS for two different contexts: stress mitigation and IR. For stress mitigation, we developed *Stress-bot*, which is a chatbot we built on the principle of active listening [100] and Motivational Interviewing (MI) skills [78]. To support IR tasks, we built an *IR-bot*, which is a chatbot that answers IR queries in the health, science, and entertainment domains. To answer **RQ1**, we varied the response time of the two chatbots based on the following geometric sequence: 2s, 4s, 8s, and 16s. We adopted this sequence from prior research in computer response theory and crowdsourcing [3, 18]. To answer **RQ3** we developed two variants of the bots each supporting a different complexity of conversational tasks: a high and a low complexity. For *Stress-bot*, the high complexity task involves a challenging scenario with multiple input stressors (any life event or past experience that causes stress) whereas the low complexity task involves a single stressor and a relatively simple situation, such as exam stress. For an *IR-bot*, the complexity of the search task was varied based on the cognitive complexity framework proposed by Kelly et al. [57] which provides different reusable sample search tasks that we used in the current study. To address **RQ2**, we made two versions of bots based on response quality: high and low. For the *Stress-bot*, a high-quality response was based on statements that were vetted by professional therapists whilst a low-quality response was simply based on small-talk. For an *IR-bot*, a high-quality response was detailed, complete and meaningful, while a low-quality response was incomplete and contained an element of hesitancy or uncertainty. In total we developed 32 experimental conditions: **2 x Context, 2 x Complexity, 2 x Quality and 4 x Delay**.

We conducted a between-subjects ($N = 478$) experiment on the Prolific crowdsourcing platform to study the influence of four factors (described above) on the waiting experience. The waiting experience was assessed by the perceived waiting time (in seconds), an appraisal of the wait (long/short judgments and affective responses), and the interaction quality. Our findings show that:

- **RQ1:** Among the four levels of delay compared in the experiment, the 16s delay resulted in more irritation and boredom in the participants. The 8s delay was considered tolerable while the 2-4s delay was perceived as robotic and unrealistically fast.
- **RQ2:** When the quality of the response content was high, participants reported a more positive waiting experience and showed a greater engagement with the CPCS.
- **RQ3:** When the complexity of the task was low, participants overestimated the elapsed time than the actual duration, and perceived quality of the response more positively.

- **RQ4:** When interacting in information retrieval tasks (IR-bot), participants underestimated the elapsed time, experienced less stress, and perceived the quality of interaction more positively than participants who performed the stress mitigation tasks (Stress-bot).

The remainder of this article is structured as follows. In Section 2, we review related work. In Section 3, we detail our methodology including the design of the conversational interfaces, the procedure and the measures used in the study. In Section 4, we detail our results including participants' feedback on long delays. In Section 5, we discuss the insights gained from our results. Finally, in Section 6, we conclude with some implications for designing future CPCs.

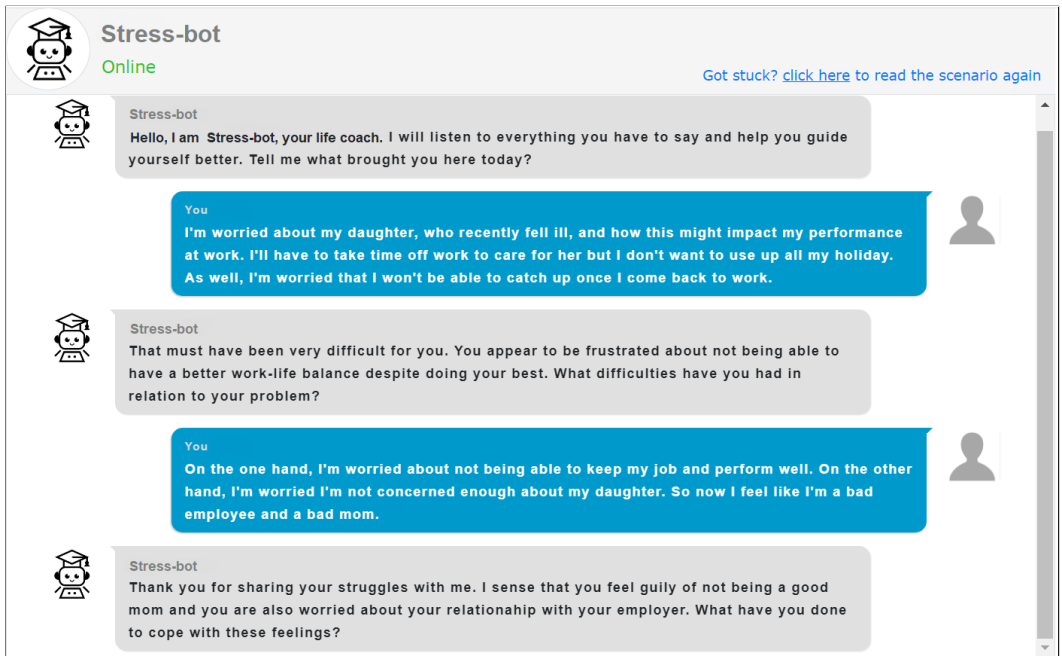


Fig. 2. User Interface of Stress-bot. This figure shows a portion of real chat between Stress-bot and one of the participants from an experimental condition where task complexity and quality of responses of Stress-bot were set to high. Stress-bot initiates the discussion and then allows users to share their ostensible worries based on the description of fictional character. Participants can access the description of the character anytime by clicking the "Got stuck?" link; after that, the description of the character is shown in a popup window.

2 RELATED WORK

2.1 Psychology of System Response Time in Human-Computer Interaction

Empirical evidence appears to confirm the notion that system response time (SRT) is an important factor that influences many aspects of the interaction, such as user satisfaction [27], performance [18], user stress [104], quality of work [53], among others. For instance, Davis and Heineke [27] found evidence that high values of perceived waiting time and actual waiting time can reduce the customer satisfaction in a service operation context. A study of response delays in browser-based applications [44] found that for every three second increase in the system response time, there is an average of .22 drop in average satisfaction and suggested that dissatisfaction may lead to discontinuing the use of the application. Toomim et al. [114] found that Mechanical Turk workers

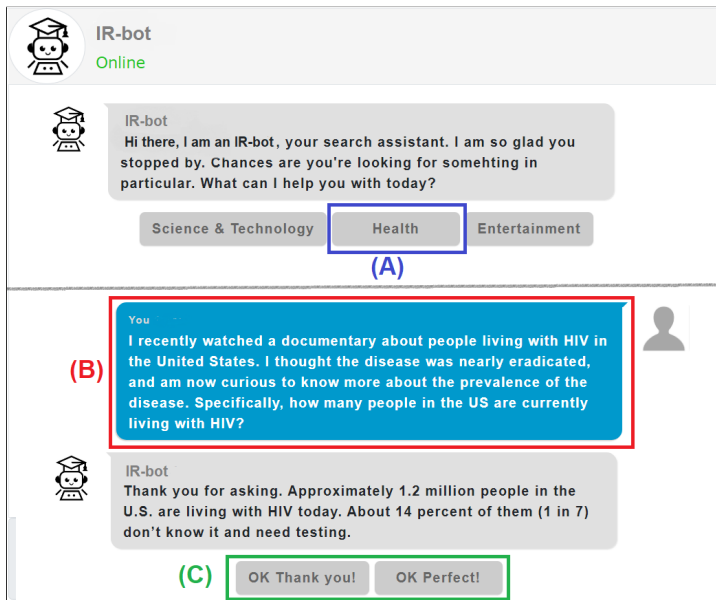


Fig. 3. User Interface of an IR-bot. (A): An IR-bot initiates the conversation by displaying some predefined domains in the form of quick replies. These quick replies disappear when the user chooses one of them.(B): After selecting a domain of interest, an IR-bot then replaces quick replies with a predefined search task, and displays it in a chat bubble. (C): Finally, users can click a quick reply to exit out of the IR task.

preferred to continue with the task when the time to complete the task was lowered as predicted by Fitts' Law. Jacko et al. [53] investigated the combined influence of network delay (short, medium, or long) and media (text or text/graphics) on usability assessment. One important finding of their experiment was that participants perceived the quality of the document to be lower when the delays were lengthier even when the text was augmented with graphics. Some researchers have studied the effect of long delays on the affective experience concerning SRT, such as frustration and stress. For example, Schleifer and Amick III [104] claimed that participants who experienced longer SRTs provided higher ratings of frustration, impatience, and stress. Other researchers studied the physiological effects (rapid heart rate and blood pressure) associated with longer SRTs. Seen in this light, Riedl and Fischer [98] found that the heart rate and blood pressure of those participants who experienced longer SRTs was elevated.

2.1.1 Acceptable System Response Time. Earlier research results pertaining to what response delays are acceptable for users diverge [75]. Shneiderman [109] stated that a system that responds in under two seconds is considered more acceptable than a system that takes longer to respond. Nielsen in his book of Usability Engineering [85] argued that "10 seconds is about the limit for keeping the user's attention focused on the dialogue." Fröhlich [35] provided empirical evidence that while interacting with a speech application, users can only tolerate silences not exceeding 8s. Another study [54] examined the effects of delayed system responses upon user stress and the mental strain they feel. Results show that 63% of their participants wanted a response delay of maximum 5 seconds, while only 20% were positively inclined towards 10-seconds delay.

2.1.2 Response Time of Chatbots. In human-chatbot interaction, only a few studies investigated the effects of variable response delays on the perceived interaction quality. Such studies mainly relate

344 the response times of chatbots with their perceived humanness and persuasiveness. Stemming
 345 from the work of persuasion in computer mediated communication, Moon [80] found a non-
 346 monotonic relationship between response latency and the persuasiveness of a chatbot's messages.
 347 Participants were invited to solve the desert survival problem together with a chatbot. Participants
 348 who interacted with a chatbot with moderate response latency (5-10s) were more willing to modify
 349 their preferences according the chatbot's suggestions, compared to those who interacted with a
 350 chatbot that responded with either long delays (13-18s) or short delays (0-1s). Gnewuch et al. [39]
 351 studied the relationship between response delays of a chatbot and its perceived humanness and
 352 social presence. The authors dynamically varied the response delays based on the complexity of the
 353 message (calculated with the Flesch-Kincaid grade level, a readability test to determine how difficult
 354 a message is to understand). Their results indicate that a chatbot that sent dynamically delayed
 355 responses was perceived to be more human-like and had more social presence than a chatbot that
 356 sent messages quickly without significant delays.

357 358 **2.2 Response Delays in Crowd-powered Conversational System**

359 *2.2.1 Real-time Crowdsourcing and its Applications.* In its early days, crowdsourcing was treated
 360 as an asynchronous mechanism by which tasks could be completed in batches, taking hours to
 361 complete. In recent years, researchers have developed methods to recruit workers on demand
 362 from crowdsourcing platforms and created techniques to reduce latency. Such methods enable
 363 real-time crowdsourcing (RTC) [66] and the resulting systems that support the RTC concept are
 364 known as real-time crowd-powered systems. Real-time crowd-powered systems use recruiting,
 365 rewarding and user interface techniques to reduce the latency of crowd input from hours to a few
 366 seconds [11]. VizWiz [12] is one of the earliest real-time crowd-powered system which helps blind
 367 users in answering visual questions about their surroundings by sending photos of objects and audio
 368 questions from their phones to crowd workers. To make the workers available on demand, VizWiz
 369 incorporates the quikTurkit algorithm that hires workers in advance and keeps them engaged
 370 with some prior tasks. Adrenaline [11] is a mobile application that leverages crowd intelligence to
 371 quickly select the best frame in a video, taking less than 2 seconds to do so in a ten-second video.
 372 Adrenaline implements retainer model where workers are paid a fixed fee for waiting and a small
 373 reward (3 cents) for quickly responding to alerts. Legion [67] allows end-users to crowdsource the
 374 real-time synchronous control of the locomotion of a robot through keyboard commands, while
 375 observing a real-time video feed. Legion:Scribe [65] is a crowd powered system that leverages the
 376 services of unskilled crowd workers to provide captions to deaf people in real time.

377
378 *2.2.2 Response Delays in CPCS.* For the reader's convenience, we have summarized the reported
 379 latencies of current CPCSs in Table 2. A pioneering example of a crowd powered conversational
 380 agent is Chorus [69], which supports information retrieval. Chorus lets workers propose responses
 381 and vote on the responses of other workers in order to support a consistent dialogue. It implements
 382 a working memory model to sustain a coherent conversation over time. Chorus was found to
 383 have a response latency of 103.4s with a single worker, which however could be reduced down to
 384 44.6s by employing multiple workers and choosing among their inputs with a voting mechanism.
 385 Since the main focus of that work was maintaining a consistent and sustained dialogue with
 386 multiple workers, authors did not examine the perceived latency of users. Chorus:view [68] is a
 387 conversational application, which was built on top of Chorus. It assists visually impaired users to
 388 ask general questions about their surroundings by engaging them in a continuous conversation
 389 with crowd workers. Chorus:View broadcasts real-time video stream and audio from the users'
 390 mobile phones to crowd workers allowing them to access the conversational feature implemented
 391 with Chorus. In a trial with blind users, Chorus:View was able to accomplish various tasks with
 392

Table 2. Summary of Response Latencies in Crowd-powered Conversational systems. None of the studies have looked into the subjective waiting experience in CPCS except reporting system response latencies.

| System | Average Response latency |
|-------------------------------|--|
| Chorus [69] | Single Worker: 103.4s Suggest Condition: 44.6s Filter Condition: 50.13s |
| Chorus:View [68] | Product Detail Task: 295s Information Finding Task: 351.2s Navigational Task: 182.3s |
| RegionSpeak [127] | 38s to 2 minutes Yelp Search: ~70.04s |
| Guardian* [47] | Rotten Tomatoes: ~48.6s Weather Underground: ~140.1s |
| Entity Extraction System [50] | 40.95s |
| InstructableCrowd+ [49] | Crowd-only: ~5 minutes Crowd voting: 20 minutes |
| Crowd of Oz [2] | Single Worker: 8.82s Group of two or four: 6.79s Group of eight: 4.12s |
| CRQA [103] | 50s |
| Evorus [48] | Not Reported |

*Total response time was inferred by adding parameter filled time and JSON filled time based on table 3 from [47]

+Total response time was inferred based on the average time users spent in conversing with the crowd and total time it took for a crowd to create a rule

an average response time of 295s, 351.2s and 182.3s for product detail, information finding and navigational tasks, respectively. However, they also did not examine how the blind users perceived the response delays. RegionSpeak [127] is an advanced tool built on top of VizWiz [12] and Chorus [69], which combines Panoramic images (by combining multiple photographic images) and parallel labeling approaches with multiple workers to label multiple objects in a scene. As a result, visually impaired users can conveniently explore the spatial layout of the objects in the space. A user test found that the average response time of RegionSpeak ranged from 38s to ~2 minutes.

Guardian [47] is a crowd powered spoken dialogue system that combines Web APIs with crowd-sourcing to enhance the scope of open dialogue systems in two phases. In the first offline phase, question and answer (QA) pairs are collected with non-expert workers for each API and then match each QA pair with the related parameter. Subsequently, in the online phase, non-expert workers extract parameter values from the running dialogue to obtain the JSON response from the web API service. After receiving the JSON response they translate it into a natural language answer understandable by the user. The performance of Guardian was described in terms of the time needed to fill a given parameter (parameter filled time) and the time needed to acquire the JSON string from the web API (JSON filled time) since the end of the user's first utterance input. Thus, the total time for receiving the answers to the users' questions is the sum of two values. For instance, the time it took to receive the correct response was found to be 70.04s for Yelp search API, 48.6s for Rotten Tomatoes and 140.1s for Weather Underground (cf. Table 3 from [47]).

Huang et al. [50] developed a crowd-powered system to extract entities from a running dialogue in real time using an ESP game. Unlike an ESP game where players are shown images to label, pairs of workers were shown a complete dialog and a description of an entity that they needed

442 to extract from the running dialog. When both workers would identify the same entity, then the
 443 task would be considered complete. An empirical evaluation showed that aggregating answers
 444 from workers yielded good quality results in less than a few seconds. The first worker arrival time
 445 since the publishing of MTurk task was 30.83s, the first response took 37.14s to arrive, while the
 446 first matched response took 40.95s on average. They did not examine the workers' self-reported
 447 perceived latency.

448 InstructableCrowd [49] is a conversational system which allows end users to communicate their
 449 requirements regarding IF-THEN rules to a group of workers. In return, crowd workers assist
 450 end-users by creating IF-THEN rules that then run on their mobile phones. Users can also review,
 451 edit, and approve these rules. In their experiment, user took on average 3 minutes and 45 seconds
 452 to converse with the crowd concerning the rules and system took another 1 minute to create a rule
 453 as per users' demands, resulting in ~5 minutes delay in total. In a condition where they merged
 454 rules from multiple workers, it took approximately 20 minutes for a system to receive rules from
 455 all workers.

456 Crowd of Oz (CoZ) [2] is a crowd powered conversational system for social robotics that can
 457 outsource conversational tasks of the Pepper robot to a synchronous group of workers. CoZ
 458 transcribes the speech of the user and forwards it to crowd workers along with audio-video (AV)
 459 feed captured through the camera and microphones of Pepper. Using a web interface that displays
 460 AV feeds and a chat box, crowd workers compose a text message, which is then spoken out by
 461 the Pepper robot. In the evaluation, authors systematically varied the number of workers who
 462 simultaneously handle the speech of the robot for a stress mitigation task, and studied its effects on
 463 the response latency and response quality. Their experiment showed response delays of 8.82s, 6.79s,
 464 6.79s and 4.12s with one, two, four and eight workers, respectively without effecting the quality of
 465 responses adversely. However, they did not examine the perceived latency for users.

466 Some researchers have combined crowdsourcing and computational intelligence to automate
 467 CPCSs. For instance, CRQA [103] is a hybrid crowd powered question answering system for
 468 informational tasks. It includes a crowdsourcing module to validate the answers generated by
 469 CRQA from existing data sources or optionally to ask workers to provide an answer if CRQA
 470 fails to generate relevant answer. Once the question is posted to CRQA, it allows workers 50s to
 471 respond. After the crowd input is received, CRQA employs a learning-to-rank model to select the
 472 final answer. Similar to CRQA, Evorus [48] also follows a hybrid approach to generate and validate
 473 candidate answers from crowd workers. Unlike CRQA that relies on existing information sources,
 474 it incorporates existing chatbots though REST APIs. Over time, it learns to select the high-quality
 475 answer from chatbots based on workers' past evaluations of answers generated by chatbots.

476 In summary, while a considerable body of research in CPCS has examined the actual latency and
 477 ways to reduce it, less attention has been paid to the latency as it is perceived by users of CPCS
 478 and the underlying factors that contribute to their waiting experience. Our study is aimed at filling
 479 this important gap in the literature.

480

481 3 METHOD

482 We conducted a between-subjects experiment on the Prolific crowdsourcing platform to study
 483 the influence of four factors of CPCS (Context, Complexity, Quality and Delay) upon the waiting
 484 experience of users (see section 3.8 for more details about the dependent variables). Specifically, we
 485 compared 32 different conditions: 2 (*context*) × 2 (*complexity*) × 2 (*quality*) × 4 (*delay*). We designed
 486 a chatbot for each context (stress and IR). The user interface of these bots was designed using
 487 the TickTalkTurk¹ library by Qiu et al. [96]. The server application was developed using Flask, a
 488

489 ¹<https://github.com/qiusihang/ticktalkturk>

490

491 Python-based web framework². The ethics review board at the [anonymous university] approved
492 the methods presented in this paper.

493 3.1 Rationale For Choosing Different Levels of Delays

495 Based on a recent study [2], which examined how varying the number of workers (N=1,2,4,8)
496 engaged concurrently with a task influences the response latency and the quality of the responses,
497 we considered the delays of 4 and 8 seconds as representative of the performance range for real
498 time crowdsourcing [2]. To account for slower crowd work and to introduce sufficient variation we
499 included a delay of 16s as the longest duration considered in this experiment.

500 Prior research in HCI indicates that the system should respond within two seconds after receiving
501 input [108]. This rule is frequently used in HCI research as a design guideline. We wanted to
502 determine if this rule still applied to hybrid intelligence driven CPCS, where responses are expected
503 to arrive faster than when only humans are involved. Furthermore, this 2s condition functioned
504 as a baseline against which contemporary hybrid CPCSs like Evorus [48] could be compared.
505 Furthermore, participants were made aware that, despite the response being received more quickly
506 in 2s condition, human corroborators were still involved to evaluate the responses, as we have
507 specified in the instructions (Cf. Section 3.7)

508 By this reasoning we arrive at the following geometric sequence of the delays: 2s, 4s, 8s, and 16s.
509 A similar geometric sequence was also used in a classic study on user interface response delays
510 [18] which examined the relationship between computer response times (2, 4, 8, 16 and 32s) and
511 user performance with simple data entry tasks.

512 3.2 Design of the Stress Management Conversational Task

514 We focus on stress management task since stress-related disorders are becoming more prevalent
515 and are connected with significant health concerns [25]. While psychotherapies have been shown to
516 reduce stress, the vast majority of the world's population continues to lack access to psychologists
517 [119]. To reach the widest possible audience, we need accessible, inexpensive, and anonymous
518 mental health consultations. Due to the inability of artificial intelligence (AI) to give answers for
519 users' particular circumstances, AI-based mental health interventions currently have a high rate of
520 attrition and low adherence [72]. As a result, it is critical to develop a crowdsourced psychological
521 intervention capable of resolving these challenges. Given its importance, some projects [55, 83]
522 integrating affective computing and crowdsourcing – dubbed affective crowdsourcing [82] – began
523 to emerge more recently. However, prolonged delays may obstruct their applications. Additionally,
524 we added the stress task because of its open-ended nature. In general, the dialogues of such
525 systems are either open-ended (there is no specific answer, and the dialogue is intrinsically useful)
526 or closed-ended (i.e., the dialogue has a specific answer). We used the IR task to address the
527 closed-ended dialogue case. Additionally, we wanted to add an open-ended case, which is why we
528 included the stress-related task. Thus, examining the delays in such dialogues would offer designers
529 with a plethora of information about the maximum permissible delay, mitigation measures for
530 reducing perceived latency, and the most efficient use of available human resources for stress-related
531 conversations in CPCS.

532 To design the stress management conversational tasks, we first collected videos clips from
533 YouTube (based on the strategy described by Pérez-Rosas et al. [90]), which demonstrate simulated
534 counselling sessions between a client and a counselor using motivational interviewing [115]. We
535 examined videos that demonstrate different situations, such as controlling anger, exam stress,
536 social anxiety, relationship problems, among others. Finally, we selected two videos that fulfill our
537

538 ²<https://flask.palletsprojects.com/en/1.1.x/>

540 criterion of complexity (see section 3.2.1 for more details). We then transcribed these videos with
 541 the assistance of Otter.ai platform³, which is an advanced AI-powered transcription platform that
 542 generates very accurate machine-generated text from voice conversations. The generated text is
 543 further augmented with speakers' tags. This feature of Otter helped us to differentiate between the
 544 client and counselor's utterances.

545 In the simulated CPCS, workers were supposed to choose the predefined or template-based user
 546 utterances rather than typing their own responses. We preferred this set up over allowing users
 547 to converse "freely" and "openly" with the simulated CPCS in order to allow a more controlled
 548 experiment. The alternatives of conversing freely with an error prone artificial agent or an actual
 549 CPCS would introduce several confounding factors regarding the direction and content of
 550 the conversation and the variability in response rates. Thus, we asked participants to role play
 551 predefined characters, and carry out dialogues with the simulated stress-based CPCS. An example
 552 of a fictional character from a simple stress task is presented in Appendix A. Full transcripts of a
 553 few dialogues from different experimental conditions and other project related details can be found
 554 at this anonymous URL: <https://bit.ly/3uPijLb>

555 The simulated CPCS was designed in the form of a chatbot that was customized to respond based
 556 on the description of the character in the narrative. One way for a conversational agent to help
 557 users resolve their negative feelings is to invite them to share thoughts that stress them with the
 558 agent, by assuming the role of an active listener. In this active listening mode, the agent lets users
 559 open up and talk about her problems without giving any specific resolution or guidance for dealing
 560 with their problems. The agent has a set of predefined utterances or questions it can present to
 561 the user. This strategy is also used in Woebot⁴, which is a very popular mental health chatbot.
 562 To simulate stress management tasks with CPCS, we developed Stress-bot. Stress-bot is simply a
 563 chatbot that we built on the Rogerian principle of active listening [100]. We will further expand on
 564 how Stress-bot works in the section 3.7. In the next two sections, we explain on how we handle the
 565 complexity and the response quality of Stress-bot.

566
 567 *3.2.1 Complexity of the Stress Task.* To investigate the influence of complexity for the stress task,
 568 we created two variants of CPCS based on Wood's pioneering work [122]. Wood [122] suggested a
 569 task model with three critical components: products, acts, and information cues. We created two
 570 levels of complexity based on the dimension of static complexity, which is related to task design. It
 571 refers to the number of distinct acts and cues required to complete the task. To create tasks with a
 572 static complexity dimension, we make our decision based on the number of stressors present in
 573 the stress task. For example, Ruscio and Holohan [102] provided a number of characteristics that
 574 can be used to characterize the complexity of a therapeutic stress case. They characterized one
 575 of the criteria as multiple, significant current stressors. Thus, the presence of additional stressors
 576 indicates that the case is more difficult to resolve. Each stressor can be mapped with a static
 577 complexity component that necessitates the execution of specific acts and information cues. For
 578 example, in a complex case, we chose a video⁵ in which the client stated two stressors: concern for
 579 her daughter, who suddenly became unwell, and concern about her performance at the company
 580 where she recently began working. Both types of stressors demand specific acts (e.g., describing
 581 stressors, feelings, causes, coping strategies) and information cues (recalling number of facts from
 582 the description of a fictional character to support a certain act) to execute the task. Thus, more
 583 stressors indicate that a task would place greater cognitive demands on the individual, thereby

584
 585 ³<https://otter.ai/>

586 ⁴<https://woebothealth.com/>

587 ⁵<https://youtu.be/osROod3Hmpg>

589 increasing its complexity. For a simple case, we chose a video⁶ in which a client has a single source
590 of stress: concerns about school performance.

592 3.3 Design of the Information Retrieval Conversational Task

593 In order to simulate IR tasks with CPCSS, we built an IR-bot. An IR-bot is a chatbot designed to
594 answer IR queries. We did not allow users to type open queries while interacting with an IR-bot.
595 Instead, we provided custom keyboard options or predefined reply options following the guidelines
596 of cognitive complexity framework [57]. We explain further how an IR-bot works in the section 3.7.

597 For the design of information retrieval tasks, we relied on the framework developed by Kelly et al.
598 [57] based on Bloom’s taxonomy which distinguishes between six types of cognitive processes:
599 *remember, understand, apply, analyze, evaluate, and create*. These cognitive processes are ordered
600 based on the amount of effort required to execute them. For instance, *remembering* tasks are the
601 simplest in information retrieval, requiring the user to just identify or recognize a fact from an
602 information source (e.g., “how many people in the US are currently living with HIV?”). On the
603 other hand, *create* tasks are the most complex tasks, requiring the user to compose a complete plan
604 by putting elements together or reorganizing elements to create something new (e.g., “Identify
605 some basic designs that you might use and create a basic plan for constructing the derby car”). So,
606 the more complex information retrieval task was based on the *create* principle while the simple
607 task was created based on the *remember* principle. Kelly et al. [57] has already developed a dataset
608 of reusable search tasks based on the cognitive complexity framework. We chose three different
609 *remember* tasks, one for each of the three different domains: *Health, Science & Technology and*
610 *Entertainment* and three *create* tasks for the same domains. Hence, our final corpus contains three
611 simple and three complex IR tasks, as shown in Table 3.

612
613 Table 3. Examples of Search Tasks, adopted from Kelly et al. [57]

| 615 Domain | Simple Case | Complex Case |
|--------------------------|--|--|
| 616 Health | I recently watched a documentary about people living with HIV in the United States. I thought the disease was nearly eradicated, and am now curious to know more about the prevalence of the disease. Specifically, how many people in the US are currently living with HIV? | My great granny’s doctor has told her that getting more exercise will increase her fitness and help her avoid injuries. I need an exercise program for her. She is 90-years old. Help me to put together two thirty-minute low impact exercise programs that she could alternate between during the week. |
| 621 Science & Technology | I recently watched a show on the Discovery Channel, about fish that can live so deep in the ocean that they’re in darkness most or all of the time. This made me more curious about the deepest point in the ocean. What is the name of the deepest point in the ocean? | After the NASCAR season opened this year, my niece became really interested in soapbox derby racing. Since her parents are both really busy, I agreed to help her build a car so that she can enter a local race. The first step is to figure out how to build a car. Identify some basic designs that you might use and help me to create a basic plan for constructing the car. |
| 626 Entertainment | I recently attended an outdoor music festival and heard a band called Wolf Parade. I really enjoyed the band and want to purchase their latest album. What is the name of their latest (full-length) album? | My local Triple-A affiliate baseball team has decided that it is time for a new mascot and are holding a contest where fans can enter suggestions. Being a loyal fan, I have decided to enter the contest. I want to suggest a mascot that will appeal to many people and will represent important qualities of a baseball team. The team is a part of the International League, so I want to avoid suggesting a mascot that is already represented in this league. Which mascot would I pick and why? |

636 ⁶<https://youtu.be/KRDH89uP8wI>

638 3.4 Operationalizing the Quality of Responses for Stress-bot and IR-bot

639 In the current study, a high-quality responses of Stress-bot and IR-bot met three criteria provided by
 640 previous research in the field of chatbot interfaces [123]: Appropriateness, empathy, and helpfulness.
 641 Appropriateness refers to the response that is more relevant to the user request; empathy means
 642 the response should provide each user personalized attention and make them feel valued; Finally,
 643 the term “helpful” refers to a response that provides relevant and specific advice in response to a
 644 user’s request and free of any ambiguity. These criteria were not met by a poor-quality response.
 645 We made sure that a high-quality response is substantially different than a low-quality response.

646 Furthermore, high quality responses of Stress-bot were constructed using motivation interviewing
 647 principles, which were vetted by psychologists. For designing the high quality templated responses
 648 of Stress-bot, we adopted a dataset based on the recent work of Park et al. [88] on designing a
 649 chatbot for a brief motivational interview on stress management. They prepared a total of 220
 650 statements which were reviewed by therapists based on the Motivational Interviewing Skills Code
 651 (MISC) [77]. For designing the templated feedback of Stress-bot, we adopted the following approach:
 652 each templated response comprised of three short statements; a) An *empathetic response* to the user’s
 653 problem (e.g., “That must have been very difficult for you.”); b) a *high-level reflection* to the user’s
 654 problem (e.g., “You appear to be frustrated about not being able to have a better work-life balance
 655 despite doing your best.”); and finally asking a follow-up *open question* (e.g., “What difficulties have
 656 you had in relation to your problem?”). On the other hand, a low-quality variant of Stress-bot
 657 asks irrelevant questions that are not directly linked to the user’s problems (e.g., “That’s okay, we
 658 can talk about it. By the way what kind of job is it?”). Furthermore, in this condition, Stress-bot
 659 did not express empathy, did not reflect or validate the user’s concerns and did not ask relevant
 660 open-ended questions. Rather, it would wrap up the discussion with an abstract advice that is not
 661 totally applicable to the problem (e.g., “Do all those activities which give you the most joy in your
 662 life and you will feel better.” or “Maybe when you get stressed you can watch something interesting
 663 on YouTube.”).

664 Next, based on the quality criteria specified by Xu et al. [123], we constructed two variants of
 665 IR-based CPCS in terms of response quality. For instance, a high-quality response is the one that
 666 was appropriate and meaningful. For instance, in response to a user query, “how many people
 667 in the US are currently living with HIV?”, a CPCS with high quality responded with “Thank you
 668 for asking. Approximately 1.2 million people in the U.S. are living with HIV today.” On the other
 669 hand, a CPCS with low-quality response was unsure about its answer and replied with “I think
 670 they are slightly over 1 million”. For an IR-bot, we crafted answers with insufficient information
 671 for low-quality responses and introduced an element of ambiguity to make participants doubt the
 672 accuracy of the answer.

674 3.5 Power and Sample Size Estimation using GPower

675 The sample size and power for the statistical tests were calculated using the Gpower tool [31, 32].
 676 It is the most widely cited analytical application for commonly used statistical tests in social and
 677 behavioral research, and it has been explained and advocated in a number of methodology studies
 678 [28, 63]. It supports a range of different power analysis techniques, the most prevalent of which are
 679 a priori and post hoc analysis. Priori analyses allow users to determine the sample size necessary to
 680 achieve a desired degree of power and effect size. In comparison, post hoc analyses are undertaken
 681 following the conclusion of the study and are used to explain non-significant results and estimate test
 682 power. GPower has been extensively employed in a variety of fields, including psychology [21, 111],
 683 transportation [117], learning and education [74, 79], human-computer interaction research [3, 20],
 684 among others.

687 A priori statistical power analysis using the G*Power application showed that to detect a small
688 effect of $\eta_p^2 = .15$ with 80% power in a 2x2x2x4 factorial between-subjects ANCOVA (32 groups,
689 alpha = .05), we would need at least 489 participants in total. We set the effect size for the current
690 study to 0.15 based on the study [73] that argue that criteria for interpreting experimentally
691 generated effect sizes differ from the frequently cited and extensively utilized guidelines proposed
692 by Cohen. Based on nearly 12,000 correlation coefficients identified in 134 meta-analyses, they
693 advised that correlation coefficients of 0.12, 0.24, and 0.41 and Cohen's ds of 0.15, 0.36, and 0.65 be
694 regarded as small, medium, and large effects in social psychology studies.

695 We also performed a post hoc power analysis to guarantee that an observed power (or post-hoc
696 power) is sufficient. The sample size of 489 was used for the statistical power analysis with alpha
697 level $p < .05$ and an effect size of 0.15. The recommended effect sizes used for this assessment were
698 as follows: small ($f_2 = .10$), medium ($f_2 = .25$), and large ($f_2 = .40$) (see [22]). The post hoc analysis
699 revealed the statistical power for this study was .80 for detecting a small effect. Thus, our sample
700 size was sufficient to obtain statistical power at the recommended level of .80 [22]. The power to
701 detect a medium-sized effect ($f_2 = .25$) was determined to be 0.99.

702 703 3.6 Participants

704 In total, 522 participants were recruited from Prolific crowdsourcing platform. Out of 522, two
705 submissions were rejected (0.38%), 29 participants returned their submissions (5.56%), and 13
706 submissions were timed out (2.49%). Hence, our remaining sample included 478 participants. We
707 restricted the experiment to only US and UK participants since our task required proficiency in
708 English. For the stress task, each worker was paid £2.00 fixed amount (£8.00/h), while in the case
709 of information retrieval task, each worker was paid £1.00 fixed amount (£8.57/h). We paid extra
710 money to participants in the stress task because it requires multiple exchanges as opposed to IR
711 task.

712 The average hourly compensation was £10.55 for the stress task and £9.55 for the IR task, based
713 on the actual amount of time spent by participants. The number of participants from the US was
714 small; there were just 16 for the stress task and 16 for the IR task across all conditions. We also
715 calculated their hourly wage in the US based on the current exchange rate of \$1.38 per British
716 pound at the time of authoring this paper. The hourly wage for the stress task was \$14.56, while
717 the hourly wage for the IR task was \$13.19 (higher than the average minimum wage in the United
718 States of \$7.25 per hour). At the time of writing this paper, this hourly wage was categorized as
719 "good" by Prolific's calculator for both US and UK participants. Participants who participated in one
720 condition were not allowed to participate in the other conditions using Prolific's built-in screening
721 feature. The detailed demographics of participants are reported in Table 4.

722 723 724 3.7 Procedure

725 In this section, we will expand on each stage of the procedure as illustrated in the Figure 4.

726
727 3.7.1 *Consent.* Upon accepting the task participants were asked to read and sign a consent form.
728 This form explained that the chatbot is powered by hybrid intelligence, which means that answers
729 are generated by a combination of computation and human input. More specifically, here are the
730 guidelines we provided to participants regarding CPCS:

731
732 Crowd-powered systems integrate computation with human intelligence, allow-
733 ing large groups of people to communicate and collaborate online. These "hybrid"
734 systems allow applications to solve problems that people and computation alone

Table 4. Participant Demographics

| Characteristics | | Stress-bot | IR bot | Total |
|-----------------|---------------------|-------------|-------------|-------------|
| Gender | Female | 165 (69.6%) | 160 (66.4%) | 325 (67.9%) |
| | Male | 70 (29.5%) | 80 (33.2%) | 150 (31.4%) |
| | Prefer not to say | 2 (.84%) | 0* (0.0%) | 2 (.42%) |
| Nationality | United Kingdom | 216 (91.1%) | 214 (88.8%) | 430 (89.9%) |
| | United States | 21 (8.9%) | 27 (11.2%) | 48 (10.0%) |
| Age | | 33.5±11.8 | 33.0±11.1 | 33.3±11.5 |
| Qualification | No Formal Education | 0 (0.0%) | 7 (2.9%) | 7 (1.5%) |
| | High School Diploma | 45 (18.9%) | 51 (21.3%) | 96 (20.2%) |
| | College Degree | 26 (10.9%) | 35 (14.6%) | 61 (12.8%) |
| | Vocational Training | 21 (8.9%) | 18 (7.5%) | 39 (8.2%) |
| | Bachelor’s Degree | 94 (39.7%) | 91 (38.0%) | 185 (38.7%) |
| | Master’s Degree | 38 (16.0%) | 24 (10.0%) | 62 (13.0%) |
| | Professional Degree | 6 (2.5%) | 9 (3.8%) | 15 (3.2%) |
| | Doctorate Degree | 4 (1.7%) | 2 (0.84%) | 6 (1.3%) |
| | Other | 3 (1.3%) | 2 (.84%) | 5 (1.0%) |

* We were unable to acquire a participant’s gender information because s/he removed his or her response to the gender question before to the study’s execution. Due to data privacy laws, Prolific did not provide that information once the participant deletes his or her response.

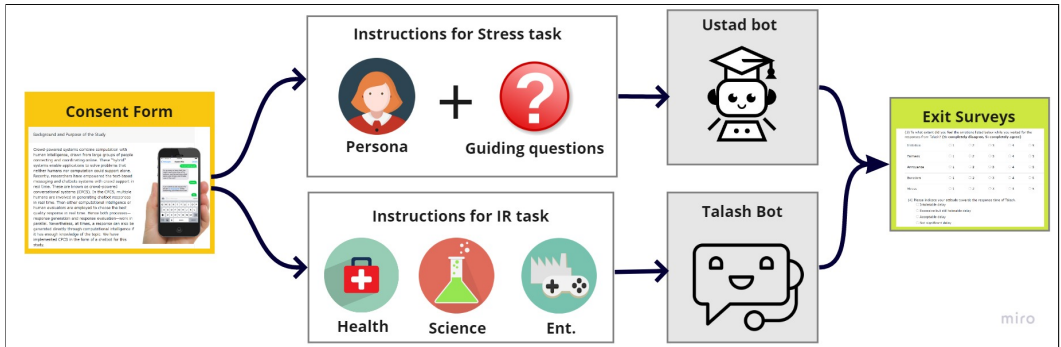


Fig. 4. Stages of Procedure

couldn’t solve. Recently, researchers have added real-time crowd support to text-based messaging and chatbot systems. These are known as crowd-powered conversational systems (CPCS). Multiple humans work together in the CPCS to generate real-time chatbot responses. The best quality response is then chosen in real time using either computational intelligence or human evaluators. Thus, both the response production and response evaluation procedures run concurrently. However, if computational intelligence has sufficient knowledge about the topic, it may sometimes generate a response directly. For this research, we used CPCS in the form of a chatbot. The purpose of the research is to investigate your perception of the chatbot, which is augmented by crowd-powered support in real time.

Participants could then interact with a simulated chatbot that replays human generated statements in a programmed sequence, while controlling the response rate. We notified participants at the end of the experiment that the chatbot’s responses were pre-programmed for ethical reasons,

785 since individuals may exhibit gullibility and embarrassment if they find they have been fooled [34].
786 Additionally, the funding agency EPSRC [15] developed ethical principles for robotics, highlighting
787 the need of users being able to “lift the curtain” and see how a robot operates on the inside.

788
789 3.7.2 *Instructions.* Participants who read and accepted the consent form were redirected to a
790 detailed instructions page. Participants assigned to groups concerning the stress management tasks
791 were presented with the description of a fictional character (see Appendix A for a simple stress
792 task). They were informed that they have to role-play a given fictional character in the conversation.
793 Below the description of the fictional character, we provided some example guiding questions that
794 the crowd-powered conversational system could ask (see Table 5) so that they get in the spirit and
795 anticipate the nature of the choices they would have to make. Because these questions were broad in
796 scope and could be applied to any stress management scenario, we feel they had no influence on the
797 study’s naturalistic setting. These generic questions were created using motivational interviewing
798 techniques used by therapists to learn about a stressed person’s inner struggles. The purpose of
799 including guiding questions was to acquaint participants with the types of queries the bot might
800 ask. Since the participants’ psychological issues were superficial, giving them guiding questions
801 enabled them to play the fictional character as accurately as possible. After they had read about the
802 fictional character and the guiding questions, we let them initiate the conversational task.

803 In case of IR tasks, participants were presented with three IR tasks (see Table 3) distinguished by
804 the complexity of the conditions in three categories or domains: *Health, Science & Technology and*
805 *Entertainment*. After they had read about the detailed description of IR search tasks, they could
806 initiate the conversational task.

807 We did not tell participants that they had to estimate the time in the instructions. We only
808 asked them to estimate the time after they engage with the task, which is called retrospective
809 paradigm [14]. The prospective paradigm, on the other hand, notifies people ahead of time that
810 they need to estimate the time after it has been elapsed. The prospective paradigm is more likely to
811 generate temporal overestimation than the retrospective paradigm [59] because people devote more
812 efforts to judge temporal information and focus less on processing non-temporal information [113].
813 Therefore, we chose the retrospective paradigm, which is less susceptible to such errors. We do
814 agree, however, that both have benefits and drawbacks. One of the drawbacks of the retrospective
815 paradigm is that people may underestimate the duration because they are unaware that they will
816 be asked to estimate the time and thus they may allocate more resources to non-temporal aspects
817 of the task. Due to these reasons, we used two ways to capture time estimation in the current study:
818 a quantitative estimate of the number of seconds participants thought the bot took, and cognitive
819 appraisal of the wait. Due to the fact that the cognitive construct is evaluative in nature and was
820 deemed to be more substantive in a previous study on time perception [40], we chose it along
821 with the quantitative estimate (seconds) to ensure completeness and to observe the actual time
822 estimates.

823 3.7.3 *Conversational Task.* For the stress management task, participants were redirected to the
824 Stress-bot. We adopted Stress-bot based on *complexity, quality, and delay* values corresponding
825 to the experimental setup – see Figure 2. The Stress-bot opens a conversation with a predefined
826 utterance: “Hello, I am Stress-bot, your life coach. I will listen to everything you have to say and
827 help you guide yourself better. Tell me what brought you here today?”. After that, the user could
828 share their ostensible problems with Stress-bot based on the description of the fictional character.
829 Participants were able to access the description of a fictional character during the conversational
830 task in case they wanted to be reminded of some details (Figure 2). In case of the high quality
831 condition, Stress-bot provides a brief reflection on the user’s problem and asks some follow up
832 questions (e.g., “Could you give me some of the details?”) to allow the user to expand on her problems
833

Table 5. Sample Guiding Questions for Stress Task

| Sample Questions | Usage |
|--|--|
| Hi, I would love to hear what that is all about. Do you have any concerns which you would like to talk with me about? | You can expect these questions to be asked at the beginning of the conversation. Explain briefly about your problem (2-3 sentences). You can specify more details later. |
| Hi, I know you had some things you wanted to talk about, so what brings you in today? | |
| I would like to understand how you see things. What brought you here? What has been the problem? | |
| What worries you about your problem? How much does that concern you? In what ways? How would you feel about that? Anything specific about your thoughts that made you feel as you do? | These questions will help you to zoom in to your problem further and focus on your most important concerns. |
| What steps can you take to avoid the frustrations? What have you done to cope with these feelings? What would your partner say about what you are doing? | These questions will help you to explore the options for a change. |

further. Finally, after 4-5 dialogue exchanges, Stress-bot provides the user with a piece of advice: *“Great! They say life is about trying things to see if they work. So, let us not give up on all the good things that you are trying, to help cope with this situation! I think you have some positive aspects to act on.”* Finally, Stress-bot also asks users about their feelings towards stress (*“How do you feel after talking about it now?”*) with some predefined replies (1: *“I am feeling relieved”*, 2: *“I am feeling neutral”*, 3: *“I am still feeling distressed”*). After selecting an option, participants were informed that this task was completed and they were redirected to the exit surveys.

In the case of the information retrieval task, we deployed an IR-bot across the experimental conditions adapting it regarding the *complexity*, *quality*, and *delay*. The conversation was initiated by an IR-bot (e.g., *“Hi There! I am an IR-bot, your search assistant. I am so glad you stopped by. Chances are you’re looking for something in particular. What can I help you with today?”*), as shown in Figure 3. Following this, users were allowed to select one domain out of the given three domains (*Health, Science & Technology and Entertainment*) – Figure 3 (A). After choosing the domain, the search task was automatically created and displayed to the users in the chat bubble based on the dataset provided by Kelly et al. [57] – Figure 3 (B). After a set delay, participants were shown a templated response. Finally, users were presented some quick replies (*“OK Thank You!”* or *“OK Perfect!”*) which they needed to select in order to exit the IR task – Figure 3 (C). We did not provide users an option to repeat the IR tasks or allow them to ask follow-up clarification questions, so as to control the complexity of the search task and the quality of the responses by an IR-bot.

We confined an IR-bot to only one turn in comparison to the Stress-bot due to the practical distinction between the bots’ natures. A task-oriented dialog system aids the user in completing specific tasks within a domain, for example, restaurant bookings, weather searches, and simple information retrieval. Typically, the emphasis of the conversation is brief and limited, and the purpose is usually reached with the fewest possible turns [42, 126]. In comparison, the primary goal of emotion-driven discussion is to enhance user engagement and rapport [2, 101, 106]; as a result, the conversation involves various turns and is typically broad and deep in scope [126]. Thus,

883 it was natural for Stress-bot’s participants to spend more time than IR-bot’s participants. However,
884 we did not impose a time limit on participants in either case.

885 After the conversational task, participants were asked to fill out the exit surveys. In the next
886 section, we expand upon the exit surveys and measures.

887 3.8 Measures & hypotheses

889 3.8.1 *Perceived Waiting Time.* The perceived waiting time was assessed by asking the participants
890 to give an estimate of the total time (in seconds) they had spent waiting between responding to the
891 system and receiving the bot’s response. More specifically, we asked: “*How did you perceive the*
892 *response time from Stress-bot/IR-bot to be during your conversation? Indicate your estimate in seconds*
893 *by simply entering the number of seconds*”

894 3.8.2 *Cognitive and Affective Appraisals of the Wait.* Attitude is considered as comprising of two
895 distinct dimensions: cognitive appraisals and affective appraisal [58]. Cognitive appraisals represent
896 the utilitarian aspect of the attitude, which is used for the products that have functional value; for
897 instance, a grammar checker software may be evaluated cognitively [116]. We measured cognitive
898 appraisal based on the perception of time spent in terms of long or short judgement [93]. It is
899 measured on a five-point scale (1: ‘very short’, 5: ‘very long’). In the survey, we asked participants:
900 “*How would you best describe the response time from Stress-bot(IR-bot) during your conversation?*”

901 On the other hand, affective appraisals refer to the hedonic aspect of the attitude, which is aimed
902 for self-fulfilling rather than functional value; for example, the evaluation of the computer game is
903 based on the affect [116]. More specifically, affective appraisals of the wait capture the emotional
904 response of users towards the response delays. This scale is sufficiently reliable (Cronbach alpha is
905 0.80) and one-dimensional [93]. It consists of five semantic items asking participants to rate the
906 extent to which they have experienced irritation, unfairness, annoyance, boredom, and stress while
907 waiting. Each item is measured on a five-point scale (1: “completely disagree” and 5: “completely
908 agree”). In a survey, we explicitly asked participants: “*To what extent did you feel the emotions listed*
909 *below while you waited for the responses from Stress-bot/IR-bot?*”. We compute the mean affective
910 experience score by averaging these items. Table 6 lists hypotheses regarding the cognitive appraisal
911 of the wait and perceived quality.

913 3.8.3 *Perceived Response Quality.* We also measured the perceived response quality of Stress-bot
914 and IR-bot on a 10-point scale (1: “Very low quality”, 10: “Very high quality”). Specifically, we asked
915 them: “*How did you perceive the response quality provided by Stress-bot (IR-bot)?*”

916 3.8.4 *Qualitative Feedback.* We also asked participants to share any additional thoughts, remarks,
917 or feedback that they may have regarding their experience interacting with Stress-bot (IR-bot). In
918 addition to that, we also asked participants to provide their suggestions on handling long delays in
919 CPCSS.

921 3.8.5 *Control variables.* In this section, we describe individual factors that can mediate the effects
922 hypothesized above. Thus, we treated them as covariates or control variables in the statistical
923 model.

924 (a) *Experience with Chatbots:* Recent work has shown that experience with chatbot technology
925 can impact the users’ self-reported preferences. For instance, a recent study shows that users
926 with more experience with chatbot technology desired more detailed features that can explain
927 the decisions made by the underlying machine learning models [9]. Another study shows that
928 more experienced users are likely to express more critical opinions about the chatbot [71]. Thus
929 we measured the self-reported prior experience of participants with chatbots using two-item scale
930 (“*I am familiar with chatbot technologies*” and “*I use chatbots frequently*”). We measured these items
931

Table 6. Hypotheses regarding appraisals of the wait and perceived quality

| RQs | Hypotheses | References in support of the hypothesis |
|-----|---|---|
| RQ1 | H1: The longer the response time, the more negative the appraisals of the wait (both cognitive and affective) | [18, 104] |
| RQ2 | H2: A poor-quality response of CPCS leads to a relatively more negative appraisal of the wait in comparison to a high-quality response. H5: The higher the quality of the bot’s responses, the greater the perceived response quality. | [30, 93, 125] |
| RQ3 | H3: The higher the complexity of the task in a CPCS, the more positive are the appraisals of the wait. | [5, 13] |
| RQ4 | H4: Longer response times leads to relatively more negative appraisals of the wait in the stress mitigation task than the information retrieval task. | [89, 107] |

on 7-point scales (1: “Completely disagree”, 7: “completely agree”) [9]. This scale showed good reliability (Cronbach alpha=0.71).

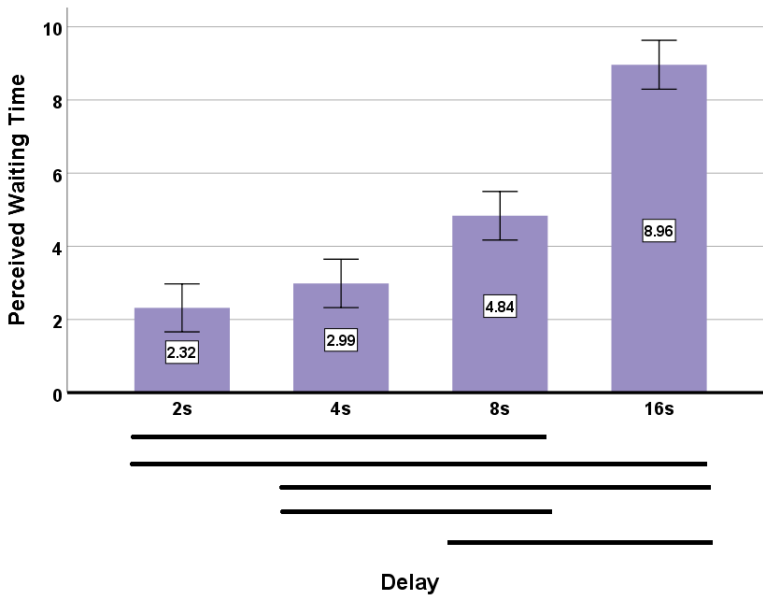
(b) *Social Orientation Towards Chatbots*: The social orientation scale [71] measures the preference of people towards realizing chatbots as “social agents”. People who have a greater tendency to treat chatbots as social entities are more likely to engage with them through natural conversation compared to those who have low tendency. Thus, it is possible that people who treat chatbots as social entities would evaluate the waiting experience more positively. The social orientation towards chatbots was measured by using a two item scale (“I think “small talks” with an AI agent or chatbot is enjoyable.” and “I like chatting casually with a chatbot.”). We measured these items on 7-point scales (1: “Completely disagree”, 2: “completely agree”). This scale has good reliability (Cronbach alpha=0.86).

(c) *Affinity for Technology Interaction* This scale can measure the users’ general technological affinity to actively engage with the technology. We were also interested to investigate whether people who have a greater inclination towards technology (e.g., “I like to occupy myself in greater detail with technical systems.”) could evaluate the waiting experience differently than those who do not have any inclination or were technophobic (e.g., “It is enough for me that a technical system works; I don’t care how or why.”). The affinity towards technology interaction (ATI) scale [33] is a 9-item scale where each item is measured on a 6-point scale (1: “completely disagree”, 6: “completely agree”). This scale has an excellent reliability (Cronbach alpha=0.88 to 0.92) supported by multiple studies with over 1500 participants. [33].

3.9 Data Analysis

We planned to study the effects of *context*, *complexity*, *quality*, and *delay* on the *perceived waiting time*, *cognitive (Long/short judgment)*, *affective waiting experience*, and *perceived response quality* using 2x2x2x4 factorial ANCOVA by treating social orientation towards chatbots, experience with chatbots and affinity for technology interaction as covariates. There was no strong correlation between the covariates, as assessed by multicollinearity test ($r < 0.3$). We also performed the homogeneity of regression slopes test before conducting ANCOVAs by investigating the interaction between highest order interaction and covariates; the assumption of homogeneity of regression

981 slopes was met for all outcome variables ($p > .005$). In a few cases, the assumptions of normality
 982 and/ or homogeneity or variances were not met. However, ANCOVA is robust with respect to
 983 deviations from homogeneity of variances and deviations from normality [10]. We removed five
 984 cases from the dataset because participants provided anomalous feedback. Except in three cases
 985 (two conditions had 13 participants and one had 14), each condition had fifteen participants in
 986 total. This happened because Prolific’s built-in algorithm immediately employs more participants
 987 to compensate for those who left the task. Due to the fact that we had just three cases in which the
 988 number of participants was slightly fewer than the needed number, we anticipate that heterogeneity
 989 of variance will not be an issue, as ANOVA is considered robust to modest deviations from this
 990 assumption [41].
 991



992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011

Fig. 5. Perceived waiting times means and standard errors for the four levels of the actual delay (independent variable). The analysis showed a main effect of delay on mean perceived waiting time. Statistically significant differences in pairwise comparisons are represented with the line segments under the horizontal axis

1016 4 RESULTS

1017 4.1 Perceived Waiting Time

1019 Before conducting an analysis, we deleted the data instances based on the three standard deviation
 1020 criteria. Factorial ANCOVA revealed no interaction effects (4-way, 3-way, and 2-way) between
 1021 all four factors for the mean perceived waiting time. However, the simple main effect of delay on
 1022 mean perceived waiting time was statistically significant ($F(3, 427) = 77.90, p < .001, \eta_p^2 = 0.354$ –
 1023 (Fig. 5). Pairwise comparisons were made for all levels of delay with a Bonferroni correction. The
 1024 mean perceived waiting time score was 2.52 (95%CI, 1.26 to 3.78) points higher for 8s delay than
 1025 2s delay ($p < .001$), 6.65 (95%CI, 5.38 to 7.91) points higher for 16s delay than 2s delay ($p < .001$),
 1026 1.85 (95%CI, .582 to 3.11) points higher for 8s delay than 4s delay ($p = .001$), 5.98 (95%CI, 4.70
 1027 to 7.25) points higher for 16s delay than 4s delay ($p < .001$) and 4.13 (95%CI, 2.86 to 5.40) points
 1028 higher for 16s delay than 8s delay ($p < .001$).
 1029

4.2 Cognitive (Long/short Judgment)

Factorial ANCOVA showed a statistically significant main effect of context (stress, IR) on the cognitive score, $F(1, 441) = 4.34, p = .038, \eta_p^2 = 0.010$). This supports our hypothesis (**H4**) that participants in the stress task judged response delays to be slightly longer ($M = 2.09$ (95%CI, 1.98 to 2.20)) than those who interacted with an IR-bot in the information retrieval tasks ($M = 1.92$ (95%CI, 1.81 to 2.03))—see Figure 6.A. Note that high scores given to the cognitive delay represents that the response delay was perceived to be longer.

We also found a statistically significant main effect of complexity (low, high) on cognitive score, $F(1, 441) = 5.73, p = .017, \eta_p^2 = 0.013$). This supports our hypothesis (**H3**) that when the complexity of the task was lower, participants perceived response delays to be slightly longer ($M = 2.10$ (95%CI, 1.99 to 2.21)) than the high complexity task ($M = 1.91$ (95%CI, 1.80 to 2.02)) – see Figure 6.B.

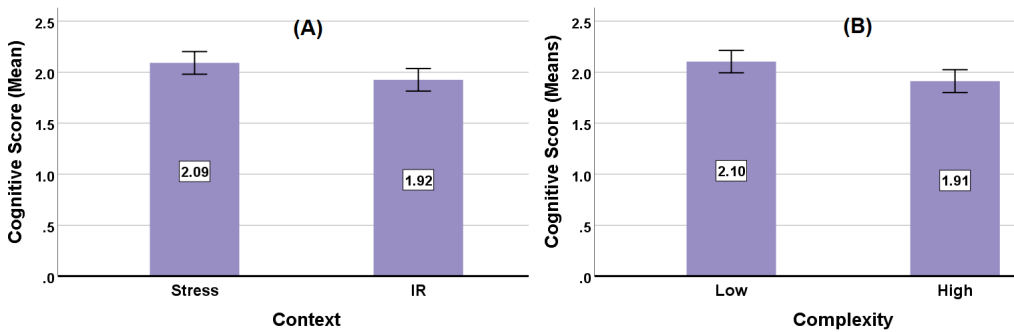


Fig. 6. (A) Result of simple main effect of *context* (stress, IR) on the mean cognitive score. Results indicate that the participants in the stress tasks perceived response delays to be slightly longer than IR tasks; (B) Simple main effect of *complexity* (low, high) on the mean cognitive score. Participants perceived response delays longer when the complexity of the task was lower

The simple main effect of delay on mean cognitive score was statistically significant ($F(3, 441) = 44.96, p < .001, \eta_p^2 = 0.234$) – Figure 7. The mean cognitive score was 0.53 (95%CI, 0.23 to 0.81) points higher for 8s delay than 2s delay ($p < .001$), 1.15 (95%CI, 0.85 to 1.44) points higher for 16s delay than 2s delay ($p < .001$), 0.48 (95%CI, 0.18 to 0.78) points higher for 8s delay than 4s delay ($p < .001$), 1.10 (95%CI, 0.80 to 1.40) points higher for 16s delay than 4s delay ($p < .001$) and 0.62 (95%CI, 0.32 to 0.92) points higher for 16s delay than 8s delay ($p < .001$) – Figure 7. This supports our hypothesis (**H1**) that when the response time was longer, participants evaluated the cognitive appraisals of the wait more negatively.

4.3 Affective Waiting Experience

Prior to presenting results, it is necessary to note that a high score of affective represents negative responses (annoying, irritating, boring, unfair, stressful). Through factorial ANCOVA, we found a statistically significant main effect of *quality* (low, high) on *affective* score, $F(1, 441) = 17.00, p < .001, \eta_p^2 = 0.037$). This shows that when the quality of the bot's responses was lower, participants experienced more negative affective emotions ($M = 2.06$ (95%CI, 1.97 to 2.16)) than in case of high quality bot responses ($M = 1.78$ (95%CI, 1.68 to 1.87)) – Figure 8. This supports our hypothesis **H2**.

Furthermore, the simple main effect of *delay* on mean *affective* score was statistically significant ($F(3, 441) = 14.51, p < .001, \eta_p^2 = 0.090$) –Figure 9. Pairwise comparisons were made with a Bonferroni correction. The mean affective score was 0.55 (95%CI, 0.287 to 0.805) points higher

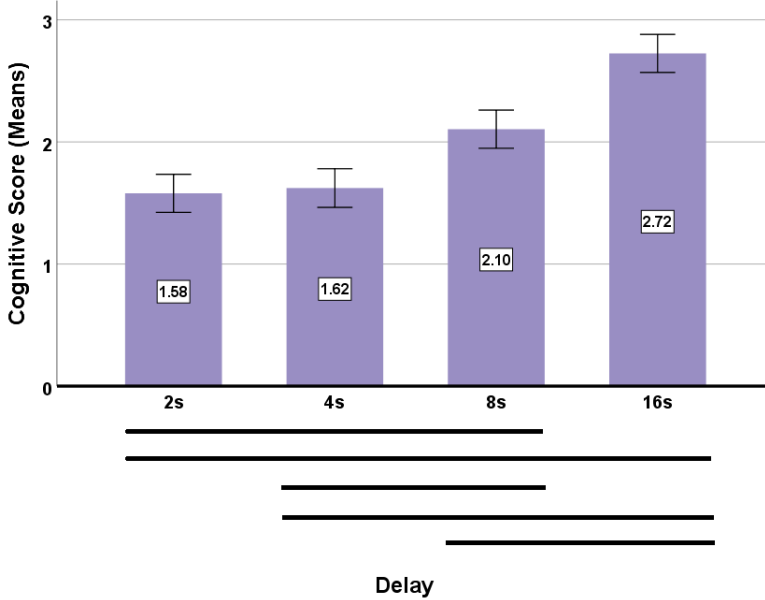


Fig. 7. Mean cognitive scores and standard errors for the four levels of the actual delay (independent variable). The analysis showed a main effect of delay on the mean cognitive score. Differences found significant with pairwise comparisons are represented with the line segments under the horizontal axis

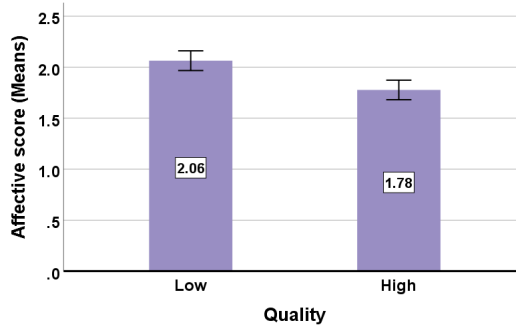


Fig. 8. The simple main effects of *quality* on the mean *affective* score. Results show that when the quality of the bot’s responses was lower, participants experienced more negative affective emotions

for 16s delay than 2s delay ($p < .001$), .53 (95%CI, 0.266 to 0.788) points higher for 16s delay than 4s delay ($p < .001$) and 0.50 (95%CI, 0.244 to 0.762) points higher for 16s delay than 8s delay ($p < .001$). This shows that participants experienced more negative affective waiting experience when they interacted with the bot that responded to them in 16s (Hypothesis **H1**). There was no significant difference in the affective waiting experience scores between 2, 4 and 8s delays.

We also found a significant main effect of the covariate *social orientation towards chatbots* on the mean *affective waiting experience* scores: ($F(1, 441) = 7.09, p = .008, \eta_p^2 = 0.016$). A negative correlation between *social orientation towards chatbots* and *affective waiting experience* was observed

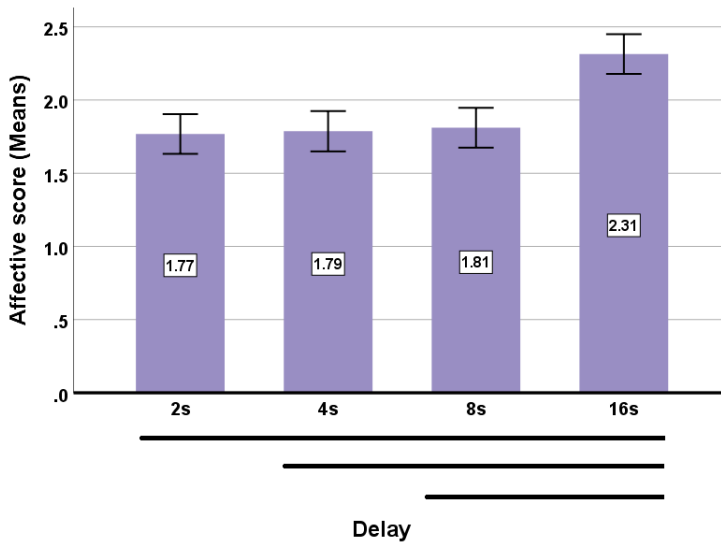


Fig. 9. Mean affective scores and standard errors for the four levels of the actual delay (independent variable). The analysis showed a main effect of delay on the mean affective score. Differences found significant with pairwise comparisons are represented with the line segments under the horizontal axis

($r = -.15, p = .001$), which means that participants who perceived chatbots as more social entities rather than machines perceived less negative affective emotions.

It should be emphasized that, while we employed two distinct task types, each of which has the ability to vary the complexity of the task, we make no comparisons of task complexity “between” the two task types. For example, we discovered that when workers were assigned to “high-complexity” tasks using the Ustad bot, their emotions were severely impacted compared to when they were assigned to IR tasks using the Talash bot. However, since each task’s complexity is defined differently, the comparisons are unfair. For example, in the case of the IR task, high complexity task was determined using the cognitive complexity framework, which defines a more complicated scenario as the one that demands extensive problem-solving skills, such as restructuring pieces into a new pattern or structure by creating, planning, or producing. While in the instance of stress management, a more complex task was one involving many stressors, which necessitates a different set of supporting conversational abilities (demonstrating reflective listening, expressing empathy, asking focusing or evoking open questions) than those required for IR tasks. Therefore, even though we found a significant interaction effect between task type and complexity on the affective score, we did not report it. We incorporate two levels of complexity into each task type, inspired by the seminal work on task complexity [122], and then compare how this may affect time perceptions. In other words, we compare complexity “within” each domain separately, namely within the IR and stress mitigation domains. We investigated task comparisons in a separate research question (RQ4), for which we constructed two distinct chatbots with distinct roles and purposes.

4.4 Perceived Response Quality

Factorial ANCOVA revealed a significant interaction effect between *complexity* and *quality* ($F(1, 441) = 10.34, p = 0.001, \eta_p^2 = 0.023$). The simple effect test shows that there was a statistically significant simple effect of *complexity* on the mean perceived response quality scores when the

quality of the bot responses was lower ($F(1, 232) = 30.15, p < .001, \eta_p^2 = 0.115$), but not for high quality bot responses ($F(1, 234) = 2.51, p = .115, \eta_p^2 = 0.011$) – Figure 10.A. Overall, when the quality of the responses was lower, participants in the low complexity tasks perceived the quality of interaction with the bots more positive ($M = 6.09$ (95%CI 5.68 to 6.50)) than when the complexity of the tasks was higher ($M = 4.45$ (95%CI 4.03 to 4.86)) – Figure 10.A. Nevertheless, complexity did not influence when the quality of the bot’s responses were high. These results show that hypothesis **H5** does not hold when the complexity of the task is higher.

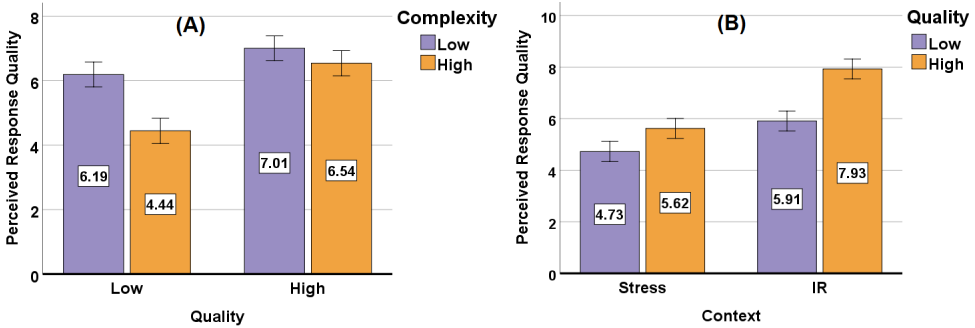


Fig. 10. (A) Interaction effects between *Quality* and *Complexity*: Results indicate that when the quality of the responses were lower, participants in the low complexity tasks perceived the quality of interaction with the bots more positive than when the complexity of the tasks were higher; (B) Interaction effects between *Quality* and *Context*: Result shows that in both contexts (stress and IR), participants significantly differentiated between the high and low quality responses.

We also found a significant interaction effect between *context* and *quality* ($F(1, 441) = 8.21, p = .004, \eta_p^2 = 0.018$). The simple effect test shows that there was a statistically significant simple effect of quality on the mean perceived response quality scores when the context was stress ($F(1, 232) = 9.44, p = .002, \eta_p^2 = 0.039$), and also when the context was IR ($F(1, 234) = 45.09, p < .001, \eta_p^2 = 0.162$). This shows that in both contexts (stress and IR), participants significantly differentiated between the high and low quality responses (supporting **H5**). Overall, participants perceived the interaction quality with an IR-bot more positive than the Stress-bot (Figure 10.B).

We also found significant main effect of the covariate *social orientation* on the mean *perceived response quality* scores: ($F(1, 441) = 49.93, p < .001, \eta_p^2 = 0.102$). A positive correlation was found between *social orientation* and *perceived response quality* ($r = .345, p < .001$) which means that participants who perceived chatbots as social entities rather than machines perceived the response quality to be higher. We have provided the summary of key findings in the Table 7.

4.5 Suggestions on Handling Long Delays

To gain insights of participants’ opinion about the design of CPCS for handling long delays, we calculated the number of suggestions that participants described and presented them in a Figure 11. More specifically, we asked participants: “Do you have any suggestions for how conversational systems can handle long delays in responses?”. We received 477 feedbacks in total from participants. The collected feedback responses were subjected to an inductive thematic analysis [17] using the Dedoose program⁷. The first iteration assessed responses and generated initial codes. The second iteration analyzed previous themes and integrated those that were redundant. After two coding

⁷<https://www.dedoose.com/>

Table 7. Summary of the key findings with respect to our research questions and hypotheses.

| RQs | Hyp. | Supported? | Summary of key findings | Reference |
|------|------|------------|---|------------------|
| RQ#1 | H1 | Yes | When the response time was longer, participants evaluated appraisals of the wait more negatively. | Section 4.2, 4.3 |
| RQ#2 | H2 | Partly | <ul style="list-style-type: none"> When the quality of the bot's utterances was lower, participants experienced more negative affective emotions. We did not find support for the cognitive component (long/short judgment). | Section 4.3 |
| | H5 | Partly | <ul style="list-style-type: none"> Participants in both contexts (stress, IR) significantly differentiated between high and low quality responses. When the quality of the bot's utterances was lower, participants perceived the quality of interaction with the bots more positive <i>only when the complexity of the task was lower</i>. | Section 4.4 |
| RQ#3 | H3 | Partly | <ul style="list-style-type: none"> When the complexity of the task was lower, participants perceived response delays to be slightly longer. When the complexity of the task was higher, participants experienced more negative emotions. | Section 4.2, 4.3 |
| RQ#4 | H4 | Yes | Participants in the stress task judged response delays to be slightly longer and experienced more negative emotions than those who were in the IR task. | Section 4.2, 4.3 |

cycles, the analysis had settled on the primary issues discussed in this paper. Now we briefly expand on each suggestion in the next sections.

4.5.1 Holding or Filler Messages. The majority of the participants recommended filler or holding messages (N = 98). Conversational fillers or holding messages are quite ubiquitous in human conversations, such as “well”, “hm”, “uh” or “let me think” [29]. They let the waiting person know that s/he is not ignored and the response is likely to come shortly. For instance, One worker (ID:385) replied: “*Perhaps they can provide a quick response first such as ‘I understand, please give me a moment to think on what you have said’ before issuing a more in-depth response...*”. Another worker (ID:370) replied: “Maybe a holding message to say, ‘please bear with me while I think about that?’”

4.5.2 Typing indicators. Typing indicators make chatbot interactions appear more natural. These indicators show that a message is being typed and thus could encourage users to wait. Participants mentioned two types of typing indicators that are very prevalent in messenger applications; 1) graphical typing indicators (N = 38); 2) text-based typing indicators (N = 11). Graphical typing indicators include three animated dots [...], progress bar or loading or thinking icon while the text-based indicators include “Person X is typing”. For example, concerning the graphical indicator, one worker (ID:264) replied as “*Maybe add ... to show that the chat bot is currently typing. There was nothing to illustrate to me that the chat bot was actually working on its next response to me which was a little annoying.*”. Example of text-based indicator was described by worker (ID:271) as follows:

1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323

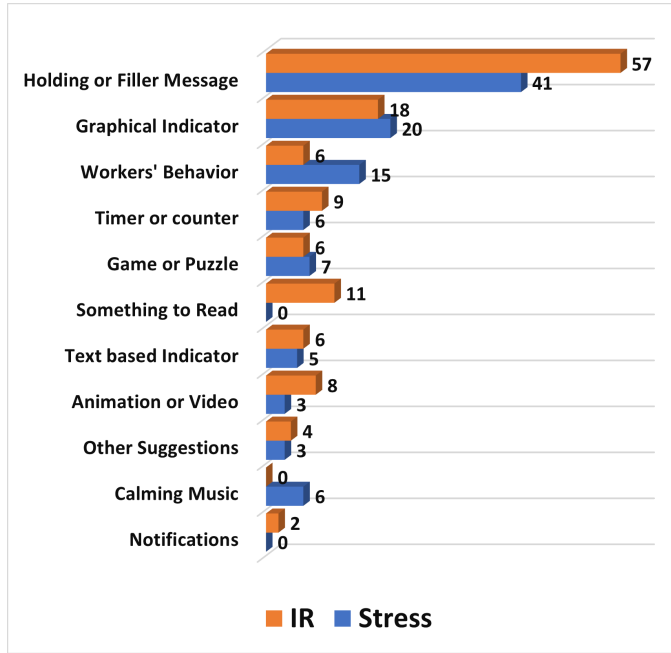


Fig. 11. participants' suggestions about handling the long delays in crowd-powered conversational agent

"I think having a notification saying 'x is typing' makes the waiting process more realistic and more reassuring"

4.5.3 Using Crowd's behaviour. Participants also suggested solutions to reducing response delays that were directly linked with the crowd's behaviour (N = 16), such as checking whether other person (worker) is still there, communicating with workers beforehand about the expected response delays, having a 'seen' function akin to messaging applications, send reminders to workers who are slow in responding. They also introduced a strategy to show partial responses that allow the user to see a response when it is being written in real time. This is similar to the IMO messaging app that allows you to see in real-time what the other person is typing. For instance, *"Check that the other person is still there..."*(ID:394) or *"Communicate with the person how long response times are expected to take at the beginning of the conversation."*(ID:278) . Regarding the partial responses, one worker (ID:268) said, *"just to display the word typing or similar so you know you will definitely get a response; the waiting was frustrating."*

4.5.4 Entertaining Activities. Participants also mentioned entertaining activities as a way to distract a waiting person. This includes playing some game or puzzle (N = 13), something interesting to read (N = 11), watching a cartoon animation or video (N = 11) and listening to some relaxing music (N = 6). [Game or Puzzle, ID:253]: *"There could be simple games like 'Pong' or 'Snake' that can keep the user occupied while waiting"* [Short readings, ID:66]: *"...Giving the user something short to read may make the delay seem less lengthy."*. [Animation or Video, ID:461]: *"Perhaps have an interesting character moving around the screen"*.

4.5.5 Notifications. Participants also mentioned that when the delays are longer, one's attention may wander off to some other tasks. Therefore an audio notification at short interval may help to

1324 retain one's attention (N = 2). Example includes (ID:63): *"Unsure as sometimes delays are unavoidable.*
 1325 *It would be good to have a ding sound when the bot responds so you can click away on other tabs, then*
 1326 *be alerted when there is a response"*

1327
 1328 4.5.6 *Other suggestions.* Other suggestions include showing emojis, engaging users with some
 1329 breathing activities and providing tools for note-taking or thought recording while waiting.

1330 4.6 Participants' Perceptions about Different Delay Levels

1331 Using the same thematic analysis approach as outlined previously, we investigated instances in
 1332 which participants expressed their emotions concerning varying levels of delay. Overall, participants
 1333 from 2s (N = 38) and 4s (N = 44) conditions believed that the responses in both conditions were
 1334 quick because it is impractical for humans to provide in-depth answers in such a short span of
 1335 time. For instance, a participant in 2s condition responded as [ID.35]: *"The response is very fast, so I*
 1336 *know that this is not a human typing a response. The response is quite meaningful and thorough. I*
 1337 *am happy with the response."* Another participant in 4s condition replied [ID.186]: *"The responses*
 1338 *from Stress-bot were very fast so sometimes it felt like he hadn't really understood but after reading*
 1339 *the responses I found they were appropriate"*.

1340 Participants in the 8s (N = 25) condition believed that the response time was not too long. This
 1341 assertion was echoed by a participant [ID.193]: *"The response time was good, and I was not kept*
 1342 *waiting a long time. If this happened on all chatboxes I don't see how anyone would complain"*. Another
 1343 participant [ID. 394] who interacted with Stress-bot commented: *"I don't feel the delay in responses*
 1344 *was excessive. As a patient, if I were waiting for that level of response from a human, I'd be happy*
 1345 *they thought about it..."*.

1346 Only a few participants (N = 13) supplied feedback on the 16s condition; those who did provided
 1347 feedback were equivocal about their feelings. Some participants (N = 5) thought the response
 1348 time was irritating, while the remainder (N = 8) thought it was reasonable. For instance, one
 1349 participant [ID.456] expressed dissatisfaction with the long delay as follows: *"I thought the response*
 1350 *was well-thought out, informative and interesting, but the delay made me think I did something wrong,*
 1351 *or something was broken at first"*. Another participant [ID. 156] replied: *"Even though I expected the*
 1352 *delay, the thought did flash across my mind that perhaps I wasn't going to get a response and my*
 1353 *internet had frozen, when I was sitting waiting..."*. Although delayed responses negatively impacted
 1354 on the user experience with CPCSS, it increased their perception that a real human was talking
 1355 instead of a bot; one participant [ID.216] replied: *"It felt like an actual person and not a chatbot"*. One
 1356 participant [ID.460] who believed that the response was reasonable commented: *"IR-bot gave a good*
 1357 *response and answered my question well within a reasonable amount of time."* Another participant
 1358 [ID.517] responded as: *"no I think the response time was fine, I have waited longer for responses in the*
 1359 *past."*

1361 5 DISCUSSION AND CONCLUSIONS

1362 5.1 RQ1: 2-4s Delay is Robotic, 8s Delay is Tolerable but 16s Delay is Captive

1364 Our findings indicate that participants perceived the waiting time to be less than it actually was
 1365 and were unable to differentiate between two and four second delays (cf. Figure 5). Thus, it is
 1366 obvious that the functional or utilitarian appraisal of time is valid until four seconds of delay. This
 1367 is owing to the engaging nature of digital interventions, particularly chatbots, which mimic human
 1368 connection through a conversational approach. This might have helped participants to divert their
 1369 focus away from the time until 4s and increased their engagement [91, 92, 95]. Theoretical models
 1370 developed by psychologists concerning time perception models can also explain this phenomenon.
 1371 For example, cognitive absorption [4] is defined as a 'state of intense involvement with software'

1372

1373 (p. 673) and is derived from Csikszentmihalyi's flow theory [24], which defines a mental state in
1374 which a person performing an activity is completely immersed in a sense of energized focus, deep
1375 involvement, and enjoyment in the process of the activity. When the system took longer than
1376 8s, however, the emotional component was affected. Given that this deep involvement is linked
1377 to subjective experiences with cognitive and emotional dimensions [91], it is likely that negative
1378 experiences concerning the appraisal of the wait started to slowly emerge after 4s but worsened
1379 after 8s. In the 2 and 4s conditions, the uncanny valley theory [81] may account for participants'
1380 assessment of the bot's utterance as robotic and unexpectedly fast. For instance, a study shows
1381 that participants perceived interaction with the chatbot as an unpleasant or eerie if they have to
1382 wait longer for the response especially when chatting with a human [110]. In our case, since the
1383 participants knew a priori that real humans are generating responses, they might have assumed
1384 that they would take some time in composing a valid response. Thus, an extremely fast response
1385 made them uncertain regarding the true nature of a CPCS as human or AI that provoked feelings
1386 of eeriness. This finding also reflects that the development of technological innovation to lower the
1387 response latency below 4s would be extravagant within the context of CPCS.

1388 Participants from the 8s condition deemed that the response time was not excessively long. This
1389 finding seems noteworthy but poses an important question: Is it conceivable to attain up to an
1390 8 second response latency for difficult conversational tasks in CPCS that require more cognitive
1391 capacity? One solution is to use speech to text (STT) services on the crowd interface. For instance,
1392 CoZ [2] was able to endow workers' web interface with STT API for quickly transliterating voices
1393 of workers into text. Furthermore, people can speak faster (compared to typing) and take shorter
1394 pauses— typically about a quarter to half a second [86]—to reflect on what they say. Therefore, an 8
1395 second duration seems sufficient for people to come up with a reasonable answer with STT. These
1396 results are also consistent with a prior finding about system response time in telephony speech
1397 applications where authors claimed that waiting time should not exceed beyond 4–8 seconds [35].
1398 Favoring this claim, Nielsen [85] suggests that 10 seconds is the upper limit for keeping the user's
1399 attention focused on the dialogue.

1400 The 16s condition elicited higher negative affective reactions from participants (cf. Figure 9).
1401 Similar to our finding, Yang et al. [125] showed that delays in the response of conversational agents
1402 were associated with negative affect (distressed, anxious, tense etc.). Likewise, Butler [18] used
1403 the same geometric sequence of response delays in simple computer-based data entry tasks. He
1404 found that when the computer took longer time to respond in printing the prompt characters, so
1405 did users in responding. Furthermore, Chatbots waits are usually open-ended, where no explicit
1406 waiting duration can be conveyed after a user inquiry. Thus, a user's attention is held captive until
1407 a chatbot response [121]. Miller [76] wrote in his report about response time in conversational
1408 transactions that "captivity of more than 15 seconds, even for information essential to him, may be
1409 more than an annoyance and disruption".

1410 In summary, within the CPCS context, researchers need to develop algorithms and workflows to
1411 keep the response latency within 4-8 seconds. For scenarios that require longer response time, such
1412 as complex information retrieval tasks, the effects of extended waiting time can be alleviated by
1413 designing and testing different time fillers [121]. Besides, the fact is that current CPCS are nowhere
1414 near the threshold of 8s; for instance, field deployments of Chorus show that 25% percent of the
1415 conversations obtained a first answer within 30 seconds [46]. Although CoZ [2] was able to keep
1416 response latency under ~8 seconds by broadcasting an audio-video stream, this is not an affordable
1417 solution for long-term deployments. Thus, our current study is the first that addresses the issue of
1418 time perceptions in CPCS and raises the urgency of developing techniques beyond already known
1419 techniques (e.g. pre-recruiting and retainer [11, 45]) to reduce the perceived latency of CPCS within
1420 the desirable range of 4-8 seconds.

1421

5.2 RQ2: Quality of Response's Content can Influence Waiting Experience

Participants perceived the affective response in terms of irritation, boredom, fairness, etc. more negatively when the quality of the bot's responses was lower, regardless of response delays (cf. Figure 8). These findings are in line with the prior work where authors studied the impact of product qualities (measured with pragmatic quality, hedonic quality, and the overall appeal) on the affective experience with conversational agents [125]. One important finding was that "quality of response content" was negatively related to the negative affect, that is, when the quality of the response content was lower, participants experienced more stress.

Given that one of the criteria for response quality in the current study was usefulness, we also analyzed the association between response times and perceived quality (Cf. section 4.4). Our findings, however, did not provide solid evidence regarding the influence of the bot's response quality on perceived delay, other than the fact that good quality responses reduce frustration. However, it appears that when the quality of the bot's answers was high, the estimation of passing time would be slightly compressed when the delays were longer; 8.8% decrease for 8s delay and 4.7% decrease for 16s delay. Thus, the high-quality of responses can be an important factor for countering the negative effects of waiting. However, further studies within the context of CPCS are needed to test this association. Additionally, because cognitive appraisal of the wait is related to time, there is ample evidence that shows that shorter delays result in greater satisfaction (see [51, 94, 112]). In other words, the more positive the wait is rated, the more satisfied the user is with the service.

To the best of our knowledge, no study in human-chatbot interaction or in CPCS has looked specifically at the influence of quality of response content on the perception of time. We are only aware of few studies in the service marketing and transit agency context that study the relationship between quality of service and perception of waiting time. For instance, Pruyn and Smidts [93] studied the attractiveness of the waiting environment in terms of comfort, spaciousness, and atmosphere on the perception of waiting time. They argued that quality of environment can serve as an element of distraction and consequently can lower the perception of elapsed time, though their study did not reveal such effects. The reason we found only a weak relationship between the response quality and perception of time could be for a variety of underlying factors; for instance, it may be the case that variation in the quality of the bot's responses in this study was not substantial enough to evoke the potential effect on the perceived waiting time.

These results also provide an interesting trade-off between content quality and latency. For instance, researchers have developed and tested a plethora of techniques to lower the objective waiting time in CPCS through pre-recruiting, queuing, and other computational techniques [11, 43, 45]. Nevertheless, lowering the objective waiting time beyond certain limit could compromise the quality of crowd output due to time pressure on workers. Our result seems to support this notion that the quality of response contents is an important factor that may influence the perception of time passing and thus, by uplifting the quality of response contents one can lower the perceptions of actual waiting time. Maintaining the high quality of crowd output is an active area of research in crowdsourcing [61]; for instance, research has shown that properly training crowd workers to support complex conversational tasks could significantly improve the quality of response content [3]. Our research has advanced the state of the art in CPCS by stating that subjective waiting time perceptions can be improved by enriching the quality of response content. As a result, this can ease the engineering efforts in real-time crowdsourcing research to improve the response latency to a desirable extent of 8 seconds based on our results.

5.3 RQ3: Perceived Time Drags and Perceived Quality Boosts Up with Low Complexity Tasks

When the tasks' complexity was lower, participants overestimated the elapsed time in comparison to high complexity tasks (Figure 6.B). Past research shows that a complex activity can lead to underestimation of elapsed time [5, 13, 59]. For instance, this assertion was supported by Khan et al. [59] who conducted an experiment to find the effect of cognitive load on time perception. The complexity of task was varied based on the cognitive load—ranging from merely paying attention to memorizing the presented items. Their results revealed a negative association between cognitive load and time perceptions, that is, when the cognitive load increased, participants underestimated the time perceptions. Likewise, Block and Gellersen [13] varied the complexity of three input techniques, including simple on-screen icon invoked by the mouse (GUI), graphical toolbar that was mounted on the keyboard to invoke commands (display keyboard) and finally a touch-sensing keyboard (touch keyboard) combined on-screen icon with corresponding touch-based keys on the keyboard. In their experiment, a touch keyboard was considered to be a cognitively more stimulating input modality than others. Their results indicate that participants who executed commands with the touch keyboard perceived the time to be shorter.

This inverse relationship between cognitive load and time perceptions can be better explained with the well-grounded *cognitive-timer* [38] and *cognitive-attentional* [113] models from psychology. For instance, the cognitive-attentional [113] model explains that a person's attentional resources are limited and are usually split between temporal and non-temporal information during an activity. Thus, if an activity has a non-temporal dimension, such as recalling more information, then the temporal estimation is affected. The cognitive-timer models also relies on similar assumptions: "(1) the existence of a cognitive timer whose purpose is to process and generate temporal information; (2) temporal information is processed by the timer by storing the number of subjective time units which have accumulated during a given interval; and (3) attentional resources are allocated continuously to enable both temporal and non-temporal information processing, with a trade-off in allocation between the cognitive timer and other cognitive modules". For instance, a high complexity search task involves more cognitive modules, such as extracting knowledge from different sources and then structuring and organizing them to create a new idea, which results in a decreased attention towards the temporal aspects; Thus, it is possible that the user may establish a belief around this notion that CPCS may take longer to respond and consequently their biological mechanism—the internal clock—[64] is adapted to underestimate the temporal judgment. On the other hand, a low complexity search task only requires extracting a relevant fact from the information sources, which requires less cognitive resources and attention is devoted towards temporal aspects. Therefore, the time appears to pass more slowly in the low complexity tasks. Thus, if CPCS takes longer to respond in low complexity tasks, it can further exacerbate users' waiting experience.

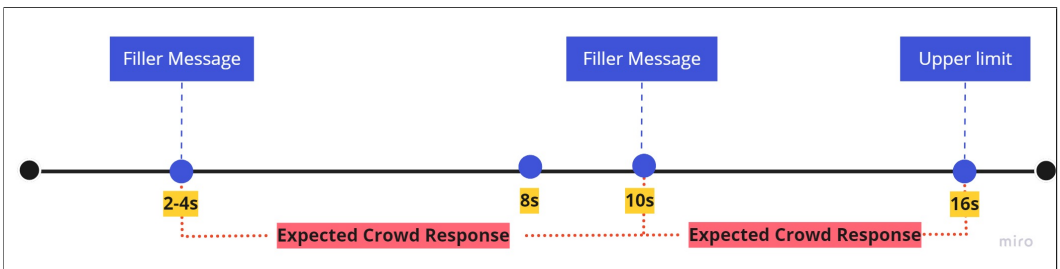
We also found a combined impact of task complexity and response quality on the perceived interaction quality with CPCS. When task complexity and response quality was low, participants perceived the interaction quality higher compared to when task complexity was high (Figure 10.A). Results obtained by Borromeo et al. [16] are consistent with our findings who studied the impact of task complexity and crowd type on work quality. They tasked workers with two types of tasks; a complex task required workers to extract more information items (e.g., title, authors, source etc.) from a given presentation, whereas a simple task required workers to only extract one information item. Their experiments indicate that the simple task generated higher quality results than the complex task by a paid crowd. Another possible interpretation of this finding is that in low complexity tasks, users may have set lower expectations regarding the response quality; as such, when the responses' quality was compromised, it did not greatly influence their

1520 perception of quality and subsequently they overestimated it. An alternative explanation could be
 1521 that participants who had less experience with chatbot technology could have overemphasized the
 1522 responses' quality. A similar idea that might explain this finding is the positive correlation that
 1523 we found between social orientation towards chatbot technology (Cf. section 4.4) and perceived
 1524 quality. Thus, it is possible that participants who had participated in these conditions perceived
 1525 chatbots more as social beings in lieu of machines, and consequently perceived the response quality
 1526 to be higher. To confirm this, we performed a t-test to examine the difference in participants'
 1527 self-reported social orientation scores between low and high-complexity tasks. We found that
 1528 when the task complexity was low, participants' self-reported social orientation scores were greater
 1529 (3.1 ± 1.65) as compared to high complexity tasks (2.6 ± 1.68) ($t(235) = 2.276, p = .024$). Thus, we
 1530 can conclude that participants who preferred more humanized social interaction with the CPCS,
 1531 perceived the quality to be higher even when it was compromised.

1532
 1533 **5.4 RQ4: Task Types Matter: Users in the Stress Task were More Affected with Waiting**

1534 Participants who interacted with an IR-bot perceived delays to be shorter (cf. Figure 6.A) and
 1535 perceived the quality of interaction to be more positive than those who interacted with the Stress-
 1536 bot (cf. Figure 10.B). Why participants, in the IR task, had a better waiting experience could be
 1537 associated with the nature of IR tasks where the only goal is to extract relevant information from
 1538 the agent as opposed to engaging in multi-turn, unconstrained and continued dialogue with the
 1539 agent to alleviate stress and maintain rapport. Thus, it is possible that participants in the IR task
 1540 were able to tolerate some delays due to its short-lived, single-turn nature. Recent research has
 1541 tended to show a similar trend where Peng et al. [89] studied the impact of different response delays
 1542 (1, 2, 4, 8s) on the users' satisfaction, in a human-robot context, for two task types: chitchat and IR.
 1543 Their results show that when the delays were longer (4s and 8s), users in the chitchat condition
 1544 showed less satisfaction than those who engaged with the IR tasks.

1545 We could relate our participants' positive perception of interaction quality with the IR task
 1546 to the way we designed the contents of responses for both tasks. For example, in the stress
 1547 task, answers were derived from templated responses that were rather generic in nature. Despite
 1548 the fact that these answers were verified by professional psychologists, it is probable that some
 1549 participants were dissatisfied with the responses due to their generic nature and lack of detailed and
 1550 contextualized feedback. In the case of IR tasks, it was comparatively easier to design a high-quality
 1551 response because the search tasks were pre-defined and were presented in the form of quick replies.
 1552 Secondly, aggregating answers from multiple information sources to tell a coherent story was a
 1553 straightforward task.
 1554



1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565 Fig. 12. Timeline: it indicates locations where the creator of CPCS can insert fillers to alleviate waiting feelings
 1566 and indicates locations where a message from the crowd is expected.
 1567
 1568

5.5 Participants' Perceptions of the Bots' Identities

Given that we did not inquire participants about their perception regarding the nature of the chatbot in this study, this raises the question of what "identity" participants perceived toward the bots when they interacted with them - Did they view them as human or AI? In previous trials involving around 1000 participants, we discovered that approximately 73% of subjects believed they interacted with a form of hybrid intelligence, whereas the other 27% believed they interacted exclusively with AI [7]. Based on this, we expect that the current experiment would yield comparable results. It is noteworthy that even if participants believed that they were interacting with AI, this is to be expected in a hybrid CPCS setting where there is a role of AI in either creating a response or evaluating the quality of human-generated output. While current research on AI and conversational agents has focused on revealing the true nature of systems in laboratory trials due to ethical issues [99], masking the identity of a bot has shown some advantages. For instance, hybrid CPCSs are comparable to chatbots often employed in customer service contexts, where concealing the true nature of a conversational agent may be helpful in certain circumstances, particularly when seamless overlapping of human (customer representative) and AI is involved [110]. Another advantage of opacity was discovered in another study, where the bot was deployed to answer users' questions on the Stack Overflow community. They found that when the bot pretended to be a person (similar to our case), it was perceived more positively than when the bot revealed its bot identity [84]. Furthermore, both bots exhibit some features that may help to perceive them as human-like. For instance, a dynamic chatbot that generated relevant follow-up questions was seen as more human-like than a static bot that did not generate such inquiries [105]. This was true in the case of the Stress-bot, which uses MI theory to ask follow-up questions through multiple turns. However, in the case of an IR-bot, we used simple, polite, and informal language to create responses, which is indicative of humanness [8]. Our future research will examine how users' perceptions change when they learn post-hoc whether the bot was controlled by humans or AI.

6 IMPLICATIONS FOR DESIGNING FUTURE CROWD-POWERED CONVERSATIONAL SYSTEMS

We found that some participants in general did not like fast responses for two reasons: 1) they thought that CPCS had not sincerely pondered and understood their problem and was only concerned about answering quickly; 2) They believed that the answer was auto generated because it is impossible for a human to read and respond to a message in such a short time. Therefore, it seemed like postponing the answer up to 4s is not strenuous in the context of CPCS. In case of longer responses, one prominent solution that was proposed by participants was using filler or holding messages (e.g., "please bear with me while I think about that")—Figure 11. Thus, the designers of CPCS can add these filler messages after certain time period to let the users believe that CPCS is still actively listening, in addition to more traditional graphical indicators (Figure 12). We propose that after a 4s period, a filler message should be displayed to acknowledge the question, and to request the user to wait until receiving an answer from the crowd. Ideally, in the next 4 to 8s, an answer should be received from the crowd. If no answer is received within 8s, then the CPCS can wait for a couple of more seconds (based on the theory proposed by Nielsen [85]) for the answer. If no answer is received within 10s, another filler message can be displayed that can explain the reasons for the delay and ask the user to further wait.

Workers in CPCS can at times take longer than 16s to respond due to variety of reasons that are outside their control, such as slow typing speed, fatigue, a poor Internet connection and other technical issues. In these scenarios, filler messages can be augmented with other options, such as those proposed by participants in our study. For instance, in case of search tasks that take longer to

1618 respond, CPCS can offer users to solve a quiz or a puzzle to sooth the effects of waiting. In case of
1619 stress tasks, it is more relevant to augment filler messages with calming music or relaxing visual
1620 images/animations that showed effectiveness in relieving stress [118]. Furthermore, since the CPCS
1621 are relying on humans, one can infer the expected waiting time based of the worker's typing speed
1622 and the average number of words that the worker has earlier written. Based on this, one can display
1623 a timer or progress bar to give a waiting estimate.

1624 To ensure and maintain high quality input from crowd workers, control mechanisms have already
1625 been proposed [6]. For instance, Chorus [69] implements an explicit voting method where workers
1626 vote on the messages proposed by their fellow workers, while Evorus [48] implements an automatic
1627 voting technique based on machine learning models to improve the quality of responses over time.
1628 Although these methods help to improve the response quality to some extent, they can be ineffective.
1629 For example, in the case there are many spammers or inexperienced workers because the majority
1630 voting would weigh the votes of all workers equally [56]. For instance, it is a challenging task for
1631 an unskilled worker to provide effective coaching to people who are distressed. Such coaching
1632 requires workers to practice a plethora of skills ranging from understanding a person's thoughts
1633 and feelings to deciding what psychological interventions to provide based on the situation. In
1634 a recent study [3], researchers explored the efficacy of a conversational user interface (CUI) for
1635 training crowd workers to deliver complex therapeutic tasks based on MI skills. Results show that
1636 workers who were trained through CUI provided better psychological interventions and felt higher
1637 self-efficacy for dealing with stress management tasks. Similar CUIs can be employed to train
1638 workers to retrieve high-quality answers in case of complex information finding tasks. Additionally,
1639 one can nudge workers to enhance their learning gain for complex conversational tasks by setting
1640 learning goals for those who have high learning goal orientation [97]. These solutions can help to
1641 improve the content quality of responses, and can consequently improve the waiting experience
1642 and engagement.

1643 Predicting complexity for complex conversational tasks could be beneficial for both workers and
1644 users of the CPCS system. For instance, if the CPCS is able to measure complexity of the utterance
1645 posed by the user, it can route the request to the workers who are fastest typists or assign the task
1646 to those who have the required skills or training to accomplish the task more quickly. Although
1647 various computation-based techniques have been developed to measure cognitive load of operators
1648 [60] or measuring the complexity of the individual messages in chatbots [39], it remains a challenge
1649 to assess complexity for diverse and complex conversational tasks. For example, in a stress task,
1650 merely measuring complexity based on the language used would not be sufficient since the task
1651 may involve a variety of stressors that would then define the complexity. Therefore, there is a need
1652 to quantify the stressful state context in stress mitigation tasks [88]. Similarly, IR tasks may also
1653 involve various dimensions that can define their complexity, such as the number of information
1654 resources needed to forge, the number of activities involved, the number of steps required and
1655 the expected workers' search behaviors in term of the query length, the number of URLs visited,
1656 etc [57]. Yang et al. [124] employed an extensive approach to quantify subjective complexity of
1657 crowdsourcing tasks based on some measurable properties of tasks, such as metadata features
1658 (e.g., title, descriptions etc.), content features (e.g., words count, links, and images etc.), and visual
1659 features (e.g., colorfulness, stylesheet etc.). Their proposed method was able to use these features
1660 to accurately predict the task complexity and task performance. Similar approaches are needed to
1661 quantify complexity of crowdsourced conversational tasks.

1662 In summary, our study highlights several key findings and design implications for CPCSs that
1663 are presented below:

1664
1665
1666

- ★ Users increasingly underestimated the waiting time when the delay was above 2s and were able to tolerate the delay up to 8s. Thus, researchers need to develop algorithms and workflows to keep the response latency within 4-8 seconds for CPCSs.
- ★ CPCSs can be equipped with different conversational fillers or acknowledgement tokens to mitigate the waiting time in CPCSs. For exceedingly longer delays, some non-temporal tasks such as solving a puzzle or an animation can be used to distract users.
- ★ Since our perception is our reality, we can aim for lowering the subjective waiting time in CPCS by uplifting the quality of response contents through workers' training, which does not require much engineering work as opposed to algorithmically lowering objective waiting time or system response time in CPCS.
- ★ An inverse relationship between task complexity and time perception holds within CPCS. Thus, computationally predicting the complexity for crowdsourced conversational tasks could be helpful in less complicated tasks to route the request to workers who are fastest typists or assign the task to those who have the required skills to accomplish the task more quickly.
- ★ People are less affected with waiting in task-centric conversations that are shallow in focus and need fewer exchanges compared to emotion-centric conversations that are deep in focus and require multiple exchanges in CPCS. In the former case, strict time constraints can be eased for crowd workers to generate answers.

7 LIMITATIONS

A limitation of the study concerns the nature of user responses. These were restricted to choosing between template-based utterances rather than typing in free responses. While this approach is adopted for reasons of ease of use in many chatbot applications, it is also not universally applicable especially in the cases where a free form conversation is intended, as is the case in cases of dealing with stress or other personal problems. Given the restricted nature of text entry in our experiments, our results may not easily generalize to cases of extended free form text input, where perhaps a largely varying time to compose an input or even the content of the system response may influence the expectations of the user regarding the response delay. Future studies could also address if the type of time-filler selected for different cases is considered suitable in different context and for different types of conversation. However, our study does inform what type of “filler” would be applicable for a certain latency.

We purposefully introduced complexity into the design to test the hypothesis that when complexity of the user input is low, participants will have less patience for extended delays than when complexity is high. As a result, we did not adjust the bot's response time based on the complexity of the user enquiry to determine how this would affect the perception of time. However, given the stress task's free-form nature, it is possible that users will occasionally enter simple responses (such as "hello") anticipating a prompt response. We controlled this concern in two ways in the Stress-bot: 1) we designed it in accordance with motivational interviewing literature, which encourages participants to ponder and enter practical responses rather than short inquiries (e.g., what concerns you about your problem?); 2) The Stress-bot initiated the conversation and greeted users, avoiding short inquiries from users, and asked a question in the beginning to break the ice (e.g., Tell me what brought you here today?). In the case of an IR-bot, the length of the user's response had no bearing on the outcome, as the user's response had been pre-programmed. We intend to do

1716 additional research to determine how varied lengths of user requests affect people’s perceptions of
 1717 CPCS delays in real-world contexts.

1718 Another limitation of our study is that we cannot ascertain how participants interpreted bots’
 1719 identities when they interacted with them. For instance, were they aware that the bots were pre-
 1720 programmed with pre-defined responses or that the responses were augmented in real time by
 1721 human intelligence? Future research will look into how the results vary between those who felt
 1722 they were interacting with AI alone versus those who thought they were interacting with a bot
 1723 powered by human intelligence.

1724

1725 **ACKNOWLEDGMENTS**

1726 This work was funded partially by Mirpur University of Science and Technology, Mirpur AJK,
 1727 Pakistan, under award number 611-19 P.D/2017, and in part by the Design@Scale Delft AI Lab, the
 1728 4TU.CEE UNCAGE project, and Eindhoven University of Technology, Netherlands. We would like
 1729 to express our gratitude to the reviewers for their helpful feedback, as well as to all participants in
 1730 our study.

1731

1732 **REFERENCES**

1733 [1] Tahir Abbas, Ujwal Gadiraju, Vassilis-Javed Khan, and Panos Markopoulos. 2021. Making Time Fly: Using Fillers to
 1734 Improve Perceived Latency in Crowd-Powered Conversational Systems. *Proceedings of the AAAI Conference on Human
 1735 Computation and Crowdsourcing* 9, 1 (Oct. 2021), 2–14. <https://ojs.aaai.org/index.php/HCOMP/article/view/18935>

1736 [2] Tahir Abbas, Vassilis-Javed Khan, Ujwal Gadiraju, Emilia Barakova, and Panos Markopoulos. 2020. Crowd of Oz: A
 1737 Crowd-Powered Social Robotics System for Stress Management. *Sensors* 20, 2 (2020). <https://doi.org/10.3390/s20020569>

1738 [3] Tahir Abbas, Vassilis-Javed Khan, Ujwal Gadiraju, and Panos Markopoulos. 2020. Trainbot: A Conversational
 1739 Interface to Train Crowd Workers for Delivering On-Demand Therapy. *Proceedings of the AAAI Conference on Human
 1740 Computation and Crowdsourcing* 8, 1 (Oct. 2020), 3–12. <https://ojs.aaai.org/index.php/HCOMP/article/view/7458>

1741 [4] Ritu Agarwal and Elena Karahanna. 2000. Time flies when you’re having fun: Cognitive absorption and beliefs about
 1742 information technology usage. *MIS quarterly* (2000), 665–694.

1743 [5] Hee-Kyung Ahn, Maggie Wenjing Liu, and Dilip Soman. 2006. Memory for Time: A Cognitive Model of Retrospective
 1744 Duration and Sequence Judgments. *Available at SSRN 897933* (2006).

1745 [6] M. Allahbakhsh, B. Benatallah, A. Ignjatovic, H. R. Motahari-Nezhad, E. Bertino, and S. Dustdar. 2013. Quality
 1746 Control in Crowdsourcing Systems: Issues and Directions. *IEEE Internet Computing* 17, 2 (2013), 76–81. <https://doi.org/10.1109/MIC.2013.20>

1747 [7] Anonymous. 2021. Making Time Fly: Using Fillers to Improve Perceived Latency in Crowd-Powered Conversational
 1748 Systems. (2021). https://osf.io/b4ket/?view_only=8e6b6ee82c764a678dba237dce265e64

1749 [8] Theo Araujo. 2018. Living up to the chatbot hype: The influence of anthropomorphic design cues and communicative
 1750 agency framing on conversational agent and company perceptions. *Computers in Human Behavior* 85 (2018), 183–189.
 1751 <https://doi.org/10.1016/j.chb.2018.03.051>

1752 [9] Zahra Ashktorab, Mohit Jain, Q. Vera Liao, and Justin D. Weisz. 2019. Resilient Chatbots: Repair Strategy Preferences
 1753 for Conversational Breakdowns. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*
 1754 (Glasgow, Scotland Uk) (CHI '19). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300484>

1755 [10] Ja Baptista-Neto, Fx Gingele, and T Leipe. [n.d.]. ZAR, JH, 1996. Biostatistical analysis, 3rd ed., New Jersey. *Environmental
 1756 Pollution* 109 ([n. d.]), 1–9.

1757 [11] Michael S. Bernstein, Joel Brandt, Robert C. Miller, and David R. Karger. 2011. Crowds in Two Seconds: Enabling
 1758 Realtime Crowd-Powered Interfaces. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software
 1759 and Technology* (Santa Barbara, California, USA) (UIST '11). Association for Computing Machinery, New York, NY,
 1760 USA, 33–42. <https://doi.org/10.1145/2047196.2047201>

1761 [12] Jeffrey P. Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C. Miller, Robin Miller, Aubrey
 1762 Tatarowicz, Brandyn White, Samuel White, and Tom Yeh. 2010. VizWiz: Nearly Real-Time Answers to Visual
 1763 Questions. In *Proceedings of the 23rd Annual ACM Symposium on User Interface Software and Technology* (New
 1764 York, New York, USA) (UIST '10). Association for Computing Machinery, New York, NY, USA, 333–342. <https://doi.org/10.1145/1866029.1866080>

[13] Florian Block and Hans Gellersen. 2010. The Impact of Cognitive Load on the Perception of Time. In *Proceedings of
 the 6th Nordic Conference on Human-Computer Interaction: Extending Boundaries* (Reykjavik, Iceland) (NordCHI '10).

- 1765 Association for Computing Machinery, New York, NY, USA, 607–610. <https://doi.org/10.1145/1868914.1868985>
- 1766 [14] Richard A Block, Peter A Hancock, and Dan Zakay. 2000. Sex differences in duration judgments: A meta-analytic
1767 review. *Memory & Cognition* 28, 8 (2000), 1333–1346.
- 1768 [15] Margaret Boden, Joanna Bryson, Darwin Caldwell, Kerstin Dautenhahn, Lilian Edwards, Sarah Kember, Paul Newman,
1769 Vivienne Parry, Geoff Pegman, Tom Rodden, Tom Sorrell, Mick Wallis, Blay Whitby, and Alan Winfield. 2017. Principles
1770 of robotics: regulating robots in the real world. *Connection Science* 29, 2 (2017), 124–129. <https://doi.org/10.1080/09540091.2016.1271400> arXiv:<https://doi.org/10.1080/09540091.2016.1271400>
- 1771 [16] Ria Mae Borromeo, Thomas Laurent, and Motomichi Toyama. 2016. The Influence of Crowd Type and Task Complex-
1772 ity on Crowdsourced Work Quality. In *Proceedings of the 20th International Database Engineering & Applications*
1773 *Symposium* (Montreal, QC, Canada) (*IDEAS '16*). Association for Computing Machinery, New York, NY, USA, 70–76.
1774 <https://doi.org/10.1145/2938503.2938511>
- 1775 [17] Virginia Braun and Victoria Clarke. 2012. Thematic Analysis. *Handbook of Research Methods in psychology* 2 (2012).
1776 <https://doi.org/10.1037/13620-004>
- 1777 [18] Thomas W Butler. 1983. Computer response time and user performance.. In *Proceedings of the SIGCHI conference on*
1778 *Human Factors in Computing Systems*. 58–62.
- 1779 [19] Jessy Ceha, Ken Jen Lee, Elizabeth Nilsen, Joslin Goh, and Edith Law. 2021. *Can a Humorous Conversational*
1780 *Agent Enhance Learning Experience and Outcomes?* Association for Computing Machinery, New York, NY, USA.
1781 <https://doi.org/10.1145/3411764.3445068>
- 1782 [20] Ki Hun Cho, Min Kyu Kim, Hwang-Jae Lee, and Wan Hee Lee. 2015. Virtual reality training with cognitive load
1783 improves walking function in chronic stroke patients. *The Tohoku journal of experimental medicine* 236, 4 (2015),
1784 273–280.
- 1785 [21] Alan J Christensen, Patricia J Moran, John S Wiebe, Shawna L Ehlers, and William J Lawton. 2002. Effect of a
1786 behavioral self-regulation intervention on patient adherence in hemodialysis. *Health Psychology* 21, 4 (2002), 393.
- 1787 [22] Jacob Cohen. 2013. *Statistical power analysis for the behavioral sciences*. Academic press.
- 1788 [23] Benilda Eleonor V Comendador, Bien Michael B Francisco, Jefferson S Medenilla, and Sharleen Mae. 2015. Pharmabot:
1789 a pediatric generic medicine consultant chatbot. *Journal of Automation and Control Engineering* 3, 2 (2015).
- 1790 [24] Mihaly Csikszentmihalyi. 1990. *Flow: The psychology of optimal experience*. New York: Harper & Row.
- 1791 [25] Tom Van Daele, Dirk Hermans, Chantal Van Audenhove, and Omer Van den Bergh. 2012. Stress Reduction Through
1792 Psychoeducation: A Meta- Analytic Review. *Health Education & Behavior* 39, 4 (2012), 474–485. <https://doi.org/10.1177/1090198111419202> PMID: 21986242.
- 1793 [26] Florian Daniel, Pavel Kucherbaev, Cinzia Cappiello, Boualem Benatallah, and Mohammad Allahbakhsh. 2018. Quality
1794 Control in Crowdsourcing: A Survey of Quality Attributes, Assessment Techniques, and Assurance Actions. *ACM*
1795 *Comput. Surv.* 51, 1, Article 7 (Jan. 2018), 40 pages. <https://doi.org/10.1145/3148148>
- 1796 [27] Mark M Davis and Janelle Heineke. 1998. How disconfirmation, perception and actual waiting times impact customer
1797 satisfaction. *international Journal of Service industry Management* (1998).
- 1798 [28] Zoltan Dienes. 2014. Using Bayes to get the most out of non-significant results. *Frontiers in psychology* 5 (2014), 781.
- 1799 [29] Kent Drummond and Robert Hopper. 1993. Acknowledgment tokens in series. *Communication Reports* 6, 1 (1993),
1800 47–53. <https://doi.org/10.1080/08934219309367561> arXiv:<https://doi.org/10.1080/08934219309367561>
- 1801 [30] Yingling Fan, Andrew Guthrie, and David Levinson. 2016. Waiting time perceptions at transit stops and stations:
1802 Effects of basic amenities, gender, and security. *Transportation Research Part A: Policy and Practice* 88 (2016), 251–264.
1803 <https://doi.org/10.1016/j.tra.2016.04.012>
- 1804 [31] Franz Faul, Edgar Erdfelder, Axel Buchner, and Albert-Georg Lang. 2009. Statistical power analyses using G* Power
1805 3.1: Tests for correlation and regression analyses. *Behavior research methods* 41, 4 (2009), 1149–1160.
- 1806 [32] Franz Faul, Edgar Erdfelder, Albert-Georg Lang, and Axel Buchner. 2007. G* Power 3: A flexible statistical power
1807 analysis program for the social, behavioral, and biomedical sciences. *Behavior research methods* 39, 2 (2007), 175–191.
- 1808 [33] Thomas Franke, Christiane Attig, and Daniel Wessel. 2019. A Personal Resource for Technology Interaction:
1809 Development and Validation of the Affinity for Technology Interaction (ATI) Scale. *International Journal*
1810 *of Human–Computer Interaction* 35, 6 (2019), 456–467. <https://doi.org/10.1080/10447318.2018.1456150>
1811 arXiv:<https://doi.org/10.1080/10447318.2018.1456150>
- 1812 [34] Norman M. Fraser and G.Nigel Gilbert. 1991. Simulating speech systems. *Computer Speech & Language* 5, 1 (1991),
1813 81–99. [https://doi.org/10.1016/0885-2308\(91\)90019-M](https://doi.org/10.1016/0885-2308(91)90019-M)
- [35] Peter Fröhlich. 2005. Dealing with system response times in interactive speech applications. In *CHI'05 Extended*
Abstracts on Human Factors in Computing Systems. 1379–1382.
- [36] Ujwal Gadiraju and Stefan Dietze. 2017. Improving Learning through Achievement Priming in Crowdsourced
Information Finding Microtasks. In *Proceedings of the Seventh International Learning Analytics & Knowledge Conference*
(Vancouver, British Columbia, Canada) (*LAK '17*). Association for Computing Machinery, New York, NY, USA, 105–114.
<https://doi.org/10.1145/3027385.3027402>

- 1814 [37] Ujwal Gadiraju, Jie Yang, and Alessandro Bozzon. 2017. Clarity is a Worthwhile Quality: On the Role of Task
1815 Clarity in Microtask Crowdsourcing. In *Proceedings of the 28th ACM Conference on Hypertext and Social Media*
1816 (Prague, Czech Republic) (*HT '17*). Association for Computing Machinery, New York, NY, USA, 5–14. <https://doi.org/10.1145/3078714.3078715>
- 1817 [38] Joseph Glicksohn. 2001. Temporal Cognition and the Phenomenology of Time: A Multiplicative Function for Apparent
1818 Duration. *Consciousness and Cognition* 10, 1 (2001), 1–25. <https://doi.org/10.1006/ccog.2000.0468>
- 1819 [39] Ulrich Gnewuch, Stefan Morana, Marc Adam, and Alexander Maedche. 2018. Faster is not always better: understanding
1820 the effect of dynamic response delays in human-chatbot interaction. (2018).
- 1821 [40] Gerald J Gorn, Amitava Chattopadhyay, Jaideep Sengupta, and Shashank Tripathi. 2004. Waiting for the web: how
1822 screen color affects time perception. *Journal of marketing research* 41, 2 (2004), 215–225.
- 1823 [41] Karen Grace-Martin. 2020. When Unequal Sample Sizes Are and Are NOT a Problem in ANOVA - The Analysis
1824 Factor. <https://www.theanalysisfactor.com/when-unequal-sample-sizes-are-and-are-not-a-problem-in-anova/>
- 1825 [42] Jonathan Grudin and Richard Jacques. 2019. Chatbots, Humbots, and the Quest for Artificial General Intelligence. In
1826 *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland UK) (*CHI '19*).
1827 Association for Computing Machinery, New York, NY, USA, 1–11. <https://doi.org/10.1145/3290605.3300439>
- 1828 [43] Daniel Haas, Jiannan Wang, Eugene Wu, and Michael J. Franklin. 2015. CLAMShell: Speeding up Crowds for
1829 Low-Latency Data Labeling. *Proc. VLDB Endow.* 9, 4 (Dec. 2015), 372–383. <https://doi.org/10.14778/2856318.2856331>
- 1830 [44] John A Hoxmeier and Chris DiCesare. 2000. System response time and user satisfaction: An experimental study of
1831 browser-based applications. *AMCIS 2000 Proceedings* (2000), 347.
- 1832 [45] Ting-Hao Huang and Jeffrey Bigham. 2017. A 10-Month-Long Deployment Study of On-Demand Recruiting for
1833 Low-Latency Crowdsourcing. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 5, 1
1834 (Sep. 2017). <https://ojs.aaai.org/index.php/HCOMP/article/view/13318>
- 1835 [46] Ting-Hao Huang, Walter Lasecki, Amos Azaria, and Jeffrey Bigham. 2016. "Is There Anything Else I Can Help You
1836 With?": Challenges in Deploying an On-Demand Crowd-Powered Conversational Agent. *Proceedings of the AAAI*
1837 *Conference on Human Computation and Crowdsourcing* 4, 1 (Sep. 2016), 79–88. <https://ojs.aaai.org/index.php/HCOMP/article/view/13292>
- 1838 [47] Ting-Hao Huang, Walter Lasecki, and Jeffrey Bigham. 2015. Guardian: A Crowd-Powered Spoken Dialog System for
1839 Web APIs. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 3, 1 (Sep. 2015).
1840 <https://ojs.aaai.org/index.php/HCOMP/article/view/13237>
- 1841 [48] Ting-Hao (Kenneth) Huang, Joseph Chee Chang, and Jeffrey P. Bigham. 2018. Evorus: A Crowd-Powered Conversational
1842 Assistant Built to Automate Itself Over Time. In *Proceedings of the 2018 CHI Conference on Human Factors in*
1843 *Computing Systems* (Montreal QC, Canada) (*CHI '18*). Association for Computing Machinery, New York, NY, USA,
1844 1–13. <https://doi.org/10.1145/3173574.3173869>
- 1845 [49] Ting-Hao K. Huang, Amos Azaria, Oscar J. Romero, and Jeffrey P. Bigham. 2019. InstructableCrowd: Creating
1846 IF-THEN Rules for Smartphones via Conversations with the Crowd. *Hum. Comput.* 6 (2019), 113–146. <https://doi.org/10.15346/hc.v6i1.7>
- 1847 [50] Ting-Hao K. Huang, Yun-Nung Chen, and Jeffrey P. Bigham. 2017. Real-time On-Demand Crowd-powered Entity
1848 Extraction. *CoRR abs/1704.03627* (2017). arXiv:1704.03627 <http://arxiv.org/abs/1704.03627>
- 1849 [51] Michael K. Hui and David K. Tse. 1996. What to Tell Consumers in Waits of Different Lengths: An Integrative
1850 Model of Service Evaluation. *Journal of Marketing* 60, 2 (1996), 81–90. <https://doi.org/10.1177/002224299606000206>
1851 arXiv:<https://doi.org/10.1177/002224299606000206>
- 1852 [52] Panagiotis G. Ipeirotis. 2010. Analyzing the Amazon Mechanical Turk Marketplace. *XRDS* 17, 2 (Dec. 2010), 16–21.
1853 <https://doi.org/10.1145/1869086.1869094>
- 1854 [53] Julie A Jacko, Andrew Sears, and Michael S Borella. 2000. The effect of network delay and media on user perceptions
1855 of web resources. *Behaviour & Information Technology* 19, 6 (2000), 427–439.
- 1856 [54] Gunn Johansson and Gunnar Aronsson. 1984. Stress reactions in computerized administrative work. *Journal of*
1857 *organizational behavior* 5, 3 (1984), 159–181.
- 1858 [55] Sepandar D Kamvar and Jonathan Harris. 2011. We feel fine and searching the emotional web. In *Proceedings of the*
1859 *fourth ACM international conference on Web search and data mining*. 117–126.
- 1860 [56] David R Karger, Sewoong Oh, and Devavrat Shah. 2011. Iterative learning for reliable crowdsourcing systems. *Neural*
1861 *Information Processing Systems*, 1953–1961.
- 1862 [57] Diane Kelly, Jaime Arguello, Ashlee Edwards, and Wan-ching Wu. 2015. Development and Evaluation of Search Tasks
for IIR Experiments Using a Cognitive Complexity Framework. In *Proceedings of the 2015 International Conference*
on *The Theory of Information Retrieval* (Northampton, Massachusetts, USA) (*ICTIR '15*). Association for Computing
Machinery, New York, NY, USA, 101–110. <https://doi.org/10.1145/2808194.2809465>
- [58] Deanna S. Kempf. 1999. Attitude formation from product trial: Distinct roles of cognition and affect for hedonic and functional products. *Psychology & Marketing* 16, 1 (1999), 35–50. [https://doi.org/10.1002/\(SICI\)1520-6793\(199901\)16:](https://doi.org/10.1002/(SICI)1520-6793(199901)16:)

- 1863 1<35::AID-MAR3>3.0.CO;2-U
- 1864 [59] Azizuddin Khan, Narendra K Sharma, and Shikha Dixit. 2006. Effect of cognitive load and paradigm on time perception. *Journal of the Indian Academy of applied Psychology* 32, 1 (2006), 37–42.
- 1865 [60] M. Asif Khawaja, Fang Chen, and Nadine Marcus. 2010. Using Language Complexity to Measure Cognitive Load for Adaptive Interaction Design. In *Proceedings of the 15th International Conference on Intelligent User Interfaces* (Hong Kong, China) (*IUI '10*). Association for Computing Machinery, New York, NY, USA, 333–336. <https://doi.org/10.1145/1719970.1720024>
- 1866 [61] Aniket Kittur, Ed H. Chi, and Bongwon Suh. 2008. Crowdsourcing User Studies with Mechanical Turk. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Florence, Italy) (*CHI '08*). Association for Computing Machinery, New York, NY, USA, 453–456. <https://doi.org/10.1145/1357054.1357127>
- 1870 [62] Kira Kretzschmar, Holly Tyroll, Gabriela Pavarini, Arianna Manzini, Iina Singh, and NeurOx Young People's Advisory Group. 2019. Can Your Phone Be Your Therapist? Young People's Ethical Perspectives on the Use of Fully Automated Conversational Agents (Chatbots) in Mental Health Support. *Biomedical Informatics Insights* 11 (2019), 1178222619829083. <https://doi.org/10.1177/1178222619829083> arXiv:<https://doi.org/10.1177/1178222619829083> PMID: 30858710.
- 1871 [63] Daniël Lakens. 2013. Calculating and reporting effect sizes to facilitate cumulative science: a practical primer for t-tests and ANOVAs. *Frontiers in psychology* 4 (2013), 863.
- 1872 [64] Carine Lallemand and Guillaume Gronier. 2012. Enhancing User EXperience during Waiting Time in HCI: Contributions of Cognitive Psychology. In *Proceedings of the Designing Interactive Systems Conference* (Newcastle Upon Tyne, United Kingdom) (*DIS '12*). Association for Computing Machinery, New York, NY, USA, 751–760. <https://doi.org/10.1145/2317956.2318069>
- 1873 [65] Walter Lasecki, Christopher Miller, Adam Sadilek, Andrew Abumoussa, Donato Borrello, Raja Kushalnagar, and Jeffrey Bigham. 2012. Real-Time Captioning by Groups of Non-Experts. In *Proceedings of the 25th Annual ACM Symposium on User Interface Software and Technology* (Cambridge, Massachusetts, USA) (*UIST '12*). Association for Computing Machinery, New York, NY, USA, 23–34. <https://doi.org/10.1145/2380116.2380122>
- 1874 [66] Walter S Lasecki, Christopher Homan, and Jeffrey P Bigham. 2014. Architecting real-time crowd-powered systems. *Human Computation* 1, 1 (2014).
- 1875 [67] Walter S. Lasecki, Kyle I. Murray, Samuel White, Robert C. Miller, and Jeffrey P. Bigham. 2011. Real-Time Crowd Control of Existing Interfaces. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology* (Santa Barbara, California, USA) (*UIST '11*). Association for Computing Machinery, New York, NY, USA, 23–32. <https://doi.org/10.1145/2047196.2047200>
- 1876 [68] Walter S. Lasecki, Phyo Thiha, Yu Zhong, Erin Brady, and Jeffrey P. Bigham. 2013. Answering Visual Questions with Conversational Crowd Assistants. In *Proceedings of the 15th International ACM SIGACCESS Conference on Computers and Accessibility* (Bellevue, Washington) (*ASSETS '13*). Association for Computing Machinery, New York, NY, USA, Article 18, 8 pages. <https://doi.org/10.1145/2513383.2517033>
- 1877 [69] Walter S. Lasecki, Rachel Wesley, Jeffrey Nichols, Anand Kulkarni, James F. Allen, and Jeffrey P. Bigham. 2013. Chorus: A Crowd-Powered Conversational Assistant. In *Proceedings of the 26th Annual ACM Symposium on User Interface Software and Technology* (St. Andrews, Scotland, United Kingdom) (*UIST '13*). Association for Computing Machinery, New York, NY, USA, 151–162. <https://doi.org/10.1145/2501988.2502057>
- 1878 [70] Younghwa Lee, Andrew N. K. Chen, and Virginia Ilie. 2012. Can Online Wait Be Managed? The Effect of Filler Interfaces and Presentation Modes on Perceived Waiting Time Online. *MIS Quarterly* 36, 2 (2012), 365–394. <http://www.jstor.org/stable/41703460>
- 1879 [71] Q. Vera Liao, Matthew Davis, Werner Geyer, Michael Muller, and N. Sadat Shami. 2016. What Can You Do? Studying Social-Agent Orientation and Agent Proactive Interactions with an Agent for Employees. In *Proceedings of the 2016 ACM Conference on Designing Interactive Systems* (Brisbane, QLD, Australia) (*DIS '16*). Association for Computing Machinery, New York, NY, USA, 264–275. <https://doi.org/10.1145/2901790.2901842>
- 1880 [72] Jake Linardon and Matthew Fuller-Tyszkiewicz. 2020. Attrition and adherence in smartphone-delivered interventions for mental health problems: A systematic and meta-analytic review. *Journal of consulting and clinical psychology* 88, 1 (2020), 1.
- 1881 [73] Andrey Lovakov and Elena Agadullina. 2017. Empirically derived guidelines for interpreting effect size in social psychology. (2017).
- 1882 [74] Kate MacFarlane and Lisa Marks Woolfson. 2013. Teacher attitudes and behavior toward the inclusion of children with social, emotional and behavioral difficulties in mainstream schools: An application of the theory of planned behavior. *Teaching and teacher education* 29 (2013), 46–52.
- 1883 [75] Robert B Miller. 1968. Response time in man-computer conversational transactions. In *Proceedings of the December 9-11, 1968, fall joint computer conference, part I*. 267–277.

- 1912 [76] Robert B. Miller. 1968. Response Time in Man-Computer Conversational Transactions. In *Proceedings of the December*
 1913 *9-11, 1968, Fall Joint Computer Conference, Part I* (San Francisco, California) (AFIPS '68 (Fall, part I)). Association for
 1914 Computing Machinery, New York, NY, USA, 267–277. <https://doi.org/10.1145/1476589.1476628>
- 1915 [77] William R Miller, Theresa B Moyers, Denise Ernst, and Paul Amrhein. 2003. Manual for the motivational interviewing
 1916 skill code (MISC). *Unpublished manuscript. Albuquerque: Center on Alcoholism, Substance Abuse and Addictions,*
University of New Mexico (2003).
- 1917 [78] William R Miller and Stephen Rollnick. 2012. *Motivational interviewing: Helping people change*. Guilford press.
- 1918 [79] Kathy Missildine, Rebecca Fountain, Lynn Summers, and Kevin Gosselin. 2013. Flipping the classroom to improve
 1919 student performance and satisfaction. *Journal of Nursing Education* 52, 10 (2013), 597–599.
- 1920 [80] Youngme Moon. 1999. The effects of physical distance and response latency on persuasion in computer-mediated
 1921 communication and human-computer communication. *Journal of Experimental Psychology: Applied* 5, 4 (1999), 379.
- 1922 [81] Masahiro Mori, Karl F MacDorman, and Norri Kageki. 2012. The uncanny valley [from the field]. *IEEE Robotics &*
Automation Magazine 19, 2 (2012), 98–100.
- 1923 [82] Robert Morris. 2011. Crowdsourcing workshop: the emergence of affective crowdsourcing. In *Proceedings of the 2011*
 1924 *annual conference extended abstracts on Human factors in computing systems*. ACM, Citeseer.
- 1925 [83] Robert Robert Randall Morris. 2015. *Crowdsourcing mental health and emotional well-being*. Ph.D. Dissertation.
 1926 Massachusetts Institute of Technology.
- 1927 [84] Alessandro Murgia, Daan Janssens, Serge Demeyer, and Bogdan Vasilescu. 2016. Among the Machines: Human-Bot
 1928 Interaction on Social Q&A Websites. In *Proceedings of the 2016 CHI Conference Extended Abstracts on Human Factors*
 1929 *in Computing Systems* (San Jose, California, USA) (CHI EA '16). Association for Computing Machinery, New York, NY,
 1930 USA, 1272–1279. <https://doi.org/10.1145/2851581.2892311>
- 1931 [85] Jakob Nielsen. 1994. *Usability engineering*. Morgan Kaufmann.
- 1932 [86] University of Gothenburg. 2015. Pauses can make or break a conversation. [www.sciencedaily.com/releases/2015/09/](http://www.sciencedaily.com/releases/2015/09/150930110555.htm)
 1933 [150930110555.htm](http://www.sciencedaily.com/releases/2015/09/150930110555.htm). Accessed: 2021-03-18.
- 1934 [87] Adam Palanica, Peter Flaschner, Anirudh Thommandram, Michael Li, and Yan Fossat. 2019. Physicians' Per-
 1935 ceptions of Chatbots in Health Care: Cross-Sectional Web-Based Survey. *J Med Internet Res* 2019;21(4):e12887
 1936 <https://www.jmir.org/2019/4/e12887> 21, 4 (apr 2019), e12887. <https://doi.org/10.2196/12887>
- 1937 [88] SoHyun Park, Jeewon Choi, Sungwoo Lee, Changhoon Oh, Changdai Kim, Soohyun La, Joonhwan Lee, and Bongwon
 1938 Suh. 2019. Designing a Chatbot for a Brief Motivational Interview on Stress Management: Qualitative Case Study. *J*
 1939 *Med Internet Res* 21, 4 (16 Apr 2019), e12231. <https://doi.org/10.2196/12231>
- 1940 [89] Zhenhui Peng, Kaixiang Mo, Xiaogang Zhu, Junlin Chen, Zhijun Chen, Qian Xu, and Xiaojuan Ma. 2020. Understanding
 1941 User Perceptions of Robot's Delay, Voice Quality-Speed Trade-off and GUI during Conversation. In *Extended Abstracts*
 1942 *of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI EA '20). Association for
 1943 Computing Machinery, New York, NY, USA, 1–8. <https://doi.org/10.1145/3334480.3382792>
- 1944 [90] Verónica Pérez-Rosas, Xinyi Wu, Kenneth Resnicow, and Rada Mihalcea. 2019. What Makes a Good Counselor?
 1945 Learning to Distinguish between High-quality and Low-quality Counseling Conversations. In *Proceedings of the 57th*
 1946 *Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence,
 1947 Italy, 926–935. <https://doi.org/10.18653/v1/P19-1088>
- 1948 [91] Olga Perski, Ann Blandford, Robert West, and Susan Michie. 2017. Conceptualising engagement with digital behaviour
 1949 change interventions: a systematic review using principles from critical interpretive synthesis. *Translational behavioral*
 1950 *medicine* 7, 2 (2017), 254–267.
- 1951 [92] Olga Perski, David Crane, Emma Beard, and Jamie Brown. 2019. Does the addition of a supportive chatbot promote
 1952 user engagement with a smoking cessation app? An experimental study. *Digital health* 5 (2019), 2055207619880676.
- 1953 [93] Ad Pruyn and Ale Smidts. 1998. Effects of waiting on the satisfaction with the service: Beyond objective time
 1954 measures1Both authors contributed equally to this article.1. *International Journal of Research in Marketing* 15, 4
 1955 (1998), 321 – 334. [https://doi.org/10.1016/S0167-8116\(98\)00008-1](https://doi.org/10.1016/S0167-8116(98)00008-1)
- 1956 [94] A Th H Pruyn and Ale Smidts. 1993. Customers' evaluations of queues: Three exploratory studies. *ACR European*
 1957 *Advances* (1993).
- 1958 [95] Sihang Qiu, Ujwal Gadiraju, and Alessandro Bozzon. 2020. Improving Worker Engagement Through Conversational
 1959 Microtask Crowdsourcing. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*
 1960 (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3313831.3376403>
- [96] Sihang Qiu, Ujwal Gadiraju, and Alessandro Bozzon. 2020. TickTalkTurk: Conversational Crowdsourcing Made Easy.
 In *Conference Companion Publication of the 2020 on Computer Supported Cooperative Work and Social Computing*. 1–5.
- [97] Amy Rechkemmer and Ming Yin. 2020. Motivating Novice Crowd Workers through Goal Setting: An Investigation
 into the Effects on Complex Crowdsourcing Task Training. *Proceedings of the AAAI Conference on Human Computation*
 and Crowdsourcing 8, 1 (Oct. 2020), 122–131. <https://ojs.aaai.org/index.php/HCOMP/article/view/7470>

- [98] René Riedl and Thomas Fischer. 2018. System response time as a stressor in a digital world: literature review and theoretical model. In *International Conference on HCI in Business, Government, and Organizations*. Springer, 175–186.
- [99] Laurel D. Riek. 2012. Wizard of Oz Studies in HRI: A Systematic Review and New Reporting Guidelines. *J. Hum.-Robot Interact.* 1, 1 (July 2012), 119–136. <https://doi.org/10.5898/JHRI.1.1.Riek>
- [100] Carl R Rogers and Richard E Farson. 1957. *Active listening*. Chicago, IL.
- [101] Emma J. Rose and Elin A. Björling. 2017. Designing for Engagement: Using Participatory Design to Develop a Social Robot to Measure Teen Stress. In *Proceedings of the 35th ACM International Conference on the Design of Communication* (Halifax, Nova Scotia, Canada) (*SIGDOC '17*). Association for Computing Machinery, New York, NY, USA, Article 7, 10 pages. <https://doi.org/10.1145/3121113.3121212>
- [102] Ayelet Meron Ruscio and Dana Rabois Holohan. 2006. Applying empirically supported treatments to complex cases: Ethical, empirical, and practical considerations. *Clinical Psychology: Science and Practice* 13, 2 (2006), 146–162.
- [103] Denis Savenkov and Eugene Agichtein. 2016. CRQA: Crowd-Powered Real-Time Automatic Question Answering System. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 4, 1 (Sep. 2016). <https://ojs.aaai.org/index.php/HCOMP/article/view/13291>
- [104] Lawrence M Schleifer and Benjamin C Amick III. 1989. System response time and method of pay: Stress effects in computer-based tasks. *International Journal of Human-Computer Interaction* 1, 1 (1989), 23–39.
- [105] Ryan M Schuetzler, Mark Grimes, Justin Scott Giboney, and Joesph Buckman. 2014. Facilitating natural conversational agent interactions: lessons from a deception experiment. (2014).
- [106] Stela H Seo, Keelin Griffin, James E Young, Andrea Bunt, Susan Prentice, and Verónica Loureiro-Rodríguez. 2018. Investigating people’s rapport building and hindering behaviors when working with a collaborative robot. *International Journal of Social Robotics* 10, 1 (2018), 147–161.
- [107] Steven C Seow. 2008. *Designing and engineering time: The psychology of time perception in software*. Addison-Wesley Professional.
- [108] T. Shiwa, T. Kanda, M. Imai, H. Ishiguro, and N. Hagita. 2008. How quickly should communication robots respond?. In *2008 3rd ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. 153–160. <https://doi.org/10.1145/1349822.1349843>
- [109] Ben Shneiderman. 1984. Response time and display rate in human performance with computers. *ACM Computing Surveys (CSUR)* 16, 3 (1984), 265–285.
- [110] Marita Skjuve, Ida Maria Haugstveit, Asbjørn Følstad, and Petter Bae Brandtzaeg. 2019. Help! Is my chatbot falling into the uncanny valley? An empirical study of user experience in human-chatbot interaction. *Human Technology* 15, 1 (2019).
- [111] T Smeets, J Leppink, M Jelcic, and HLGJ Merckelbach. 2009. Shortened versions of the Gudjonsson Suggestibility Scale meet the standards. *Legal and criminological Psychology* 14, 1 (2009), 149–155.
- [112] Shirley Taylor. 1994. The Effects of Filled Waiting Time and Service Provider Control over the Delay on Evaluations of Service. *Journal of the Academy of Marketing Science* 23, 1 (1994), 38–48. <https://doi.org/10.1177/0092070395231005> arXiv:<https://doi.org/10.1177/0092070395231005>
- [113] Ewart AC Thomas and Wanda B Weaver. 1975. Cognitive processing and time perception. *Perception & psychophysics* 17, 4 (1975), 363–367.
- [114] Michael Toomim, Travis Kriplean, Claus Pörtner, and James Landay. 2011. Utility of human-computer interactions: Toward a science of preference measurement. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 2275–2284.
- [115] Janet Treasure. 2004. Motivational interviewing. *Advances in Psychiatric Treatment* 10, 5 (2004), 331–337.
- [116] Hans Van Der Heijden. 2002. On the cognitive-affective structure of attitudes toward information systems. *ICIS 2002 Proceedings* (2002).
- [117] Andrea E Waylen, Mark S Horswill, Jane L Alexander, and Frank P McKenna. 2004. Do expert drivers have a reduced illusion of superiority? *Transportation research part F: traffic psychology and behaviour* 7, 4-5 (2004), 323–331.
- [118] Taelyr Weekly, Nicole Walker, Jill Beck, Sean Akers, and Meaghann Weaver. 2018. A Review of Apps for Calming, Relaxation, and Mindfulness Interventions for Pediatric Palliative Care Patients. *Children* 5, 2 (2018). <https://doi.org/10.3390/children5020016>
- [119] World Health Organization (WHO). 2017. Mental health: massive scale-up of resources needed if global targets are to be met. <https://www.who.int/news/item/06-06-2018-mental-health-massive-scale-up-of-resources-needed-if-global-targets-are-to-be-met>
- [120] Alex C. Williams, Harmanpreet Kaur, Gloria Mark, Anne Loomis Thompson, Shamsi T. Iqbal, and Jaime Teevan. 2018. *Supporting Workplace Detachment and Reattachment with Conversational Intelligence*. Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi-org.tudelft.idm.oclc.org/10.1145/3173574.3173662>

- 2010 [121] Philipp Wintersberger, Tobias Klotz, and Andreas Riener. 2020. Tell Me More: Transparency and Time-Fillers to
2011 Optimize Chatbots' Waiting Time Experience. In *Proceedings of the 11th Nordic Conference on Human-Computer Inter-*
2012 *action: Shaping Experiences, Shaping Society* (Tallinn, Estonia) (NordiCHI '20). Association for Computing Machinery,
2013 New York, NY, USA, Article 76, 6 pages. <https://doi.org/10.1145/3419249.3420170>
2014 [122] Robert E Wood. 1986. Task complexity: Definition of the construct. *Organizational behavior and human decision*
2015 *processes* 37, 1 (1986), 60–82.
2016 [123] Anbang Xu, Zhe Liu, Yufan Guo, Vibha Sinha, and Rama Akkiraju. 2017. A New Chatbot for Customer Service on Social
2017 Media. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (Denver, Colorado, USA) (CHI
2018 '17). Association for Computing Machinery, New York, NY, USA, 3506–3510. <https://doi.org/10.1145/3025453.3025496>
2019 [124] Jie Yang, Judith Redi, Gianluca Demartini, and Alessandro Bozzon. 2016. Modeling Task Complexity in Crowdsourcing.
2020 *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 4, 1 (Sep. 2016). [https://ojs.aaai.org/](https://ojs.aaai.org/index.php/HCOMP/article/view/13283)
2021 [index.php/HCOMP/article/view/13283](https://ojs.aaai.org/index.php/HCOMP/article/view/13283)
2022 [125] Xi Yang, Marco Aurisicchio, and Weston Baxter. 2019. Understanding Affective Experiences with Conversational
2023 Agents. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk)
2024 (CHI '19). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300772>
2025 [126] Zheng Zhang, Ryuichi Takanobu, Qi Zhu, MinLie Huang, and XiaoYan Zhu. 2020. Recent advances and challenges in
2026 task-oriented dialog systems. *Science China Technological Sciences* (2020), 1–17.
2027 [127] Yu Zhong, Walter S. Lasecki, Erin Brady, and Jeffrey P. Bigham. 2015. RegionSpeak: Quick Comprehensive Spatial
2028 Descriptions of Complex Images for Blind Users. In *Proceedings of the 33rd Annual ACM Conference on Human Factors*
2029 *in Computing Systems* (Seoul, Republic of Korea) (CHI '15). Association for Computing Machinery, New York, NY,
2030 USA, 2353–2362. <https://doi.org/10.1145/2702123.2702437>

A SAMPLE STRESS MANAGEMENT TASK

2030 Carla is a 21-year-old unmarried girl. She is a final year bachelor student studying Innovation
2031 Sciences. She is a very hard-working student and has received an 'A' grade on all of her
2032 assignments. Her teachers believe that she is one of the best students in the class. A new
2033 student named Hanna, has recently moved from another class to Carla's class. Hanna is
2034 also a very talented student and received 'A' grades on all of her assignments. The presence
2035 of Hanna is stressing Carla out. Prior to Hanna's arrival, Carla considered herself to be
2036 the best student in class, with the best grades. Now she feels like she needs to compete
2037 with Hanna and is wary of losing out on the 'Best Student Award', an award presented
2038 by the University to the student with the most stellar academic record at the end of the
2039 study program. Previously, Carla found assignments to be easy and straightforward but
2040 now she is obsessed about her performance in the assignments and tends to worry about
2041 unnecessary details. For instance, she is constantly sending emails to her teachers to clarify
2042 the assignments and to make sure that she is doing well according to the instructions. She
2043 believes that if she fails to earn that award, then she would be a failure. She finds this
2044 situation increasingly overwhelming, and this makes her feel distressed and worried. When
2045 Carla is stressed out, she sits in a quiet place and tells herself that this award is something
2046 confined to her university, and nobody would know or care about it in the grand scheme of
2047 things. So, it would not really impact her future employment prospects or what her peers
2048 may think of her. This strategy is called self-talk, and this relieves her stress to some degree.
2049

2050
2051
2052
2053
2054
2055
2056
2057
2058