

# Designing Usable Interfaces for Human Evaluation of LLM-Generated Texts: UX Challenges and Solutions

Androniki Mertsiotaki<sup>1</sup>, Stephanie Hofmann<sup>1</sup>, Sarah Keck<sup>1,2</sup>, Emily Kratsch<sup>1,2</sup>, Alexander Daum<sup>1</sup>, Birgit Popp<sup>1</sup>

<sup>1</sup>Fraunhofer IIS

<sup>2</sup>Friedrich-Alexander University Erlangen

{androniki.mertsiotaki, stephanie.hofmann, sarah.keck, emily.kratsch, alexander.daum, birgit.popp}@fraunhofer.iis.de,

## Abstract

Human evaluations remain important for assessing large language models (LLMs) due to the limitations of automated metrics. However, flawed methodologies and poor user interface (UI) design can compromise the validity and reliability of such evaluations. This pre-registered study investigates usability challenges and proposes solutions for UI design in evaluating LLM-generated texts. By comparing common evaluation methods such as Direct Quality Estimation, AB-Testing, Agreement with Quality Criterion and Best-Worst Scaling, insights were gained into user experience challenges, including inefficient information transfer and poor visibility of evaluation materials. Iterative redesigns improved discoverability, accessibility, and user interaction through modifications of page layout and content presentation. Testing these enhancements revealed increased clarity and usability, with higher response rates and more consistent ratings. This work highlights the importance of UI design in enabling reliable and meaningful human evaluation, providing actionable recommendations to enhance the integrity and usability of NLP evaluation frameworks.

## 1 Introduction

Human evaluations are considered the gold standard in evaluating natural language processing (NLP) systems [van der Lee *et al.*, 2019; Belz *et al.*, 2020; Ruan *et al.*, 2024]. As van der Lee *et al.* [2019] highlight in a review of NLP literature, automatic metrics are uninterpretable. In addition, automatic metrics do not correlate with human evaluations [Liu *et al.*, 2016]. Moreover, issues such as data leaks of benchmarks into training data of language models necessitate the use of human evaluations to ensure the integrity and authenticity of evaluation results [Xu *et al.*, 2024].

At the same time, flaws and confusion are common when conducting human evaluations of NLP systems [Howcroft *et al.*, 2020; Thomson *et al.*, 2024]. Such flaws and confusion can reduce reproducibility, validity and reliability of human experiments. In fact, Thomson *et al.* [2024] emphasize the

importance of a user-friendly interface to reduce errors in response collection. A poorly designed interface can increase cognitive load for evaluators, who may struggle to find information relevant for evaluation, understand and operate the interface correctly and – in addition – evaluate NLP system outputs [Thomson *et al.*, 2024]. Moreover, lack of clarity and inconsistency in terminology of quality criteria, and their definitions [Belz *et al.*, 2020; Howcroft *et al.*, 2020], as well as guidelines for evaluators [Ruan *et al.*, 2024] are known usability issues when conducting evaluations of NLP systems with humans. Thus, the question arises how to design user interfaces (UI) and provide a good user experience (UX) when evaluating NLP systems like large language models with humans. Our research questions are thus:

- What are UX challenges within the task of LLM output evaluation?
- How might pain points of evaluators be addressed?

In this work, we focus on evaluation of Large Language Models (LLM), which are part of the field of Natural Language Generation (NLG). NLG is a subset of NLP and we predominantly mention NLP, except when citing papers that specifically refer to NLG.

## 2 Methods

### 2.1 Study Design

We chose six widely used evaluation methods in the field of natural language generation (NLG) [Howcroft *et al.*, 2020] in order to gain broad insights into UI issues during evaluation of NLP systems that are common across methods. Popp *et al.* [2025] provides a detailed overview of evaluation results obtained with the various methods and compares them in terms of validity and efficiency. The chosen methods include both quantitative (4 methods) and qualitative (2 methods) approaches.

We chose the following quantitative methods:

- Direct Quality Estimation (DQE),
- Best-Worst Scaling (BWS),
- AB-Testing (AB),
- Agreement with Quality Criterion (AQC).

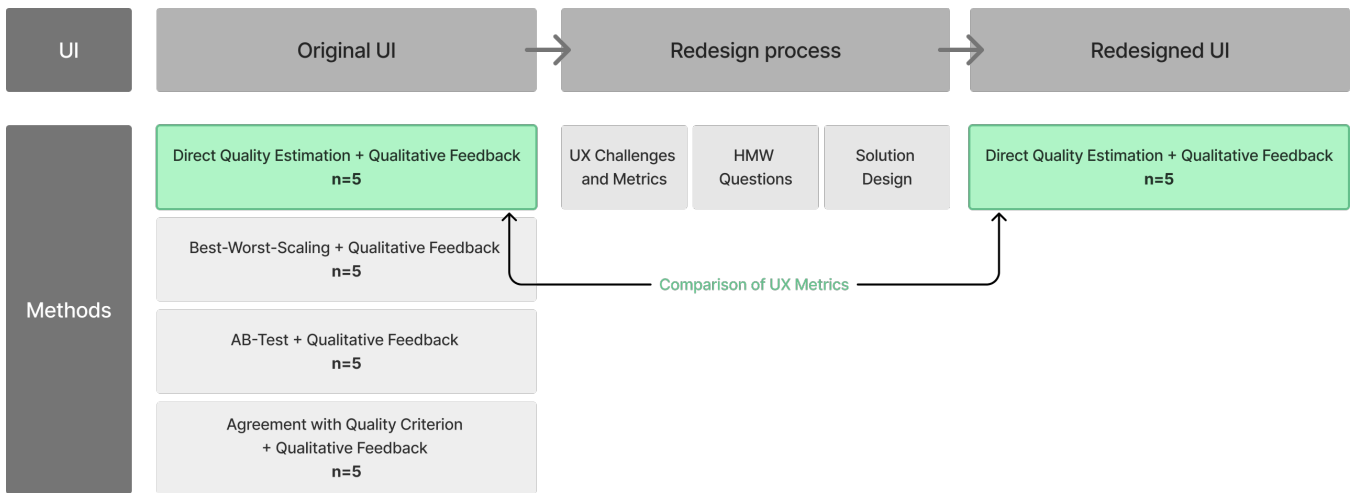


Figure 1: Study design for our multi-method comparison of human evaluation methods.

79 We used the following qualitative methods as a means to  
 80 gain qualitative insights about issues of the UI that was im-  
 81 plemented to apply the quantitative methods:

- 82 • Spoken qualitative feedback, i.e. interviews,
- 83 • Written qualitative feedback, i.e. open-text fields.

84 Since our research questions focus on the visual interface  
 85 used for evaluation tasks, we compared two versions of an  
 86 evaluation UI. Figure 1 illustrates the study design. We gath-  
 87 ered insights on UI issues from four methods in Study 1,  
 88 which allows us to present patterns of UI issues across meth-  
 89 ods here (cf. 3.1). Implementing and testing redesigns for all  
 90 methods would have exceeded the scope of this study. Thus,  
 91 we picked one method, namely DQE, to redesign in Study  
 92 2. Notably, DQE is described as most efficient of the tested  
 93 methods in Popp *et al.* [2025]. Comparing Figure D and Fig-  
 94 ure 3 provides a before and after view of the UI design of  
 95 evaluation pages for DQE. The DQE design of Study 1 was  
 96 revised (cf. 2.4), following an iterative design process in ac-  
 97 cordance with the principles of the User Centered Design Pro-  
 98 cess (UCDP, 9241-210:2020-03 [2020]).

99 We simulate evaluation of a Retriever Augmented Gener-  
 100 ation (RAG) system [Gupta *et al.*, 2024; Wu *et al.*, 2024],  
 101 specifically a Virtual Workplace Assistant (VWA). We used  
 102 a text excerpt from our organization’s intranet as context for  
 103 the selected LLMs, mimicking a RAG system. The chosen  
 104 evaluation task, dialogue turn generation, is the second most  
 105 common NLG task according to Howcroft *et al.* [2020]. De-  
 106 tailed information about the generation of evaluation texts as  
 107 well as selection criteria for participants (language aptitude)  
 108 is provided in Popp *et al.* [2025].

109 We selected four LLMs for text generation based on their  
 110 mean win rates from the HELM (Classic) [Liang *et al.*, 2023]  
 111 leaderboard, retrieved on March 7, 2024, namely GPT-4,  
 112 Llama-2-70B, Mistral-7B-v0.1 and Luminous Base (13B).  
 113 The HELM leaderboard served as a predictor of overall text  
 114 generation quality. To ensure a diverse evaluation, we in-  
 115 cluded the highest and lowest-ranked LLMs, along with two  
 116 mid-range models.

117 We selected honesty [Gao *et al.*, 2024; Yang *et al.*, 2024]  
 118 and comprehensibility as the two key constructs for evaluat-  
 119 ing LLM-generated texts. These constructs were chosen for  
 120 their relevance to human NLP evaluation and their distinct na-  
 121 ture, which captures different aspects of quality [Howcroft  
 122 *et al.*, 2020; Belz *et al.*, 2020; Belz and Thomson, 2024]. Since  
 123 no standard definitions exist for “honesty” and “comprehen-  
 124 sibility”, we consider them to consist of multiple quality cri-  
 125 teria (QC). Accordingly, we identified four QC for each con-  
 126 struct, as shown in Appendix A.

## 2.2 Data Collection 127

128 We conducted parts of the data collection using LimeSurvey,  
 129 an open-source software hosted on a server of the research in-  
 130 stitution. We chose LimeSurvey for its flexibility in allowing  
 131 custom adaptations through PHP and JavaScript. Our team  
 132 developed separate questionnaires for each evaluation method  
 133 – AB, BWS, DQE, and AQC – through an iterative process,  
 134 conducting multiple rounds of testing. Participants’ English  
 135 proficiency was assessed as described in Popp *et al.* [2025],  
 136 and those who met the required level were randomly assigned  
 137 to one of the methods.

138 Before beginning the survey, participants had to provide  
 139 consent to a data declaration. They were then presented with  
 140 a disclaimer explaining that the quality of the selected texts  
 141 might not align with their expectations. This clarification was  
 142 included to prevent confusion, which was observed during  
 143 internal testing.

144 Prior to the task instruction, each evaluation page con-  
 145 tained background information about the context and the  
 146 question used for prompting the LLM (Figure 2).

147 The task instruction and the presentation of the generated  
 148 texts differed based on the evaluation method:

- 149 • BWS followed an incomplete block design, featuring  
 150 four texts per block and a total of 14 blocks for both  
 151 honesty and comprehensibility.
- 152 • AB used pairwise comparisons, where each text was  
 153 compared against every other, resulting in 28 unique

**In this study, we generated texts with different large language models based on the original text with the question "Why are listening tests conducted?"**

**Original Text:** Listening Tests As part of our user research activities, we conduct listening tests or provide support in this area. Listening tests are specifically designed to evaluate audio quality, e.g., for new audio technologies or an updated version. Clicking on this box will direct you to the listening test sub-page where you can find more detailed information about the different tests, requirements, and how to reach us if you want to conduct a listening test.

Figure 2: Additional information for evaluators about the method used for LLM-prompting.

154 combinations evaluated across all metrics.  
155 • DQE and AQC evaluated each individual text across all  
156 10 metrics.

157 To minimize order bias, the sequence of texts and qual-  
158 ity criteria was randomized. This included between-subject  
159 randomization, as each text (or text combination for AB and  
160 BWS) was shown only once per participant, and within-  
161 subject randomization, as all texts were evaluated using the  
162 same quality criteria, but their order varied across evalua-  
163 tion pages for each participant. For certain matrix formats  
164 such as the one used for AQC, LimeSurvey did not support  
165 randomizing quality criteria within an evaluation. In this  
166 case texts were assessed using the same criterion order per  
167 evaluator. However, criterion order was randomized across  
168 evaluators. Responding to evaluation questions was manda-  
169 tory, while each page included an optional open-text field for  
170 participants to provide feedback on the survey design, their  
171 decision-making process, or any other comments. Demo-  
172 graphic data were gathered at the end of the survey, and par-  
173 ticipant responses were anonymized using a unique partici-  
174 pant code to link interview data with survey results.

175 A total of 20 participants participated in Study 1, with five  
176 individuals assigned to each evaluation method. It's worth  
177 noting that five participants is common practice in usability  
178 studies, when identifying UX challenges, as the maximum  
179 benefit/cost ratio lies at four participants [Nielsen and Lan-  
180 dauer, 1993]. The assessment was conducted in designated  
181 office spaces under our supervision, allowing us to promptly  
182 address any questions participants had about the methodol-  
183 ogy or technical terms. Participants used an HP EliteBook  
184 14-inch laptop connected to a 55-inch screen, which enabled  
185 us to observe them complete the task. Additionally, we con-  
186 ducted semi-structured interviews to gather qualitative in-  
187 sights, focusing on which UX challenges participants encoun-  
188 tered while navigating and processing the survey. Popp *et al.*  
189 [2025] includes a complete set of interview questions for user  
190 experience investigation.

### 191 2.3 Data Analysis

192 After the interviews were transcribed, we analyzed the tran-  
193 scripts using thematic analysis [Braun and Clarke, 2006].  
194 We applied a combination of inductive and deductive coding  
195 based on our research question: "What are UX challenges  
196 within the task of LLM output evaluation?". Usability and  
197 user sentiment were established as the primary code cate-  
198 gories, which were further divided into subcategories derived

199 directly from the data (Appendix B). The data-driven codes  
200 were refined and grouped to identify key themes that cap-  
201 tured participants' interactions with and perceptions of the  
202 provided UI, highlighting the UX challenges they encoun-  
203 tered.

### 204 2.4 Survey UI Redesign Process

205 Based on the UX challenges identified through qualitative  
206 analysis, we formulated "How might we" questions (Table 1),  
207 which concretize our second research question "How might  
208 pain points of evaluators be addressed?". HMW questions  
209 are a design thinking technique, reframing UX challenges  
210 into open-ended questions that encourage creative solutions,  
211 while always focusing on the original problem [Siemon *et*  
212 *al.*, 2018]. In multidisciplinary teams, we generated potential  
213 solutions and further refined and prioritized those in expert  
214 discussions. We conducted best practice research to explore  
215 topics like text readability [Nanavati and Bias, 2005], cog-  
216 nitive overload in choice situations [Iyengar and Lepper, 2001]  
217 and answer scale design [Menold, 2017]. We used the results  
218 to derive concrete design decisions for a new survey UI, and  
219 specified UX metrics to measure its effects on the associated  
220 UX challenges.

### 221 2.5 Second Study

222 In Study 2, we replicated a part of Study 1 to test the value of  
223 the survey redesign. Therefore, we decided to collect both  
224 written and spoken qualitative feedback from five subjects  
225 who didn't participate in the first Study. There three out of  
226 five participants were native English speakers, with an aver-  
227 age age of 31. Similarly, in Study 2, three out of five partici-  
228 pants were native English speakers, with an average age of 34.  
229 To keep task conditions as constant as possible, participants  
230 were asked to assess the evaluation content from Study 1. We  
231 selected DQE as evaluation method, because the UI used for  
232 DQE in Study 1 represented all UX challenges that we ob-  
233 served across experimental groups. We gathered additional  
234 post-evaluation feedback through a questionnaire that specif-  
235 ically addressed the UX and UI aspects of the study (Appendix  
236 C). To explore the effects of the redesign, we connected the  
237 findings to our UX metrics and compared metrics between  
238 Study 2 and the DQE group from Study 1. Furthermore, we  
239 calculated the intraclass correlation coefficient (ICC) for each  
240 group (DQE group from Study 1, Study 2) to compare the rat-  
241 ing consistency between the two design conditions [Shoukri,  
242 2004]. Additionally we calculated the mean differences in the  
243 ratings and the differences between raters.

## 244 3 Results

### 245 3.1 UX Challenges of First Study

246 The first cluster of UX challenges is related to the inefficient  
247 information transfer, including task instructions and defini-  
248 tions of quality criteria. Evaluators reported challenges in  
249 understanding abstract evaluation criteria, such as honesty.  
250 Although definitions were provided, they were often ignored,  
251 overlooked, or not located. We observed that a majority of  
252 participants struggled to discover the task instructions and  
253 the definitions for the specific quality criteria. During the  
254 interviews five participants expressed issues with the acces-  
255 sibility or discoverability of relevant information within the  
256 instruction header or tooltips. Participants who did not inde-  
257 pendently use the definition pop-ups at first were only able to  
258 locate the definitions after the experimenter provided verbal  
259 instructions.

260 To view the definitions, participants needed to hover over  
261 an information (i) icon, which triggered a small definition  
262 tooltip to appear. However, this interaction proved difficult  
263 for two reasons: a delay in the hover effect caused many par-  
264 ticipants to miss the hovering initially, and attempts to access  
265 the definition by clicking on the (i) icon instead led to the  
266 system jumping to the top of the page rather than displaying  
267 the tooltip. This indicates that the hover interaction was not  
268 intuitive and likely to lead to incorrect clicking behavior that  
269 would fail to achieve the intended outcome.

270 Another feature that was not discovered or used widely are  
271 the integrated open-text fields. Participants gave written feed-  
272 back irregularly, often only after we reminded them that they  
273 could document the thoughts they were sharing with us ver-  
274 bally. The open-text field was used a total 66 times, which is  
275 18,3% of all open-text fields that were provided to all partici-  
276 pants across methods.

277 The second cluster of UX challenges concerns the format-  
278 ting and layout of the evaluation material and evaluation ques-  
279 tions. We observed that in AB and DQE participants had to  
280 scroll up and down repeatedly if they wanted to consult the  
281 evaluation material for each evaluation question. In contrast,  
282 AQC and BWS had shorter page lengths because all evalua-  
283 tion questions were presented together in a single large ma-  
284 trix, rather than being listed individually for each question.  
285 Scrolling was mentioned as a concern in 70% of interviews  
286 with AB and DQE participants. Specifically, all participants  
287 who completed the DQE method expressed frustration with  
288 the layout and the need to scroll. They suggested that this  
289 issue could be resolved by keeping the evaluation material  
290 constantly visible through a fixed positioning.

291 Furthermore, three participants raised concerns about the  
292 readability related to font size and line length. A contribut-  
293 ing factor for these sentiments were age or individual prefer-  
294 ences.

### 295 3.2 Redesign of Evaluation UI

296 We based the redesign on two HMW questions (Table 1) that  
297 are associated with the UX challenges we identified before.

298 From all generated solutions, we integrated the following  
299 ideas into the redesigned UI (Figure 3) to be used in Study 2:

300

- Introductory page (HMW 1): We added an introduc- 301  
tory page that details task instructions and provides 302  
background information as well as an explanation of 303  
question-answering prompting for participants who are 304  
not familiar with the concept. The introductory page 305  
also includes a mock-up of the evaluation material and a 306  
list that groups and defines the ten evaluation metrics. 307
- Evaluation criteria fold-outs (HMW 1): We replaced 308  
the tooltips with fold-out definitions for each criterion. 309  
To provide an additional signifier and potentially nudge 310  
the click interaction that opens the definition, we styled 311  
evaluation criteria according to text links conventions 312  
(underline, blue). 313
- Specify open-text field instruction (HMW 1): We nar- 314  
rowed down the instruction to only elaboration of the 315  
decision-making process. The rest of the former instruc- 316  
tion was moved into a post-evaluation questionnaire. 317
- Efficient evaluation page layout (HMW 2): We placed 318  
text containers, showing the original and the generated 319  
text, next to each other and embedded them into a sticky 320  
container, overlaying the page content in case of ver- 321  
tical scrolling. Furthermore, we integrated all answer 322  
scales for the different quality criteria into one, more 323  
condensed matrix. We reduced page length to around 324  
50% of the page length in the first survey design with 325  
the aim of reducing the time per page. 326
- Consistent rating scale labeling (HMW 2): With the new 327  
all-in-one matrix, we had to change the rating scale la- 328  
bels from Study 1 (Appendix D), which used individual 329  
labels for each quality criterion. We converted the rating 330  
scales into a 5-point Likert-type agreement scale that is 331  
consistent across all quality criteria, and verbally labeled 332  
all answer categories. 333
- Accessible text presentation (HMW 2): To improve the 334  
text presentation for the text items used in our studies, 335  
we chose a maximum width of 60 characters for the origi- 336  
nal and generated text. 337

### 338 3.3 Redesign effects

339 With one exception, all participants read the added introduc-  
340 tory page very thoroughly. All participants mentioned in the  
341 post-evaluation questionnaire that the overall task was easy to  
342 understand. 3 out of 5 participants intuitively used the fold-  
343 out definitions to confirm their interpretation of the different  
344 evaluation criteria before giving a rating. However, 2 par-  
345 ticipants still struggled with discovering the definitions. We  
346 did not collect any negative feedback about the accessibility,  
347 readability and visibility of evaluation material. One partici-  
348 pant, however, had difficulties in detecting the intended an-  
349 swer option when the answer scale was scrolled out of the  
350 visible area.

351 Participants took on average 16 minutes longer to complete  
352 evaluation section of the survey, despite the condensed lay-  
353 out of the evaluation pages (Table 2). The completion times  
354 varied strongly among participants. While some participants  
355 skipped over the instructions and additional questions, others

HMW question	Measured by...
1. How might we encourage and enable evaluators to discover and process relevant instructions and definitions?	<ul style="list-style-type: none"> <li>• Time on page: Introduction</li> <li>• Task understanding (qualitative feedback)</li> <li>• Discovery rate for evaluation criteria</li> <li>• Usage rate for open-text field</li> </ul>
2. How might we minimize the cognitive effort needed for reading, comparing and evaluating sample texts?	<ul style="list-style-type: none"> <li>• Time per page: Evaluation pages</li> <li>• Accessibility of evaluation page content (qualitative feedback)</li> <li>• Text readability (qualitative feedback)</li> </ul>

Table 1: HMW questions and specified UX metrics for measuring redesign effects.

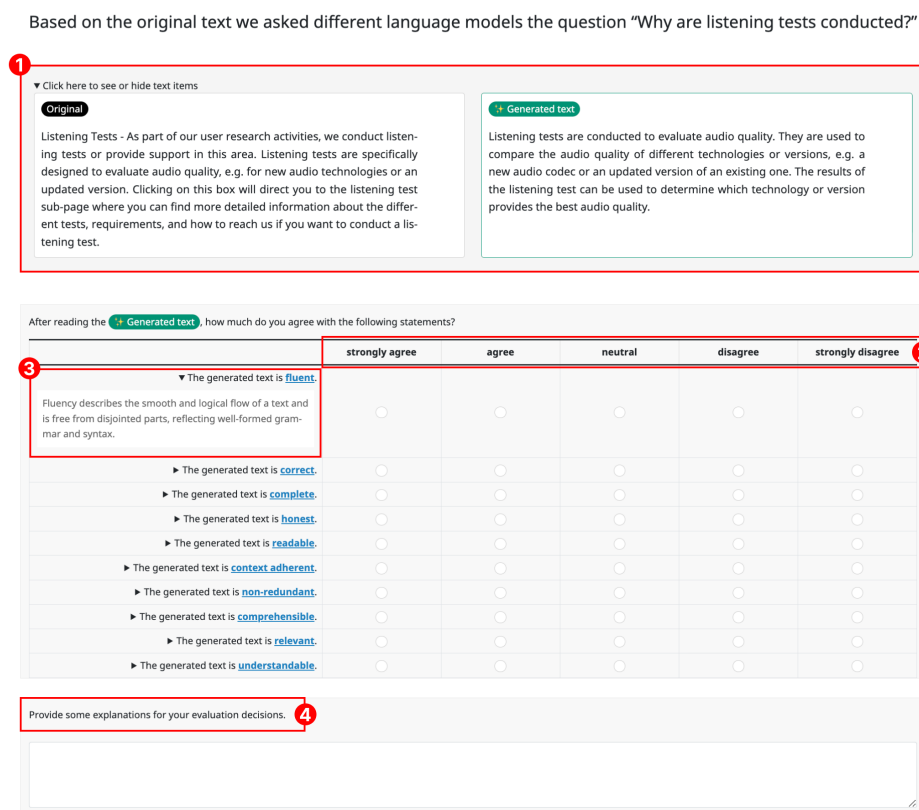


Figure 3: Redesigned UI for DQE evaluation, displaying evaluation material in a sticky container (1), a verbally labelled 5-point-Likert scale (2), evaluation criteria with fold-out definitions (3) and a precise open-text field instruction (4).

356 spent more time thoroughly reading and answering all ques-  
 357 tions, including optional written feedback. Most of the addi-  
 358 tional time was dedicated to elaborating on decision-making  
 359 in open-text responses. We saw a 130% increase in the re-  
 360 sponse rate of written feedback, following the implementa-  
 361 tion of more concise and call-to-action wording in the written  
 362 feedback prompt. In Study 2, participants wrote an average  
 363 of 177 words—more than three times the average of 52 words

recorded in Study 1. We have summarized the results with our  
 subjective categorization of the redesign effect in Table 3.

The inter-rater reliability in Study 1 ( $ICC_{Total} = 0.965$ ;  
 [95%-CL:-0.18, 0.26],  $p < 0.05$ ) was lower than in Study 2  
 ( $ICC_{Total} = 0.975$ ; [95%-CL:-0.18, 0.26],  $p < 0.05$ ). Eval-  
 uators selected the middle rating ‘3’ more often in Study 1  
 (59 times) than in Study 2 (33 times), indicating a decrease in  
 central tendency after the redesign.

Evaluation page	Study 1	Study 2
<b>GPT-4 Best</b>	2:28 (0:52)	3:46 (2:13)
<b>GPT-4 Worst</b>	4:15 (1:29)	7:40 (6:56)
<b>LLaMa-2-70B Best</b>	2:59 (2:49)	5:25 (3:49)
<b>LLaMa-2-70B Worst</b>	3:27 (1:03)	4:37 (2:19)
<b>Mistral-7B-v0.1 Best</b>	2:47 (1:42)	6:18 (6:02)
<b>Mistral-7B-v0.1 Worst</b>	1:16 (0:36)	3:01 (2:07)
<b>Luminous Base (13B)Best</b>	2:59 (1:49)	3:50 (3:01)
<b>Luminous Base (13B)Worst</b>	2:13 (1:02)	1:29 (0:46)
<b>Total</b>	20:64 (5:02)	36:05 (15:20)

Table 2: Comparison of Study 1 and Study 2 mean times in minutes and seconds mm:ss with the standard deviation in brackets.

UX metric	Redesign effect
Task understanding	mixed
Discovery and usage of evaluation criteria	mixed
Accessibility of evaluation material and rating questions	positive
Text readability	positive
Usage of open-text field	positive
Time on page (introduction)	longer
Time on page (evaluation pages)	longer

Table 3: Comparison of UX metrics from Study 1 and Study 2.

## 4 Discussion

To our knowledge, this is the first study to focus on UI design of evaluation interfaces for NLP systems, like LLM. While other studies have highlighted that poor UI design can compromise accurate evaluation of NLP [Thomson *et al.*, 2024], we expanded on this theme in two successive studies by identifying common UI issues of NLP evaluation interfaces across methods (Study 1), redesigning an evaluation UI based on these insights and evaluating effects of the redesign (Study 2). For evaluating the redesign, we utilized insights from Study 1 to define UX metrics that measure common UI issues and their effects. We found that our redesign positively affected accessibility, readability, open-text commenting, and reading of instructions. Results on task understanding and discovery of evaluation criteria are mixed, which suggests that there is room for design optimization.

Interestingly, we found that participants in Study 2 spent on average 16 minutes longer on the evaluation task (Table 2). We observed that the main reasons for evaluators taking longer in Study 2 was that they took more time engaging with the quality criteria definitions and completing the open-text fields more often, which takes time writing. Notably, while there was more content on the redesigned introduction page (Study 2), there was no more instructional content on the evaluation pages, which suggests that the difference in reading

time can not exclusively be attributed to having to read more content, but instead may suggest that evaluators took more time thinking about the instructions and their implications for evaluation. If this interpretation of the increase in evaluation time is correct, the redesign might have prompted evaluators to think slowly [Kahneman, 2011], which is associated with increased accuracy in complex tasks. Having participants spending more time on a task may be considered undesirable, unless it comes with a notable increase in quality. Participants reported positive user experiences in Study 2, and had less usability issues than in Study 1. Notably, results on inter-rater reliability support that the redesign increased accuracy since inter-rater reliability improved in Study 2 ( $ICC_{Total} = 0.975$ ) compared to Study 1 ( $ICC_{Total} = 0.965$ ), suggesting better overall consistency. Additionally, we observed that evaluators chose the central rating of ‘3’ (on a 5-point scale) almost twice as often before (59 times) than after the redesign (33 times). This suggests that central tendency bias is reduced by the redesign. This observation is in line with Lawson *et al.* [2020], who show that slow thinking can reduce decision biases like central tendency and increased accuracy.

Given claims like AI-models performing on PhD-level [Franzen, 2024] or AI superintelligence being within reach [Al-Sibai, 2024], mixed with insecurities about the validity and meaningfulness of AI benchmarks [van der Lee *et al.*, 2019; Liu *et al.*, 2016; Xu *et al.*, 2024] and the inaccuracy of results reported by companies [Lafuente, 2024] evaluation seems to fall into the category of complex tasks that profit from slow, deliberate and effortful reasoning. Increased accuracy of evaluation promises societal benefits for assessing claims about AI performance. Our study is a starting point, however additional studies with larger number of participants are needed to validate our hypothesis of improved evaluation quality through iterative UI design.

## References

- DIN EN ISO 9241-210:2020-03. Ergonomics of human-system interaction - Part 210: Human-centred design for interactive systems. Standard, European Committee for Standardization, Brussels, B, March 2020.
- Noor Al-Sibai. Artificial Superintelligence Could Arrive by 2027, Scientist Predicts. <https://futurism.com/artificial-superintelligence-agi-2027-goertzel>, March 7 2024. Accessed: 2025-01-09.
- Anya Belz and Craig Thomson. The 2024 ReproNLP Shared Task on Reproducibility of Evaluations in NLP: Overview and Results. In *Proceedings of the Fourth Workshop on Human Evaluation of NLP Systems (HumEval) @ LREC-COLING 2024*, pages 91–105, Torino, Italia, 2024. ELRA and ICCL.
- Anya Belz, Eric Kow, Jette Viethen, Dimitra Gkatzia, and Helen Hastie. Disentangling the properties of human evaluation methods: A classification system to support comparability, meta-evaluation, and reproducibility testing. In *Proceedings of the 13th International Conference on Natural Language Generation (INLG 2020)*, pages 182–194, 2020.

- 453 Virginia Braun and Victoria Clarke. Using thematic anal- 508  
454 ysis in psychology. *Qualitative Research in Psychology*, 509  
455 3(2):77–101, 2006.
- 456 Carl Franzen. Forget GPT-5! OpenAI launches new AI 510  
457 model family o1 claiming PhD-level performance. <https://venturebeat.com/ai/forget-gpt-5-openai-launches-new-ai-model-family-o1-claiming-phd-level-performance/>, 511  
458 September 12 2024. Accessed: 2025-01-09. 512  
459 513
- 461 Chujie Gao, Siyuan Wu, Yue Huang, Dongping Chen, Qi- 514  
462 hui Zhang, Zhengyan Fu, Yao Wan, Lichao Sun, and Xian- 515  
463 gliang Zhang. HonestLLM: Toward an Honest and Helpful 516  
464 Large Language Model. In *The Thirty-eighth Annual Con- 517  
465 ference on Neural Information Processing Systems*, 2024. 518
- 466 Shailja Gupta, Rajesh Ranjan, and Surya Narayan Singh. A 519  
467 Comprehensive Survey of Retrieval-Augmented Genera- 520  
468 tion (RAG): Evolution, Current Landscape and Future Di- 521  
469 rections, 2024. 522
- 470 David M. Howcroft, Verena Rieser, Vera Demberg, and 523  
471 Michael White. Twenty years of confusion in human eval- 524  
472 uation: NLG needs evaluation sheets and standardized de- 525  
473 finitions. In *Proceedings of the 13th International Confer- 526  
474 ence on Natural Language Generation (INLG 2020)*, pages 527  
475 169–181, 2020. 528
- 476 Sheena Iyengar and Mark Lepper. When Choice is Demoti- 529  
477 vating: Can One Desire Too Much of a Good Thing? *Jour- 530  
478 nal of personality and social psychology*, 79:995–1006, 01 531  
479 2001. 532
- 480 Daniel Kahneman. *Thinking, Fast and Slow*. Farrar, Straus 533  
481 and Giroux, New York, 2011. 534
- 482 Alejandro Cuadron Lafuente. SWE-Bench-Verified-O1- 535  
483 reasoning-high-results. [https://huggingface.co/datasets/](https://huggingface.co/datasets/AlexCuadron/SWE-Bench-Verified-O1-reasoning-high-results) 536  
484 AlexCuadron/SWE-Bench-Verified-O1-reasoning-high- 537  
485 results, 2024. Accessed: 2025-01-09. 538
- 486 M. Asher Lawson, Richard P. Larrick, and Jack B. Soll. Com- 539  
487 paring fast thinking and slow thinking: The relative ben- 540  
488 efits of interventions, individual differences, and inferen- 541  
489 tial rules. *Judgment and Decision Making*, 15(5):660–684, 542  
490 2020. 543
- 491 Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, 544  
492 Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak 545  
493 Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic 546  
494 Evaluation of Language Models. *Transactions on Ma- 547  
495 chine Learning Research*, 2023. [https://openreview.net/](https://openreview.net/forum?id=iO4LZibEqW) 548  
496 forum?id=iO4LZibEqW. 549
- 497 Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, 550  
498 Laurent Charlin, and Joelle Pineau. How NOT to evalu- 551  
499 ate your dialogue system: An empirical study of unsuper-  
500 vised evaluation metrics for dialogue response generation.  
501 In *Proceedings of the 2016 Conference on Empirical Meth-  
502 ods in Natural Language Processing*, 2016.
- 503 Natalja Menold. Rating-Scale Labeling in Online Surveys: 508  
504 An Experimental Comparison of Verbal and Numeric Rat- 509  
505 ing Scales with Respect to Measurement Quality and Re- 510  
506 spondents’ Cognitive Processes. *Sociological Methods &  
507 Research*, 49:004912411772969, 10 2017.
- Anuj Nanavati and Randolph Bias. Optimal Line Length in 508  
Reading—A Literature Review. *Visible Language*, 01 2005. 509
- Jakob Nielsen and Thomas K. Landauer. A mathematical 510  
model of the finding of usability problems. In *Proceed- 511  
ings of ACM INTERCHI ’93 Conference*, pages 206–213, 512  
1993. 513
- Birgit Popp, Sarah Keck, Androniki Mertsiotaki, Emily 514  
Kratsch, and Alexander Daum. Which Method(s) to Pick 515  
When Evaluating Large Language Models with Humans? 516  
A Comparison of Six Methods, 2025. Preprint, to be sub- 517  
mitted for publication in February 2025. 518
- Jie Ruan, Wenqing Wang, and Xiaojun Wan. Defining and 519  
Detecting Vulnerability in Human Evaluation Guidelines: 520  
A Preliminary Study Towards Reliable NLG Evaluation. In 521  
*Proceedings of the 2024 Conference of the North American 522  
Chapter of the Association for Computational Linguistics: 523  
Human Language Technologies (Volume 1: Long Papers)*, 524  
pages 7965–7989, Mexico City, Mexico, 2024. Associa- 525  
tion for Computational Linguistics. 526
- Mohamed M. Shoukri. *Measures of Interobserver Agree- 527  
ment*. Chapman and Hall/CRC, Boca Raton, FL, 2nd edi- 528  
tion, 2004. 529
- Dominik Siemon, Felix Becker, and Susanne Robra-Bissantz. 530  
How Might We? From Design Challenges to Business In- 531  
novation. *Journal of Creativity and Business Innovation*, 532  
4:96–110, 12 2018. 533
- Craig Thomson, Anya Belz, and Helen Hastie. Common 534  
flaws in running human evaluation experiments in NLP. 535  
*Computational Linguistics*, 50(2):123–135, 2024. 536
- Chris van der Lee, Albert Gatt, Emiel van Miltenburg, Sander 537  
Wubben, and Emiel Kraemer. Human evaluation of auto- 538  
matically generated text: Current trends and best practice 539  
guidelines. In *Proceedings of the 12th International Con- 540  
ference on Natural Language Generation (INLG 2019)*, 541  
pages 92–101, 2019. 542
- Shangyu Wu, Ying Xiong, Yufei Cui, Haolun Wu, Can Chen, 543  
Ye Yuan, Lianming Huang, Xue Liu, Tei-Wei Kuo, Nan 544  
Guan, and Chun Jason Xue. Retrieval-Augmented Gener- 545  
ation for Natural Language Processing: A Survey, 2024. 546
- Ruibao Xu, Zhenhailong Wang, Rongze Fan, and Pengfei Liu. 547  
Benchmarking benchmark leakage in large language mod- 548  
els. *arXiv preprint arXiv:2404.18824*, 2024. 549
- Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and 550  
Pengfei Liu. Alignment for Honesty, 2024. 551

## 552 A Appendix: Quality Criteria

Honesty	Comprehensibility
Correctness	Fluency
Context adherence	Non-redundancy
Relevancy	Understandability
Completeness	Readability

Table 4: This table shows the metrics used in our present study. We assume a hierarchical structure of the metrics, with honesty and comprehensibility being the overarching constructs that encompass the listed sub-constructs.

553 The definitions of the quality criteria for honesty and com-  
554 prehensibility are described in Popp *et al.* [2025].

## 555 B Taxonomy of Qualitative Data Analysis

### 556 1. Usability

- 557 (a) Definition discoverability: Includes challenges users  
558 face in locating the definitions within the interface.
- 559 (b) Fixed text positioning: Participants prefer text el-  
560 ements to remain visible while scrolling for easier  
561 access.
- 562 (c) Hover interaction: Challenges with the hover inter-  
563 action of definition tooltips.
- 564 (d) Instruction presentation: Comments on discover-  
565 ability and information richness of provided in-  
566 structions.
- 567 (e) Introductory page: Importance of accessible and  
568 understandable introductory page.
- 569 (f) Metrics presentation: Improving the clarity and or-  
570 ganization of metrics presentation.
- 571 (g) Navigation: User experiences and preferences re-  
572 lated to navigating the interface.
- 573 (h) Placement of generated text/ layout: This focuses  
574 on how generated text is displayed in relation to  
575 original text.
- 576 (i) Readability: Concerns about text size and layout  
577 affecting reading comfort.

### 578 2. User Sentiment

- 579 (a) Approach: This encompasses various methods and  
580 strategies participants employ to evaluate and com-  
581 pare texts.
- 582 (b) Assumption: This reflects participants' precon-  
583 ceived notions and expectations about features or  
584 outcomes.
- 585 (c) Confusion: Participants expressed uncertainty and  
586 misunderstanding regarding features of the study.
- 587 (d) Content: Participants express satisfaction with the  
588 contents or navigation of the interface.
- 589 (e) Curiosity: This reflects a desire for deeper under-  
590 standing and exploration of information.
- 591 (f) Dislikes: Participants expressed dissatisfaction  
592 with various features.
- 593 (g) Doubt: Participants express uncertainty and skepti-  
594 cism regarding the evaluation process.

- (h) Likes: Participants express positive sentiments 595  
about features of the study. 596
- (i) Struggle: Participants experience challenges either 597  
in terms of understanding the methodology or in 598  
regards to UI aspects. 599
- (j) Surprise: Participants expressed astonishment at 600  
the quality of generated texts. 601
- (k) Wishes: This reflects participants' desires for im- 602  
provements in their experience. 603

## C Appendix: Post-Evaluation Questionnaire 604

Statements rated on a 5-point agreement scale: 605

### 1. General Usability 606

- (a) I think that I would need the support of a technical 607  
person to be able to complete the survey. 608
- (b) I found the overall task easy to master. 609
- (c) I found it easy to complete the survey. 610
- (d) I found it easy to navigate through the survey. 611
- (e) I found it easy to understand the survey. 612

### 2. Navigation 613

- (a) There was sufficient guidance or instructions on 614  
how to complete each section of the survey. 615
- (b) The survey navigation was intuitive and easy to fol- 616  
low. 617
- (c) Overall, I am satisfied with the support information 618  
when completing the tasks. 619

### 3. Clarity of Questions and Response Options 620

- (a) The answer options were appropriately formatted 621  
for the questions. 622
- (b) It was easy/intuitive to find the word definitions. 623
- (c) It was easy to choose an appropriate response. 624
- (d) The questions were clearly worded and understand- 625  
able. 626

### 4. Visual Design and Layout 627

- (a) The color contrast was comfortable to read. 628
- (b) The visual design helped you stay focused on the 629  
questions. 630
- (c) Overall, I liked the survey design and layout. 631
- (d) The font size was easy to read. 632



## D Appendix: Evaluation Page UI Study 1

In this study we generated texts with different large language models based on the original text with the question "Why are listening tests conducted?".

**Original Text:** Listening Tests - As part of our user research activities, we conduct listening tests or provide support in this area. Listening tests are specifically designed to evaluate audio quality, e.g. for new audio technologies or an updated version. Clicking on this box will direct you to the listening test sub-page where you can find more detailed information about the different tests, requirements, and how to reach us if you want to conduct a listening test.

Read the generated text and answer the questions.

---

**Generated Text:** Listening tests are conducted to evaluate audio quality, e.g. for new audio technologies or an updated version. They help us to understand the user's perception of audio quality, to find out which audio quality is preferred by the user and to identify areas for improvement.

---

**\*Fluency**

	5 (very fluent)	4	3	2	1 (not fluent at all)
How do you rate the text regarding its fluency?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Correctness**

	5 (very correct)	4	3	2	1 (not correct at all)
How do you rate the text regarding its correctness?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Comprehensibility**

	5 (very comprehensible)	4	3	2	1 (not comprehensible at all)
How do you rate the text regarding its comprehensibility?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Readability**

	5 (very readable)	4	3	2	1 (not readable at all)
How do you rate the text regarding its readability?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Understandability**

	5 (very helpful)	4	3	2	1 (not helpful at all)
How do you rate the text regarding its understandability?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Honesty**

	5 (very honest)	4	3	2	1 (not honest at all)
How do you rate the text regarding its honesty?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Answer Relevancy**

	5 (very relevant)	4	3	2	1 (not relevant at all)
How do you rate the text regarding its answer relevancy?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Context Adherence**

	5 (very adherent)	4	3	2	1 (not adherent at all)
How do you rate the text regarding its context adherence?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Non-Redundancy**

	5 (non-redundant)	4	3	2	1 (very redundant)
How do you rate the text regarding its nonredundancy?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

**\*Completeness**

	5 (very complete)	4	3	2	1 (not complete at all)
How do you rate the text regarding its completeness?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

---

Do you have any comments, for example on why you made a decision, or on the task, the method or the layout of the experiment?

Figure 4: Evaluation page UI for DQE of Study 1.