

1 **Highlights**

2 **The Way You Assess Matters: User Interaction Design of Survey Chatbots**  
3 **for Mental Health**

4 Yucheng Jin, Li Chen, Xianglin Zhao, Wanling Cai

- 5 • This work investigates how the interaction design of psychological assess-  
6 ment with closed-ended questions could influence user responses to open-  
7 ended questions in a survey chatbot for mental health.
- 8 • An empirical study shows the significant effects of *interaction style* (form-  
9 based vs. conversation-based) on user-perceived assessment credibility and  
10 self-awareness.
- 11 • A structural equation model illustrates the mediating role of perceived as-  
12 sessment credibility in the effects of psychological assessment design on  
13 user responses to the subsequent open-ended questions.

# The Way You Assess Matters: User Interaction Design of Survey Chatbots for Mental Health

Yucheng Jin<sup>a</sup>, Li Chen<sup>a</sup>, Xianglin Zhao<sup>a</sup>, Wanling Cai<sup>a</sup>

<sup>a</sup>Department of Computer Science, Hong Kong Baptist University, Hong Kong, China

---

## Abstract

The global pandemic has pushed human society into a mental health crisis, prompting the development of various chatbots to supplement the limited mental health workforce. Several organizations have employed mental health survey chatbots for public mental status assessments. These survey chatbots typically ask closed-ended questions (Closed-EQs) to assess specific psychological issues like anxiety, depression, and loneliness, followed by open-ended questions (Open-EQs) for deeper insights. While Open-EQs are naturally presented conversationally in a survey chatbot, Closed-EQs can be delivered as embedded forms or within conversations, with the length of the questionnaire varying according to the psychological assessment. This study investigates how the *interaction style* of Closed-EQs and the *questionnaire length* affect user perceptions regarding survey credibility, enjoyment, and self-awareness, as well as their responses to Open-EQs in terms of quality and self-disclosure in a survey chatbot. We conducted a 2 (*interaction style*: form-based vs. conversation-based)  $\times$  3 (*questionnaire length*: short vs. middle vs. long) between-subjects study (N=213) with a loneliness survey chatbot. The results indicate that the form-based interaction significantly enhances the perceived credibility of the assessment, thereby improving response quality and self-disclosure in subsequent Open-EQs and fostering self-awareness. We discuss our findings for the interaction design of psychological assessment in a survey chatbot for mental health.

**Keywords:** Chatbots, survey design, open-ended questions, psychological assessment, self-disclosure, mental health, loneliness

---

## 1. Introduction

The rise of mental health issues among young adults has become a significant public health challenge [1, 2, 3, 4], further intensified by the global pandemic's

43 impact on various aspects of life [5, 6, 7]. Early detection and intervention are cru-  
44 cial for providing targeted support and treatments [8, 9]. With the rapid advance-  
45 ment in artificial intelligence (AI), several organizations, including universities,  
46 hospitals, and public sectors, have begun utilizing mental health survey chatbots  
47 for conducting psychological assessments to determine individuals' mental states  
48 and needs [10, 11, 12]. Compared with traditional web-based surveys, chatbot  
49 surveys have demonstrated advantages in response rate, user engagement, and re-  
50 sponse quality due to the natural conversation and interactive features [13, 14].

51 Mental health surveys typically contain two primary types of questions [15,  
52 16]: *closed-ended questions* (Closed-EQs) often based on psychological scales  
53 like the UCLA Loneliness Scale consisting of twenty Closed-EQs [17], and open-  
54 ended questions (Open-EQs) that delve into deeper individual insights [16], pro-  
55 moting spontaneous and less biased responses [15]. Research in web surveys  
56 has revealed correlations between responses to Closed-EQs and subsequent Open-  
57 EQs. For example, participants dissatisfied with job or e-services through Closed-  
58 EQs tended to disclose more details about negative feelings in subsequent Open-  
59 EQs [18, 19, 20]. However, little work has investigated *if* and *how* the design  
60 choices of Closed-EQs influence user responses to Open-EQs, particularly in men-  
61 tal health survey chatbots. Existing work has primarily investigated leveraging  
62 a chatbot to respectively improve the response quality of Closed-EQs or Open-  
63 EQs [14, 21]. Our research aims to bridge the gap by exploring the effects of  
64 two prominent design factors (i.e., *interaction style* and *questionnaire length*) of  
65 a psychological assessment with Closed-EQs on user responses to the follow-up  
66 Open-EQs in a mental health survey chatbot.

67 The *interaction style* and *questionnaire length* are two crucial design factors  
68 of Closed-EQs [22, 23]. Prior studies have shown that, compared to conventional  
69 form-based interactions on webpages, employing conversation-based interactions  
70 can enhance the quality of responses to Closed-EQs [14]. Additionally, research  
71 has demonstrated that the *questionnaire length* can influence participation and  
72 completion rate [23, 24], as well as the response quality [25]. In our study, we  
73 experimented with both form-based and conversation-based interactions in our  
74 chatbot's psychological assessment. The manipulation of questionnaire length is  
75 based on the three validated versions of the UCLA loneliness scale [17], including  
76 short (three items), middle (ten items), and long (twenty items), respectively. This  
77 led to a 2 (interaction style: form-based vs. conversation-based)  $\times$  3 (question-  
78 naire length: short vs. middle vs. long) between-subjects study, enabling us to  
79 address the following four research questions with empirical evidence.

80 **RQ1:** How does the *interaction style* of an assessment influence the users'

81 perceptions of a mental health survey chatbot (i.e., w.r.t. enjoyment, assessment  
82 credibility, and self-awareness)?

83 **RQ2:** How does the *interaction style* of an assessment influence user re-  
84 sponses to the follow-up Open-EQs (i.e., w.r.t. response quality and self-disclosure)  
85 in a mental health survey chatbot?

86 **RQ3:** How does the *questionnaire length* of an assessment influence the users’  
87 perceptions of a mental health survey chatbot?

88 **RQ4:** How does the *questionnaire length* of an assessment influence user  
89 responses to the follow-up Open-EQs in a mental health survey chatbot?

90 Our study provides practical design implications to designers of survey chat-  
91 bots for mental health. To the best of our knowledge, this is the first study that em-  
92 pirically analyzes how psychological assessment design influences user responses  
93 to Open-EQs within a mental health survey chatbot. Consequently, the contribu-  
94 tions of our work are three-fold:

- 95 **1. Empirical evidence of the effects of psychological assessment (with Closed-  
96 EQs) design on user responses to the follow-up Open-EQs in a survey  
97 chatbot for mental health.** Our findings reveal the effective design choices  
98 for the psychological assessment that could motivate respondents to provide  
99 quality responses and stimulate deep self-disclosure in Open-EQ.
- 100 **2. Analysis of the causal relationship between the design factors of psycho-  
101 logical assessment and the measures for user responses to Open-EQs.**  
102 We employed a structural equation model (SEM) to identify how users’ *per-  
103 ceived assessment credibility*, as a mediator, links psychological assessment  
104 design factors to the critical metrics of user responses to Open-EQs such as  
105 response quality and self-disclosure.
- 106 **3. Design recommendations of psychological assessment in a survey chat-  
107 bot for mental health.** Based on our findings, we present several practical  
108 design recommendations. For instance, form-based interaction is preferable  
109 for psychological assessments, as it leads to a higher perceived assessment  
110 credibility compared to conversation-based interaction.

## 111 **2. Related work**

### 112 *2.1. Loneliness Among Young Adults and Its Measurement*

113 Loneliness is a common distressing feeling that is closely associated with ad-  
114 verse mental health states, such as depression and anxiety [26, 27, 28, 29]. Young

115 people are more susceptible to loneliness compared to other age groups, due to a  
116 dramatic increase in socioemotional demands at their unique life stage [30, 31].  
117 The social restrictions imposed to control the spread of COVID-19 have notably  
118 diminished social contact for the youth, exacerbating their feelings of loneliness  
119 and leading to increased psychological distress [32, 33, 34]. For example, fol-  
120 lowing the outbreak of COVID-19, up to 60% of young adults in America have  
121 reported symptoms indicative of psychological distress [35].

122 Early detection and intervention of loneliness are crucial for young adults, as  
123 these steps can help them mitigate its long-term effects on their mental health  
124 and support them in establishing healthier social connections and networks [34].  
125 When measuring loneliness, the UCLA Loneliness Scale and its related shorter  
126 forms are widely acknowledged and recommended as the primary tools for as-  
127 sessing loneliness [36]. As for intervention strategies, recent studies highlight the  
128 effectiveness of chatbots as an innovative method to offer essential social sup-  
129 port. They serve as a valuable tool in fostering users' reflection on their emotional  
130 self-awareness, social awareness, and interpersonal relationships, which will be  
131 described in detail in the following section. Considering the context of our study  
132 and the prevalence of loneliness among young adults, particularly in the era of the  
133 COVID-19 pandemic, our study has focused on loneliness in our psychometric  
134 assessments.

## 135 2.2. *Chatbots for Mental Health*

136 Chatbots have great potential to promote mental health by conversing with  
137 users to provide psychological assessment, training, and therapy [10]. For ex-  
138 ample, Woebot <sup>1</sup> and Wysa <sup>2</sup> are representative chatbots for mental health; and  
139 their efficacy has been proven by clinical research [37, 38]. To assess users' emo-  
140 tional state or the severity of a specific mental health issue, some chatbots ask  
141 questions based on some well-known psychological scales, such as PHQ-9 De-  
142 pression Test Questionnaire [39] and Generalized Anxiety Disorder Assessment  
143 (GAD-7) [40]. Performing assessment in a chatbot tends to be an effective way to  
144 collect mental health data, comparable to physical interviews in terms of response  
145 rate [11]. Based on users' responses to the assessment questions, chatbots pro-  
146 vide empathetic responses, emotion diary, mindfulness exercises, and goal setting  
147 to help users cope with mental health issues [41, 42]. Existing Human-Computer

---

<sup>1</sup><https://woebothealth.com/>

<sup>2</sup><https://www.wysa.io/>

148 Interaction (HCI) research in mental health chatbots focuses on improving conver-  
149 sation skills to demonstrate compassion and empathy [43, 44] and promote user  
150 self-disclosure [45, 46], and integrating various practices for mental health (e.g.,  
151 expressive writing [47], motivational interview [48], and social support [49]) into  
152 chatbots. However, little work has studied the psychological assessment design  
153 and its impacts on user responses in a survey chatbot for mental health.

### 154 2.3. Design for Online Psychological Assessment

155 The computer-based psychological assessment allows users to employ valid  
156 psychological scales to quickly gauge a specific mental health aspect such as  
157 loneliness, anxiety, and depression [50]. The psychological assessment is often  
158 performed by asking users to answer a set of closed-ended questions, similar to  
159 the questionnaire. Interaction style and questionnaire length are two major design  
160 factors that could influence the participation rate and response quality of a ques-  
161 tionnaire [51, 52, 53, 54, 55]. Therefore, we mainly review the related work of  
162 *interaction style* and *questionnaire length* that we have manipulated in our study.

#### 163 2.3.1. Interaction Style

164 Prior work shows mixed effects of the interaction style on user responses to  
165 questionnaires. The ways of showing the questions (multiple short pages vs. a  
166 long scrollable page) and adding more interactive elements (i.e., pop-up menus,  
167 button scales, and numerical labeling) do not yield a significant difference in user  
168 response behavior [56, 57]. In contrast, compared to the item-by-item questions,  
169 showing questions in a matrix may increase non-response items [58]. Addition-  
170 ally, interaction style could affect users' perceived credibility of information on  
171 the web [59]. Within a chatbot, some social characteristics (e.g., proactivity and  
172 conscientiousness) could also influence users' perceived credibility [60]. As such,  
173 we hypothesize that *interaction type* of psychological assessment would influence  
174 the assessment credibility (**H1**).

175 Previous studies show that adding interactive elements (e.g., interactive prob-  
176 ing and interactive feedback) to the questionnaire could improve the response  
177 quality for the follow-up open-ended questions [51, 61]. Compared with the  
178 form-based questionnaire, the conversation-based survey behaves as a virtual in-  
179 terviewer and intrinsically enriches interactivity through conversation, enhancing  
180 the response quality [14] and enjoyment [62]. Therefore, we hypothesize that  
181 the conversation-based psychological assessment would lead to higher enjoyment  
182 (**H2**) and higher response quality in open-ended questions (**H3**).

### 183 2.3.2. Questionnaire Length

184 Numerous studies have investigated the effects of questionnaire length on a  
185 variety of indicators of a questionnaire, such as participation rates [53], dropout  
186 rates [54, 55], and response quality [25, 24]. Although longer questionnaires  
187 may discourage initial participation due to a higher response burden, no empir-  
188 ical evidence indicates “shorter is better” [63]. The short questionnaires are often  
189 criticized due to lower reliability [63]. As such, we hypothesize that the shorter  
190 questionnaire would negatively influence assessment credibility (**H4**). Moreover,  
191 participating in a psychological assessment can enhance self-awareness [64], and  
192 a longer assessment requires users to spend more time reflecting on their mental  
193 status, which may increase mental health awareness. Thus, we hypothesize that  
194 a longer questionnaire could lead to a higher self-awareness of loneliness in our  
195 study (**H5**).

196 According to a meta-analysis of response rates in web surveys [65], the length  
197 is not always associated with response rates. Nevertheless, adopting a longer  
198 questionnaire generally tends to decrease the response rate and cause a higher  
199 dropout rate [54, 23]. However, the quality of the responses does not necessarily  
200 deteriorate with a lengthy questionnaire as long as participants’ motivation can be  
201 maintained [25].

### 202 2.4. Closed-EQs versus Open-EQs

203 The *closed-ended questions* (Closed-EQs) and *open-ended questions* (Open-  
204 EQs) are two major types of questions in web surveys. Closed-EQs are more ef-  
205 fective for gathering quantitative data [66], and Open-EQs perform better at mea-  
206 suring knowledge and obtaining more reliable and in-depth information [67, 16].  
207 However, Open-EQs may increase the burden of the respondents [68] and the  
208 non-response rate due to more required cognitive efforts [69, 70]. Prior work  
209 showed the correlation between the responses to Closed-EQs and those to Open-  
210 EQs in web surveys for job satisfaction and user experience of e-service websites.  
211 Precisely, the dissatisfied employees, as measured via Closed-EQs about job sat-  
212 isfaction, were more likely to provide negative responses to Open-EQs [20] and  
213 disclose more content of negative feelings in Open-EQs [19]. Likewise, users with  
214 negative experiences of the e-service measured by Likert scale questions (a kind  
215 of Closed-EQs) tended to respond more to the comment-specific Open-EQs than  
216 those with positive experiences [18].

217 A mental health survey chatbot may ask users to answer Closed-EQs for a  
218 psychological assessment and Open-EQs for additional or detailed information  
219 regarding the assessment results. However, it is unclear how the psychological

220 assessment design could influence user responses to Open-EQs in a survey chatbot  
221 for mental health. Previous studies have mainly revealed the relationship between  
222 Closed-EQs and Open-EQs based on user responses [18, 20, 19], while our work  
223 aims to investigate how the design aspects of Closed-EQs (i.e., interaction style  
224 and questionnaire length) influence users' responses to Open-EQs for collecting  
225 more in-depth data about mental health.

## 226 2.5. Perceptions of Mental Health Survey

227 Our study measures user perceptions of the mental health survey in terms of  
228 assessment credibility, self-awareness, and enjoyment.

### 229 2.5.1. Assessment Credibility

230 The users' perception of the psychological assessment results [71] (named *as-*  
231 *essment credibility* in this work) is crucial as it could affect their health-related  
232 behaviors and decisions [72, 73]. Broadly speaking, the psychological assessment  
233 result is a type of health information. Previous studies have revealed several fac-  
234 tors that could influence the perceived credibility of online health information, in-  
235 cluding source expertise [74, 75, 76] (i.e., the rating of the source), website design  
236 (e.g., layout, interactivity, visual design) [77, 76], the language used online [75],  
237 and ease of use [77].

### 238 2.5.2. Self-Awareness

239 Self-awareness refers to being conscious of users' own feelings, thoughts, be-  
240 liefs, and behaviors, which is key to effective counseling and psychotherapy [78].  
241 In the context of mental health, self-awareness is more about emotional self-  
242 awareness that can be gauged from four aspects: identifying emotions, empathy,  
243 managing emotions, and social skills [79]. Psychological assessment provides  
244 users with early problem detection and feedback, which in turn increases their  
245 self-awareness and general knowledge [64]. Thus, the design of these assess-  
246 ments is fundamental in fostering users' self-awareness regarding their mental  
247 health status.

### 248 2.5.3. Enjoyment

249 Enjoyment is a hedonic experience with which users deeply engage in an en-  
250 joyable activity [80]. Lin et al. [81] proposed a scale to measure enjoyment of  
251 the web experience based on three dimensions: engagement, positive affect, and  
252 fulfillment. Several studies have demonstrated the positive effects of chatbots on



253 the effectiveness of surveys [62, 82] and the persuasion of health insurance rec-  
254 ommendations [83], which are mediated by perceived enjoyment. Furthermore,  
255 enabling chatbot self-disclosure [45] or anthropomorphic cues [84, 85] can im-  
256 prove enjoyment, in turn promoting behavioral intentions (e.g., intention to use).

## 257 2.6. Evaluation of User Responses to Open-EQs

258 The main goal of asking Open-EQs is to collect richer data logically concern-  
259 ing response quality and self-disclosure [86]. Previous studies on survey chatbots  
260 evaluate user responses to Open-EQs mainly from response quality and the degree  
261 of self-disclosure [87, 21].

### 262 2.6.1. Response Quality

263 Compared to the responses to Closed-EQs, the responses to Open-EQs are  
264 free-form answers in an open text format, the quality of which can be gauged by  
265 some objective metrics such as response length, number of themes, response time,  
266 and item non-response [88]. For the Open-EQs in a chatbot, researchers employ  
267 Gricean Maxims (i.e., informativeness, specificity, relevance, and clarity) [21],  
268 readability [89], and sentiment intensity [90] to measure response quality.

### 269 2.6.2. Self-Disclosure

270 As an indicator of user engagement in chatbots, self-disclosure measures to  
271 what extent users would like to share their personal information, thoughts, and  
272 feelings [91], which is particularly important for the chatbot to understand the  
273 users' mental status [46]. Various self-reported instruments, such as Jourard Self-  
274 Disclosure Questionnaire (JSDQ) [92], Distress Disclosure Index (DDI) [93], and  
275 Self-Disclosure Index (SDI) [94], have been developed to measure self-disclosure  
276 by asking participants to rate their tendency to disclose information about their  
277 attitudes, opinions, and feelings on a Likert scale. Besides, the self-disclosure  
278 can also be rated by assessors from breadth (i.e., the range of discussed topics)  
279 and depth (i.e., the level of details discussed for a specific topic ) [95]. Our study  
280 adopts both *subjective* and *objective* measurements to gauge self-disclosure in the  
281 user responses to Open-EQs. As the level of self-awareness is found to be pos-  
282 itively related to self-disclosure during computer-mediated communication [96],  
283 we, therefore, hypothesize that users' self-disclosure is positively associated with  
284 self-awareness (**H6**). Additionally, the credibility of health information could in-  
285 fluence the self-disclosure of personal health information [97, 98]. As such, we  
286 hypothesize that a higher level of assessment credibility would lead to a higher  
287 degree of self-disclosure (**H7**) for Open-EQs.

### 288 3. Method

289 We employed a mixed method of qualitative and quantitative approaches to  
290 study how two design features of the psychological assessment (i.e., interaction  
291 style and questionnaire length) influence user perceptions of the assessment and  
292 user responses to Open-EQs.

#### 293 3.1. Study Background

294 To address our raised research questions in a real-world setting of mental  
295 health service, we designed and developed a chatbot (called Percy) to help college  
296 students cope with loneliness during COVID-19 in collaboration with the Coun-  
297 seling and Development Center (CDC) of Hong Kong Baptist University (HKBU)  
298 that provides free and confidential counseling to students as well as consultation  
299 and referral services for staff. Participants were recruited through email invitations  
300 sent by the CDC of the university. We took precautions to minimize potential bi-  
301 ases and priming effects by providing clear instructions and ensuring participants  
302 understood the purpose of the study without explicitly influencing their responses  
303 toward loneliness. Percy bot has three distinct functions: 1) psychological assess-  
304 ment of loneliness and overall mood [Figures 1(a-d)], 2) asking Open-EQ to get  
305 additional information about the feeling of loneliness [Figure 1(e)], and 3) offer-  
306 ing some practical suggestions for managing loneliness [Figure 1(f)], for example,  
307 “*Call a friend or join an online group.*”

#### 308 3.2. Participants

309 The study targets college students who experience loneliness during the COVID-  
310 19 pandemic. The Research Ethics Committee of Hong Kong Baptist University  
311 granted ethics [human (non-clinical)] clearance approval for this study. We re-  
312 cruited 330 participants using mailing lists and public bulletin boards for three  
313 weeks. As a result, 266 participants successfully finished the entire study. To en-  
314 sure the quality of data, we filtered participants by four criteria: 1) the detected  
315 outliers (N=14) having extraordinarily long or short completion time based on the  
316 interquartile range (IQR), 2) the participants (N=10) who failed in two attention  
317 check questions, 3) the participants (N=7) who gave the meaningless responses  
318 (e.g., “nono” and “xxx”) to all the Open-EQs, 4) the participants (N=22) who  
319 gave the same answers to all the questions asked in the post-study. Finally, we  
320 kept 213 valid participants for further analyses. Among those 213 valid partici-  
321 pants, 80.28% of them (N=171) are female (because HKBU has a 1.7 : 1 ratio of

322 female students to male students <sup>3</sup>), 89.67% of them (N=191) are 18 to 25 years  
323 old, 7.98% of them (N=17) are aged 25 to 30, and 2.35% (N=5) are older than  
324 30. In addition, 78.87% of participants (N=168) are Hong Kong locals, and the  
325 rest are international students. To thank participants for supporting our research,  
326 30 participants who completed the study were drawn to receive a supermarket  
327 coupon valued at 200 HKD ( $\approx$ 25.7 USD).

### 328 3.3. Design Manipulations

#### 329 3.3.1. Manipulation of Interaction Style

330 We offered two interaction styles for answering the questions in the psycho-  
331 logical assessment: *form-based* and *conversation-based*. The choice of the two  
332 alternative interaction styles for the psychological assessment is based on review-  
333 ing the user interface design guidelines of several major conversational platforms  
334 such as Messenger<sup>4</sup> and WhatsApp<sup>5</sup>. For example, the form-based interaction is  
335 proposed based on the Webview in Messenger.

336 *Form-based*. The Percy bot offered an alternative way to present the questions  
337 of a psychological assessment in which all questions are embedded in a web form  
338 (see Figure 1(b)). We think the form-based interaction could increase psycho-  
339 logical assessment efficiency while maintaining the interactivity of assessing their  
340 mental health in the chatbot.

---

<sup>3</sup>[https://intl.hkbu.edu.hk/student-exchange/incoming-students/why-hkbu/  
fast-facts](https://intl.hkbu.edu.hk/student-exchange/incoming-students/why-hkbu/fast-facts)

<sup>4</sup><https://developers.facebook.com/docs/messenger-platform>

<sup>5</sup><https://www.facebook.com/brand/resources/whatsapp/user-interface>

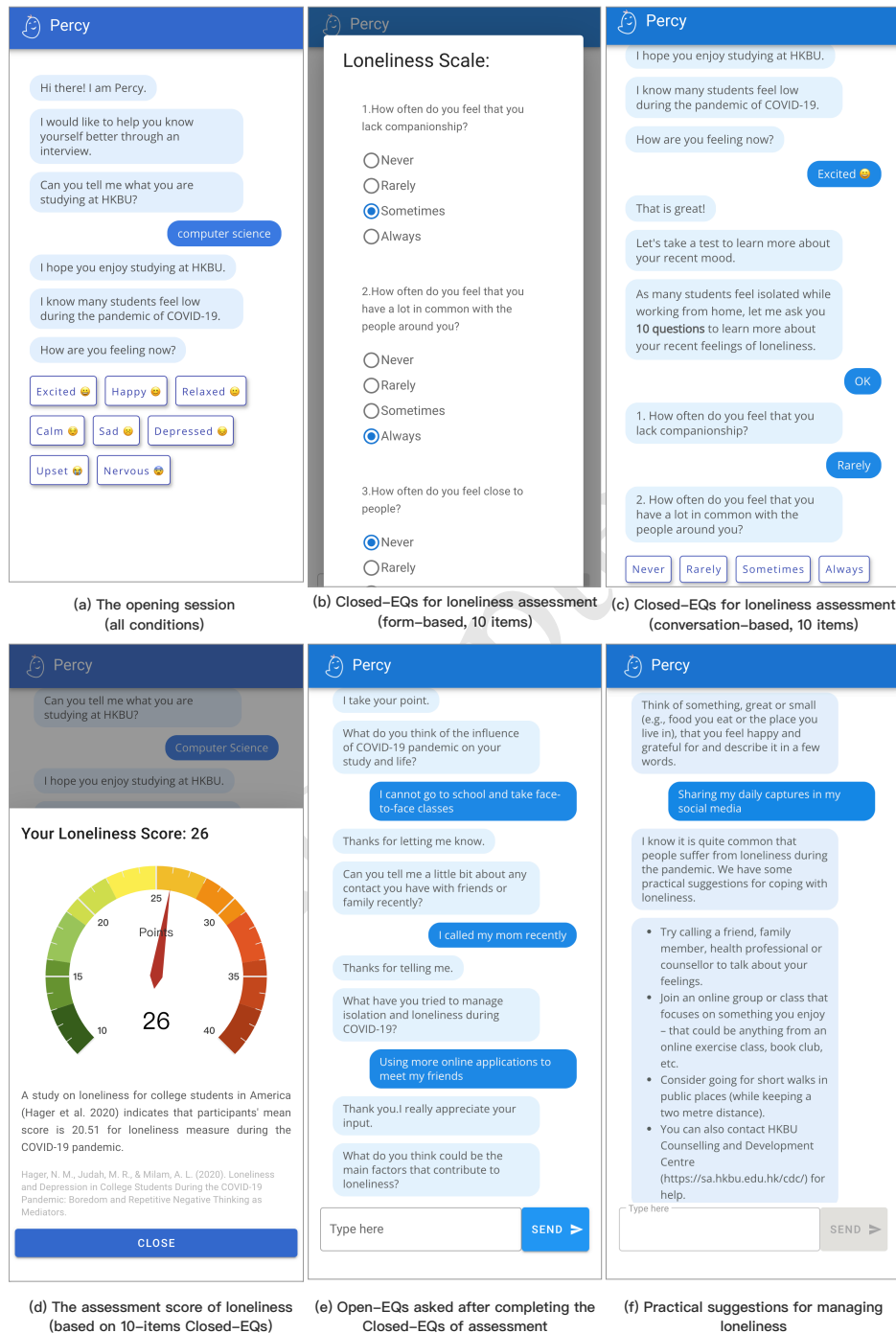


Figure 1: Screenshots of Percy bot: (a) the opening session of conversation and mood recording, (b) the loneliness assessment with the web form, (c) the loneliness assessment in the conversation, (d) the result of loneliness assessment, (e) Open-EQ for getting additional information about the feeling of loneliness, and (f) practical suggestions for coping with loneliness.

341 *Conversation-based.* In this condition, all the loneliness psychological as-  
342 sessment questions were presented in the conversational style. Users can an-  
343 swer a question by clicking one of the buttons under the dialog in conversation  
344 that contains, for instance, selecting one from four options: “Never”, “Rarely”,  
345 “Sometimes”, and “Always” (see Figure 1(c)). The transformation from a web  
346 survey to a conversational survey could improve response quality and user en-  
347 gagement [14, 62].

### 348 3.3.2. *Manipulation of Questionnaire Length*

349 The longer questionnaire can result in a “straight-line” response pattern, which  
350 means more identical answers to most Closed-EQ [25]. Thus, we think the ques-  
351 tionnaire length could influence users’ patience and carefulness in the psycho-  
352 logical assessment. Moreover, the increased response burden caused by a long  
353 questionnaire may influence response quality and response length for Open-EQs.

354 In this study, our chatbot specializes in surveying university students’ lone-  
355 liness during the pandemic of COVID-19. UCLA loneliness scale is the most  
356 widely used instrument for assessing loneliness [17], and it has three validated  
357 length versions, including three items, ten items, and twenty items, respectively  
358 [99, 17]. Based on the three versions, we determined three questionnaire lengths:  
359 short (three items), middle (ten items), and long (twenty items). The questions  
360 in the short version are measured on a three-point scale (1 = Hardly Ever; 2 =  
361 Some of the Time; 3 = Often) [99], while the questions in the middle and long  
362 versions are rated on a four-point scale (1 = Never; 2 = Rarely; 3 = Sometimes; 4  
363 = Always) [17].

### 364 3.4. *User Study Design and Procedure*

365 Based on our two independent variables, *interaction style* and *questionnaire*  
366 *length*, we designed a 2 (interaction style: form-based vs. conversation-based)  
367  $\times$  3 (questionnaire length: short vs. middle vs. long) between-subjects study.  
368 Figure 2 shows an overview of the study design, including the following three  
369 major phases:

370 ***Pre-study.*** First, we asked all participants to sign a consent form and read an  
371 information page describing Percy’s main features and explaining the steps they  
372 should follow to finish the study. After that, we asked participants to answer three  
373 questions about their demographics, including age, gender, and nationality.

374 Moreover, we asked participants to indicate their current mood from eight  
375 options based on two dimensions of core-affect [100], including excited, happy,  
376 relaxed, calm, sad, depressed, upset, and nervous (Figure 1(a)).

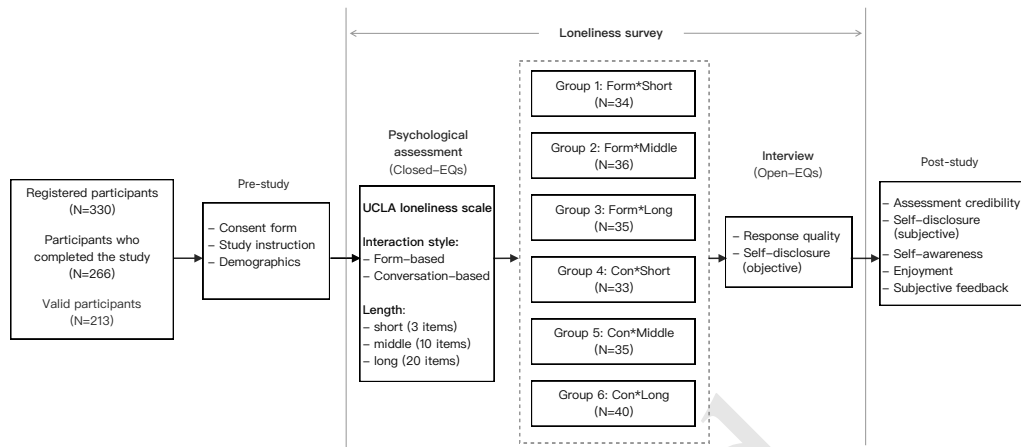


Figure 2: User study design and procedure.

377 **Loneliness survey.** The loneliness survey contains a psychological assess-  
 378 ment (measured by Closed-EQs) and an interview (measured by Open-EQs). The  
 379 psychological assessment has six variants combining two design manipulations:  
 380 *interaction style* and *questionnaire length*. Following the between-subjects de-  
 381 sign, we randomly assigned participants to one of six conditions. When users  
 382 finished the psychological assessment, a result page popped up, showing a loneli-  
 383 ness score, a semicircle meter with color gradients for the score level, and an ex-  
 384 planation with a reference for the score (Figure 1(d)). The participants were then  
 385 guided to the interview session after closing the result page. During the interview,  
 386 the chatbot asked seven Open-EQ (see Table A.4 in Appendix A) to understand  
 387 the participants’ feelings of loneliness during COVID-19 deeply. As the chatbot’s  
 388 responses may likely influence how users chat with it [21], our chatbot only gen-  
 389 erated some general responses to users’ answers to avoid such interference. These  
 390 responses vary and depend on the content of users’ answers, for example, “*Thank*  
 391 *you. I appreciate your input.*” or “*Thank you for your thoughtful input.*” are  
 392 possible responses for the user answers of rich content, e.g., “*I wish to be around*  
 393 *my family more often where I can be myself more. I also think exercising regu-*  
 394 *larly can help.*”, while “*Got it.*” or “*I understand!*” are for simple and brief user  
 395 answers, e.g., “*It’s fine.*” or “*nothing*”.

396 **Post-study.** Participants were required to complete a questionnaire containing  
 397 sixteen five-point Likert scale questions (Table 1) to indicate their perceived as-  
 398 sessment credibility, self-awareness, enjoyment, and self-disclosure. In addition,  
 399 we asked participants to answer five Open-EQs (see Table B.5 in Appendix B)

400 to understand their in-depth opinions on Percy.

401 *3.5. Measurement and Analysis*

402 This study measured users’ perceptions of the loneliness survey based on as-  
 403 sessment credibility, self-awareness, and enjoyment. Moreover, we adopted sev-  
 404 eral metrics for response quality and subjective and objective measures for self-  
 405 disclosure in user responses to Open-EQs.

Table 1: Post-Study Questionnaire for Measuring User Perceptions of the Survey and Self-Disclosure

Construct	Item	Loading
<b>Assessment Credibility</b> (Cronbach alpha: 0.894; AVE: 0.741)		
	I am convinced that the score can indicate my feelings of loneliness.	0.709
	I am confident I will trust my loneliness score.	0.770
	The loneliness score calculated by the Percy bot can be trusted	0.674
<b>Self-Awareness</b> (Cronbach alpha: 0.818; AVE: 0.607)		
	I have insight into myself.	
	I recognize the stress and worry in my current life.	0.696
	I understand myself well.	
	I generally feel positive about self-awareness.	0.581
	The Percy bot made me aware of my loneliness.	0.754
<b>Enjoyment</b> (Cronbach alpha: 0.841; AVE: 0.649)		
	I enjoy talking with the Percy bot.	0.716
	I feel enjoyable when I converse with the Percy bot.	0.798
	I would like to answer survey questions with the Percy bot.	0.612
<b>Self-Disclosure (subjective)</b> (Cronbach alpha: 0.758; AVE: 0.610)		
	I think I have told my real feelings to the Percy bot.	0.605
	I think I have provided sufficient information to the Percy bot.	0.578
	The design of the interview Percy bot made me think longer about my responses compared to traditional surveys.	
	If time allows, I would like to spend more time elaborating my responses to let the Percy bot understand me better.	
	I am not willing to reveal my feelings to the Percy bot. (reversed)	

*Note:* The items marked in gray were dropped due to a poor loading value (< 0.5) or high cross-loading value (> 12) measured by modification index [101].

406 *3.5.1. Perceptions of Loneliness Survey*

407 Perceptions of the loneliness survey refer to participants' feelings and attitudes  
408 towards the loneliness assessment (Closed-EQs) and the interview (Open-EQs).  
409 We employed a set of questions (see Table 1) to measure three constructs: assess-  
410 ment credibility, self-awareness, and enjoyment. All these questions were mea-  
411 sured on a five-point Likert scale. We run a confirmatory factor analysis (CFA)  
412 to establish the validity of these question items. Commonly accepted cutoff val-  
413 ues for convergent validity are 0.7 for Cronbach's alpha, 0.5 for average variance  
414 extracted (AVE) [102], and 0.5 for factor loading.

- 415 • *Assessment credibility.* It measures to what extent the psychological assess-  
416 ment result can be trusted and believed. According to Hilligoss and Rieh's  
417 credibility framework consisting of three levels of credibility judgments:  
418 construct, heuristics, and interaction [103], We composed three questions  
419 to measure participants' perceived credibility of their loneliness assessment  
420 (Cronbach alpha: 0.894; AVE: 0.741).
- 421 • *Self-awareness.* Self-awareness is the participant's ability to know and un-  
422 derstand their feelings and behaviors. We measured self-awareness based  
423 on the three validated questions of a Self-Awareness Outcomes Question-  
424 naire (SAOQ) [104] (Cronbach alpha: 0.818; AVE: 0.607).
- 425 • *Enjoyment.* It gauges how much the participants enjoyed chatting with  
426 Percy. We used three validated questions from a questionnaire for evalu-  
427 ating recommendations in a mental health app [105] to measure enjoyment  
428 (Cronbach alpha: 0.841; AVE: 0.649).

429 *3.5.2. Response Quality*

430 In this study, we did not measure the response quality of Closed-EQs using  
431 methods such as differentiation response index (i.e., satisficing behavior of choos-  
432 ing the same response every time) [106] because these metrics are usually applied  
433 to assessing whether participants are serious and attentive for answering the ques-  
434 tions in general surveys such as internet usage behavior [14] and course satisfac-  
435 tion [62]. In our opinion, the motivation for completing a mental health survey  
436 differs from answering a general survey. The participants are more motivated by a  
437 need to understand their mental health status more accurately. Moreover, choosing  
438 the same response to all the questions in a short psychological assessment (e.g.,  
439 the short loneliness assessment with five Closed-EQs) does not necessarily mean  
440 satisfying behavior.



441 For Open-EQs, we measured the response quality based on Gricean Maxims  
 442 theory [107] that has often been used to evaluate the quality of users’ responses  
 443 in chatbots [87, 21]. Gricean Maxims was developed based on the cooperative  
 444 principle for enabling effective conversational communication by concretely con-  
 445 sidering four aspects: quantity, quality, relevance, and manner [108]. According  
 446 to the definition of Gricean Maxims, the aspect of “quality” refers to being truth-  
 447 ful in communication. Due to the general difficulty in assessing the truthfulness  
 448 of user responses [21], we did not measure this aspect. In our study, we concretely  
 449 adopted four quality metrics (i.e., informativeness, specificity, relevance, clarity)  
 450 used to evaluate user responses to Open-EQs in a chatbot [21], which were pro-  
 451 posed based on three Gricean Maxims aspects: quantity, relevance, and manner  
 452 (see Table 2). We measured these metrics based on user responses to all Open-EQ  
 453 asked by our Percy bot.

Table 2: Quality Metrics Defined Based on Gricean Maxims [21]

Gricean Maxims	Definition	Quality Metric	Definition
Quantity	One should be as informa- tive as possible.	<i>Informativeness</i>	A participant’s response should be as informative as possible.
		<i>Specificity</i>	A participant’s response should give as much infor- mation as needed.
Relevance	One should provide relevant information.	<i>Relevance</i>	A participant’s response should be relevant to a question.
Manner	One should communicate in a clear and orderly manner.	<i>Clarity</i>	A participant’s response should be clear.

- 454 • *Informativeness*. Per the maxim of quantity, the communication should be  
 455 as informative as possible. The measure of informativeness in users’ re-  
 456 sponses based on Formula (1) [21] that calculates the sum of a word’s sur-  
 457 prisal based on the inverse of its occurrence frequency in four major English  
 458 corpora, including British National Corpus [109], the Brown Corpus [110],  
 459 Webtext <sup>6</sup>, and the NPS Chat Corpus [111].

<sup>6</sup><https://github.com/teropa/nlp/tree/master/resources/corpora/webtext>

$$I(\text{Response}) = \sum \log_2 \frac{1}{F(\text{word}_n)} \quad (1)$$

460 • *Response quality index.* We measured the overall response quality by re-  
 461 sponse quality index (RQI) [21] that combines three quality metrics: speci-  
 462 ficity, relevance, and clarity, as shown in Formula (2) and respectively de-  
 463 fined in Table 2. The measures of the three quality metrics follow a man-  
 464 ual assessment method, and we defined three levels (0,1,2) for each met-  
 465 ric. In total, we collected 1,491 text responses from 213 participants. We  
 466 followed a standard coding protocol to code each response. First, we ran-  
 467 domly selected 10% of responses and then asked two researchers to finish  
 468 the coding independently. After that, they discussed the differences in cod-  
 469 ing, and a third researcher was involved in voting for the irreconcilable dif-  
 470 ferences. The coding criteria became more consistent after the discussion.  
 471 Finally, they finished coding for the rest of the responses. The Cohen’s  
 472 kappa of each set of coding (Specificity:  $\kappa=0.73$ , Relevance:  $\kappa=0.81$ , Clar-  
 473 ity:  $\kappa=0.89$ ) indicates good inter-rater reliability of the coded items <sup>7</sup>.

$$RQI = \sum_{n=1}^N \text{specificity}[i] * \text{relevance}[i] * \text{clarity}[i] \quad (2)$$

(N is the number of responses in a completed assessment)

474 Table 3 shows some examples of our coded responses. *Specificity* refers to  
 475 the level of details the response provides, and a specific response should  
 476 convey meaningful insights (0 – generic description only, 1 – specific con-  
 477 cepts, and 2 – specific concepts with detailed examples). *Relevance* mea-  
 478 sures to which extent the answer is relevant to the question asked during the  
 479 interview (0 – irrelevant, 1 – somewhat relevant, and 2 – relevant). *Clarity* is  
 480 measured based on the human effort of understanding the text (0 – illegible  
 481 text, 1 – incomplete sentences, and 2 – clearly articulated response).

### 482 3.5.3. Self-Disclosure

483 Self-disclosure involves sharing personal thoughts, feelings, or experiences  
 484 about oneself with others [113]. The quality of user responses to Open-EQs in a  
 485 survey is linked to the extent of self-disclosure [86], signifying the extent to which

---

<sup>7</sup>Slight: 0.0-0.2; Fair: 0.21-0.4; Moderate: 0.41-0.6; Substantial: 0.61-0.8; Almost Perfect: 0.81-1 [112].

Table 3: Examples of Coded Responses to the Open-Ended Question Open-EQ7 (“Think of something that you feel happy and grateful for, great or small (e.g., *the food you eat or the place you live in*).”)

Response Example	Rating
<i>“my family, including my father, even though he had passed away. Also, my husband. All about love; I know they love me even though I don’t know how to express the gratitude.”</i>	Specificity:2, Relevance:2, Clarity:2, Self-disclosure:2
<i>“Money”</i>	Specificity:2, Relevance:1, Clarity:0, Self-disclosure:0
<i>“Listening to my favorite music and watching my favorite reality show .”</i>	Specificity:2, Relevance:2, Clarity:1, Self-disclosure:1
<i>“Everything will be fine.”</i>	Specificity:1, Relevance:2, Clarity:0, Self-disclosure:0

486 users are willing to share information with the chatbot. In Open-EQs, we assessed  
 487 self-disclosure based on users’ subjective feelings and objective metrics of user  
 488 responses, such as the breadth and depth of content.

- 489 • *Self-disclosure (subjective)*. It assesses participants’ subjective perspec-  
 490 tives on sharing their feelings and thoughts about loneliness. The questions  
 491 for measuring subjective self-disclosure, as depicted in Table 1, have been  
 492 adapted from those used to evaluate user responses in a survey chatbot [62]  
 493 (Cronbach’s alpha: 0.758; AVE: 0.610).
- 494 • *Self-disclosure (objective)*. It gauged the extent to which participants shared  
 495 their personal feelings and thoughts with the chatbot. We manually evalu-  
 496 ated the level of self-disclosure based on the breadth and depth of topics  
 497 conveyed in user responses to the seven Open-EQs (0 – a brief description  
 498 with no specific topic, 1 – a brief description with a specific topic, and 2 –  
 499 a detailed description with one specific topic / a description with multiple  
 500 topics) [91]. The self-disclosure coding demonstrated substantial inter-rater  
 501 reliability, as evidenced by Cohen’s kappa score of 0.69. As illustrated in  
 502 the example (the first example in Table 3), higher levels of self-disclosure  
 503 may encompass more detailed and private topics.

### 504 3.6. Interaction Behavior

505 We also recorded response length for Open-EQs and engagement duration to  
 506 understand better how much users would like to interact with the chatbot.

- 507 • *Response length.* Response length was counted by the number of words  
508 in each participant’s responses to all seven Open-EQs during the interview.  
509 The response length is usually proportional to the engagement duration.
- 510 • *Engagement duration.* Engagement duration measured the time a partici-  
511 pant spent answering all the Open-EQs in the interview session of the lone-  
512 liness survey. A longer engagement duration could mean the participant  
513 invests more effort in thinking and answering the Open-EQ.

## 514 4. Results

515 This section presents the main results related to each research question. For  
516 the convenience of illustration, we use an expression of **interaction\*length** to  
517 denote each experimental condition in the remaining parts of this manuscript.  
518 In this expression, interaction can be “Con” or “Form”, respectively standing  
519 for *conversation-based* and *form-based*, and length can be “Short”, “Middle”,  
520 or “Long”. For example, Con\*Middle refers to the condition where participants  
521 assessed their loneliness by completing the middle-length UCLA loneliness scale  
522 (ten items) through conversation-based interaction for Closed-EQs.

523 To investigate two design factors (i.e., interaction style and questionnaire length),  
524 we employed a 2x3 factorial design in our study. Additionally, we need to run  
525 multiple regression analyses to test our research hypotheses. To achieve this, we  
526 have opted to use structural equation modeling (SEM) to analyze our results, given  
527 its capacity to evaluate multivariate causal relationships simultaneously within a  
528 statistical estimation procedure [114]. Table C.6 in Appendix C presents the  
529 descriptive statistics of the dependent variables (DVs) for six experimental condi-  
530 tions derived from a 2x3 factorial design.

### 531 4.1. Structural Equation Modeling

532 We use *lavaan*,<sup>8</sup> an R package to build our SEM model. Some dependent vari-  
533 ables (DVs), such as informativeness, engagement duration, and response length,  
534 were measured differently from the five-point Likert scale for measuring the DVs  
535 related to user perceptions, resulting in much larger values. Therefore, we nor-  
536 malized the values of these dependent variables by using the *scale()* function in  
537 R, which scales the data based on the mean value and the standard deviation. In  
538 addition, as our data do not conform to the normal distribution, we chose a more

---

<sup>8</sup><https://lavaan.ugent.be/>

539 robust estimator, “MLR,” in our SEM analysis. The sample size of our study  
 540 meets a CFA/SEM rule of thumb that 10:1 is the recommended ratio of subjects  
 541 to observable variables (N:q) [115] and the recommended sufficient sample size  
 542 (N = 200) for structural equation modeling [116, 117]. Following the procedure  
 543 of trimming non-significant paths in SEM model [118], we obtain our resulting  
 544 model (see Figure 3) showing a good fit <sup>9</sup>:  $\chi^2(149) = 209.323, p = .003$  <sup>10</sup>; root  
 545 mean squared error of approximation (RMSEA) = 0.044; 90% CI: [0.029, 0.057];  
 546 Comparative Fit Index (CFI) = 0.969; Tucker-Lewis Index (TLI) = 0.963. In addition,  
 547 we utilized the R package, *semPower*, <sup>11</sup> to execute a post-hoc power analysis  
 548 for our obtained model. The analysis revealed a high power level (power > .98)  
 549 with a sample size of N = 213 to identify misspecifications of a model (involving  
 550 df = 149 degrees of freedom) corresponding to RMSEA  $\geq$  .05 at an alpha error  
 551 level of .05.

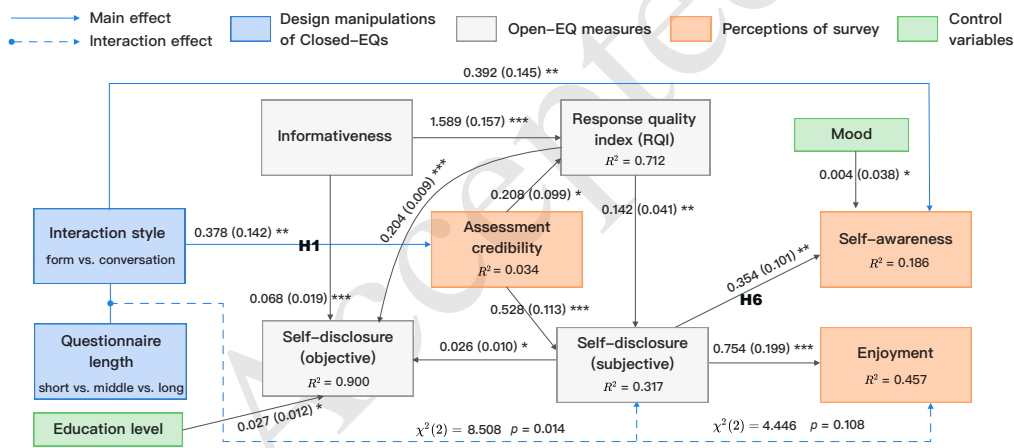


Figure 3: The structural equation model for our user study’s data. Significance levels: \*\*\*  $p < .001$ , \*\*  $p < .01$ , \*  $p < .05$ . The numbers on the edges refer to the  $\beta$  coefficient and standard error (in parentheses) of the causal relationship.  $R^2$  is the proportion of variance explained by the model. Factors are scaled to have an SD of 1. The paths labeled with H1 and H6 indicate these two paths support hypotheses H1 and H6.

<sup>9</sup>Hu and Bentler [119] proposed cutoff values for several fit indices to be: CFI > .96, TLI > .95, and RMSEA < .05, with the upper bound of its 90% CI below 0.10.

<sup>10</sup>A model should not have a non-significant  $\chi^2$ , but this statistic is regarded as too sensitive [120].

<sup>11</sup><https://github.com/moshagen/semPower>

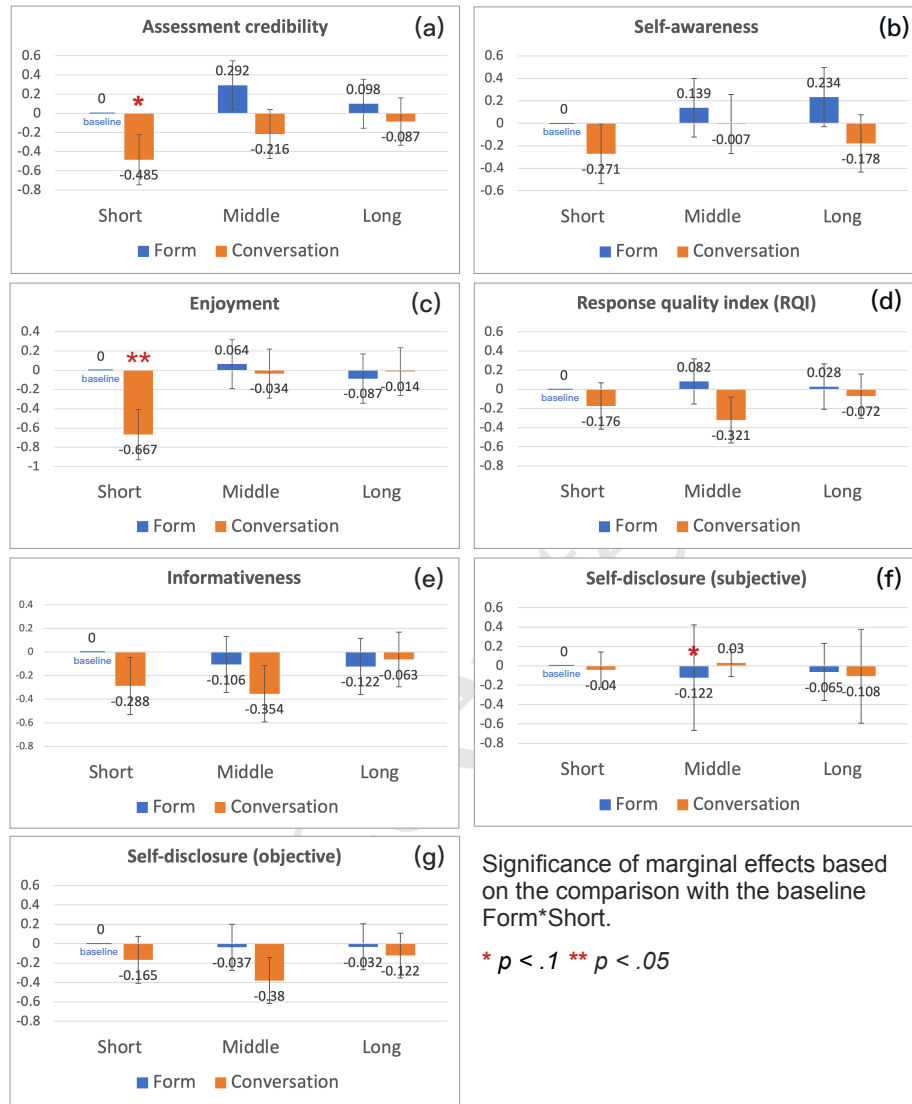


Figure 4: Marginal effects of interaction style and questionnaire length on different DVs. The effects of the baseline Form\*Short are set to zero, and the y-axis is scaled by the sample standard deviation. Significance levels: \*\* $p < .05$ , \*  $p < .1$ .

552 In addition, to understand how the values of a dependent variable (e.g.,  
 553 assessment credibility) change with variation of the independent variable (IV) (e.g.,  
 554 interaction style), we analyzed the marginal effects of the two IVs (i.e., interac-  
 555 tion style and questionnaire length) on each DV, assuming other covariates to be

556 fixed [121]. Figure 4 shows the marginal effects of dependent variables that are  
557 associated with significant main effects or interaction effects of two design factors.  
558 In order to effectively gauge or test our hypothesis, we also consider the potential  
559 influence of control variables (such as age, gender, education level, and mood)  
560 on the dependent variable. The findings indicate that education level significantly  
561 impacts self-disclosure, while mood significantly affects self-awareness.

#### 562 4.2. *The Effects of Interaction Style with Closed-EQs on Perceptions of the Survey* 563 *(RQ1)*

564 The SEM model (Figure 3) shows a direct positive effect of interaction style on  
565 assessment credibility ( $\beta = 0.378, p < .01$ ). Moreover, as depicted in Figure 4(a),  
566 the conversation-based design appears to compromise user-perceived assessment  
567 credibility, particularly when combined with the short questionnaire. Con\*Short  
568 was lower than the baseline with marginal significance ( $p < .1$ ). Thus, we can  
569 accept the hypothesis **H1**: the form-based psychological assessment would lead  
570 to higher assessment credibility. Moreover, the model shows no other significant  
571 effects of interaction style on enjoyment and response quality in Open-EQs. Thus,  
572 we cannot accept the hypothesis **H2**: the conversation-based psychological assess-  
573 ment leads to higher enjoyment, and the hypothesis **H3**: the conversation-based  
574 psychological assessment leads to higher response quality in Open-EQs. The  
575 marginal effects on enjoyment (Figure 4(c)) indicate that combining conversation-  
576 based interaction and a short questionnaire could lower enjoyment, and Con\*Short  
577 is significantly lower than the baseline in terms of enjoyment ( $p < .05$ ). In addi-  
578 tion to testing our hypothesized effects, the model shows a significant effect of  
579 interaction style on self-awareness ( $\beta = 0.392, p < .01$ ). The marginal effects  
580 on self-awareness (Figure 4(b)) show that form-based interaction leads to higher  
581 self-awareness than conversation-based interaction regardless of the questionnaire  
582 length.

#### 583 4.3. *The Effects of Questionnaire Length on Perceptions of the Survey (RQ2)*

584 Manipulating questionnaire length does not directly affect any investigated  
585 measures for users' perceptions of the survey. Thus, we could not accept the hy-  
586 pothesis **H4**: a shorter questionnaire leads to lower assessment credibility, and  
587 the hypothesis **H5**: a longer questionnaire leads to higher self-awareness. Even  
588 though not statistically significant, users seem to perceive higher assessment credi-  
589 bility with the form-based design when completing a middle questionnaire (refer  
590 to Figure 4(a)), and they attain increased self-awareness by completing a longer  
591 questionnaire (as seen in Figure 4(b)). Furthermore, we find an interaction effect

592 of interaction style and questionnaire length on enjoyment, which is marginally  
593 significant,  $\chi^2(2) = 4.446$ ,  $p = .108$ . In other words, the effects of questionnaire  
594 length on enjoyment depend on the interaction style. Specifically, the distinc-  
595 tion between the short questionnaire and questionnaires of other lengths is more  
596 pronounced with conversation-based interaction than with form-based interaction  
597 (see Figure 4(c)).

#### 598 4.4. *The Effects of Interaction Style with Closed-EQs on User Responses to Open-* 599 *EQs (RQ3)*

600 The SEM model (Figure 3) does not show any direct effect of interaction  
601 style on response quality and self-disclosure measures. Despite no significant  
602 direct main effects of interaction style on response quality, the form-based de-  
603 sign could positively influence self-disclosure (subjective and objective) and RQI  
604 through assessment credibility. The assessment credibility positively influences  
605 self-disclosure (subjective) ( $\beta = 0.528$ ,  $p < .001$ ) and RQI ( $\beta = 0.208$ ,  $p < .05$ ),  
606 which in turn positively influences self-disclosure (objective). Thus, the signif-  
607 icant effects of assessment credibility on self-disclosure (subjective) and self-  
608 disclosure (objective) allow us to accept the hypothesis **H7**: higher credibility  
609 leads to more self-disclosure in Open-EQs.

610 Specifically, the significant paths (*P1*: Interaction style  $\rightarrow$  Assessment credi-  
611 bility  $\rightarrow$  Self-disclosure (subjective)  $\rightarrow$  Self-disclosure (objective)) and (*P2*: In-  
612 teraction style  $\rightarrow$  Assessment credibility  $\rightarrow$  RQI  $\rightarrow$  Self-disclosure (objective))  
613 indicate a *mediating role* of assessment credibility in the effects of interaction style  
614 on self-disclosure (objective) in Open-EQs. Figure 4(g) shows that regardless of  
615 the questionnaire length, conversation-based interaction results in lower levels of  
616 self-disclosure (objective) compared to form-based interaction. However, the total  
617 indirect effect of assessment credibility on self-disclosure (objective) is minimal  
618 ( $\beta = 0.057$ ).

#### 619 4.5. *The Effects of Questionnaire Length on User Responses to Open-EQs (RQ4)*

620 The model does not show any main effects of questionnaire length on response  
621 quality. The marginal effects of questionnaire length on RQI and informative-  
622 ness illustrate the non-significant difference caused by the manipulation of ques-  
623 tionnaire length (see Figure 4(d and e)). Compared with the baseline condition  
624 (Form\*Short), the short and middle questionnaires lead to lower response quality  
625 with the conversation-based design.

626 Despite no main effect of questionnaire length on self-disclosure measures, we  
627 find a significant interaction effect of interaction style and questionnaire length on



628 self-disclosure (subjective),  $\chi^2(2) = 8.508$ ,  $p < .05$ , indicating that the effect of  
629 questionnaire length on self-disclosure (subjective) depends on interaction style.  
630 For instance, the marginal effect on subjective self-disclosure (Figure 4(f)) indi-  
631 cates that the middle questionnaire results in the highest subjective self-disclosure,  
632 with marginal significance ( $p < .01$ ) when combined with conversation-based in-  
633 teraction, whereas it leads to the lowest subjective self-disclosure when combined  
634 with form-based interaction.

#### 635 4.6. *Relations Between User Responses to Open-EQs and Perceptions of the Sur-* 636 *vey*

637 The model also reveals the relationships between the perceptions of the survey  
638 (i.e., enjoyment and self-awareness) and user responses to Open-EQs. Specifi-  
639 cally, the significant path (P3: Informativeness  $\rightarrow$  RQI  $\rightarrow$  Self-disclosure (subjec-  
640 tive)  $\rightarrow$  Self-Awareness & Enjoyment) confirms the mediated effects of informa-  
641 tiveness and response quality on self-awareness and enjoyment. As self-disclosure  
642 (subjective) positively influences self-awareness ( $\beta = 0.354$ ,  $p < .01$ ), we could  
643 accept the hypothesis **H6**: higher self-disclosure is positively associated with self-  
644 awareness. Interestingly, self-disclosure (subjective) has a strong positive effect  
645 on enjoyment ( $\beta = 0.754$ ,  $p < .001$ ), indicating that participants who are willing  
646 to disclose their personal feelings and experiences are more likely to perceive en-  
647 joyment while interacting with the survey chatbot. Moreover, the significant path  
648 (P4: Assessment credibility  $\rightarrow$  Self-disclosure (subjective)  $\rightarrow$  Self-Awareness &  
649 Enjoyment) suggests that participants who perceive higher assessment credibility  
650 tend to disclose their feelings and thoughts about loneliness with the chatbot and  
651 then perceive higher self-awareness and enjoyment.

#### 652 4.7. *Interaction Behavior*

653 We recorded the number of words in each participant's responses to all Open-  
654 EQs (response length) and the total time they spent answering them (engage-  
655 ment duration). Design manipulations do not directly affect response length and  
656 engagement duration. Nevertheless, the conversation-based interaction leads to  
657 shorter responses than the form-based interaction, and the condition of Form\*Middle  
658 has the longest response on average ( $M=60.9$  words,  $SD=44.6$ ). Furthermore, the  
659 questionnaire length positively influences engagement duration when adopting the  
660 conversation-based interaction, and the condition of Form\*Middle has the longest  
661 engagement duration ( $M=339.8$  seconds,  $SD=277.7$ ).

662 4.8. *Subjective Feedback*

663 To better understand participants' subjective experiences of two design manip-  
664 ulations in our survey chatbot, we performed a thematic analysis [122] based on  
665 participants' responses to the five Open-EQs in the post-study (Table B.5). Two  
666 authors independently finished half of the responses and addressed the conflicts  
667 in coding through additional discussion, resulting in an almost perfect inter-rater  
668 agreement among coding tested by Cohen's kappa ( $\kappa = 0.85$ )<sup>12</sup>. One author fin-  
669 ished coding the remaining responses and discussed them with another author to  
670 reach a consensus on the codes.

671 **The Length of Questionnaire.** Using a short questionnaire could potentially  
672 diminish the credibility of the assessment. Although the questionnaire length  
673 does not significantly influence assessment credibility according to the quantita-  
674 tive analysis, a short questionnaire seems to decrease users' perceived assessment  
675 credibility. Certain participants trying with the short version of the assessment  
676 believed that incorporating more question items could enhance the credibility of  
677 the test, as two participants noted,

678 *"I think there could be more questions to indicate my loneliness score*  
679 *better."* (P7, Form\*Short)

680 *"I don't think people's loneliness can be scored when people just an-*  
681 *swer three questions."* (P170, Con\*Short)

682 **Interaction Style of Psychological Assessment.** Moreover, compared with  
683 the form-based interaction, the conversation-based interaction offers a more casual  
684 way for users to answer the questions measured by the Likert scale. However, it  
685 may also make the questionnaire perceived as less formal, aligning with the result  
686 of quantitative analysis. One participant called,

687 *"It is just like chatting. But I don't really agree with the score, and it*  
688 *may need an adjustment to have more options. Maybe 0 to 10."* (P40,  
689 Con\*Middle)

690 It seems that presenting the questions of a psychological assessment via the  
691 conversational style decreases the questionnaire's formality [60], which in turn  
692 influences users' perceived assessment credibility. However, some participants

---

<sup>12</sup>Slight: 0.0-0.2; Fair: 0.21-0.4; Moderate: 0.41-0.6; Substantial: 0.61-0.8; Almost Perfect: 0.81-1 [112].

693 doubted the assessment’s credibility because of the ambiguous measurement stan-  
694 dard for loneliness; for example,

695 *“...these are some general questions, cannot be sure if the score is*  
696 *trustworthy cause people have different standards.”* (P206, Con\*Middle)

697 Additionally, a few participants also complained about the increased interac-  
698 tion time caused by the conversation. For example, one participant stated,

699 *“Filling in an online form can be boring if there are too many ques-*  
700 *tions. Chatting with the Percy bot is interesting, at least with more*  
701 *interaction. But chatting with a bot can be time-consuming.”* (P137,  
702 Form\*Middle)

703 **Psychological Assessment Result.** The assessment score is key to self-awareness.  
704 Many participants claimed that they became more aware of their loneliness status  
705 by finishing the psychological assessment. One participant noted,

706 *“I think the questions asked were relevant for calculating the loneli-*  
707 *ness score. I am aware of what my feelings are during the pandemic.”*  
708 (P108, Form\*Middle)

709 Some participants thought the reference on the result page (see Figure 2(d))  
710 showing the mean score of others who completed this loneliness assessment helped  
711 them better understand their loneliness status.

712 *“Comparing to the mean score, I know more about my status among*  
713 *people.”* (P137, Form\*Short)

## 714 **5. Discussion**

715 Prior research has highlighted the benefits of using a survey chatbot as com-  
716 pared to a conventional survey delivered through web forms. This study delves  
717 deeper into the refined design aspects of a survey chatbot within the scope of  
718 mental health. More specifically, we explore the impact of the interaction style  
719 and length of psychological assessments featuring Closed-EQs on the quality of  
720 responses to subsequent Open-EQs within a survey chatbot. Thus, the findings  
721 from this investigation are contextualized within a survey chatbot environment  
722 that presents both Closed-EQs and Open-EQs.

723 Before discussing the results of our study, we first briefly summarize our re-  
724 search findings based on quantitative and qualitative results.

- 725 1. **The interaction style of psychological assessment significantly affects**  
726 **the assessment credibility and self-awareness. The influenced assess-**  
727 **ment credibility could influence response quality and self-disclosure for**  
728 **Open-EQs.** The participants who completed the psychological assessment  
729 via the form-based interaction were more convinced by the assessment,  
730 thereby being more engaged in responding to the follow-up Open-EQs and  
731 being more aware of their feelings.
- 732 2. **The questionnaire length does not significantly impact the assessment**  
733 **credibility and user responses to Open-EQs.** Although there is an in-  
734 teraction effect between interaction style and questionnaire length on self-  
735 disclosure (subjective) and enjoyment, questionnaire length has no signifi-  
736 cant main effect on any dependent variables.
- 737 3. **The assessment credibility mediates the effects of psychological assess-**  
738 **ment design on users' responses to Open-EQs.** The psychological assess-  
739 ment design has *indirect* positive impacts on users' self-disclosure (objec-  
740 tive) and response quality index (RQI) through the assessment credibility.

### 741 5.1. *Psychological Assessment Design*

742 The psychological assessment is vital for monitoring mental health status and  
743 delivering timely adaptive interventions in a mental health survey chatbot [123].  
744 This is especially crucial when access to mental health services is limited, as seen  
745 during events like the COVID-19 pandemic [124]. With this in mind, our inves-  
746 tigation focuses on how the design of the psychological assessment with Closed-  
747 EQs could impact users' perceptions of the assessment and their responses to  
748 Open-EQs in a survey chatbot.

#### 749 5.1.1. *Interaction Style of Closed-EQs*

750 Our study investigated two interaction styles of psychological assessment with  
751 closed-ended questions in a survey chatbot: form-based and conversation-based.  
752 Previous studies have demonstrated the benefits of conversation-based design over  
753 form-based design for the entire survey in terms of response quality [14, 62, 21],  
754 without making a distinction between Closed-EQs and Open-EQs. However, we  
755 found that within a survey chatbot, the form-based interaction leads to higher  
756 assessment credibility with Closed-EQs, which in turn leads to higher response  
757 quality in Open-EQs. We argue that survey design for psychological assessments  
758 is different from surveying course satisfaction [62], gamers' opinions [21], and

759 Internet usage behavior [14] in previous studies. In contrast to traditional surveys,  
760 the psychological assessment is frequently succeeded by an assessment score or  
761 report, aiming to provide users with an understanding of their health status and en-  
762 courage positive health behavior changes [50]. This process may lead participants  
763 to take the assessment questions more seriously, as inaccurate self-assessments  
764 could potentially impact mental health [125].

765 Despite the benefits of casual communication (e.g., more communicative [126],  
766 or a strong feeling of being involved [127]), formal communication has been  
767 proven to be associated with high information credibility [128]. Furthermore, a  
768 prior study showed that with a task-oriented chatbot, users are more likely to feel  
769 like performing a task in a natural, casual, informal conversation rather than in  
770 goal-directed settings [129]. Therefore, we speculate that the casual communica-  
771 tion conveyed by the conversation-based design may decrease the users' perceived  
772 formality of assessment and weaken their perceived assessment credibility.

773 Moreover, our study shows that the conversation-based interaction signifi-  
774 cantly increases interaction time than the form-based one while adopting a long  
775 questionnaire (Figure 4(i)), which aligns the findings of a previous study on a  
776 survey chatbot with Closed-EQs [14]. Unlike the responses to Open-EQs, which  
777 could be diverse free-text inputs, the responses to Closed-EQs are based on pre-  
778 defined content, such as the Likert scale or multiple-choice. We think that the  
779 increased response time of the psychological assessment may imply a lower ef-  
780 ficiency of assessment rather than higher user engagement. The conversational  
781 interaction may especially cause users' displeasure at the slow pace of complet-  
782 ing a long questionnaire. Therefore, we wonder how we may make a trade-off  
783 between the advantages of the conversation-based design (e.g., natural interac-  
784 tion, less non-differentiation in a rating task, aka a "straight-line" response [14])  
785 and its disadvantages (e.g., low efficiency). For example, one participant (P137,  
786 Form\*Short) stated, "*Chatting with the Percy bot is quite interesting, at least with  
787 more interaction. However, chatting with a bot can be time-consuming.*" Thus,  
788 a form-based design could be more suitable for presenting a questionnaire in a  
789 chatbot because it maintains the formality and efficiency of the questionnaire and  
790 does not influence users' perceived interactivity of responding to the follow-up  
791 Open-EQs in the survey chatbot.

792 Therefore, we suggest **adopting a form-based design for the psychological**  
793 **assessment in a survey chatbot for mental health.** Although the conversation-  
794 based design has distinct advantages over the form-based design, such as inter-  
795 active content [14], reciprocity [45], and human-like communication [44, 21], it  
796 also imposes more interaction time on users [14, 21]. More notably, the form-

797 based design makes participants perceive higher assessment credibility than the  
798 conversation-based. Therefore, chatbot designers could embed a form-based psy-  
799 chological assessment into the chatbot before asking Open-EQs through conver-  
800 sation. This hybrid design may also combat the survey-taking fatigue in case the  
801 participants are expected to be more engaged in responding to the Open EQs [21].  
802 On the one hand, users may feel they are still answering questions in the chatbot;  
803 on the other hand, they may focus more on questionnaire content, which is less  
804 tedious than following the humdrum conversation pattern to answer Closed-EQs.

### 805 *5.1.2. The Length of Questionnaire with Closed-EQs*

806 Information completeness is a major factor that influences the perceived cred-  
807 ibility of health information [130]. The length of the questionnaire reflects how  
808 much information is collected for assessment, which could affect the complete-  
809 ness of the assessment information. Thus, we investigated how the questionnaire  
810 length influences the assessment credibility. However, we did not find a signif-  
811 icant main effect of questionnaire length on users' perceived assessment credi-  
812 bility, probably because the participants did not perceive significantly different  
813 assessment results regarding information completeness with three different ques-  
814 tionnaire lengths (short, middle, and long). Moreover, our results also indicate  
815 that questionnaire length does not have a significant main effect on the response  
816 quality and self-disclosure in Open-EQs, which echoes the findings of prior work  
817 that the response quality of Open-EQs is not associated with the survey length  
818 [63, 24]. Thus, keeping the assessment as short as possible is unnecessary, but the  
819 content (questions) of the psychological assessment should satisfy the users' as-  
820 sessment needs [63]. Additionally, the significant interaction effects of interaction  
821 style and questionnaire length on enjoyment and subjective self-disclosure in the  
822 follow-up Open-EQs suggest that the determination of questionnaire length might  
823 also depend on the questionnaire's interaction style. Therefore, we suggest that  
824 **designers may determine the questionnaire length based on user needs and**  
825 **the interaction style of the questionnaire.**

826 Moreover, according to a recent literature survey on the instruments used in  
827 the psychological assessment of mental health and health behavior [50], among  
828 21 surveyed questionnaires (e.g., GAD-7 for anxiety [40], PHQ-7 for depres-  
829 sion [39], PHQ-15 for physical symptoms [131]), the questionnaire length varies  
830 from 2 to 28 items, similar to the range used in our study. Consequently, our find-  
831 ings regarding the impact of questionnaire length could potentially be applied to  
832 scenarios utilizing other psychological assessments.

### 833 5.1.3. Assessment Credibility

834 Users' perceived credibility of health information significantly impacts their  
835 behavioral intention of using the health informatics service [132]. In our study, the  
836 structural equation model (Figure 3) demonstrates a mediating role of assessment  
837 credibility in the effects of the interaction style of psychological assessment on the  
838 metrics evaluating users' responses to Open-EQs. The users' perceived credibility  
839 of assessment is critical to the mental health survey, as it could influence user  
840 engagement in the activities at a later stage [132], for example, answering Open-  
841 EQs in a mental health survey.

842 Online health information can be categorized mainly into scientific and expe-  
843 riential information [133]. The results of the psychological assessment provided  
844 by the agent belong to scientific information, the credibility of which is mainly as-  
845 sessed based on reference credibility [133]. Thus, our psychological assessment  
846 result (score) page also shows an academic reference (Figure 1(d)) to justify the  
847 interpretation of the assessment score (Figure 1(d)). However, we wonder if par-  
848 ticipants could notice the study's reference and how much it may help them justify  
849 the result. Our qualitative results indicate that although we provide a descriptive  
850 explanation of the psychological assessment results based on a reference (Fig-  
851 ure 1(d)), some participants still do not trust the assessment score due to the am-  
852 biguous measurement standard for loneliness, for example, "...cannot sure if the  
853 score is trustworthy cause people have different standards." (P206, Con\*Middle)  
854 Therefore, the future design may allow users to ask for further explanations of  
855 the psychological assessment results through conversation. When addressing user  
856 inquiries about assessment results, the conversational explanation may be consid-  
857 ered more convincing by users due to the persuasive potential of the chatbot [83].

858 In general, the credibility of information on the web can also be influenced by  
859 multiple aspects of the information medium, such as content format, design of user  
860 interface, and interactivity [59]. With the evolution of human-computer interac-  
861 tion, virtual agents' simulated human-human interaction is increasingly popular  
862 for mental health because of greater interactivity that supports therapeutic con-  
863 versation [134]. However, should we deliver all the services in a mental health  
864 chatbot through conversation? For the psychological assessment, our study re-  
865 sults suggest that the participants perceived higher assessment credibility with  
866 the form-based assessment questionnaire than with the conversation-based ques-  
867 tionnaire. As most mental health surveys still adopt form-based questionnaires,  
868 the conversation-based interaction style probably does not conform to the partici-  
869 pants' mental model of taking a psychological assessment.

## 870 5.2. *User Responses to Open-EQs*

871 We evaluated user responses to Open-EQs in our survey chatbot from multiple  
872 aspects, among which self-disclosure and response quality have more often been  
873 emphasized in the previous studies [44, 21, 45].

874 Self-disclosure refers to revealing personal and even sensitive information to  
875 others [135]. Prior work has identified its important role in building trust [136] and  
876 intimacy [137] for communication. In our study, users' subjective self-disclosure  
877 is satisfying (above 3.8 out of 5) in all the experimental conditions. Still, their ob-  
878 jective self-disclosure (below 1 out of 2) is not as good as the subjective measure.  
879 The discrepancy between the two measures might be due to our chatbot's limited  
880 social skills. Since our study has aimed to investigate the effects of psychological  
881 assessment design on users' self-disclosure in Open-EQs, we did not incorporate  
882 the social characteristics into the chatbot design, such as proactivity (e.g., ac-  
883 tive listening [44]) and emotional intelligence (e.g., empathetic responses [138]),  
884 which, however, could encourage honest self-disclosure during the communica-  
885 tion [139].

886 The interaction style indirectly influences subjective and objective self-disclosure  
887 through assessment credibility, while questionnaire length does not (Figure 3).  
888 Despite the lack of a main effect of questionnaire length, questionnaire length  
889 seems to influence the effect of interaction style on self-disclosure (subjective).  
890 Although participants thought the design manipulations of psychological assess-  
891 ment did not significantly influence their willingness to disclose themselves (sub-  
892 jective self-disclosure) for Open-EQs, in practice, they showed more self-disclosure  
893 in form-based conditions than in conversation-based conditions. This may imply  
894 that the form-based interaction is more favorable than the conversation-based in-  
895 teraction regarding users' self-disclosure in their responses to Open-EQs.

896 We measured the response quality of Open-EQs from multiple dimensions,  
897 and the Form\*Middle design leads to the highest response quality index (RQI),  
898 and the Form\*Short design has the highest informativeness. We argue that per-  
899 ceiving higher assessment credibility in the form-based questionnaire motivates  
900 participants who feel lonely to talk with the survey chatbot. Furthermore, the re-  
901 sponse quality of Open-EQs is highly associated with objective self-disclosure,  
902 which aligns with the findings of existing work [21, 140].

## 903 6. **Limitations**

904 Our study has several limitations that need to be mentioned while interpreting  
905 our research findings, including the unbalanced gender distribution, narrow scope



906 of mental health, and limited social communication skills of our chatbot.

907 *First*, our primary target group is university students who may suffer from  
908 loneliness. To reach a broad audience, we have collaborated with the Counsel-  
909 ing and Development Center (CDC) of Hong Kong Baptist University (HKBU)  
910 to recruit participants within the university. However, we encountered an imbal-  
911 ance in the gender distribution of our participants, primarily because HKBU has a  
912 higher ratio of female students. In addition, existing research suggests that lone-  
913 liness is more commonly experienced by males than females [141]. However, the  
914 analysis of gender as a control variable on all dependent variables did not yield  
915 significance. Therefore, the gender imbalance should not significantly impact the  
916 generalizability of our findings.

917 *Second*, we investigated the design of the psychological assessment only for  
918 loneliness because the loneliness scale has three validated length versions, which  
919 meets our requirement of manipulating questionnaire lengths as short, middle,  
920 and long. Strictly speaking, loneliness is not a mental health issue, but it is closely  
921 related to various mental health issues such as anxiety, stress, and depression [142,  
922 143]. Lonely people may behave differently from those who suffer from mental  
923 health issues regarding self-disclosure intentions. For example, lonely people are  
924 more willing to disclose private information than those connected [144], while  
925 individuals with depression and anxiety are associated with lessened emotional  
926 self-disclosure [145]. Therefore, further study is needed to validate to what extent  
927 our findings on the psychological assessment design can be generalized to a survey  
928 chatbot for screening other mental health issues.

929 *Third*, our current survey design is that Open-EQs were positioned immedi-  
930 ately after Closed-EQs. While this sequential arrangement is common in mental  
931 health survey design, there are some alternative methods to mix Open-EQs and  
932 Closed-EQs. For example, participants could explain their choices of a Closed-  
933 EQ through the following Open-EQ. This highlights the need for further research  
934 to explore diverse psychological measurement design approaches in survey chat-  
935 bots.

936 *Fourth*, since we have focused on investigating the impacts of the psychologi-  
937 cal assessment design on user responses to Open-EQs, our survey chatbot provides  
938 relatively unified responses according to the length of users' responses. For ex-  
939 ample, "*I understand.*" or "*Thank you. I really appreciate your input.*". However,  
940 some participants expected to receive more meaningful and personalized feedback  
941 while conversing with the chatbot. For example, "*...the bot response does not re-*  
942 *ply authentically according to my response.*" (P56, Con\*Middle) In the future, we  
943 plan to incorporate sophisticated social communication skills, such as active lis-

944 tening [44] and bot self-disclosure [46, 45] into a survey chatbot for mental health.  
945 Besides, the chatbot powered by large language models (LLMs) [146], e.g., Chat-  
946 GPT,<sup>13</sup> has demonstrated an impressive ability to understand and generate natural  
947 language in conversation. Therefore, we will consider leveraging the LLMs to  
948 generate engaging and empathetic responses so as to improve user engagement in  
949 the survey chatbot for mental health.

## 950 7. Conclusions

951 We conducted a field study (N=213) that investigated how two prominent de-  
952 sign factors of the psychological assessment (i.e., *interaction style* and *question-*  
953 *naire length*) influence user responses to the open-ended questions (Open-EQs)  
954 in a survey chatbot for mental health. The results indicate that the form-based  
955 interaction is more favored than the conversation-based interaction for the psy-  
956 chological assessment regarding users' perceived assessment credibility and self-  
957 awareness. The increased assessment credibility could further stimulate more  
958 self-disclosure and quality responses in Open-EQs. Moreover, although the ques-  
959 tionnaire length has a limited impact on user responses to Open-EQs, we suggest  
960 that the questionnaire length could be adapted to the assessment purpose and con-  
961 tent or be determined based on participants' time pressure. To the best of our  
962 knowledge, most existing works on mental health chatbots focus on enhancing  
963 chatbots' communication skills to increase user engagement and response qual-  
964 ity [44, 21, 46, 45]. However, little work has investigated the potential effect of  
965 the psychological assessment design in a survey chatbot for mental health. Fi-  
966 nally, we explain our findings through an SEM model containing all design fac-  
967 tors, response quality and self-disclosure in Open-EQs, and the users' perceptions  
968 of the survey. By investigating two prominent design factors of the psychologi-  
969 cal assessment in a survey chatbot for mental health, we believe that the findings  
970 could be suggestive for researchers and practitioners to better leverage the chatbot  
971 technology for improving the quality and user experience of their mental health  
972 survey.

## 973 Acknowledgements

974 This work was supported by the Hong Kong Research Grants Council (RGC)  
975 GRF project (RGC/HKBU12201620), the Hong Kong Baptist University IG-FNRA

---

<sup>13</sup><https://chat.openai.com/chat>

976 project (RC-FNRA-IG/21-22/SCI/01), and the Hong Kong Baptist University Start-  
977 up Grant (RC-STARTUP/21-22/23). In addition, we want to thank our collabo-  
978 rators working in the Counseling and Development Center (CDC) of HKBU, Ms.  
979 Vicki Kwan, Ms. Chloe Li, and Ms. Wendy Cheung, and all the participants who  
980 participated in our study.

## 981 **References**

- 982 [1] D. Eisenberg, Countering the troubling increase in mental health symp-  
983 toms among us college students, *Journal of Adolescent Health* 65 (5) (2019)  
984 573–574.
- 985 [2] R. Mojtabai, M. Olfson, B. Han, National trends in the prevalence and  
986 treatment of depression in adolescents and young adults, *Pediatrics* 138 (6)  
987 (2016).
- 988 [3] R. D. Goodwin, A. H. Weinberger, J. H. Kim, M. Wu, S. Galea, Trends  
989 in anxiety among adults in the united states, 2008–2018: Rapid increases  
990 among young adults, *Journal of psychiatric research* 130 (2020) 441–446.
- 991 [4] T. Gagné, I. Schoon, A. Sacker, Trends in young adults’ mental distress and  
992 its association with employment: Evidence from the behavioral risk factor  
993 surveillance system, 1993–2019, *Preventive Medicine* 150 (2021) 106691.
- 994 [5] N. R. Magson, J. Y. Freeman, R. M. Rapee, C. E. Richardson, E. L. Oar,  
995 J. Fardouly, Risk and protective factors for prospective changes in ado-  
996 lescent mental health during the covid-19 pandemic, *Journal of youth and*  
997 *adolescence* 50 (1) (2021) 44–57.
- 998 [6] L. Liang, H. Ren, R. Cao, Y. Hu, Z. Qin, C. Li, S. Mei, The effect of covid-  
999 19 on youth mental health, *Psychiatric quarterly* 91 (3) (2020) 841–852.
- 1000 [7] D. Courtney, P. Watson, M. Battaglia, B. H. Mulsant, P. Szatmari, Covid-19  
1001 impacts on child and youth anxiety and depression: challenges and oppor-  
1002 tunities, *The Canadian Journal of Psychiatry* 65 (10) (2020) 688–691.
- 1003 [8] E. J. Costello, Early detection and prevention of mental health problems:  
1004 developmental epidemiology and systems of support, *Journal of Clinical*  
1005 *Child & Adolescent Psychology* 45 (6) (2016) 710–717.

- 1006 [9] J. M. Levitt, N. Saka, L. H. Romanelli, K. Hoagwood, Early identifica-  
1007 tion of mental health problems in schools: The status of instrumentation,  
1008 *Journal of School Psychology* 45 (2) (2007) 163–191.
- 1009 [10] A. A. Abd-Alrazaq, M. Alajlani, A. A. Alalwan, B. M. Bewick, P. Gardner,  
1010 M. Househ, An overview of the features of chatbots in mental health: A  
1011 scoping review, *International Journal of Medical Informatics* 132 (2019)  
1012 103978.
- 1013 [11] I. Hungerbuehler, K. Daley, K. Cavanagh, H. G. Claro, M. Kapps, et al.,  
1014 Chatbot-based assessment of employees’ mental health: Design process  
1015 and pilot implementation, *JMIR Formative Research* 5 (4) (2021) e21678.
- 1016 [12] A. Schick, J. Feine, S. Morana, A. Maedche, U. Reininghaus, et al., Validity  
1017 of chatbot use for mental health assessment: Experimental study, *JMIR*  
1018 *mHealth and uHealth* 10 (10) (2022) e28082.
- 1019 [13] M. E. Te Pas, W. G. Rutten, R. A. Bouwman, M. P. Buise, User experi-  
1020 ence of a chatbot questionnaire versus a regular computer questionnaire:  
1021 prospective comparative study, *JMIR Medical Informatics* 8 (12) (2020)  
1022 e21982.
- 1023 [14] S. Kim, J. Lee, G. Gweon, Comparing data from chatbot and web surveys:  
1024 Effects of platform and conversational style on survey response quality, in:  
1025 *Proceedings of the 2019 CHI conference on human factors in computing*  
1026 *systems*, 2019, pp. 1–12.
- 1027 [15] U. Reja, K. L. Manfreda, V. Hlebec, V. Vehovar, Open-ended vs. close-  
1028 ended questions in web questionnaires, *Developments in applied statistics*  
1029 19 (1) (2003) 159–177.
- 1030 [16] O. Friberg, J. H. Rosenvinge, A comparison of open-ended and closed  
1031 questions in the prediction of mental health, *Quality & Quantity* 47 (3)  
1032 (2013) 1397–1411.
- 1033 [17] D. W. Russell, UCLA loneliness scale (version 3): Reliability, validity, and  
1034 factor structure, *Journal of personality assessment* 66 (1) (1996) 20–40.
- 1035 [18] R. Zhou, X. Wang, L. Zhang, H. Guo, Who tends to answer open-ended  
1036 questions in an e-service survey? the contribution of closed-ended answers,  
1037 *Behaviour & Information Technology* 36 (12) (2017) 1274–1284.

- 1038 [19] I. Borg, C. Zuell, Write-in comments in employee surveys, *International*  
1039 *Journal of Manpower* (2012).
- 1040 [20] R. M. Poncheri, J. T. Lindberg, L. F. Thompson, E. A. Surface, A comment  
1041 on employee surveys: Negativity bias in open-ended responses, *Organiza-*  
1042 *tional Research Methods* 11 (3) (2008) 614–630.
- 1043 [21] Z. Xiao, M. X. Zhou, Q. V. Liao, G. Mark, C. Chi, W. Chen, H. Yang, Tell  
1044 me about yourself: Using an ai-powered chatbot to conduct conversational  
1045 surveys with open-ended questions, *ACM Transactions on Computer-*  
1046 *Human Interaction (TOCHI)* 27 (3) (2020) 1–37.
- 1047 [22] S. Ganassali, The influence of the design of web survey questionnaires on  
1048 the quality of responses, in: *Survey research methods*, Vol. 2, 2008, pp.  
1049 21–32.
- 1050 [23] M. Galesic, M. Bosnjak, Effects of questionnaire length on participation  
1051 and indicators of response quality in a web survey, *Public opinion quarterly*  
1052 73 (2) (2009) 349–360.
- 1053 [24] B. Burchell, C. Marsh, The effect of questionnaire length on survey re-  
1054 sponse, *Quality and quantity* 26 (3) (1992) 233–244.
- 1055 [25] A. R. Herzog, J. G. Bachman, Effects of questionnaire length on response  
1056 quality, *Public opinion quarterly* 45 (4) (1981) 549–559.
- 1057 [26] L. A. Peplau, D. Perlman, *Loneliness: A sourcebook of current theory,*  
1058 *research, and therapy*, (No Title) (1982).
- 1059 [27] E. Hards, M. E. Loades, N. Higson-Sweeney, R. Shafran, T. Serafimova,  
1060 A. Brigden, S. Reynolds, E. Crawley, E. Chatburn, C. Linney, et al., Lone-  
1061 liness and mental health in children and adolescents with pre-existing men-  
1062 tal health problems: A rapid systematic review, *British Journal of Clinical*  
1063 *Psychology* 61 (2) (2022) 313–334.
- 1064 [28] J. Christiansen, P. Qualter, K. Friis, S. Pedersen, R. Lund, C. Andersen,  
1065 M. Bekker-Jeppesen, M. Lasgaard, Associations of loneliness and social  
1066 isolation with physical and mental health among adolescents and young  
1067 adults, *Perspectives in public health* 141 (4) (2021) 226–236.

- 1068 [29] J. T. Cacioppo, L. C. Hawkley, R. A. Thisted, Perceived social isolation  
1069 makes me sad: 5-year cross-lagged analyses of loneliness and depressive  
1070 symptomatology in the Chicago Health, Aging, and Social Relations Study.,  
1071 *Psychology and Aging* 25 (2) (2010) 453.
- 1072 [30] M. Lasgaard, K. Friis, M. Shevlin, “where are all the lonely people?” a  
1073 population-based study of high-risk groups across the life span, *Social Psy-  
1074 chiatry and Psychiatric Epidemiology* 51 (2016) 1373–1384.
- 1075 [31] T. E. Keller, M. Perry, R. Spencer, Reducing social isolation through formal  
1076 youth mentoring: opportunities and potential pitfalls, *Clinical Social Work  
1077 Journal* 48 (2020) 35–45.
- 1078 [32] N. A. Mayorga, T. Smit, L. Garey, A. K. Gold, M. W. Otto, M. J. Zvolensky,  
1079 Evaluating the interactive effect of COVID-19 worry and loneliness on mental  
1080 health among young adults, *Cognitive Therapy and Research* (2022) 1–9.
- 1081 [33] K. Cooper, E. Hards, B. Moltrecht, S. Reynolds, A. Shum, E. McElroy,  
1082 M. Loades, Loneliness, social relationships, and mental health in adoles-  
1083 cents during the COVID-19 pandemic, *Journal of Affective Disorders* 289  
1084 (2021) 98–104.
- 1085 [34] S. Marchini, E. Zaurino, J. Bouziotis, N. Brondino, V. Delvenne, M. Del-  
1086 haye, Study of resilience and loneliness in youth (18–25 years old) during  
1087 the COVID-19 pandemic lockdown measures, *Journal of Community Psy-  
1088 chology* 49 (2) (2021) 468–480.
- 1089 [35] A. Martinez, S. Nguyen, The impact of COVID-19 on college student well-  
1090 being (2020).
- 1091 [36] M. Panayiotou, J. C. Badcock, M. H. Lim, M. J. Banissy, P. Qualter, Mea-  
1092 suring loneliness in different age groups: The measurement invariance of  
1093 the UCLA Loneliness Scale, *Assessment* 30 (5) (2023) 1688–1715.
- 1094 [37] K. K. Fitzpatrick, A. Darcy, M. Vierhile, Delivering cognitive behavior  
1095 therapy to young adults with symptoms of depression and anxiety using  
1096 a fully automated conversational agent (Woebot): a randomized controlled  
1097 trial, *JMIR Mental Health* 4 (2) (2017) e7785.
- 1098 [38] B. Inkster, S. Sarda, V. Subramanian, An empathy-driven, conversational  
1099 artificial intelligence agent (Wysa) for digital mental well-being: real-world

- 1100 data evaluation mixed-methods study, *JMIR mHealth and uHealth* 6 (11)  
1101 (2018) e12106.
- 1102 [39] K. Kroenke, R. L. Spitzer, J. B. Williams, The phq-9: validity of a brief  
1103 depression severity measure, *Journal of general internal medicine* 16 (9)  
1104 (2001) 606–613.
- 1105 [40] R. L. Spitzer, K. Kroenke, J. B. Williams, B. Löwe, A brief measure for  
1106 assessing generalized anxiety disorder: the gad-7, *Archives of internal  
1107 medicine* 166 (10) (2006) 1092–1097.
- 1108 [41] K. Denecke, S. Vaaheesan, A. Arulnathan, A mental health chatbot for reg-  
1109 ulating emotions (sermo)-concept and usability test, *IEEE Transactions on  
1110 Emerging Topics in Computing* 9 (3) (2020) 1170–1182.
- 1111 [42] R. R. Morris, K. Kouddous, R. Kshirsagar, S. M. Schueller, Towards an  
1112 artificially empathic conversational agent for mental health applications:  
1113 system design and user perceptions, *Journal of medical Internet research*  
1114 20 (6) (2018) e10148.
- 1115 [43] M. Lee, S. Ackermans, N. Van As, H. Chang, E. Lucas, W. IJsselsteijn,  
1116 Caring for vincent: a chatbot for self-compassion, in: *Proceedings of the  
1117 2019 CHI Conference on Human Factors in Computing Systems*, 2019, pp.  
1118 1–13.
- 1119 [44] Z. Xiao, M. X. Zhou, W. Chen, H. Yang, C. Chi, If i hear you correctly:  
1120 Building and evaluating interview chatbots with active listening skills, in:  
1121 *Proceedings of the 2020 CHI Conference on Human Factors in Computing  
1122 Systems*, 2020, pp. 1–14.
- 1123 [45] Y.-C. Lee, N. Yamashita, Y. Huang, W. Fu, ” i hear you, i feel you”: En-  
1124 couraging deep self-disclosure through a chatbot, in: *Proceedings of the  
1125 2020 CHI conference on human factors in computing systems*, 2020, pp.  
1126 1–12.
- 1127 [46] Y.-C. Lee, N. Yamashita, Y. Huang, Designing a chatbot as a mediator for  
1128 promoting deep self-disclosure to a real mental health professional, *Pro-  
1129 ceedings of the ACM on Human-Computer Interaction* 4 (CSCW1) (2020)  
1130 1–27.

- 1131 [47] S. Park, A. Thieme, J. Han, S. Lee, W. Rhee, B. Suh, “i wrote as if i were  
1132 telling a story to someone i knew.”: Designing chatbot interactions for ex-  
1133 pressive writing in mental health, in: Designing Interactive Systems Con-  
1134 ference 2021, 2021, pp. 926–941.
- 1135 [48] S. Park, J. Choi, S. Lee, C. Oh, C. Kim, S. La, J. Lee, B. Suh, Designing a  
1136 chatbot for a brief motivational interview on stress management: Qualita-  
1137 tive case study, *Journal of medical Internet research* 21 (4) (2019) e12231.
- 1138 [49] P. B. Bae Brandtzæg, M. Skjuve, K. K. Kristoffer Dysthe, A. Følstad, When  
1139 the social becomes non-human: Young people’s perception of social sup-  
1140 port in chatbots, in: Proceedings of the 2021 CHI Conference on Human  
1141 Factors in Computing Systems, 2021, pp. 1–13.
- 1142 [50] R. G. Maunder, J. J. Hunter, An internet resource for self-assessment of  
1143 mental health and health behavior: Development and implementation of  
1144 the self-assessment kiosk, *JMIR mental health* 5 (2) (2018) e9768.
- 1145 [51] J. L. Holland, L. M. Christian, The influence of topic interest and interac-  
1146 tive probing on responses to open-ended questions in web surveys, *Social  
1147 Science Computer Review* 27 (2) (2009) 196–212.
- 1148 [52] F. G. Conrad, M. P. Couper, R. Tourangeau, M. Galesic, T. Yan, Interactive  
1149 feedback can improve the quality of responses in web surveys, in: *Proc.  
1150 Surv. Res. Meth. Sect. Am. Statist. Ass.*, 2005, pp. 3835–3840.
- 1151 [53] M. Revilla, C. Ochoa, Ideal and maximum length for a web survey., *Inter-  
1152 national Journal of Market Research* 59 (5) (2017).
- 1153 [54] P. Edwards, I. Roberts, P. Sandercock, C. Frost, Follow-up by mail in clini-  
1154 cal trials: does questionnaire length matter?, *Controlled clinical trials* 25 (1)  
1155 (2004) 31–52.
- 1156 [55] E. Deutskens, K. De Ruyter, M. Wetzels, P. Oosterveld, Response rate and  
1157 response quality of internet-based surveys: an experimental study, *Market-  
1158 ing letters* 15 (1) (2004) 21–36.
- 1159 [56] A. Peytchev, M. P. Couper, S. E. McCabe, S. D. Crawford, Web survey  
1160 design: Paging versus scrolling, *International Journal of Public Opinion  
1161 Quarterly* 70 (4) (2006) 596–607.



- 1162 [57] U.-D. Reips, Context effects in web surveys, *Online social sciences* (2002)  
1163 69–80.
- 1164 [58] M. Liu, A. Cernat, Item-by-item versus matrix questions: A web survey  
1165 experiment, *Social Science Computer Review* 36 (6) (2018) 690–706.
- 1166 [59] C. N. Wathen, J. Burkell, Believe it or not: Factors influencing credibility  
1167 on the web, *Journal of the American society for information science and  
1168 technology* 53 (2) (2002) 134–144.
- 1169 [60] A. P. Chaves, M. A. Gerosa, How should my chatbot interact? a survey  
1170 on social characteristics in human–chatbot interaction design, *International  
1171 Journal of Human–Computer Interaction* 37 (8) (2021) 729–758.
- 1172 [61] M. Oudejans, L. M. Christian, Using interactive features to motivate and  
1173 probe responses to open-ended questions, *Social and behavioral research  
1174 and the internet: Advances in applied methods and research strategies  
1175* (2010) 304–332.
- 1176 [62] T. Wambsganss, R. Winkler, M. Söllner, J. M. Leimeister, A conversational  
1177 agent to improve response quality in course evaluations, in: *Extended Ab-  
1178 stracts of the 2020 CHI Conference on Human Factors in Computing Sys-  
1179 tems, 2020*, pp. 1–9.
- 1180 [63] S. Rolstad, J. Adler, A. Rydén, Response burden and questionnaire length:  
1181 is shorter better? a review and meta-analysis, *Value in Health* 14 (8) (2011)  
1182 1101–1108.
- 1183 [64] D. W. Eby, L. J. Molnar, J. T. Shope, J. M. Vivoda, T. A. Fordyce, Improv-  
1184 ing older driver knowledge and self-awareness through self-assessment:  
1185 The driving decisions workbook, *Journal of safety research* 34 (4) (2003)  
1186 371–381.
- 1187 [65] C. Cook, F. Heath, R. L. Thompson, A meta-analysis of response rates  
1188 in web-or internet-based surveys, *Educational and psychological measure-  
1189 ment* 60 (6) (2000) 821–836.
- 1190 [66] Y. Guo, R. W. Proctor, G. Salvendy, What do users want to see? a content  
1191 preparation study for consumer electronics, in: *International Conference  
1192 on Human-Computer Interaction, Springer, 2009*, pp. 413–420.

- 1193 [67] J. A. Krosnick, Questionnaire design, in: *The Palgrave handbook of survey*  
1194 *research*, Springer, 2018, pp. 439–455.
- 1195 [68] P. F. Lazarsfeld, The controversy over detailed interviews—an offer for ne-  
1196 *gotiation*, *Public opinion quarterly* 8 (1) (1944) 38–60.
- 1197 [69] M. Denscombe, The length of responses to open-ended questions: A com-  
1198 *parison of online and paper questionnaires in terms of a mode effect*, *Social*  
1199 *Science Computer Review* 26 (3) (2008) 359–368.
- 1200 [70] M. Emde, M. Fuchs, Using adaptive questionnaire design in open-ended  
1201 *questions: A field experiment*, in: *American Association for Public Opin-*  
1202 *ion Research (AAPOR) 67th Annual Conference*, San Diego, USA, 2012,  
1203 *pp. 1–13*.
- 1204 [71] R. H. Dana, T. A. Hoffmann, Health assessment domains: Credibility and  
1205 *legitimization*, *Clinical Psychology Review* 7 (5) (1987) 539–555.
- 1206 [72] B. Kitchens, C. A. Harle, S. Li, Quality of health-related online search  
1207 *results*, *Decision Support Systems* 57 (2014) 454–462.
- 1208 [73] Y. Zhang, Searching for specific health-related information in m edline p  
1209 *lus: Behavioral patterns and user experience*, *Journal of the Association for*  
1210 *Information Science and Technology* 65 (1) (2014) 53–68.
- 1211 [74] M. S. Eastin, Credibility assessments of online health information: The  
1212 *effects of source expertise and knowledge of content*, *Journal of Computer-*  
1213 *Mediated Communication* 6 (4) (2001) JCMC643.
- 1214 [75] E. L. Jenkins, J. Ilicic, A. M. Barklamb, T. A. McCaffrey, Assessing  
1215 *the credibility and authenticity of social media content for applications*  
1216 *in health communication: scoping review*, *Journal of medical Internet re-*  
1217 *search* 22 (7) (2020) e17296.
- 1218 [76] X. Zhao, L. Chen, Y. Jin, X. Zhang, Comparing button-based chatbots with  
1219 *webpages for presenting fact-checking results: A case study of health in-*  
1220 *formation*, *Information Processing & Management* 60 (2) (2023) 103203.
- 1221 [77] L. Sbaffi, J. Rowley, Trust and credibility in web-based health informa-  
1222 *tion: a review and agenda for future research*, *Journal of medical Internet*  
1223 *research* 19 (6) (2017) e218.

- 1224 [78] A. L. Pieterse, M. Lee, A. Ritmeester, N. M. Collins, Towards a model  
1225 of self-awareness development for counselling and psychotherapy training,  
1226 *Counselling Psychology Quarterly* 26 (2) (2013) 190–207.
- 1227 [79] K. D. Killian, Development and validation of the emotional self-awareness  
1228 questionnaire: A measure of emotional intelligence, *Journal of Marital and*  
1229 *Family Therapy* 38 (3) (2012) 502–514.
- 1230 [80] M. Csikszentmihalyi, M. Csikzentmihaly, *Flow: The psychology of opti-*  
1231 *mal experience*, Vol. 1990, Harper & Row New York, 1990.
- 1232 [81] A. Lin, S. Gregor, M. Ewing, Developing a scale to measure the enjoyment  
1233 of web experiences, *Journal of Interactive Marketing* 22 (4) (2008) 40–57.
- 1234 [82] N. Abbas, T. Pickard, E. Atwell, A. Walker, University student surveys us-  
1235 ing chatbots: Artificial intelligence conversational agents, in: *International*  
1236 *Conference on Human-Computer Interaction*, Springer, 2021, pp. 155–169.
- 1237 [83] C. Ischen, T. Araujo, G. van Noort, H. Voorveld, E. Smit, “i am here to  
1238 assist you today”: The role of entity, interactivity and experiential percep-  
1239 tions in chatbot persuasion, *Journal of Broadcasting & Electronic Media*  
1240 64 (4) (2020) 615–639.
- 1241 [84] M. C. Han, The impact of anthropomorphism on consumers’ purchase de-  
1242 cision in chatbot commerce, *Journal of Internet Commerce* 20 (1) (2021)  
1243 46–65.
- 1244 [85] T. Rietz, I. Benke, A. Maedche, The impact of anthropomorphic and func-  
1245 tional chatbot design features in enterprise collaboration systems on user  
1246 acceptance, in: *14th International Conference on Wirtschaftsinformatik*,  
1247 2019, pp. 1642–1656.
- 1248 [86] R. C. Hanna, B. Weinberg, R. P. Dant, P. D. Berger, Do internet-based  
1249 surveys increase personal self-disclosure?, *Journal of Database Marketing*  
1250 *& Customer Strategy Management* 12 (2005) 342–356.
- 1251 [87] B. Jacquet, A. Hullin, J. Baratgin, F. Jamet, The impact of the gricean max-  
1252 ims of quality, quantity and manner in chatbots, in: *2019 international con-*  
1253 *ference on information and digital technologies (idt)*, IEEE, 2019, pp. 180–  
1254 189.

- 1255 [88] J. D. Smyth, D. A. Dillman, L. M. Christian, M. McBride, Open-ended  
1256 questions in web surveys: Can increasing the size of answer boxes and pro-  
1257 viding extra verbal instructions improve response quality?, *Public Opinion*  
1258 *Quarterly* 73 (2) (2009) 325–337.
- 1259 [89] R. Flesch, *Marks of readable style; a study in adult education.*, Teachers  
1260 College Contributions to Education (1943).
- 1261 [90] B. Pang, L. Lee, et al., Opinion mining and sentiment analysis, *Foundations*  
1262 *and Trends® in information retrieval* 2 (1–2) (2008) 1–135.
- 1263 [91] A. Barak, O. Gluck-Ofri, Degree and reciprocity of self-disclosure in online  
1264 forums, *CyberPsychology & Behavior* 10 (3) (2007) 407–417.
- 1265 [92] S. M. Jourard, *Self-disclosure, An experimental analysis of the transparent*  
1266 *self* (1971).
- 1267 [93] J. H. Kahn, R. M. Hessling, Measuring the tendency to conceal versus  
1268 disclose psychological distress, *Journal of Social and Clinical Psychology*  
1269 20 (1) (2001) 41–65.
- 1270 [94] L. C. Miller, J. H. Berg, R. L. Archer, Openers: Individuals who elicit in-  
1271 timate self-disclosure., *Journal of personality and social psychology* 44 (6)  
1272 (1983) 1234.
- 1273 [95] M. Nguyen, Y. S. Bin, A. Campbell, Comparing online and offline self-  
1274 disclosure: A systematic review, *Cyberpsychology, Behavior, and Social*  
1275 *Networking* 15 (2) (2012) 103–111.
- 1276 [96] A. N. Joinson, Self-disclosure in computer-mediated communication: The  
1277 role of self-awareness and visual anonymity, *European journal of social*  
1278 *psychology* 31 (2) (2001) 177–192.
- 1279 [97] R. Garrett, J. Chiu, L. Zhang, S. D. Young, A literature review: website  
1280 design and user engagement, *Online journal of communication and media*  
1281 *technologies* 6 (3) (2016) 1.
- 1282 [98] L. E. Kam, W. G. Chismar, Online self-disclosure: model for the use  
1283 of internet-based technologies in collecting sensitive health information,  
1284 *International journal of healthcare technology and management* 7 (3-4)  
1285 (2006) 218–232.

- 1286 [99] M. E. Hughes, L. J. Waite, L. C. Hawkey, J. T. Cacioppo, A short scale for  
1287 measuring loneliness in large surveys: Results from two population-based  
1288 studies, *Research on aging* 26 (6) (2004) 655–672.
- 1289 [100] P. M. Desmet, M. H. Vastenburger, N. Romero, Mood measurement with  
1290 pick-a-mood: review of current methods and design of a pictorial self-  
1291 report scale, *Journal of Design Research* 14 (3) (2016) 241–279.
- 1292 [101] Y. Rosseel, The lavaan tutorial, Department of Data Analysis: Ghent Uni-  
1293 versity (2014).
- 1294 [102] e. a. Hair, Joseph F., *Multivariate Data Analysis: A Global Perspective*. 7th  
1295 ed., Upper Saddle River: Prentice Hall, 2009.
- 1296 [103] B. Hilligoss, S. Y. Rieh, Developing a unifying framework of credibility  
1297 assessment: Construct, heuristics, and interaction in context, *Information*  
1298 *Processing & Management* 44 (4) (2008) 1467–1484.
- 1299 [104] A. Sutton, Measuring the effects of self-awareness: Construction of the  
1300 self-awareness outcomes questionnaire, *Europe’s journal of psychology*  
1301 12 (4) (2016) 645.
- 1302 [105] S. Pieritz, M. Khwaja, A. A. Faisal, A. Matic, Personalised recommenda-  
1303 tions in mental health apps: The impact of autonomy and data sharing, in:  
1304 *Proceedings of the 2021 CHI Conference on Human Factors in Computing*  
1305 *Systems*, 2021, pp. 1–12.
- 1306 [106] J. A. McCarty, L. J. Shrum, The measurement of personal values in survey  
1307 research: A test of alternative rating procedures, *Public Opinion Quarterly*  
1308 64 (3) (2000) 271–298.
- 1309 [107] H. P. Grice, *Logic and conversation*, in: *Speech acts*, Brill, 1975, pp. 41–58.
- 1310 [108] L. Dybkjær, N. O. Bernsen, H. Dybkjær, Grice incorporated: cooperativity  
1311 in spoken dialogue, in: *Proceedings of the 16th conference on Computa-*  
1312 *tional linguistics-Volume 1*, 1996, pp. 328–333.
- 1313 [109] G. N. Leech, *100 million words of english: the british national corpus*  
1314 *(bnc)*, Language Research (1992).
- 1315 [110] K. Hofland, S. Johansson, *Word frequencies in british and american en-*  
1316 *glish*, Norwegian computing centre for the Humanities, 1982.

- 1317 [111] E. N. Forsythand, C. H. Martell, Lexical and discourse analysis of online  
1318 chat dialog, in: International Conference on Semantic Computing (ICSC  
1319 2007), IEEE, 2007, pp. 19–26.
- 1320 [112] J. R. Landis, G. G. Koch, The measurement of observer agreement for cat-  
1321 egorical data, *biometrics* (1977) 159–174.
- 1322 [113] K. Kays, K. Gathercoal, W. Buhrow, Does survey format influence self-  
1323 disclosure on sensitive question items?, *Computers in Human Behavior*  
1324 28 (1) (2012) 251–256.
- 1325 [114] R. C. MacCallum, J. T. Austin, Applications of structural equation model-  
1326 ing in psychological research, *Annual review of psychology* 51 (1) (2000)  
1327 201–226.
- 1328 [115] T. A. Kyriazos, et al., Applied psychometrics: sample size and sample  
1329 power considerations in factor analysis (efa, cfa) and sem in general, *Psy-  
1330 chology* 9 (08) (2018) 2207.
- 1331 [116] J. Wang, X. Wang, Sample size for structural equation modeling, *Structural  
1332 equation modeling: Applications using Mplus* (2012) 391–428.
- 1333 [117] J. J. Hoogland, A. Boomsma, Robustness studies in covariance structure  
1334 modeling: An overview and a meta-analysis, *Sociological Methods & Re-  
1335 search* 26 (3) (1998) 329–367.
- 1336 [118] K. A. Bollen, R. H. Hoyle, Latent variables in structural equation modeling,  
1337 *Handbook of structural equation modeling* (2012) 56–67.
- 1338 [119] L.-t. Hu, P. M. Bentler, Cutoff criteria for fit indexes in covariance structure  
1339 analysis: Conventional criteria versus new alternatives, *Structural equation  
1340 modeling: a multidisciplinary journal* 6 (1) (1999) 1–55.
- 1341 [120] P. M. Bentler, D. G. Bonett, Significance tests and goodness of fit in the  
1342 analysis of covariance structures., *Psychological bulletin* 88 (3) (1980) 588.
- 1343 [121] E. C. Norton, B. E. Dowd, M. L. Maciejewski, Marginal ef-  
1344 fects—quantifying the effect of changes in risk factors in logistic regression  
1345 models, *Jama* 321 (13) (2019) 1304–1305.
- 1346 [122] V. Braun, V. Clarke, Using thematic analysis in psychology, *Qualitative  
1347 research in psychology* 3 (2) (2006) 77.

- 1348 [123] I. Nahum-Shani, S. N. Smith, B. J. Spring, L. M. Collins, K. Witkiewitz,  
1349 A. Tewari, S. A. Murphy, Just-in-time adaptive interventions (jitais) in mo-  
1350 bile health: key components and design principles for ongoing health be-  
1351 havior support, *Annals of Behavioral Medicine* 52 (6) (2018) 446–462.
- 1352 [124] J. P. Chung, W.-s. Yeung, Staff mental health self-assessment during the  
1353 covid-19 outbreak, *East Asian Archives of Psychiatry* 30 (1) (2020) 34.
- 1354 [125] C. R. Colvin, J. Block, D. C. Funder, Overly positive self-evaluations and  
1355 personality: negative implications for mental health., *Journal of personality  
1356 and social psychology* 68 (6) (1995) 1152.
- 1357 [126] F. Heylighen, J.-M. Dewaele, Formality of language: definition, measure-  
1358 ment and behavioral determinants, *Interner Bericht, Center “Leo Apostel”,  
1359 Vrije Universiteit Brussel* 4 (1999).
- 1360 [127] K. M. Rennekamp, P. Witz, Linguistic formality and perceived  
1361 engagement–investors’ reactions to two unique characteristics of social me-  
1362 dia disclosures, *SSRN* (2017).
- 1363 [128] M. A. Hamilton, Message variables that mediate and moderate the effect of  
1364 equivocal language on source credibility, *Journal of Language and Social  
1365 Psychology* 17 (1) (1998) 109–143.
- 1366 [129] P. Thomas, M. Czerwinski, D. McDuff, N. Craswell, G. Mark, Style and  
1367 alignment in information-seeking conversation, in: *Proceedings of the 2018  
1368 Conference on Human Information Interaction & Retrieval, 2018*, pp. 42–  
1369 51.
- 1370 [130] G. Eysenbach, J. Powell, O. Kuss, E.-R. Sa, Empirical studies assessing  
1371 the quality of health information for consumers on the world wide web: a  
1372 systematic review, *Jama* 287 (20) (2002) 2691–2700.
- 1373 [131] K. Kroenke, R. L. Spitzer, J. B. Williams, B. Löwe, The patient health  
1374 questionnaire somatic, anxiety, and depressive symptom scales: a system-  
1375 atic review, *General hospital psychiatry* 32 (4) (2010) 345–359.
- 1376 [132] D.-H. Shin, S. Lee, Y. Hwang, How do credibility and utility play in the user  
1377 experience of health informatics services?, *Computers in Human Behavior*  
1378 67 (2017) 292–302.

- 1379 [133] R. Lederman, H. Fan, S. Smith, S. Chang, Who can you trust? credibility  
1380 assessment in online health forums, *Health Policy and Technology* 3 (1)  
1381 (2014) 13–25.
- 1382 [134] H. Gaffney, W. Mansell, S. Tai, et al., Conversational agents in the treatment  
1383 of mental health problems: mixed-method systematic review, *JMIR mental*  
1384 *health* 6 (10) (2019) e14166.
- 1385 [135] I. Altman, D. A. Taylor, *Social penetration: The development of interper-*  
1386 *sonal relationships.*, Holt, Rinehart & Winston, 1973.
- 1387 [136] L. R. Wheeless, J. Grotz, The measurement of trust and its relationship to  
1388 self-disclosure, *Human Communication Research* 3 (3) (1977) 250–257.
- 1389 [137] P. C. Cozby, Self-disclosure: a literature review., *Psychological bulletin*  
1390 79 (2) (1973) 73.
- 1391 [138] S. Devaram, Empathic chatbot: Emotional intelligence for empathic chat-  
1392 bot: Emotional intelligence for mental health well-being, arXiv preprint  
1393 arXiv:2012.09130 (2020).
- 1394 [139] G. M. Lucas, J. Gratch, A. King, L.-P. Morency, It’s only a computer: Vir-  
1395 tual humans increase willingness to disclose, *Computers in Human Behav-*  
1396 *ior* 37 (2014) 94–100.
- 1397 [140] A. J. Bush, A. Parasuraman, Assessing response quality. a self-disclosure  
1398 approach to assessing response quality in mall intercept and telephone in-  
1399 terviews, *Psychology & Marketing* 1 (3-4) (1984) 57–71.
- 1400 [141] M. Barreto, C. Victor, C. Hammond, A. Eccles, M. T. Richins, P. Qual-  
1401 ter, Loneliness around the world: Age, gender, and cultural differences in  
1402 loneliness, *Personality and Individual Differences* 169 (2021) 110066.
- 1403 [142] W. D. Killgore, S. A. Cloonan, E. C. Taylor, N. S. Dailey, Loneliness: A  
1404 signature mental health concern in the era of covid-19, *Psychiatry research*  
1405 290 (2020) 113117.
- 1406 [143] T. Richardson, P. Elliott, R. Roberts, Relationship between loneliness and  
1407 mental health in students, *Journal of Public Mental Health* (2017).



- 1408 [144] Y. Al-Saggaf, S. Nielsen, Self-disclosure on facebook among female users  
 1409 and its relationship to feelings of loneliness, *Computers in Human Behavior*  
 1410 36 (2014) 460–468.
- 1411 [145] J. H. Kahn, A. M. Garrison, Emotional self-disclosure and emotional avoid-  
 1412 ance: Relations with symptoms of depression and anxiety., *Journal of coun-  
 1413 seling psychology* 56 (4) (2009) 573.
- 1414 [146] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal,  
 1415 A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., Language models  
 1416 are few-shot learners, *Advances in neural information processing systems*  
 1417 33 (2020) 1877–1901.

1418 **Appendix A. Open-ended Questions**

Table A.4: The Open-Ended Questions Asked During the Interview Session

ID	Question
Open-EQ1	<i>In general, how would you describe your current mood?</i>
Open-EQ2	<i>What do you think of the influence of the COVID-19 pandemic on your study and life?</i>
Open-EQ3	<i>Can you tell me a little bit about any contact you have had with friends or family recently?</i>
Open-EQ4	<i>What have you tried to manage isolation and loneliness during COVID-19?</i>
Open-EQ5	<i>What do you think could be the main factors contributing to loneliness?</i>
Open-EQ6	<i>What would it take for you to feel happier or more at peace?</i>
Open-EQ7	<i>Think of something that you feel happy and grateful for, great or small (e.g., the food you eat or the place you live in).</i>

1419 **Appendix B. Post-study Questions**

Table B.5: The Questions Asked in the Post-Study

ID	Question
Post-Q1	<i>What do you think of answering the questions to know your loneliness score?</i>
Post-Q2	<i>What do you think of knowing your mental status by chatting with such a bot?</i>
Post-Q3	<i>What do you think of answering the questions in conversation with the Percy bot instead of filling in an online form?</i>
Post-Q4	<i>What do you think of describing your feelings through talking with the Percy bot?</i>
Post-Q5	<i>What questions that the Percy bot asked may make you feel concerned about?</i>

1420 **Appendix C. Descriptive Statistics of Dependent Variables**

Table C.6: Descriptive Statistics of Dependent Variables

<b>Dependent variable</b>	<b>Form*Short</b> (N=34) Mean (SD)	<b>Form*Middle</b> (N=36) Mean (SD)	<b>Form*Long</b> (N=35) Mean (SD)	<b>Con*Short</b> (N=33) Mean (SD)	<b>Con*Middle</b> (N=35) Mean (SD)	<b>Con*Long</b> (N=40) Mean (SD)
<b>Subjective Experiences</b>						
Assessment Credibility	3.69 (0.66)	<b>3.89</b> (0.56)	3.74 (0.67)	3.34 (0.95)	3.53 (0.84)	3.62 (0.77)
Self-Awareness	3.58 (0.66)	3.75 (0.65)	<b>3.78</b> (0.66)	3.50 (0.72)	3.67 (0.82)	3.57 (0.67)
Enjoyment	3.83 (0.68)	<b>3.91</b> (0.80)	3.78 (0.71)	3.40 (0.73)	3.84 (0.95)	3.83 (0.69)
<b>Response Quality</b>						
Informativeness	<b>671.0</b> (458.3)	625.8 (414.4)	618.6(406.3)	547.5 (329.4)	519.3 (359.0)	644.2 (554.5)
Specificity	1.08 (0.47)	<b>1.10</b> (0.40)	1.08 (0.41)	0.97 (0.38)	0.94 (0.40)	1.01 (0.40)
Relevance	1.85 (0.21)	<b>1.90</b> (0.17)	1.87 (0.18)	1.86 (0.27)	1.82 (0.21)	1.87 (0.23)
Clarity	1.53 (0.32)	1.58 (0.29)	<b>1.60</b> (0.27)	1.54 (0.28)	1.47 (0.30)	1.55 (0.33)
RQI	3.41 (2.19)	<b>3.56</b> (1.86)	3.46 (1.88)	3.08 (1.68)	2.80 (1.84)	3.27 (1.87)
<b>Self-Disclosure</b>						
Self-Disclosure (sub.)	4.18 (0.68)	3.82 (0.75)	4.00 (0.66)	3.92 (0.68)	<b>4.24</b> (0.69)	4.16 (0.57)
Self-Disclosure (obj.)	<b>0.96</b> (0.53)	0.94 (0.47)	0.95 (0.50)	0.88 (0.42)	0.78 (0.51)	0.90 (0.46)
Response Length	<b>60.9</b> (44.6)	56.2 (39.5)	55.8 (39.4)	48.5 (31.8)	45.8 (34.6)	57.6 (51.5)
Engagement Duration	267.4 (119.4)	<b>339.8</b> (277.7)	278.9 (162.8)	261.9 (104.5)	291.0 (184.1)	333.6 (226.0)

Note: 1. RQI is calculated based on specificity, relevance, and clarity by using Formula (2). 2. The highest value of each dependent variable is marked in bold.