

Article

INTERACTIVE HESITATION SYNTHESIS AND ITS EVALUATION

Simon Betz ^{1,2,4*}, Birte Carlmeyer ^{1,3,4}, Petra Wagner ^{1,2} and Britta Wrede ^{1,3}

¹ Cognitive Interaction Technology Center (CITEC), Bielefeld University

² Phonetics and Phonology Workgroup, Faculty of Linguistics and Literary Studies, Bielefeld University

³ Applied Informatics Group, Faculty of Technology, Bielefeld University

⁴ Dialogue Systems Group, Faculty of Linguistics and Literary Studies, Bielefeld University

* Correspondence: simon.betz@uni-bielefeld.de; Tel.: +49-521-106-3518

Academic Editor: name

Version December 8, 2017 submitted to *Multimodal Technologies and Interact.*; Typeset by L^AT_EX using class file mdpi.cls

Abstract: Conversational spoken dialogue systems that interact with the user rather than merely reading text can be equipped with hesitations to manage the dialogue flow and the users' attention. Based on a series of empirical studies, we built an elaborated hesitation synthesis strategy for dialogue systems that inserts hesitations of scalable extent wherever needed in the ongoing utterance. So far, evaluations of hesitating systems have shown that synthesis quality is affected negatively by hesitations, but that there is improvement in interaction quality. We argue that due to its conversational nature, hesitation synthesis needs interactive evaluation rather than traditional MOS-based questionnaires. To prove this point, we dually evaluate our system's speech synthesis component: on the one hand, linked to the dialogue system evaluation, on the other hand, in the traditional MOS way. This way we are able to analyze and discuss differences that arise due to the evaluation methodology. Our results suggest that MOS scales are not sufficient to assess speech synthesis quality, which has implications for future research that are discussed in this paper. Furthermore, our results indicate that hesitations work well to increase task performance and that an elaborated strategy is necessary to avoid likability issues.

Keywords: Speech Synthesis; Evaluation; Hesitation; Virtual Agents; Interaction.

1. Introduction

1.1. Motivation and aims of this study

Synthetic speech is applied in various fields and it has entered the realm of everyday life, be it in public transportation announcements, telephone customer services, smartphone speech output or smart-home environments. Despite the interactive nature of many of these applications, the speech output remains to be rather static, typically reading out pre-defined texts or often responding with an awkward delay.

Also, a special feature of synthetic speech is its "fluency", i.e. it does not contain the hesitations, reformulations or filled pauses typical for human spontaneous speech production. Rather, speech output, once generated, is produced in a single, non-interrupted fashion. The study we are presenting in this paper rests on the assumption that this is suboptimal for many human-machine interactions where listeners need to actually process information that is synthetically generated, and where a human speaker would try to deliver the information in a way which is suited to the listener's attention level, to enable him or her to follow and process what is being said (previously explored in [1,2]). In

30 order to test this assumption, we will explore the space of possible improvements of speech synthesis
31 for interactive purposes using synthesized hesitations.

32 Our assumption rests on the finding that the hesitations produced in spontaneous speech
33 communication are not merely disturbances or "errors" of human speech production. Rather, they
34 serve an important role in dialogue: They allow the speaker extra time in situations where this
35 is needed, e.g. when searching for the right thing to say, and to signal this to the listener. That
36 way, hesitations help to keep the metaphorical right to speak. It has been shown previously that
37 spoken dialogue systems can utilize hesitations to bridge gaps in dialogue, and to successfully handle
38 interruptions and attention shifts, e.g. [1–3].

39 In this study we explore the applicability of an elaborated hesitation synthesis strategy that is
40 based on observations of human hesitations. Upon an event of hesitation, a hierarchical hesitation
41 insertion is triggered that continues "buying time" as long as possible or until it receives a signal
42 for ending the hesitation mode. The start and stop signals for hesitation insertion are inferred from
43 user's attention: When users look away, the system will enter hesitation mode until users re-focus.
44 Furthermore, we test a model of optimal hesitation placement. Compared to previous hesitating
45 systems, the approach presented here allows for dynamic hesitation insertion in the middle of an
46 utterance and for flexible, scalable hesitation clusters tailored for hesitation events of various extents.

47 Due to the intrinsically interactive fashion of our hesitation strategy, its evaluation is not
48 straightforward. While the system as a whole can be evaluated with interactive measures such as task
49 performance, speech synthesis components are usually evaluated using non-interactive measures,
50 in which listeners are asked to rate the quality of synthetic speech, typically individual utterances,
51 using Mean Opinion Scores (MOS). Despite numerous criticisms for this method, alternatives have
52 seldom been proposed [4–8]. Also, to our knowledge, there exists no previous study that actually
53 verifies these critical voices. We therefore test our hesitation synthesis twice: First, we evaluate it in
54 direct connection with the dialogue system evaluation in interaction, and interpret objective measures
55 like task performance and efficiency alongside subjective user ratings of system features such as
56 synthesis quality. Second, we assess the subjective speech synthesis quality in a non-interactive,
57 crowdsourcing-based parallel study that uses the same stimuli. That way, we can compare user task
58 performance and their subjective impression of a system with subjective ratings where the interaction
59 quality is not part of the evaluation strategy. Ultimately, we hope to not only be able to evaluate our
60 own synthesis approach, but also shed light on the issue of what traditional approaches to speech
61 synthesis evaluation actually reveal.

62 1.2. Structure of this paper

63 In the following chapter, we provide further background for this study. First, we define the
64 term *hesitation* as we use it and give a brief overview of its research history (2.1). We continue
65 with the description of a model of incremental speech production, which serves as a foundation for
66 defining and discussing hesitations, incremental spoken dialogue systems and synthesis strategies
67 (2.2). We continue with a brief introduction to dialogue systems with a special focus on systems with
68 incremental processing, which we will work on in this study (2.3). With this foundation set up, we
69 turn towards our model for a hesitation synthesis strategy for incremental spoken dialogue systems
70 based on studies of human speech production (3). In the empirical part of this paper we present two
71 experiments that make up this study. First, we describe the methods and results of an interaction
72 study (4), continuing with a crowdsourcing-based parallel study for evaluation purposes (5). The
73 experiments sections are each concluded with short discussions. The main study is then concluded
74 with a general discussion of both experiments and their implications (6).

75 2. Background and Related Work

76 2.1. A brief introduction to hesitations

77 Hesitations are lexical and non-lexical elements that delay information delivery in speech. The
78 most common hesitations include fillers, silences, lengthenings and repetitions, cf. Example 1.¹

"Take these ... uhm ... the, the red line to the university"

Example 1: Different surface forms of hesitation: lengthening, silence, filler, repetition

79 It has been noted that hesitations do not only buy time, but that they are a useful strategy for both
80 speaker and listener to manage the conversation. [11] suggested that speakers intentionally decide
81 to produce a filler as either "uh" or "uhm", the former denoting only a small delay, the latter a major
82 problem accompanied by a longer pause in speech. This leads to the assumption that this difference
83 in form is a listener-oriented strategy, a means to ensure that the interlocutor is informed about the
84 dialogue state and does not attempt to grab the conversational floor too early.

85 It has further been observed that hesitations, with their property to control information timing
86 in dialogue, are linked to users visual attention. This relationship may be bilateral: [12] found that
87 speakers hesitate when the listener is apparently distracted, and [13,14] found that listeners may
88 heighten attention when a hesitation is uttered.

89 While it is highly controversial whether hesitations and disfluencies are produced in order to
90 signal something to the listener (see [15] for an overview), or if it is merely the fact that the listener can
91 do something with the information, it can be claimed that listeners can make use of at least the extra
92 time that hesitations grant in dialogue, an effect that is replicable for human-machine interactions,
93 e.g. [1–3,16].

94 Shifting the focus back to the speaker, with the aim in mind to adapt speaker strategies for
95 dialogue systems, we encounter several common reasons for hesitating. Speakers might have trouble
96 retrieving the correct, or the most appropriate item (cf. Example 2). They might run out of things
97 to say before having conveyed the intended message (cf. Example 3). The dialogical situation might
98 change, causing a change in speech plan, that needs time (cf. Example 4).

"The capital of Serbia is ... uhm ... Belgrade."

Example 2: Difficulty retrieving an infrequent lexical item.

"There is no direct flight to Sydney ... uhm ... today or tomorrow..."

Example 3: Travel agent giving information, but the database query takes time.

"You can take a seat ... in the living room."

Example 4: Originally, the plan was to offer a seat in the kitchen, but as the interlocutor apparently shifted her attention to the living room during the dialogue, a new speech plan was realized.

99 The above three are all fictional examples, but they shed light on the various usages of
100 hesitations. The surface forms might be indefinitely complex for every situation, with any

¹ Traditionally, hesitations are often associated with disfluencies. In this study we only consider hesitation phenomena. For an excellent overview on the historical entanglement of hesitations and disfluencies, see [9]. For the most influential descriptive work on disfluencies in general, see [10].

101 combination of the elements suggested in the introductory Example 1. The challenge in this study
102 will be to model plausible surface forms of hesitations for a dialogue system that can use them on the
103 fly whenever the situation requires it.

104 2.2. Incremental speech production

105 Hesitations are closely related to the way humans speak. When initiating an utterance, speakers
106 have not fully pre-planned what to say and how to say it. Instead, they plan and produce speech
107 *incrementally*, in a piece-meal fashion unfolding over time [17]. Doing so, speakers use and interpret
108 information from interlocutors rapidly and simultaneously formulate their own speech plan([17,18],
109 summarized in [3]). Despite the lack of a complete speech plan, human speech requires ahead
110 planning of a certain degree. Psycholinguistic studies suggest that speakers plan at least one word
111 ahead, usually more [3]. Evidence for the concept of *incremental processing* comes from several
112 observable phenomena of spontaneous speech, many of those closely related to hesitations:

- 113 • **Anticipatory speech production errors.** (e.g. "a cuff of coffee") where parts of the utterance are
114 produced in advance, or metathetically switched around, anticipating upcoming phonemes or
115 syllables.
- 116 • **Hesitation lengthening form in English.** ("Theee:" vs. "the") The lengthened form has a
117 different vowel quality (i: vs. ə), so the speaker must be aware of upcoming challenges in
118 the speech plan, cf. [19].
- 119 • **Different types of fillers.** ("uh" vs. "uhm") The former appears to denote minor, the latter major
120 problems in the speech plan, both requiring ahead planning [20].
- 121 • **Hesitation occurrence probability.** Hesitations are more likely to occur before longer
122 utterances [21].

123 Models of incremental speech production inspire the design of incremental spoken dialogue
124 systems, which will be further described in section 2.3. In this study, we investigate whether
125 human-like features that are typical of incremental processing, such as hesitation phenomena, are
126 suitable for dialogue systems as well. Special attention will be paid to the concept of the articulatory
127 buffer, which provides insights how to commence hesitation in incremental spoken dialogue systems.

128 The concept of the articulatory buffer was introduced in Levelt's model of speech production [18]
129 (p. 414) to describe the lookahead of several words that speakers have access to when speaking. It
130 describes a temporary storage for words that have been planned, but have not yet been articulated.
131 The content of the buffer can be amended when the speech plan changes. Based on [18] and [22], Li
132 and Tilsen [23] hypothesize that the material in the articulatory buffer can be lengthened by speakers
133 in order to buy time for solving word retrieval problems. We assume that this might not only be
134 the case for word retrieval issues, but make the proposition that this may hold as a general strategy
135 for phonetically producing hesitations. Based on this assumption, we present in this study a general
136 model for hesitation insertion in conversational dialogue systems.

137 2.3. Dialogue systems

138 Dialogue Systems are programs that communicate with users in text and/or speech form. They
139 are generally distinguished into task-oriented dialogue agents and chatbots. The latter are designed
140 for extensive conversations, for entertainment or practical application, traditionally in text form. The
141 former are designed to interact with the user in a limited domain in short task-oriented conversations,
142 for example to give directions or control home appliances. Well-known present-day examples would
143 be Siri, Alexa or Google Home. These current task-oriented dialogue systems are based on speech in-
144 and output. The scope of application is limited to small domains, but the interaction has become more
145 like spoken conversation between humans as more computational power and better speech synthesis
146 became available. One major shortcoming of these systems is their lack of adaptivity that contrasts
147 their field of application. They can only produce static responses, but are incapable of interpreting

148 user feedback or handling interruptions. It could thus be stated that these systems are less interactive
149 than they should be. They perform their tasks, but cannot do anything conversational beyond that.

150 Addressing the adaptivity and interactivity issue, a strand of research evolved that aims to
151 develop *conversational dialogue systems* that are capable of *talking* instead of merely *reading* out
152 pre-defined responses. One key feature on the way to more interactivity is incrementality.² As
153 described in section 2.2, human dialogue does not work like a ball-tossing game, but rather
154 simultaneously: Responses are planned while the interlocutor is speaking. It can be shown that
155 limited-domain dialogue systems can make use of incremental processing to achieve human-like
156 interaction speed [24].

157 Hesitations are a useful feature for incremental spoken dialogue systems. On the one hand, these
158 systems might need to buy time for re-planning and can use hesitations to do so. On the other hand,
159 the incrementality enables the system to hesitate immediately and flexibly. To develop conversational
160 dialogue systems, various approaches have been proposed, with incremental processing, with
161 various forms and functions of hesitation and with both incrementality and hesitations.

162 [3] built an incremental system based on general, abstract model for incremental processing [25]
163 that employs turn-initial hesitations ("eh...", "well...", "wait a minute...") to buy time to generate a
164 response (or in this case: time for the wizard to type the answer). This system exploits the fact
165 that hesitations do not commit content to the conversation, they can literally be used as fillers to
166 bridge gaps in dialogue. [26] conducted an experiment in a driving simulator, during which a virtual
167 assistant told the driver about appointments on that day. It was shown that a system that hesitates
168 by means of silences whenever a difficult situation occurs, improves both the participants driving
169 performance as well as their recall of information presented during the task. [27] uses hesitations
170 in human-robot interaction as a disengagement strategy. A directions-giving robot produces lexical
171 hesitations ("so...", "let's see...") after own speaking turns to bridge the awkward silence during which
172 the user has to decide whether she wants to continue the interaction or not. Interestingly, this usage
173 of hesitations is contrary to many other studies that highlight the usefulness of hesitations to gain
174 attention and to continue interacting.

175 [1,2] use hesitations (silence) as a user-oriented strategy, based on observations of the human
176 interaction partner. They investigated the effect of self-interruptions as a strategy to regain the visual
177 attention of distracted users in a smart-home setting with a virtual agent. They showed that insertion
178 of silence whenever the attention of the users shifts away, has a positive effect of the attention of
179 the user, but at the cost of less positive subjective ratings. In a similar scenario the authors could
180 show that incremental information presentation leads to a better task performance [28]. Whereas the
181 authors could show that listener-oriented insertion of hesitations (in this case: silences) has a positive
182 effect on the interaction, the self-interrupting agent was perceived less friendly in all three studies.
183 [16] found that hesitation lengthenings, as long as they are shorter than 800ms, have a positive effect
184 on users' task performance in a game setting.

185 All systems presented here reported positive effects on the interactivity. Not all systems
186 evaluated speech synthesis quality, but those that do report negative effects. This hints at a
187 shortcoming, a trade-off between interactivity and sound quality that is a key issue for current and
188 future research in this field. An off-line evaluation study [29] suggests, that different hesitations
189 strategies differ inherently with regard to sound quality: while lengthenings and silences get
190 relatively good user feedback (stimuli with lengthening got even better user feedback than fluent
191 baseline stimuli), fillers, and other disfluencies like mid-word cutoffs are dispreferred. The same
192 authors investigated the reasons for the good performance of lengthening and found a paradox
193 situation: the reason for the good rating of synthetic lengthening might be that they are barely

² In this study, we explore incremental spoken dialogue systems. It is worthwhile noting that it was recently demonstrated that an interactive system capable of handling interruptions can be built without incremental processing [7].

194 perceivable. In a follow-up study [30] showed that even corpus annotators with the task to label
 195 disfluencies miss up to 80% of lengthening instances that can be identified with semi-automatic
 196 classification. This makes lengthening a promising candidate for application in conversational
 197 dialogue systems.

198 Based on the assumption that the underlying reasons for hesitations are similar in dialogue
 199 systems and humans, and in the light of the positive effect hesitations have on the interactive
 200 capacities of dialogue systems, we will explore a hesitation strategy for dialogue systems that
 201 generates a suitable hesitation initiation, overall duration and phonetic structure, and is based on
 202 observations of hesitation strategies in conversations among humans. Doing so, we hope to improve
 203 our system regarding subjective ratings compared to [1,2,28], by using a smoother hesitation insertion
 204 strategy that will not, as we hope, evoke a notion of rudeness.

205 3. Towards a hesitation synthesis strategy for incremental spoken dialogue systems

206 3.1. A model for hesitation insertion in incremental spoken dialogue systems

207 Given the insights summarized in section 2.3, we now propose an elaborated and dynamic
 208 hesitation insertion strategy, that can be evoked (1) while a dialogue system is speaking, (2) and
 209 that determines the best entry point, given an event of hesitation, (3) and the best temporal extension
 210 of a hesitation. In this section, we walk through the details of the algorithm that can be seen as our
 211 general model for hesitation insertion in dialogue systems. In section 3.2, we give details on how we
 212 realized the implementation of it for this study.

213 The aim of the strategy proposed here is to buy as much time as possible for the speaker, by
 214 lengthening words in the articulatory buffer and inserting silences. Only in severe cases, where even
 215 more time is needed, will other measures, such as fillers, be employed (cf. Figure 1). This approach is
 216 governed by technical constraints. The choice to prioritize lengthening and silence is due to the simple
 217 fact that they can be synthesized with better sound quality [29], the absence of which is a weakness
 218 of many incremental systems. Moreover, this strategy is motivated by the general assumption stated
 219 in 2.2 that suggests that a hesitation is always initiated by lengthening the phonetic material available
 220 in the articulatory buffer.

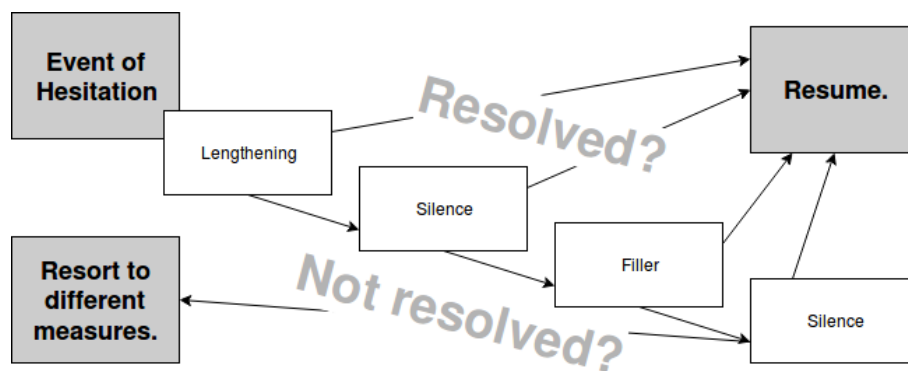


Figure 1. Hesitation insertion strategy

221 The strategy depicted in Figure 1 can be summarized as follows: While an event of hesitation is
 222 active, execute the following steps:

- 223 1. Apply lengthening to best target.
- 224 2. Insert first silence.
- 225 3. Insert filler.
- 226 4. Insert second silence.

227 When the hesitation ends during any of these steps, the original speech plan is resumed. If all steps
 228 have been run through without the event of hesitation ending, resort to different measures. In the
 229 following, we walk through the individual steps in more detail.

230 **While event of hesitation is active.** As described in section 2.1, there are various reasons for
 231 hesitating. Any of these reasons could be accounted for in a dialogue system. It could also be a
 232 wizard-of-oz setting, where there is a "start" and a "stop" button to delimit the event.

233 1. **Apply lengthening to best target.** Hesitation lengthening does not occur arbitrarily. Given the
 234 concept of the articulatory buffer, speakers start hesitating as soon as possible, which means, at
 235 the next appropriate syllable. Several linguistic and phonetic factors determine which syllable
 236 that is, and how much that syllable can be stretched in duration. To summarize findings of [31]
 237 and [32]:

- 238 • Lengthening prefers closed-class ("function") words.
- 239 • Lengthening prefers, in this order, nasals, long vowels and diphthongs, short vowels, other
 240 non-plosive sounds.³
- 241 • The extent of the lengthening is governed by the elasticity of the phone in question.

242 The lengthening continues until the phone has been stretched to its maximum, or until
 243 hesitation mode ends, whichever occurs first.

244 2. **Insert first silence.** If the lengthening has not bought enough time to resolve the event of
 245 hesitation, silence can be added. Following the suggestion of a standard maximum silence of 1
 246 second in conversation [33], this silence will continue for maximally 1000ms, or until hesitation
 247 mode ends.⁴

248 3. **Insert filler.** If the previous steps did not buy enough time, as a more severe measure of
 249 hesitation, fillers ("uhm") can be added. Short fillers ("uh") denote minor pauses and are thus
 250 not adequate for long hesitation loops [20].

251 4. **Insert second silence.** If after the filler the hesitation mode is still not resolved, a second silence
 252 can be added, with the same rules as the first silence.

253 5. **Resort to different measures.** Systems need a strategy to continue when the above steps do not
 254 suffice to buy enough time to resolve the event of hesitation. This strategy is depending on the
 255 architecture. Some examples of how a system could proceed:

- 256 • Wait for hesitation event to end.
- 257 • Re-enter the loop or parts of it to buy more time.
- 258 • Repeat parts of previously uttered speech to buy more time (cf. Example 1).
- 259 • Resume own speech plan if possible, despite event of hesitation is not over.

260 3.2. Implementing the algorithm

261 In the following, we describe how the individual concepts of the model described in the previous
 262 section 3 are realized in this study.

263 3.2.1. Event of hesitation

264 In this study, we define an event of hesitation as the time interval a user does not maintain
 265 eye-contact to our virtual agent. This is based on one of the reasons from section 2.1 - change in
 266 dialogue environment. We deploy hesitations as a user-oriented strategy (cf. [2]), as a response to
 267 visual attention shifts. The goal is to assist users in their task by only giving them information while
 268 they are paying attention.

³ The latter is language-specific. In some languages, plosives can be lengthened (e.g. Swedish) in others not (e.g. German).

⁴ For a more elaborated analysis of pauses and their duration, see [34].

269 3.2.2. Different measures

270 This definition for events of hesitation also governs the strategy for continuation. In this case, it
271 is simply waiting for the hesitation to end, i.e. the user looking back.

272 3.2.3. Lengthening

273 Lengthening is the starting point for hesitations. The appropriate target syllable is selected from
274 the words in the buffer. We included a lookahead with a 5-word limit, in order for the hesitation
275 not to start too late after an attention shift. That means that the best target is selected from the
276 upcoming words, but no later than 5 words after the trigger. Based on the preference hierarchy for
277 lengthening targets described in the previous section 3, our system iterates over the buffer, searching
278 for the optimal syllable (i.e. a nasal in a function word), increasing the tolerance for less appropriate
279 targets with each iteration.

280 The duration of the lengthening is inferred from mean duration values from previous corpus
281 studies, from which a so-called stretch factor is deducted. This factor is calculated by generating
282 Gaussian random numbers with the mean duration and standard deviation for each phoneme. The
283 highest number from 10,000 samples is selected and divided by the mean duration. This factor reflects
284 how much a given phoneme needs to be stretched in duration to achieve its average maximum. This
285 factor was additionally multiplied by 1.5 for this study, because, as it is the nature of lengthening, the
286 original duration increase was barely audible.

287 3.2.4. Fillers

288 Due to technical problems, fillers are not included in our main study. Four participants were
289 recorded in a condition with fillers, but it became apparent, that the negative impact on sound quality
290 is too great for the time being. This issue will be addressed in future studies. As will be described in
291 section 4.2, we explored the usability of data with this preliminary "full hesitation" version.

292 3.2.5. Silences

293 As fillers are left out, the main study operates with only the first silence. In the general model,
294 it is designed to last 1000ms. In our implementation, the duration is variable as we wait for the user
295 to re-focus. (In the exploratory condition with fillers, the first silence lasts for 1000ms and the second
296 silence lasts until the users re-focus.)

297 3.2.6. Technical implementation

298 From the technical side, the hesitation algorithm is integrated as separate module into an
299 existing incremental spoken dialogue system [35], which uses a toolkit for incremental dialogue
300 processing [36] and MaryTTS [37] as speech synthesis back-end.

301 4. Experiment 1: interaction study

302 To evaluate the effect of hesitation in a human-agent interaction, we conducted an interaction
303 study in *the Cognitive Service Robotics Apartment*⁵ (CSRA) [38]. The apartment consists of three rooms
304 (kitchen, living room and hallway) which are equipped with various sensors for visual tracking and
305 recording.

306 The strategy for hesitation synthesis described in section 3 is evaluated by means of a task in
307 which the participants have to perform a memorization task. A virtual agent provides a background
308 story and instructs the participants to look for hidden treats at seven different places in the apartment.
309 The dialogue system underlying the virtual agent is implemented in two different versions: one

⁵ <https://cit-ec.de/en/csra>



Figure 2. Person being instructed by virtual agent on a screen.

310 *baseline* version without hesitations or adaptations of any sort, and a *hesitating* version that monitors
 311 participant's attention shifts via gaze tracking and that enters hesitation mode whenever participants
 312 look away from the virtual agent.

313 Our hypotheses for this experiment are:

- 314 1. We expect memory task performance to benefit from the presence of hesitations.
- 315 2. We expect that presence of hesitations influences user ratings of perceived synthesis quality.
 316 (Undirected)
- 317 3. We expect no negative impact of the presence or absence of hesitation on the system's likability.

318 4.1. Methods

319 We use a between-subjects design, i.e. each participant interacts with the system in either the
 320 baseline condition or in the hesitation condition. Before the main study starts, participants are asked
 321 to fill out a declaration of consent to be recorded. In addition, they must complete a short memory
 322 test, in which they are presented a pre-constructed audio file containing ten words produced by a
 323 synthetic voice. The voice is MaryTTS's [37] German female HHM voice with no further modification.
 324 The words are German nouns that fall into five categories (professions, food, sports, buildings, cities),
 325 two in each category. Each participant is presented with the same words and order of words. They
 326 are then asked to say aloud as many of the words as they can remember. The resulting *memory test*
 327 *score* is surveyed with a checklist for later comparison to the recall rates in the main study, in order to
 328 calculate task efficiency (i.e. how good did participants perform relative to their memory capacity).

329 The main study is set in the kitchen and living room of the smart home. As a platform we use the
 330 simulation of the anthropomorphic head Flobi [39] (cf. Figure 2) displayed on a screen in the kitchen
 331 area of the smart apartment. With a web-cam on top of the screen, the agent is able to detect faces and
 332 estimate the current visual focus of attention of the human interaction partner [40]. Flobi is also able
 333 to show facial expressions and to pay attention to the current focus of discourse by looking at it.

334 As soon as a participant appears in front of Flobi, it starts talking (cf. figure 2). It first introduces
 335 itself and the apartment and then instructs participants about the task they are to perform: Each
 336 participant is asked to search for treats that have allegedly been hidden in various places in the
 337 apartment (cf. figure 3). The agent lists all potential hiding places, asking the participant to memorize

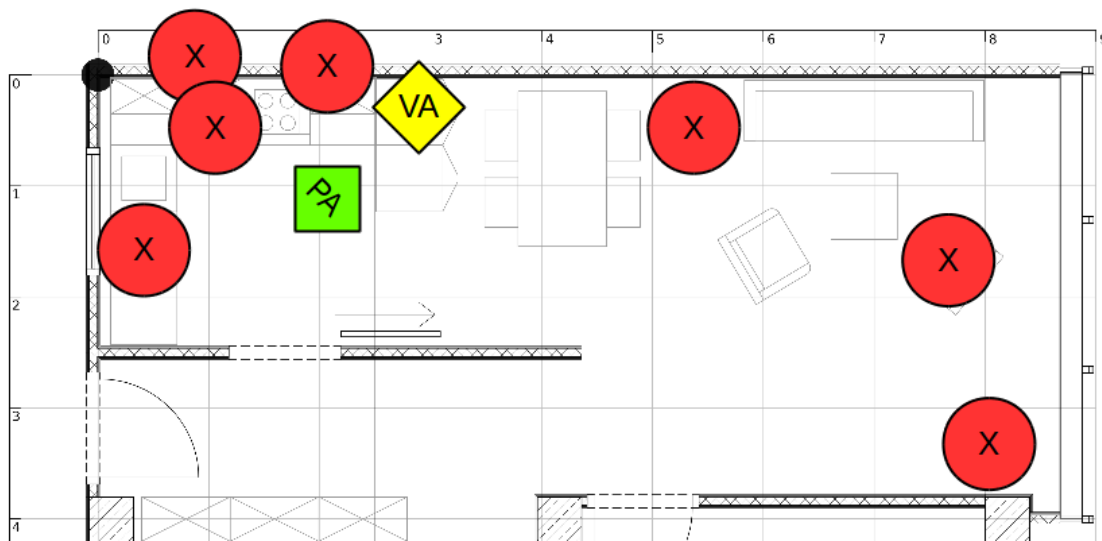


Figure 3. 2D map of the smart-home environment. (X) denotes hiding places of treats, [VA] the position of the screen with the virtual agent, [PA] the initial position of the participant.

338 and later investigate these. The task is embedded in a story about construction workers that have just
 339 left the apartment and caused confusion in the agent’s sensors, due to the dust they stirred.⁶ This
 340 creates a plausible pre-text for the agent to list all possible hiding places for the participant later to
 341 remember, with the hint that it is not sure whether it got all places correctly. During the instruction
 342 phase, there is intentional audiovisual distraction at three fixed points in time. This is included to
 343 ensure some degree of distraction and gaze shift for each participant. The distractions are: (1) Lighting
 344 up a door handle in the participants’ field of vision, (2) The experimenter entering the room to bring
 345 a code for use in the questionnaire, (3) A music beat being played for two seconds.

346 As soon as the agent has finished the instruction, the participants start investigating the possible
 347 hiding places. They are asked to call out each place before looking at it, to ensure that they remember
 348 the places and that they do not search the entire place and find things by chance. The interaction is
 349 monitored audiovisually in an adjacent room. The number of treats thus retrieved is taken down for
 350 each participant as *finding rate*. After the interaction, participants fill out a questionnaire that assesses
 351 their subjective impression of the system quality on 24 dimensions using 7-point Likert scales (based
 352 on the Godspeed questionnaire [41]), and in which they also rate their impression of speech synthesis
 353 quality using a 5-point MOS scale. Additionally, demographic data and previous experiences with
 354 robotic systems, the agent Flobi and speech synthesis systems in general are surveyed. Finally,
 355 participants are asked one question in a follow-up interview regarding the interaction, namely, if
 356 they felt that the agent adapted to their behavior in any way. All participants receive monetary
 357 compensation.

358 The entire interaction is recorded via four cameras mounted on the ceiling of the apartment.
 359 In addition, various system events for later analysis are collected (for further information about this
 360 process refer to [42]).

361 The collected data were entered into a generalized linear model (glm) with *finding rate* as
 362 dependent variable, *hesitation condition* as fixed factor, *memory test score*, *gender* and *age* as control
 363 variables. To include individual memory performance in participants’ retrieval performance, we
 364 calculated an efficiency measure: $efficiency = \frac{MemoryScore(\%)}{FindingRate(\%)}$. This is to take into account the users’
 365 individual memory capacities and to normalize results accordingly. As efficiency scores are not

⁶ There was actual visible construction work in the apartment at the time of the study, which inspired this narrative.

366 normally distributed, we used a Mann-Whitney-U test to check for effects on *efficiency* by *hesitation*
 367 *condition*. The same test was then used to analyze users' feedback on synthesis quality with regard to
 368 *hesitation condition*.

369 To evaluate the questionnaires regarding the user's perception of the agent, based on [41],
 370 the responses are grouped into five key concepts (*anthropomorphism, animacy, likability, perceived*
 371 *intelligence* and *safety*). Using Shapiro-Wilk and Bartlett tests, we found the data of all five concepts
 372 to be normally distributed and to show equal variances, qualifying the data for a t-test of *key concept*
 373 and *hesitation condition*.

374 4.2. Results and discussion

375 We recorded 37 trials with 24 female and 13 male participants in total. Participants were recruited
 376 on the university campus and via campus-related social media. Mean age was 24.6 ($SD = 4.2$). Two
 377 participants had to be excluded from the analysis because their language competence did not suffice
 378 to follow the instructions correctly. 17 participants interacted with the baseline system (ten female and
 379 seven male), and 14 with the hesitation system (ten female and four male). These 31 trials provide the
 380 core for our analysis. In addition, four participants (three female and one male) were recorded in the
 381 full hesitation condition for exploratory purposes, cf. section 3.2. The participants are balanced with
 382 regards to the their prior experience with robotic systems, the virtual agent Flobi (mostly no or very
 383 little experience) and speech systems in general (little experience).

384 4.2.1. Finding rate

385 On average, the number of items found is higher in the hesitation condition ($M = 6.36, SD =$
 386 0.84) than in the baseline condition ($M = 5.71, SD = 1.21$), (cf. Figure 4, left panel). The glm analysis
 387 shows that the effect is not significant ($\beta = 0.8, SE = 0.44, z = 1.84, p = 0.065$).

388 4.2.2. Efficiency

389 Efficiency increases in the hesitation condition ($M = 1.22, SD = 0.3$) compared to the baseline
 390 ($M = 1.5, SD = 0.58$), (cf. Figure 4, 3rd panel from the left). The Mann-Whitney-U test shows no
 significant effect of *hesitation condition* on *efficiency* ($W = 79, p = 0.11$)

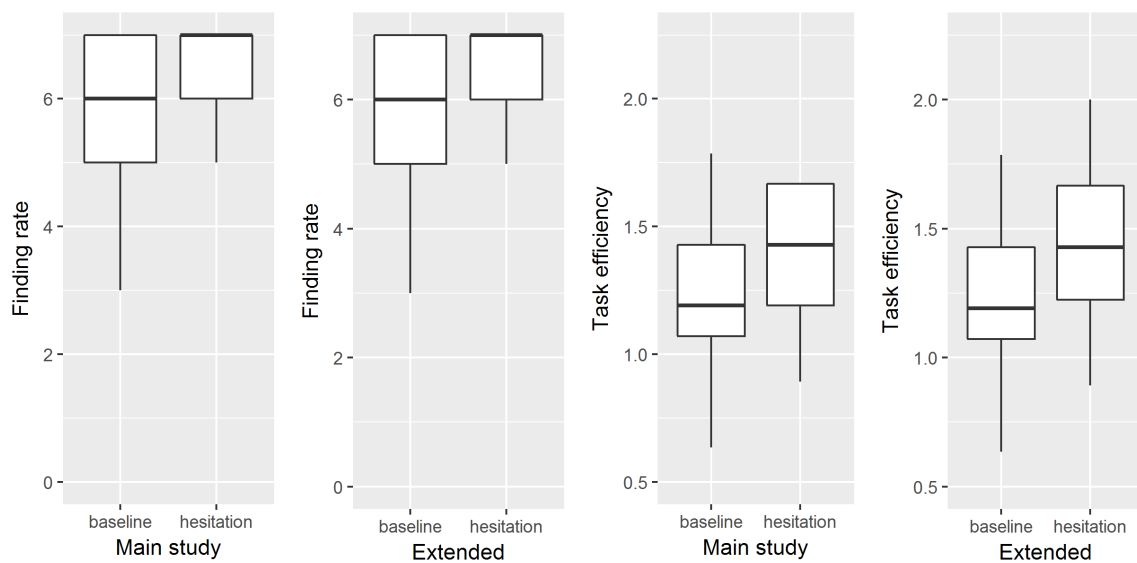


Figure 4. Task performance and efficiency.

4.2.3. Subjective speech synthesis quality.

On average, using a 5-point MOS scale (1 = "very bad", 5 = "very good") users rate synthesis quality worse in the hesitation condition ($M = 1.36, SD = 0.84$) compared to the baseline condition ($M = 2.53, SD = 0.62$), cf. Figure 6, left panel. The Mann-Whitney-U test shows that there is a significant effect of *hesitation condition* on users' perception of synthesis quality ($W = 203, p = 0.0004$).

4.2.4. Subjective rating of the agent.

We conducted t-tests for an effect of *hesitation condition* on each subjective ratings of the five key concepts *anthropomorphism*, *animacy*, *likability*, *perceived intelligence* and *safety*. The factor *hesitation condition* had no significant influence on any of the user feedbacks regarding these concepts, cf. Figure 5. Aside from the questionnaire results, participants were encouraged to give free-text

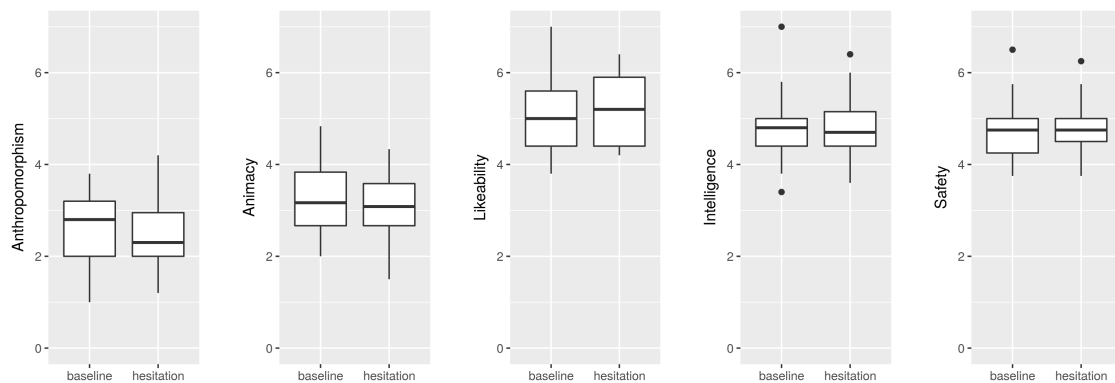


Figure 5. Subjective ratings for the five key concepts.

feedback in a comments box in the questionnaire, and they were asked regarding their perception of adaptivity after the study. In previous studies, a system that employed silence rather than hesitation to adapt to participant's level of attention, increased the attention of distracted users [2], but was perceived as less likable [2] and rude [1]. This effect appears to be lost in this study, as participants reported that they rather liked the system, which is also reflected in the questionnaire data in both conditions (cf. Figure 5).

Regarding the adaptivity, most people did not report anything in the baseline condition; some people had the impression that the agent followed their gaze (which is not the case, but the agent looks into the directions of the places he talks about, and users are likely to look in the same direction.) In the hesitation condition, many participants noticed the hesitations, but could not figure out what triggers them. Some reported that they like this feature, as it grants more time for searching, but most others were put off by the disfluent delivery: In total we have negative sound quality feedback from 13 out of 18 participants that were recorded in the hesitation conditions. In the following interview, however, the notion was rather that the adaptivity is positive and promising for the future, given improvements in the technical realization.

4.2.5. Exploratory extension of analysis.

As the tendencies observed for finding rate and efficiency failed to reach the 0.05 significance level by only a small margin, we hypothesized that the effect might reach significance if more trials were recorded. As we have at our disposal four recordings with the full hesitation condition (cf. section 3.2), we re-did the analyses with the same 17 trials for the baseline condition and with all 18 hesitation trials combined as the hesitation condition. The effect on finding rate still does not reach significance, however by a very small margin ($\beta = 1.03, SE = 0.53, z = 1.96, p = 0.0504$). The effect on efficiency becomes significant, when all trials are considered ($W = 83.5, p = 0.02$), (cf. Figure 4).

425 This suggests that there is indeed an impact of hesitations that needs to be considered. We assume
 426 that these effects will be confirmed in a follow-up study with a bug-fixed version of the system and
 427 with more participants.

428 4.2.6. Summary

429 The results gathered here point in expected directions: Speech synthesis quality suffers from the
 430 presence of hesitation, but task performance appears to benefit from it. The evaluation of subjective
 431 ratings on the five key concepts as well as qualitative evaluation of user feedback suggests that the
 432 hesitation algorithm tested in this study is acceptable. Thus, for the first study we can state that
 433 hypotheses (1) and (3) can be accepted for now, and with respect to hypothesis (2), the results suggest
 a negative impact of hesitations on user's perception of synthesis quality.

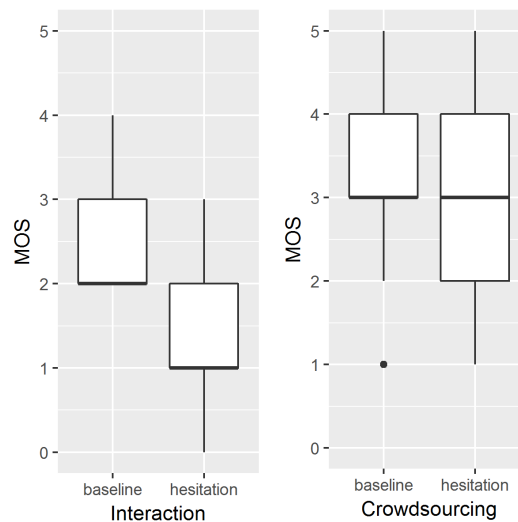


Figure 6. 5-point MOS scale user feedback on synthesis quality.

434

435 5. Experiment 2: crowdsourcing-based evaluation of hesitation synthesis

436 In order to assess the quality of the hesitation synthesis in a non-interactive setting, we conducted
 437 a parallel online crowdsourcing study. In this evaluation, we used a more traditional approach to
 438 speech synthesis evaluation, namely a classic MOS-scale rating task without any interaction between
 439 participants and system. This is done in order to shed light on our underlying assumption that an
 440 interactive approach to synthesis evaluation indeed may lead to different conclusions with respect to
 441 synthesis quality. Our main hypothesis for this experiment is undirected, i.e. we do expect a different
 442 outcome in terms of speech synthesis quality than we achieved in experiment 1. We do not make any
 443 claims about the direction of this hypothesis, as the non-interactive setting may have unforeseeable
 444 effects. So far, our only expectation is that the result will differ from the interaction study.

445 5.1. Methods

446 Participants listened to a series of 14 synthetic audio stimuli and rated them individually for their
 447 overall quality on a 5-point MOS scale (1 = "very bad", 5 = "very good"). Participants were recruited
 448 using mailing lists and social media, and the evaluation builds on a web-based, crowdsourcing
 449 approach. The listening test was set up using the platform PERCY [43], specially designed for online
 450 audio-based perception studies. Unlike experiment 1, but very much like standard MOS-based
 451 synthesis evaluations, participants rated the synthesis quality of each individual stimulus. The
 452 participants were not compensated for their participation.

453 For maximal comparison with the interaction study, we again chose a between-subjects design
454 with a single controlled independent variable *hesitation condition*, which has the two levels *hesitation*
455 and *baseline*. That is, participants listened to either stimuli containing hesitations only, or to stimuli
456 not containing any hesitations. This may create a deviation between our two experiments, as in the
457 interactive study, the presence, absence and length of a hesitation was determined by the participant's
458 individual behavior, and was not necessarily present or absent in each stimulus. Demographic data
459 and information about the output device and individual listening situation is surveyed as well, but
460 not analyzed further.

461 Before the actual listening tests, participants received some background information of what
462 is being tested (a synthetic voice for usage in an intelligent apartment). They also received some
463 instructions on the procedure of the experiment, i.e. how to use the scale and how long the experiment
464 is likely to last. In both conditions, participants were presented with 14 stimuli which were based
465 upon the text input given to the virtual agent in experiment 1. That way, participants get the same
466 background story (and text) as in the first experiment. Stimuli are divided into 6 introductory, 7
467 instructive and 1 concluding utterance. They are presented in the same order for each participant,
468 to generate a coherent story, and to ensure maximal similarity with experiment 1. In the baseline
469 condition (non-hesitation), the stimuli are produced with MaryTTS's [37] female German HMM
470 voice, with no further modification. For the hesitation condition, lengthenings and silent pauses are
471 woven into each stimulus: In the instructive stimuli, the silent pauses are set to 2000ms, in all other
472 stimuli, silences are set to 1000ms. This difference in duration is motivated by experiment 1, which by
473 design leads to longer pause intervals in the instructions, because participants tend to look around the
474 apartment when possible hiding places are mentioned, these gaze shifts triggering hesitation mode.
475 Lengthenings are applied to syllables preceding the silence with the same durational parameters as
476 in the first study. A list of the stimuli used in this experiment can be found in appendix A.

477 The collected data were entered into a linear mixed effects model with *MOS ratings* as dependent
478 variable, *hesitation condition* as fixed factor, and *stimulus*, *gender* and *age* as random factors (random
479 intercepts). This model was compared to a less complex model, leaving out the fixed factor *hesitation*
480 *condition* using a likelihood ratio test. All statistical tests were carried out in R, using the R-package
481 *lme4* (version 1.1-12).

482 5.2. Results and discussion

483 We collected ratings from 44 participants (29 female, 15 male) with an age range between 18
484 and 46 years (median: 24.5). With one exception, all participants reported to have entered school in
485 Germany, so we expect them to have a native competence in German. No participant reported any
486 hearing problems. Most participants were raised in the vicinity of Bielefeld, a few in Bavaria. The
487 listening tests typically lasted less than 5 minutes, including the time needed to provide demographic
488 background data. For subsequent analyses, we pooled all participants' data, independent of listening
489 situation, and including one participant who reported to have entered school out of Germany, as the
490 fact that s/he managed to follow the instructions is an indicator of a sufficiently high competence in
491 German.

492 On average, MOS-ratings were slightly higher in the baseline condition ($M = 3.28, SD = 0.93$)
493 as compared to the hesitation condition ($M = 2.96, SD = 0.93$) (cf. Figure 6). In the LMER-model
494 containing the fixed factor *hesitation*, the absence of hesitation has a slightly positive, but no significant
495 effect on MOS-ratings ($\beta = 0.31, SE = 0.18, t = 1.78, p = 0.08$). This lack of an effect is further
496 confirmed by the model comparison (likelihood ratio test between models with and without the factor
497 *hesitation*), which does not reveal a significant difference either.

498 These results are perhaps surprising insofar, as there were a reasonable number of participants
499 for both conditions (> 20), as the test gave listeners a chance to rate each stimulus without being
500 distracted by an ancillary task as in experiment 1, and since participants were confronted with
501 hesitations in each stimulus in the *hesitation condition*. Still, it can only be concluded that even though

502 there is a tendency for stimuli to be rated as slightly less pleasant when hesitations are present, this
503 detrimental effect is not perceived to be significantly strong by listeners in the classic non-interactive
504 approach to speech synthesis evaluation. Of course, most MOS-type analyses rely on within-subjects
505 designs. It is possible, that participants would have rated the stimuli containing hesitations as
506 less good when given a chance for a direct comparison with a stimulus not containing hesitations.
507 However, our aim was to test the influence of an interactive task on speech synthesis ratings. A
508 within-subjects approach would have made such a comparison impossible.

509 6. General Discussion

510 This study yields several insights that demand discussion. We improve the conversational
511 capabilities of a dialogue system by integrating a strategy for dynamic insertion of synthesized
512 hesitations. The experimental results suggest that hesitations are a useful and viable strategy in
513 interaction with users, as they increase task efficiency. Our evaluation is, however, not limited to
514 objective assessments of the system as a whole, rather, we also assessed subjective system ratings via
515 participant feedback.

516 Of special interest in this study is the feedback on speech synthesis quality. In addition to the
517 interaction study, we conducted a parallel crowdsourcing experiment with comparable stimuli in
518 order to compare ratings gathered within and without interactive settings. Regarding evaluations in
519 dialogue system and speech synthesis research, we observe that: (1) In dialogue system evaluation,
520 the speech synthesis quality is often not assessed. (2) In speech synthesis evaluation, user ratings
521 are surveyed in MOS-based questionnaires regarding stimuli presented without interaction with the
522 system. The results gathered in this study support a claim that has often been uttered in the speech
523 synthesis community lately: Non-interactive evaluation of speech synthesis does not work, or at least,
524 it assesses aspects of quality that differ from those gathered in interactive settings. Even if it could be
525 guaranteed that what is being assessed really is the "pure" synthesis quality, then it is totally unclear
526 what to do with this information. Speech synthesis is not used in the void, there is always some
527 application or interaction associated with it.

528 Our study highlights this point. As can be seen in Figure 6, there are two main differences
529 between MOS-ratings after interaction and after the non-interactive crowdsourcing evaluation: First,
530 stimuli are generally rated better without prior interaction, second, the presence of hesitation only
531 makes a significant difference in the interaction study. The reason for this discrepancy lies in the
532 nature of the two experimental settings. The crowdsourcing experiment uses neatly pre-constructed
533 stimuli, the interaction study adapts and enhances the stimuli on the fly with spontaneous speech
534 phenomena. The latter will cause artifacts that detriment the synthesis quality, which will be noticed
535 by users and reflected in their feedback. This is the general problem with synthesis evaluation:
536 Experimental results from MOS-based questionnaires are not the same as those gathered in interaction
537 studies (And, while being closer to in-the-wild application, interaction studies are still not the reality
538 of application.)

539 An important question that arises is: how to gather quality measures that do account for the
540 interactive nature of speech synthesis applications? In general, there are two possible starting points:
541 use the dialogue system evaluation to infer something for speech synthesis quality, or make offline
542 evaluations more interactive. There is no obvious way to get precise first-hand user feedback on
543 synthesis quality from an interaction study, as the interaction cannot be interrupted in between to
544 ask for feedback. Neither can task performance measures from the study be used to directly infer the
545 impact of the speech synthesis. One conceivable option would be to have external evaluators review
546 the recorded interactions and give feedback on the synthesis quality every given time interval. It
547 thus appears more fruitful to enrich offline evaluations. If the stimuli that participants have to rate
548 would be embedded in small-scale interactive scenarios, interactive measures like reaction time, task
549 completion time or task performance in general could be surveyed in addition to the MOS feedback,

550 helping to analyze and interpret the results. Preliminary tests with relative task completion time for
551 instructive stimuli in connection with MOS-feedback were explored in [16].

552 Speech synthesis evaluation as of now is an unsolved problem. Speech synthesis does not exist
553 without interaction, thus it makes no sense to evaluate it without. If any given speech synthesis
554 system achieved good MOS scale ratings, it would at least be necessary to test the system in
555 interaction to see if the results can be justified. If the system cannot reach the same quality level in
556 interaction due to technical limitations, as observed in this study, then the off-line version could serve
557 as a gold standard to be reached in interaction via further development of the system. Non-interactive
558 MOS-based evaluation, however, maximally reflects the opinion of a user testing it in a disembodied
559 way without the application it may be designed for.

560 Turning to other objectives of this study, it is to be asked what our evaluation results tell us about
561 the actual system that we tested.

562 It is in general satisfying that there is a tendency towards more task performance and efficiency.
563 The detrimental effect observed for synthesis quality, in turn, highlights the need for improvement.
564 The fact that some of the effects can be attributed to the fact that the technical realization of our
565 hesitation model yielded some audible artifacts, gives rise to the question if a simpler strategy could
566 not have achieved the same thing. It may appear unnecessary to develop and implement a complex
567 model that yields technical problems that could have been avoided by simply being silent. In a
568 previous study that used silence only as an attention-driven hesitation strategy [2], an increase the
569 visual attention and hesitations in terms of silence increase the task performance was noticed at well
570 [28], but the hesitating system was perceived as comparably less friendly. This is an effect that we
571 cannot observe in our study - the presence of hesitation has no detrimental or beneficial effect on
572 perceived friendliness. Also, feedback gathered in the comments section of the questionnaire and in
573 the short interview after the study suggests that people regard the adaptive strategy of the system
574 positively, despite the fact that many are rather put off by the disfluent speech delivery. This suggests
575 that the general approach to overtly indicate system hesitation is a promising extension for (virtual)
576 agents' dialogue systems, and doing so with more sophisticated methods than plainly being silent is
577 credited by users. In a follow-up study we will explore further the applicability of our model with
578 some extensions regarding the realization of hesitations in order to minimize the irritating effects
579 reported for this first prototype.

580 To conclude, given some necessary improvements on the technical side, we expect the hesitation
581 model to have future application and we will explore that in follow-up studies. The evaluation itself
582 also needs improvements; synthesis designed for interaction needs to be evaluated in interaction.
583 It is, as of now, one of the greatest challenges for the speech synthesis community to develop and
584 establish evaluation paradigms that allow to go beyond pure MOS scales.

585 **Acknowledgments:** This research was carried out as part of the CITEC Large Scale Project "Computational
586 Service Robotics Apartment" (CSRA) and was supported by the Cluster of Excellence Cognitive Interaction
587 Technology 'CITEC' (EXC 277) at Bielefeld University, which is funded by the German Research Foundation
588 (DFG). We warmly thank our participants; Christoph Draxler, who invested much of his time and helped setting
589 up experiment 2 in virtually no time; Timo Baumann and Soledad Lopez-Gambino, who gave advice regarding
590 several issues with the synthesis in InProTK; Monika Chromik and Ayla Canpolat for their massive support
591 during experiment 1 and all the people, who invested a lot of their time to set up the CSRA as a platform for
592 interaction research.

593 **Author Contributions:** All authors conceived and designed the experiments, and analyzed the data jointly. Birte
594 Carlmeyer and Simon Betz conducted experiment 1. Simon Betz constructed the stimuli for experiment 2. Petra
595 Wagner conducted experiment 2. Simon Betz, Birte Carlmeyer and Petra Wagner wrote the paper.

596 **Conflicts of Interest:** The authors do not declare any conflict of interests.

597 Appendix A. Stimuli for crowdsourcing study

598 The following stimuli are used for the crowdsourcing experiment described in section 5.
599 Lengthened syllables are indicated by appended colons. Pauses are indicated by seconds in brackets.

600 Lengthening durations are determined as described in section 3.2.3. Stimuli for the baseline condition
601 are the same, except without lengthenings and pauses.

602 Introduction

- 603 1. "Hallo, schön, dass du an: (1.0) dieser Studie teilnimmst."
- 604 2. "Ich werde dir heute ein wenig über dieses Apartment erzählen, und (1.0) dann habe ich eine
605 kleine Aufgabe für dich."
- 606 3. "Du könntest mir nämlich beim Suchen helfen. Hier sind eben ein paar: (1.0) Sachen verloren
607 gegangen."
- 608 4. "Einige Handwerker waren hier im Apartment und (1.0) haben die Küche umgebaut."
- 609 5. "Ich konnte wegen des Staubs leider nicht genau erkennen, wo die: (1.0) Sachen versteckt
610 wurden."

611 Instruction

- 612 1. "Jemand hat die Waschmaschine bedient und (2.0) das Waschpulverfach geöffnet."
- 613 2. "Und ich habe gesehen, wie jemand zur Pflanze im Wohnzimmer gegangen ist, und (2.0) etwas
614 am Blumentopf gemacht hat."
- 615 3. "Danach hat jemand die Bescheckschublade geöffnet und (2.0) hat dort rumgewühlt."
- 616 4. "Und dann habe ich beobachtet dass jemand den Schrank über der: (2.0) Mikrowelle aufgemacht
617 hat."
- 618 5. "Dann wurde einer der Stühle im: (2.0) Wohnzimmer bewegt."
- 619 6. "Irgend etwas ist mit den Kaffeetassen auf dem Tisch im: (2.0) Wohnzimmer passiert."
- 620 7. "Zu guter Letzt war noch jemand am Bescheckfach der: (2.0) Spülmaschine."

621 Conclusion

- 622 1. "Schau in beliebiger Reihenfolge an: (1.0) den Orten nach, die ich dir genannt habe."

623 Bibliography

- 624 1. Carlmeyer, B.; Schlangen, D.; Wrede, B. Exploring self-interruptions as a strategy for regaining the
625 attention of distracted users. Proceedings of the 1st Workshop on Embodied Interaction with Smart
626 Environments - EISE '16. Association for Computing Machinery (ACM), 2016.
- 627 2. Carlmeyer, B.; Schlangen, D.; Wrede, B. "Look at Me!": Self-Interruptions as Attention Booster?
628 Proceedings of the Fourth International Conference on Human Agent Interaction - HAI '16. Association
629 for Computing Machinery (ACM), 2016.
- 630 3. Skantze, G.; Hjalmarsson, A. Towards incremental speech generation in conversational systems. *Computer
631 Speech and Language* 27 2013.
- 632 4. King, S. What speech synthesis can do for you (and what you can do for speech synthesis). Proceedings
633 of the 18th International Congress of the Phonetic Sciences (ICPhS 2015).
- 634 5. Mendelson, J.; Aylett, M. Beyond the Listening Test: An Interactive Approach to TTS Evaluation.
635 Proceedings of the 18th Annual Conference of the International Speech Communication Association
636 (Interspeech 2017, Stockholm), 2017, pp. 249–253.
- 637 6. Rosenberg, A.; Ramabhadran, B. Bias and Statistical Significance in Evaluating Speech Synthesis
638 with Mean Opinion Scores. Proceedings of the 18th Annual Conference of the International Speech
639 Communication Association (Interspeech 2017, Stockholm), 2017, pp. 3976–3980.
- 640 7. Wester, M.; Braude, D.A.; Potard, B.; Aylett, M.; Shaw, F. Real-Time Reactive Speech Synthesis:
641 Incorporating Interruptions. Proceedings of the 18th Annual Conference of the International Speech
642 Communication Association (Interspeech 2017, Stockholm), 2017, pp. 3996–4000.
- 643 8. Wagner, P.; Betz, S. Speech Synthesis Evaluation – Realizing a Social Turn. Tagungsband Elektronische
644 Sprachsignalverarbeitung (ESSV), 2017, p. 167–172.
- 645 9. Eklund, R. Disfluency in Swedish human–human and human–machine travel booking dialogues. PhD
646 thesis, Linköping University Electronic Press, 2004.
- 647 10. Shriberg, E. Preliminaries to a Theory of Speech Disfluencies. *Ph D. thesis University of California* 1994.

- 648 11. Clark, H.H.; Tree, J.E.F. Using uh and um in spontaneous speaking. *Cognition* **2002**, *84*, 73–111.
- 649 12. Goodwin, C. Conversational organization. *Interaction between speakers and hearers* **1981**.
- 650 13. Tree, J.E.F. Listeners' uses of um and uh in speech comprehension. *Memory & cognition* **2001**, *29*, 320–326.
- 651 14. Collard, P. Disfluency and listeners' attention: An investigation of the immediate and lasting effects of
652 hesitations in speech. PhD thesis, University of Edinburgh, 2009.
- 653 15. Corley, M.; Stewart, O.W. Hesitation disfluencies in spontaneous speech: The meaning of um. *Language
654 and Linguistics Compass* **2008**, *2*, 589–602.
- 655 16. Betz, S.; Zariß, S.; Wagner, P. Synthesized lengthening of function words - The fuzzy boundary between
656 fluency and disfluency. Proceedings of the International Conference Fluency and Disfluency, 2017.
- 657 17. Kempen, G.; Hoenkamp, E. Incremental sentence generation: Implications for the structure of a syntactic
658 processor. Proceedings of the 9th conference on Computational linguistics-Volume 1. Academia Praha,
659 1982, pp. 151–156.
- 660 18. Levelt, W.J.M. *Speaking: From Intention to Articulation*; MIT Press, 1989.
- 661 19. Shriberg, E. To 'errrr'is human: ecology and acoustics of speech disfluencies. *Journal of the International
662 Phonetic Association* **2001**, *31*, 153–169.
- 663 20. Clark, H. Speaking in Time. *Speech Communication* **36** **2002**.
- 664 21. Shriberg, E. Disfluencies in switchboard. Proceedings of International Conference on Spoken Language
665 Processing, 1996, Vol. 96, pp. 11–14.
- 666 22. Shriberg, E. Toerrrr'is human: ecology and acoustics of speech disfluencies. *Journal of the International
667 Phonetic Association* **2001**, *31*, 153–164.
- 668 23. Li, J.; Tilsen, S. Phonetic evidence for two types of disfluency. Proceedings of ICPhS 2015, 2015.
- 669 24. Skantze, G.; Schlangen, D. Incremental Dialogue Processing in a Micro-Domain. Proceedings of the 12th
670 Conference of the European Chapter of the Association for Computational Linguistics (EACL 2009), 2009,
671 pp. 745–753.
- 672 25. Schlangen, D.; Skantze, G. A General, Abstract Model of Incremental Dialogue Processing. *Dialogue and
673 Discourse* **2011**, *2*, 83–111.
- 674 26. Kousidis, S.; Kennington, C.; Baumann, T.; Buschmeier, H.; Kopp, S.; Schlangen, D. Situationally
675 Aware In-Car Information Presentation Using Incremental Speech Generation: Safer, and More Effective.
676 Proceedings of the EACL 2014 Workshop on Dialogue in Motion, 2014, pp. 68–72.
- 677 27. Bohus, D.; Horvitz, E. Managing Human-Robot Engagement with Forecasts and... Um... Hesitations.
678 Proc. of the 16th International Conference on Multimodal Interaction; ACM: New York, USA, 2014; pp.
679 2–9.
- 680 28. Chromik, M.; Carlmeyer, B.; Wrede, B. Ready for the Next Step?: Investigating the Effect of Incremental
681 Information Presentation in an Object Fetching Task. Proc. of the Companion of the HRI 2017 ACM/IEEE.
682 ACM, 2017.
- 683 29. Betz, S.; Wagner, P.; Schlangen, D. Micro-Structure of Disfluencies: Basics for Conversational Speech
684 Synthesis. Proceedings of the 16th Annual Conference of the International Speech Communication
685 Association (Interspeech 2015, Dresden), 2015, pp. 2222–2226.
- 686 30. Betz, S.; Voße, J.; Zariß, S.; Wagner, P. Increasing Recall of Lengthening Detection via Semi-Automatic
687 Classification. Proceedings of the 18th Annual Conference of the International Speech Communication
688 Association (Interspeech 2017, Stockholm), 2017, pp. 1084–1088.
- 689 31. Betz, S.; Wagner, P.; Vosse, J. Deriving a strategy for synthesizing lengthening disfluencies based on
690 spontaneous conversational speech data. *Phonetik und Phonologie* **12**, 2016.
- 691 32. Betz, S.; Voße, J.; Wagner, P. Phone Elasticity in Disfluent Contexts. *Fortschritte der Akustik - DAGA
692 2017*, 2017.
- 693 33. Jefferson, G. Preliminary notes on a possible metric which provides for a "standard maximum" silence
694 of approximately one second in conversation. In *Conversation: An Interdisciplinary Perspective.*; Roger, D.;
695 Bull, P., Eds.; 1989.
- 696 34. Lundholm Fors, K. Production and Perception of Pauses in Speech. PhD thesis, 2015.
- 697 35. Carlmeyer, B.; Schlangen, D.; Wrede, B. Towards Closed Feedback Loops in HRI: Integrating InproTK
698 and PaMini. Proceedings of the 2014 Workshop on Multimodal, Multi-Party, Real-World Human-Robot
699 Interaction. ACM, 2014, ICMI-MMRWHRI '14, pp. 1–6.

- 700 36. Baumann, T.; Schlangen, D. The InproTK 2012 Release. NAACL-HLT Workshop on Future Directions
701 and Needs in the Spoken Dialog Community: Tools and Data; Association for Computational Linguistics:
702 Stroudsburg, PA, USA, 2012; SDCTD '12, pp. 29–32.
- 703 37. Schroeder, M.; Trouvain, J. The German text-to-speech synthesis system MARY: A tool for research,
704 development and teaching. *International Journal of Speech Technology*, 6:365-377. **2003**.
- 705 38. Wrede, S.; Leichsenring, C.; Holthaus, P.; Hermann, T.; Wachsmuth, S. The Cognitive Service Robotics
706 Apartment: A Versatile Environment for Human-Machine Interaction Research. *KI - Kuenstliche Intelligenz*
707 (*Special Issue Smart Environments*) **2017**.
- 708 39. Lütkebohle, I.; Hegel, F.; Schulz, S.; Hackel, M.; Wrede, B.; Wachsmuth, S.; Sagerer, G. The Bielefeld
709 Anthropomorphic Robot Head "Flobi". 2010 IEEE International Conference on Robotics and Automation.
710 IEEE, 2010, pp. 3384–3391.
- 711 40. Schillingmann, L.; Nagai, Y. Yet another gaze detector: An embodied calibration free system for the iCub
712 robot. 2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids), 2015, pp. 8–13.
- 713 41. Bartneck, C.; Kulić, D.; Croft, E.; Zoghbi, S. Measurement Instruments for the Anthropomorphism,
714 Animacy, Likeability, Perceived Intelligence, and Perceived Safety of Robots. *International Journal of Social*
715 *Robotics* **2009**, 1, 71–81.
- 716 42. Holthaus, P.; Leichsenring, C.; Bernotat, J.; Richter, V.; Pohling, M.; Carlmeyer, B.; Köster, N.; zu Borgsen,
717 S.M.; Zorn, R.; Schiffhauer, B.; Engelmann, K.F.; Lier, F.; Schulz, S.; Cimiano, P.; Eyssel, F.; Hermann, T.;
718 Kummert, F.; Schlangen, D.; Wachsmuth, S.; Wagner, P.; Wrede, B.; Wrede, S. How to Address Smart
719 Homes with a Social Robot? A Multi-modal Corpus of User Interactions with an Intelligent Environment.
720 Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016);
721 European Language Resources Association: Paris, France, 2016.
- 722 43. Draxler, C. Online Experiments with the Percy Software Framework - Experiences and some Early
723 Results. Proceedings of the Ninth International Conference on Language Resources and Evaluation
724 (LREC'14); Chair), N.C.C.; Choukri, K.; Declerck, T.; Loftsson, H.; Maegaard, B.; Mariani, J.; Moreno,
725 A.; Odijk, J.; Piperidis, S., Eds.; European Language Resources Association (ELRA): Reykjavik, Iceland,
726 2014.

727 © 2017 by the authors. Submitted to *Multimodal Technologies and Interact.* for possible open
728 access publication under the terms and conditions of the Creative Commons Attribution license
729 (<http://creativecommons.org/licenses/by/4.0/>)