

Is your LLM Ageist?

Anonymous ACL submission

Abstract

As society’s reliance on LLMs becomes stronger with each passing day, the concerns over its treatment of less-tech-savvy and other discriminated demographics grow and grow. We examine the effects of age on LLM use on both the input and output sides, over two tasks representing an assisted creativity scenario and an information request scenario. The results of our human study suggest that LLMs do not significantly discriminate against younger or older adults, despite their clear stylistic differences. We do, however, find that assisting creativity is easier for models than providing information. We also find that the practice of LLM-as-a-judge is a reliable proxy for self-evaluation questionnaires on these scenarios.

1 Introduction

Human language undergoes notable changes during healthy aging, affecting both spoken and written communication (Biran et al., 2023; Cho et al., 2021). These changes, driven by age-related sensory and cognitive alterations, are reflected across multiple language domains, including lexical-semantic processing, phonology, grammar, comprehension, and discourse abilities (Burke and Shafto, 2011; Baciú and Roger, 2024). Recent research indicates that age affects vocabulary and syntax in digital communication, with younger adults using more diverse and descriptive language than their older counterparts (Alshahrani, 2023). In addition, older adults make more typing errors and provide fewer corrections (Kalman et al., 2015). These findings highlight the distinct linguistic styles that different age groups may employ in written communication and how they interact with AI systems.

On the side of artificial intelligence, the development of large language models (LLMs), such as ChatGPT (Brown et al., 2020; Naveed et al., 2025),

have made AI tools widely accessible across various domains (Pack and Maloney, 2023; Burger et al., 2023). However, despite their impressive capabilities, these systems are not neutral. LLMs are known to exhibit demographic biases, including systematic under-representation and age bias (Zhang et al., 2021; Chu et al., 2023). Beyond these structural biases, LLM behavior can also vary substantially depending on how prompts are phrased (Sclar et al., 2024; Hackmann et al., 2024).

In this study, we focus on different age groups interacting with LLMs. We suspect that, given the way wording and style of user input directly shape model responses, differences in how younger and older adults formulate prompts for LLMs may influence the quality and effectiveness of the output. While prior work shows that sociodemographic cues in prompts can influence how LLMs simulate demographic groups (Lutz et al., 2025) and that implicit personalization may occur when models infer user traits from conversational signals (Neplenbroek et al., 2025), no studies to date have examined how naturally-produced prompts by different demographic groups differ in structure, nor how such differences influence model outputs. This leaves open the question of how user characteristics, prompt formulation, and LLM performance interact in real-world settings. We address this problem in a human-subject study, comparing older and younger adults in their interaction with an LLM in two distinct use modes: an *assisted creativity* scenario and an *information request* scenario.

Along three axes of comparison, we find somewhat surprising results. Although the language used by older and younger adults addressing the LLM varies qualitatively, there are few quantitative differences in the model’s responses, either in length or in quality, according to statistical measures and judgments by both the users themselves and external LLMs. We do, however, find con-

082 sistent differences between the two task scenarios,
083 indicating that models struggle with providing in-
084 formation but excel at assisting creativity.¹

085 2 Method

086 The study employed a 2×2 mixed design, with
087 Age Group (younger adults, ages 20–45, vs. older
088 adults, ages 65+) as a between-subjects factor and
089 Task Type (assisted creativity vs. information re-
090 quest) as a within-subjects factor.

091 2.1 Participants

092 Our sample consisted of 90 participants, of whom
093 45 were older adults ($M = 70.53$, $SD = 4.69$,
094 range = 65–84; 31 female, 14 male) and 45 were
095 younger adults ($M = 31.89$, $SD = 5.20$, range =
096 21–40; 29 female, 14 male, two participants who
097 selected “prefer not to say”). All participants were
098 native English speakers residing in the UK with no
099 prior experience with AI systems such as ChatGPT
100 or Gemini. Participants were recruited via Prolific²
101 across two independent trials, one for each age
102 group. Participants completed the study in an aver-
103 age of 22:13 minutes ($SD = 5 : 46$, range 13:46–
104 48:02). See full details for recruitment, eligibility
105 criteria, and screening procedures in Appendix A.1.

106 2.2 Design and Procedure

107 Participants first received an explanation of GPT
108 and completed a brief training phase to become
109 familiar with interactions with the AI system. They
110 then completed two experimental tasks, which were
111 presented in counterbalanced order to control for
112 sequence effects. All interactions were open-ended,
113 unconstrained in both length and structure, and
114 routed through the G4R system (Kim, 2025), using
115 the GPT-4o-mini model. Complete instructions for
116 each stage are provided in Appendix A.2.

117 **Information request task (trip planning).** Par-
118 ticipants were instructed to interact with the system
119 to obtain a complete plan for a multi-day trip to
120 Andorra. This interaction was intended to elicit
121 information-seeking use of the model.

122 **Assisted creativity task (greeting writing).** Par-
123 ticipants were instructed to interact with the AI to
124 generate a personalized greeting-card message for
125 a close friend or family member. This interaction
126 was intended to elicit assisted creative use of the
127 model.

¹Data and code will be shared with publication.

²<https://prolific.com/>

Dimension	Question Form
Relevance	To what extent did the response address what you asked for?
Usefulness	How useful was the response for your purpose?
Clarity	How clear and understandable was the response?
Sufficiency	If this were the only answer you received, would it be sufficient?

Table 1: Human evaluation questionnaire items used to assess model performance.

128 2.3 Measures

129 **Surface quantities.** We extracted from the con-
130 versations the response times for both model and
131 participant, measured in seconds; as well as text
132 characteristics from user input and model output:
133 total word count and spelling error rate. These mea-
134 sures were used to index response length, verbosity,
135 and linguistic quality.

136 **Lexical features.** We measured lexical differ-
137 ences between the groups in terms of word fre-
138 quency counts (raw and Older/Younger ratio), as
139 well as Positive Pointwise Mutual Information
140 (PPMI). Full details in Appendix B.1).

141 **Self-evaluation.** After each task, participants
142 completed a questionnaire assessing their satisfac-
143 tion with the model’s output across four dimen-
144 sions: relevance, usefulness, clarity, and sufficiency
145 (see Table 1 for the exact formulation). Each dimen-
146 sion was rated on a 5-point Likert scale (1 = Not at
147 all, 5 = Completely), and then analyzed separately.

148 **External quality evaluations.** We solicited other
149 LLMs (Google Gemini, Meta LLama, and Anthropic Claude) for ratings mirroring those of the self-evaluation. Full details of the evaluation procedure are provided in Appendix B.2.

153 3 Results

154 **Surface quantities.** The main effects of age and
155 task on statistical measures in the conversations
156 are shown in Table 2, and comprehensive statisti-
157 cal analyses are presented in Appendix C. Overall,
158 we found that task type was the primary source
159 of variation across all measures in both user in-
160 put and model output, whereas age group effects
161 were limited and specific. Participants in both age
162 groups produced longer inputs, and GPT generated

Measure	Task	Age	Task × Age	Direction
<i>User Input Characteristics</i>				
# of messages	–	–	–	–
Word count	*	–	–	Trip > Greet
Error rate	*	–	–	Trip > Greet
<i>GPT Output Characteristics</i>				
Response time	***	–	–	Trip > Greet
Output length	***	–	*	↑ Trip (YA)
Error rate	***	–	–	Trip > Greet

Table 2: Summary of significant effects for user input and GPT output characteristics. Asterisks indicate statistically significant effects (* $p < .05$, ** $p < .01$, *** $p < .001$). ↑ indicates a larger task-related increase in a task for an age group.

163 substantially longer, slower responses during the in-
164 formation request (trip planning) task compared to
165 the assisted creativity (greeting) task. This pattern
166 was consistent across word counts, spelling error
167 rates, and response times, reflecting the increased
168 cognitive and informational demands of informa-
169 tion request interactions. Nonetheless, the informa-
170 tion request task did present some age-group out-
171 put variation, with longer responses generated for
172 younger adults. Additional analyses that included
173 participant input length as a covariate supported
174 these results, demonstrating that task-dependent
175 age differences in GPT output length are not ex-
176 plained by variations in participant input length or
177 input–output coupling.

178 Taken together, the findings indicate that dif-
179 ferences in both user input and GPT output were
180 driven primarily by task demands. However, age-
181 related differences observed in the information re-
182 quest task point to a potential tendency for GPT to
183 generate longer responses for younger adults under
184 more complicated demands, suggesting bias con-
185 gruent with the assumption that AI systems target
186 more dominant user demographics.

187 **Lexical features.** We analyzed the conversation
188 texts based on relative frequency and Positive Point-
189 wise Mutual Information (PPMI), identifying sys-
190 tematic age-related differences in both participant
191 prompts and model responses. When prompts
192 were combined across tasks, core vocabulary (i.e.,
193 the most frequent words) was similar across age
194 groups. However, as demonstrated in Figure 1,
195 PPMI revealed consistent differences in lexical em-
196 phasis. Prompts from younger adults contained
197 more action- and production-oriented language,

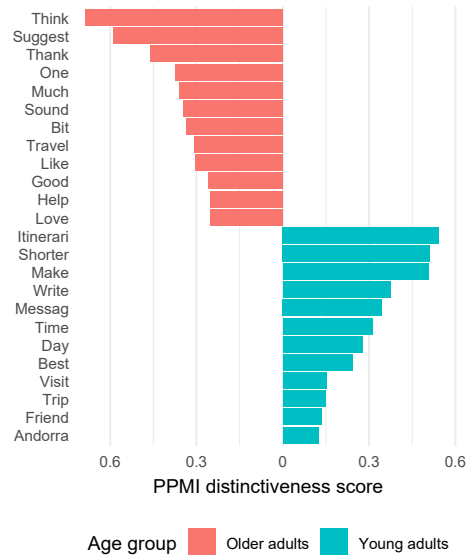


Figure 1: Top word stems per age group in participants' prompts to the language model by age group, ranked by PPMI, on both tasks. Higher absolute PPMI values indicate stronger-than-expected associations between a stem and a given age group.

198 whereas those from older adults more frequently
199 included reflective, evaluative, and interpersonal
200 terms. These patterns persisted after stemming, in-
201 dicated stable interaction styles (full results and
202 analysis in Appendix D.1).

203 PPMI analyses of the model's responses revealed
204 complementary age-related differences in its out-
205 puts, particularly in the information request task,
206 as seen in Figure 2. Responses directed to older
207 adults emphasized transportation, infrastructure,
208 and procedural travel information, whereas those
209 to younger adults highlighted experiential, lifestyle-
210 oriented, and contextual content. In the assisted
211 creativity task, responses to older adults were more
212 strongly associated with relational and life-course
213 framing, while those to younger adults emphasized
214 immediacy, affective tone, and action-oriented lan-
215 guage (see Appendix D.2).

216 These findings suggest a reciprocal adaptation
217 process. The systematic alignment between user
218 prompts and model responses may reflect age-
219 related communication norms: interactions with
220 older users prioritize stability and relational con-
221 text, whereas interactions with younger users em-
222 phasize exploration and immediacy.

223 **Self-evaluation.** Participants' rated the model's
224 output across four dimensions: relevance, useful-
225 ness, clarity, and sufficiency. Comprehensive sta-

Measure	Assisted Creativity		Information Request		Age Group Difference
	Younger Adults	Older Adults	Younger Adults	Older Adults	
Relevance	4.53 (0.11)	4.36 (0.11)	4.11 (0.12)	4.40 (0.12)	*
Usefulness	4.18 (0.14)	4.22 (0.14)	4.02 (0.14)	4.31 (0.14)	–
Clarity	4.53 (0.09)	4.73 (0.09)	3.96 (0.14)	4.47 (0.14)	**
Sufficiency	4.13 (0.16)	4.16 (0.16)	3.69 (0.17)	3.91 (0.17)	–

Table 3: Mean self-evaluation ratings (and standard errors) by task type and age group. The final column indicates significant differences between age groups within a task condition (* $p < .05$, ** $p < .01$).

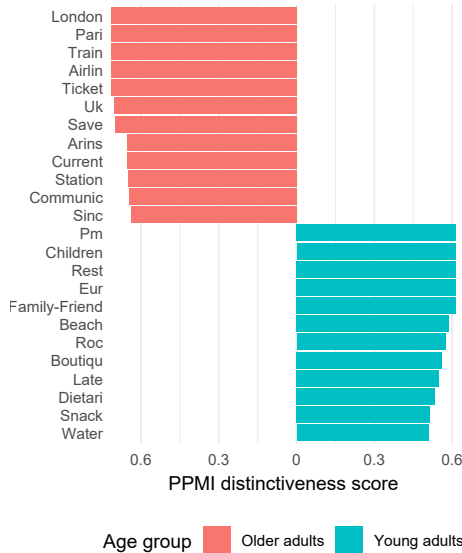


Figure 2: Top word stems per age group in the language model’s responses, ranked by PPMI, on the information request task (trip planning).

tistical analyses for each measure are provided in Appendix E. In Table 3, we present the aggregated ratings across all dimensions. We find that task type and age group influenced self-evaluations selectively rather than uniformly: responses in the greeting task received higher ratings for clarity and sufficiency for task completion compared to those in the trip planning task, but not much difference was found on relevance and usefulness. This may suggest that users’ expectations for a creativity use-case are lower than for obtaining information with respect to specificity of the answer, or that the model is less successful in producing reliable information than free text.

Age-related differences were most evident in clarity and relevance: older adults consistently rated the model’s responses as clearer than those of younger adults across all tasks. For relevance, younger adults were more sensitive to task type, whereas older adults provided consistently high ratings across contexts.

In summary, participants’ subjective evaluations of the model’s responses were primarily shaped by task characteristics rather than by age-related evaluation patterns. The observed age-related differences seem to reflect selective sensitivity to task structure rather than overall evaluation tendencies.

LLM as a judge. We prompted non-GPT LLMs to assess the model’s responses across six dimensions: relevance, usefulness, clarity, sufficiency, accuracy, and conversational flow. Our results (provided in Appendix F) suggest that task type was the primary source of systematic variation in ratings judged by LLM. Specifically, across the four dimensions shared with self-evaluation, responses in the greeting task received higher ratings than those in the trip planning task. This difference in ratings between tasks is consistent with trends in participant self-evaluations and aligns with previous research that demonstrated similarities between human and model-based judgments (Ho et al., 2025; Wang et al., 2025). Such findings may indicate that external LLM judges, like participants themselves, tend to perceive the model as more successful in creative contexts than in tasks requiring specific information retrieval. Age group did not produce main effects or interactions across any evaluation dimension in the LLM judgments, again echoing self-evaluation results.

4 Conclusion

Within the context of ongoing concern over the “technological gap” exacerbated by AI systems and their biased behavior towards already-discriminated groups, we find our results point towards a cautiously optimistic prospect: a popular system with no knowledge of the user’s age did not, as a whole, present age-related bias when faced with users picked to not have prior experience with such possible bias. We are also encouraged by the fact that independent LLMs served as reliable judges of model outputs, facilitating further (much needed) exploration of these questions.

288 **Limitations**

289 Subjective evaluation of model performance has
290 only been collected from the participants them-
291 selves and from external LLMs. We have not in-
292 cluded subjective evaluations from external par-
293 ticipants, although our own cursory glances over
294 conversations suggested that the self-evaluations
295 are mostly reliable. In any event, they teach us
296 about people’s expectations from models at differ-
297 ent ages.

298 This study investigated only two types of tasks,
299 each requiring distinct cognitive and communica-
300 tive skills. Therefore, the observed age-related
301 language differences may be influenced by task-
302 specific demands and may not be applicable to
303 a broader range of everyday activities. Further
304 research examining language use across a wider
305 range of tasks could clarify which age-related pat-
306 terns are consistent across contexts and which are
307 specific to particular tasks.

308 **Ethical Considerations**

309 The study was approved by the ethics committee
310 of the authors’ institution.

311 **References**

312 Abdulaziz Alshahrani. 2023. [Analyzing the impact](#)
313 [of generational variations on vocabulary and syn-](#)
314 [tax in instagram posts of older and younger adults.](#)
315 *Eurasian Journal of Applied Linguistics*, 9(1):202–
316 212.

317 Monica Baciu and Elise Roger. 2024. [Finding the](#)
318 [words: How does the aging brain process language?](#)
319 [a focused review of brain connectivity and compen-](#)
320 [satory pathways.](#) *Topics in Cognitive Science*.

321 Kenneth Benoit, Kohei Watanabe, Haiyan Wang, Paul
322 Nulty, Adam Obeng, Stefan Müller, and Akitaka Mat-
323 suo. 2018. [quanteda: An r package for the quantita-](#)
324 [tive analysis of textual data.](#) *Journal of Open Source*
325 *Software*, 3(30):774.

326 Michal Biran, Aviah Gvion, and Shira Shmueli-Samuel.
327 2023. [Language in healthy ageing: A comparison](#)
328 [across language domains.](#) *Folia Phoniatrica et Lo-*
329 *gopaedica*, 75(2):90–103.

330 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
331 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
332 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
333 Askeel, and 1 others. 2020. [Language models are](#)
334 [few-shot learners.](#) *Advances in neural information*
335 *processing systems*, 33:1877–1901.

336 Bastian Burger, Dominik K Kanbach, Sascha Kraus,
337 Matthias Breier, and Vincenzo Corvello. 2023. On

the use of ai-based tools like chatgpt to support man- 338
agement research. *European journal of innovation 339*
management, 26(7):233–241. 340

Deborah M. Burke and Meredith A. Shafto. 2011. 341
Language and aging. In Fergus I. M. Craik and Timo- 342
thy A. Salthouse, editors, *The Handbook of Aging 343*
and Cognition, 3 edition, pages 373–443. Psychology 344
Press. 345

Sunghye Cho, Naomi Nevler, Sanjana Shellikeri, Na- 346
talia Parjane, David J Irwin, Neville Ryant, Sharon 347
Ash, Christopher Cieri, Mark Liberman, and Murray 348
Grossman. 2021. [Lexical and acoustic characteristics](#)
349 [of young and older healthy adults.](#) *Journal of Speech,*
350 *Language, and Hearing Research*, 64(2):302–314. 351

Charlene H. Chu, Simon Donato-Woodger, Shehroz S. 352
Khan, Rune Nyrup, Kathleen Leslie, Alexandra 353
Lyn, Tianyu Shi, Andria Bianchi, Samira Abbas- 354
gholizadeh Rahimi, and Amanda Grenier. 2023. [Age-](#)
355 [related bias and artificial intelligence: a scoping re-](#)
356 [view.](#) *Humanities and Social Sciences Communica-*
357 *tions*, 10(1):510. 358

Franz Faul, Edgar Erdfelder, Albert-Georg Lang, and 359
Axel Buchner. 2007. [G* power 3: A flexible statisti-](#)
360 [cal power analysis program for the social, behavioral,](#)
361 [and biomedical sciences.](#) *Behavior research methods,*
362 39(2):175–191. 363

Stefan Hackmann, Haniyeh Mahmoudian, Mark Stead- 364
man, and Michael Schmidt. 2024. [Word importance](#)
365 [explains how prompts affect language model outputs.](#)
366 *Preprint*, arXiv:2403.03028. 367

Xanh Ho, Jiahao Huang, Florian Boudin, and Akiko 368
Aizawa. 2025. [Llm-as-a-judge: Reassessing the](#)
369 [performance of llms in extractive qa.](#) *Preprint*,
370 arXiv:2504.11972. 371

Yoram M Kalman, Gitit Kavé, and Daniil Umanski. 372
2015. [Writing in a digital world: self-correction](#)
373 [while typing in younger and older adults.](#) *Interna-*
374 [tional Journal of Environmental Research and Public](#)
375 [Health](#), 12(10):12723–12734. 376

Jin Kim. 2025. [How to capture and study conversations](#)
377 [between research participants and chatgpt: Gpt for](#)
378 [researchers \(g4r.org\).](#) *Preprint*, arXiv:2503.18303. 379

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, 380
Ruo Chen Xu, and Chenguang Zhu. 2023. [G-eval:](#)
381 [Nlg evaluation using gpt-4 with better human align-](#)
382 [ment.](#) In *Proceedings of the 2023 Conference on*
383 *Empirical Methods in Natural Language Processing*,
384 pages 2511–2522. Association for Computational
385 Linguistics. 386

Marlene Lutz, Indira Sen, Georg Ahnert, Elisa Rogers, 387
and Markus Strohmaier. 2025. [The prompt makes](#)
388 [the person\(a\): A systematic evaluation of sociodemo-](#)
389 [graphic persona prompting for large language models.](#)
390 In *Findings of the Association for Computational Lin-*
391 *guistics: EMNLP 2025*, pages 23212–23237, Suzhou,
392 China. Association for Computational Linguistics. 393

394 Burt L Monroe, Michael P Colaresi, and Kevin M Quinn.
395 2008. Fightin’ words: Lexical feature selection and
396 evaluation for identifying the content of political con-
397 flict. *Political Analysis*, 16(4):372–403.

398 Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad
399 Saqib, Saeed Anwar, Muhammad Usman, Naveed
400 Akhtar, Nick Barnes, and Ajmal Mian. 2025. A com-
401 prehensive overview of large language models. *ACM*
402 *Transactions on Intelligent Systems and Technology*,
403 16(5):1–72.

404 Vera Neplenbroek, Arianna Bisazza, and Raquel Fer-
405 nández. 2025. [Reading between the prompts: How](#)
406 [stereotypes shape LLM’s implicit personalization.](#)
407 In *Proceedings of the 2025 Conference on Empirical*
408 *Methods in Natural Language Processing*, pages
409 20378–20411, Suzhou, China. Association for Com-
410 putational Linguistics.

411 Austin Pack and Jeffrey Maloney. 2023. [Using genera-](#)
412 [tive artificial intelligence for language education re-](#)
413 [search: Insights from using openai’s chatgpt.](#) *TESOL*
414 *Quarterly*, 57(4):1571–1582.

415 Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane
416 Suhr. 2024. [Quantifying language models’ sensitiv-](#)
417 [ity to spurious features in prompt design or: How i](#)
418 [learned to start worrying about prompt formatting.](#)
419 In *Proceedings of the International Conference on*
420 *Learning Representations (ICLR)*.

421 Peter D. Turney and Patrick Pantel. 2010. From fre-
422 quency to meaning: Vector space models of se-
423 mantics. *Journal of Artificial Intelligence Research*,
424 37:141–188.

425 Ruiqi Wang, Jiyu Guo, Cuiyun Gao, Guodong Fan,
426 Chun Yong Chong, and Xin Xia. 2025. [Can llms](#)
427 [replace human evaluators? an empirical study of llm-](#)
428 [as-a-judge in software engineering.](#) *Proceedings of*
429 *the ACM on Software Engineering*, 2(ISSTA):1955–
430 1977.

431 Sheng Zhang, Xin Zhang, Weiming Zhang, and Anders
432 Sjøgaard. 2021. [Sociolectal analysis of pretrained](#)
433 [language models.](#) In *Proceedings of the 2021 Con-*
434 *ference on Empirical Methods in Natural Language*
435 *Processing*, pages 4581–4588. Association for Com-
436 putational Linguistics.

437 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
438 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
439 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,
440 Joseph E. Gonzalez, and Ion Stoica. 2023. Judging
441 LLM-as-a-judge with MT-bench and chatbot arena.
442 *Advances in Neural Information Processing Systems*,
443 36:46595–46623.

A Detailed Experimental Setup

A.1 Participants

The final sample consisted of 90 participants, including 45 younger adults (ages 20–45) and 45 older adults (ages 65+), all of whom were native English speakers residing in the UK. Participants were recruited via Prolific³ across two independent trials, one for each age group. Participants were required to have completed at least 50 previous studies on Prolific with an approval rate of 98% or higher, to report regular internet usage, and to have no prior experience with AI systems such as ChatGPT, Gemini, or Claude. Since slightly more older adults were initially recruited, the younger adults were subsequently matched to ensure equal group sizes. The planned sample size ($n = 84$) was determined through an a priori power analysis using G*Power (Faul et al., 2007) ($\alpha = .05$, power = .95, $f = .20$). A total of 880 individuals initiated the screening process. Of these, 790 were excluded—785 due to prior experience with AI and 5 due to a technical issue—resulting in a final sample of 90 participants. The older adult group ($n = 45$) had a mean age of 70.53 years ($SD = 4.69$; range = 65–84) and included 31 females and 14 males, whereas the younger adult group ($n = 45$) had a mean age of 31.89 years ($SD = 5.20$; range = 21–40) and included 29 females, 14 males, and two participants who selected “prefer not to say.” Participants completed the study in an average of 22 min 13 s ($SD = 5$ min 46 s), with completion times ranging from 13 min 46 s to 48 min 02 s.

All participants provided informed consent. They were informed of their right to withdraw from the trial at any time without penalty. Participants were also notified that they would receive credit for their participation time and that their data would be used exclusively for research purposes. Compensation for all participants was provided in accordance with Prolific’s standard payment policies. Participants who completed the study received £3.75 for an average duration of 25 minutes (approximately £12/hour). Participants who were screened out after only completing the initial screening question received £0.10 for that stage.

A.2 Instructions Given to Participants

Below is the full text of the instructions provided to participants for each of the three tasks: GPT

³<https://prolific.com/>

explanation, the training phase, the information request task, and the assisted creativity task.

GPT explanation: Introduction to ChatGPT

Thank you for agreeing to take part in this study. In the study, you will be asked to complete two different tasks that involve interacting with a tool called ChatGPT — a type of Artificial Intelligence (AI) that can understand and generate human-like language. You do not need any previous experience with AI. The tasks will guide you through how to use GPT, and you’ll be asked to type a few questions or instructions and observe how it responds.

What is ChatGPT?

ChatGPT stands for Generative Pre-trained Transformer. It is a powerful AI system that was exposed to a large amount of text from the internet, including books, articles, and conversations, which allows it to respond in a natural, human-like way.

When you write a message (called a prompt), ChatGPT uses what it has learned to generate a reply. But it’s not just a question-and-answer system. GPT is designed to carry out a conversation, so you can continue the interaction, clarify what you mean, refine your request, or build on what it said previously — much like you would in a real discussion with a person.

People use ChatGPT for many different purposes, such as:

- Asking for information or explanations (e.g., “What causes climate change?”)
- Getting ideas (e.g., “What can I cook for dinner with only pasta and spinach?”)
- Planning (e.g., “Help me plan a 3-day trip to Barcelona.”)
- Improving writing (e.g., “Can you make this message sound more polite?”)
- Creating things (e.g., poems, stories, recipes, or even computer programs)

In this study, we are interested in how people interact with GPT. There are no right or wrong answers — just be yourself and try your best to communicate naturally.

Training Phase: first short task Before moving on to the main exercises, you’ll now have a chance to explore how ChatGPT works briefly. Start by asking ChatGPT for something simple — for example, a recipe. Once you receive the response, ask

ChatGPT to adjust it based on a specific preference or need — for instance, removing one ingredient because you want a version without it. Try to notice how the chatbot responds and adapts. This task is designed to help you become familiar with interacting with ChatGPT.

You have 4 minutes for this short task. You may complete it in less time if you feel you have reached a satisfactory result. After a few minutes, a “continue” button will appear, allowing you to move on at any point. However, it is recommended to use the full 4 minutes to work on the task. If you need more time to finalize your answer, it is also fine to continue a little beyond the 4 minutes. On the next screen, you will see the task interface. At the bottom of the page, a timer will show the remaining time.

Information Request Task [Trip Planning]

Trip to Andorra Task: Traveling can be a wonderful way to discover new places, relax and learn about new cultures. But planning a trip can often be challenging, especially if you’re traveling to a new country or if you’re unsure of everything you want to do and perhaps traveling alone for the first time. In such cases, it can be helpful to turn to the assistance of chatbots. We’ll now look at whether and how ChatGPT can help you plan trips. These sample prompts (input or instructions you type for the AI) might help you with your trip planning:

- “You are a tour guide. I’m very interested in theater in Naples, please tell me more about what places and buildings I should visit and in what order.”
- “What is the cheapest destination for a 3-day city trip in Europe? My budget is around 1000 euros.”
- “List me free museums in Amsterdam. I am primarily interested in modern art.”

Your task now is to plan a 4-day trip to Andorra in September. Whether you travel alone or with others, where you stay, whether you travel around, what activities you do, etc., are entirely up to you. Please plan your trip as concretely as possible. Use your own language, do not copy-paste from the example prompts above. Ask the chatbot to plan the itinerary as if you had no prior knowledge, relying entirely on its suggestions. Try to create the itinerary as “automated” as possible using (almost exclusively) the chatbot.

You have 7 minutes for the task “Travel to Andorra.” You may complete it in less time if you feel you have reached a satisfactory result. After a few minutes, a “continue” button will appear, allowing you to move on at any point. However, it is recommended to use the full 7 minutes to work on the task. If you need more time to finalize your answer, it is also fine to continue a little beyond the 7 minutes. On the next screen, you will see the task interface. At the bottom of the page, a timer will show the remaining time.

Assisted Creativity Task [Greeting Writing]

Birthday Greeting Task: Sending a thoughtful greeting can be a wonderful way to show someone you care. But finding the right words can often be difficult, especially when you’re trying to strike the right tone: warm, kind, and appropriate for the occasion. Whether it’s a birthday, a wedding, the birth of a child, or a retirement celebration, many people turn to chatbots like ChatGPT to help them write meaningful and well-phrased messages. We’ll now look at whether and how ChatGPT can help you write greetings. These sample prompts (input or instructions you type for the AI) might help you with your greeting:

- “Write a warm and heartfelt message to congratulate a friend on the birth of their first child.”
- “Create a short and cheerful wedding message suitable for a greeting card.”
- “Generate a respectful and appreciative note for someone retiring after many years of work.”

Your task now is to write an instruction for ChatGPT to generate a birthday message for someone close to you. Whether it’s a family member, a friend, or someone else you care about is entirely up to you. The tone, length, and level of formality are also yours to decide. Please write your prompt as clearly and concretely as possible. Use your own language, do not copy-paste from the example prompts above. However, avoid “unnecessary” personal contributions in the form of your own wording or personal details. Try to formulate the instruction as “automated” as possible using (almost exclusively) the chatbot.

You have 7 minutes for the task “Birthday Greeting.” You may complete it in less time if you feel you have reached a satisfactory result. After a few

minutes, a “continue” button will appear, allowing you to move on at any point. However, it is recommended to use the full 7 minutes to work on the task. If you need more time to finalize your answer, it is also fine to continue a little beyond the 7 minutes. On the next screen, you will see the task interface. At the bottom of the page, a timer will show the remaining time.

A.3 Self-Evaluation Questionnaire

After each task, participants completed a self-evaluation questionnaire assessing their satisfaction with the model’s output across four dimensions: Relevance, Usefulness, Clarity, and Sufficiency. Each dimension was rated on a 5-point Likert scale (1 = Not at all, 5 = Completely).

Instructions and Items Participants were presented with the following prompt:

“Based on the most recent response you received from ChatGPT, please evaluate the performance of the model by answering the following questions:”

The questionnaire consisted of the following four items:

1. **Relevance:** To what extent did the response address what you asked for?
2. **Usefulness:** How useful was the response for your purpose?
3. **Clarity:** How clear and understandable was the response?
4. **Sufficiency:** If this were the only answer you received, would it be sufficient?

All items were answered using the following 5-point scale: (1) Not at all, (2) Slightly, (3) Somewhat, (4) Mostly, and (5) Completely.

B Detailed Analysis Setups

B.1 Lexical Analysis Method Details

To examine systematic differences in language use between age groups and to characterize the linguistic distinctiveness of the model’s responses, we conducted a quantitative text analysis using R and the *quanteda* package (Benoit et al., 2018).S The analysis was applied to two distinct corpora: (1) participants’ prompts to the model, and (2) the model’s generated responses.

Data Aggregation and Preprocessing. To capture stable interaction styles and reduce noise from idiosyncratic utterance-level variations, textual data were aggregated at the participant level. This process yielded one concatenated document per participant for each analysis target. Before analysis, all texts underwent a standard preprocessing pipeline. First, the text was tokenized into individual tokens, with punctuation, numbers, and symbols removed. Second, standard English stopwords were removed. Third, tokens shorter than 2 characters were filtered out to focus on content-bearing vocabulary. Finally, we applied the Snowball stemming algorithm implemented in the *quanteda* package (Benoit et al., 2018) to reduce words to their root forms (e.g., merging “planning” and “planned” into “plan”). This ensured that morphological variations did not obscure broader lexical patterns.

Relative Frequency Ratio. To visualize broad differences in word usage, we constructed a Document-Feature Matrix (DFM). To ensure robustness, features were filtered to retain only words that appeared in at least 2 participants’ documents ($docfreq \geq 2$). For each retained word, we calculated its relative frequency within each age group (count divided by the group’s total token count). We then computed the ratio of these frequencies (Older Adults / Younger Adults) to identify words disproportionately used by one group relative to the other, following standard approaches for lexical comparison (Monroe et al., 2008).

Lexical Distinctiveness (PPMI). To identify specific vocabulary markers that statistically distinguish between age groups, we computed Positive Pointwise Mutual Information (PPMI) (Turney and Pantel, 2010). For this analysis, the participant-level DFM was aggregated by Age Group. To mitigate the known PMI bias toward low-frequency terms, we applied a secondary frequency threshold, retaining only words with a minimum total frequency of 30 across the entire corpus.

PMI quantifies the discrepancy between the observed joint probability of a word (w) and a group (g) and their expected probability under the assumption of independence. The metric was calculated as:

$$PMI(w, g) = \log_2 \left(\frac{P(w, g)}{P(w) \times P(g)} \right) \quad (1)$$

We retained only positive values to focus on char-

acteristic associations:

$$\text{PPMI}(w, g) = \max(\text{PMI}(w, g), 0) \quad (2)$$

Higher PPMI scores indicate that a specific word stem appears in a group’s text significantly more often than would be expected by chance, given its overall frequency in the corpus.

B.2 LLM-as-a-Judge: Method Details

Recent studies have demonstrated that large language models can serve as evaluators of natural language generation quality when given structured prompts and explicit scoring criteria (Zheng et al., 2023; Liu et al., 2023). Building on this approach, the present study implemented an LLM-as-a-Judge framework to assess the quality of AI-generated responses in the Assisted Creativity task (Greeting writing) and the Information Request task (Trip planning). This framework uses advanced large language models as impartial evaluators to score complete conversational transcripts against predefined evaluation criteria.

Data Preprocessing and Age Masking. To eliminate potential biases related to participants’ ages, a regex-based masking algorithm was applied to redact age-related information from the conversation logs prior to evaluation. The algorithm identified and masked various age patterns, such as explicit numerical declarations (e.g., “I am 70 years old”), phrases indicating upcoming milestones (e.g., “turning 30”), and context-dependent descriptors (e.g., “in my twenties”). All identified instances were replaced with a neutral token (###). The LLM judges received only these masked transcripts as input. By removing explicit age indicators, the scoring process focused exclusively on the content and quality of the interactions. As a result, the evaluation procedure was fully blinded to participants’ age groups.

Models and Configuration. To ensure robust and diverse evaluation, three distinct large language models from different model families were selected as judges. The models included:

- **Google Gemini:** Specifically, the gemini-2.5-pro version.
- **Anthropic Claude:** Using the claude-sonnet-4-5-20250929 model.
- **Meta Llama (via Groq):** Using the llama-3.3-70b-versatile model.

This selection facilitates comparisons across different architectures and providers, thereby reducing reliance on the inductive biases of any single model. To maximize reproducibility and minimize variability in evaluation scores, the temperature parameter for all models was set to 0.0.

Evaluation Metrics and Prompting. Each judge received the complete, age-masked transcript of the conversation along with task-specific instructions. A comprehensive system prompt was developed to guide the evaluation process, directing the LLMs to act as “expert evaluators of human–AI interactions.” The models were instructed to produce a structured JSON response containing both reasoning and a numerical score for each evaluation dimension, as defined below. We defined six specific evaluation dimensions to capture different aspects of the interaction quality. The judges were provided with the following definitions for each metric:

1. **Task Relevance:** Assesses the extent to which the AI’s responses directly addressed the specific requests made by the user throughout the conversation.
2. **Usefulness:** Evaluates the practical value of the responses, specifically considering how useful the AI’s suggestions and content were for the user’s stated purpose.
3. **Clarity:** Measures how understandable, coherent, and clear the AI’s responses were to read and comprehend.
4. **Sufficiency:** Determines whether the responses provided enough information to complete the task (e.g., plan the trip or write the greeting) without the user needing to seek further external input.
5. **Accuracy:** Checks the factual correctness of the information provided and ensures the responses were free from errors or misleading statements.
6. **Conversational Flow:** Evaluates the smoothness and naturalness of the interaction, assessing how well the AI maintained coherence across multiple turns of the conversation.

Scoring Scale. The models were instructed to score each dimension using a 5-point Likert scale, ranging from 1 (“Not at all”) to 5 (“Completely”). To ensure rigorous evaluation, the prompt explicitly

directed the models to interpret the scale conservatively and reserve high scores for responses that clearly exceeded the standard of quality.

Aggregation of Judge Scores. To reduce bias from individual models, ratings from three judges were used. For each evaluation dimension, the scores from Gemini, Claude, and Llama were averaged, yielding a single judge-averaged score for each participant and task. All six dimensions were analyzed separately, and no overall score was computed. These judge-averaged scores were used as dependent variables in subsequent statistical analyses to examine the effects of age group and task type.

C Surface Quantities

We begin by reporting the results of the **surface quantities analysis**. This analysis examines the structural and quantitative properties of the interaction, providing a baseline characterization of the dialogue flow. We present findings regarding volume (word count, message count), technical accuracy (spelling errors), and timing (latency) for both participants' prompts and the model's responses.

C.1 User Input Characteristics

Number of messages to GPT A 2×2 mixed-design ANOVA was conducted to examine whether the number of messages sent to GPT differed by age group (younger adults vs. older adults; between-subjects) and task type (assisted creativity task [greeting writing] vs. information request task [trip planning]; within-subjects) (see Figure 3)

The main effect of age group was not significant, $F(1, 88) = 0.60, p = .44, \eta_g^2 = .005$, indicating that older adults ($M = 4.84, SE = 0.27$) and younger adults ($M = 5.14, SE = 0.27$) sent a similar number of messages. The main effect of task type was also not significant, $F(1, 88) = 0.38, p = .54, \eta_g^2 = .000$, with a similar number of messages for the assisted creativity task ($M = 5.06, SE = 0.21$) and the information request task ($M = 4.93, SE = 0.23$). The interaction between age group and task type was not significant, $F(1, 88) = 0.25, p = .62, \eta_g^2 = .000$, indicating that system interaction was similar across age groups and task types. These results indicate that the number of messages sent to the system was similar across age groups and task types, suggesting that both younger and older adults engaged with

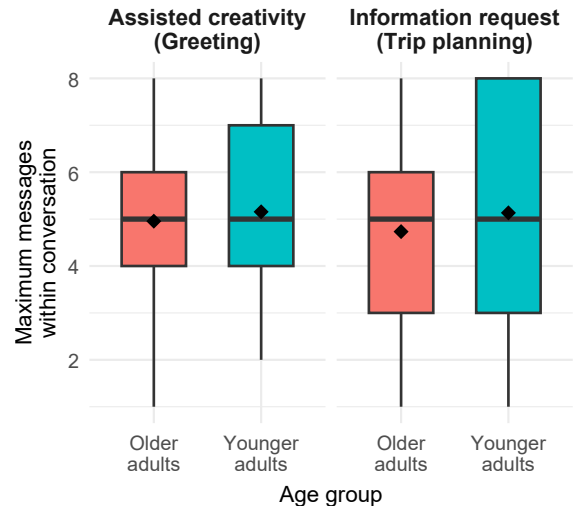


Figure 3: Maximum number of messages sent to GPT by age group and task type. Boxes represent interquartile ranges, horizontal lines indicate medians, and diamonds denote means.

the model to a comparable extent regardless of the interaction context.

Number of words sent to GPT A 2×2 mixed-design ANOVA was conducted to examine whether the total number of words sent to GPT differed by age group (younger adults vs. older adults; between-subjects) and task type (assisted creativity task [greeting writing] vs. information request task [trip planning]; within-subjects) (see Figure 4)

The main effect of age group was not significant, $F(1, 88) = 0.61, p = .44, \eta_g^2 = .005$, indicating that older adults ($M = 74.6, SE = 5.0$) and younger adults ($M = 69.1, SE = 5.0$) generated a similar number of words in their prompts. A significant main effect of task type was found, $F(1, 88) = 9.23, p = .003, \eta_g^2 = .013$, with participants in both age groups producing more words in the information request task ($M = 76.0, SE = 3.4$) than in the assisted creativity task ($M = 67.6, SE = 4.1$). The interaction between age group and task type was not significant, $F(1, 88) = 0.07, p = .79, \eta_g^2 = .000$, indicating that the amount of text sent to GPT was similar across both age groups and task types. These results indicate that word production was similar across age groups, while varying by task type, suggesting that both younger and older adults produced comparable amounts of text, with longer prompts generated in the information request task regardless of interaction context.

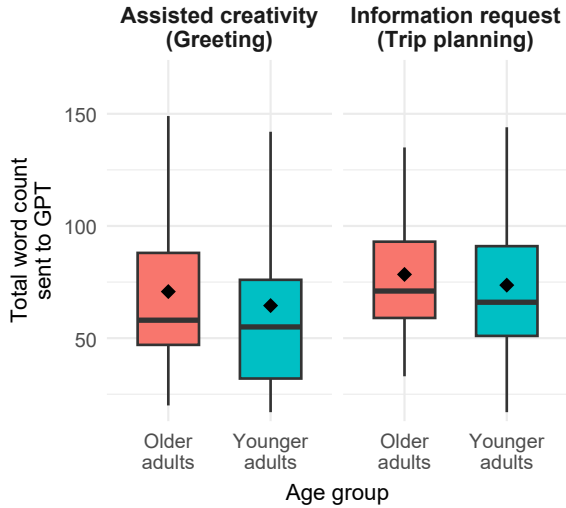


Figure 4: Total word count sent to GPT by age group and task type. Boxes represent interquartile ranges, horizontal lines indicate medians, and diamonds denote means.

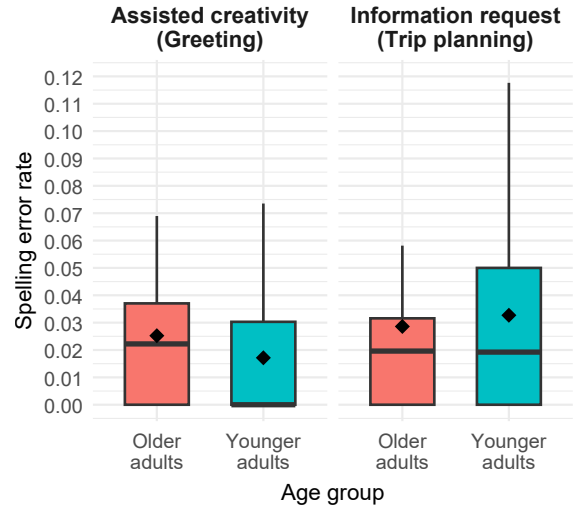


Figure 5: Spelling error rate in messages sent to GPT by age group and task type. Boxes represent interquartile ranges, horizontal lines indicate medians, and diamonds denote means.

Spelling Error Rate in Messages Sent to GPT
 Spelling error rate was operationalized as the number of misspellings per message, identified automatically using the English LanguageTool model (language-tool-python v3.2.1). Only matches explicitly labeled as misspellings were included, excluding suggestions related to grammar, style, or punctuation. For each message, the total number of detected misspellings was computed and used as a raw count measure of spelling errors. Messages containing no textual content were assigned a spelling error rate of zero.

A 2×2 mixed-design ANOVA was conducted to examine whether the spelling error rates (misspellings per message) sent to GPT differed by age group (younger adults vs. older adults; between-subjects) and task type (assisted creativity [greeting writing] vs. information request [trip planning]; within-subjects) (see Figure 5)

The main effect of age group was not significant, $F(1, 88) = 0.12, p = .73, \eta_g^2 = .000$, indicating that older adults ($M = 0.027, SE = 0.004$) and younger adults ($M = 0.025, SE = 0.004$) had a similar spelling error rate. A significant main effect of task type was found, $F(1, 88) = 6.91, p = .010, \eta_g^2 = .022$, indicating higher spelling error rates in the information request task ($M = 0.0307, SE = 0.00372$) than in the assisted creativity task ($M = 0.0212, SE = 0.00304$), averaged across age groups. The interaction between age group and task type did not reach statistical significance,

$F(1, 88) = 2.86, p = .094, \eta_g^2 = .009$, indicating that the spelling error rate sent to GPT was similar across both age groups and task types. These results indicate that spelling error rates were similar across age groups, while varying by task type, suggesting that both younger and older adults exhibited comparable levels of spelling errors, with higher error rates observed in the information request task regardless of interaction context.

C.2 GPT Output Characteristics

GPT response time A 2×2 mixed-design ANOVA was conducted to examine whether GPT response time differed by age group (younger adults vs. older adults; between-subjects) and task type (assisted creativity [greeting writing] vs. information request [trip planning]; within-subjects) (see Figure 6)

The main effect of age group was not significant, $F(1, 88) = 0.32, p = .58, \eta_g^2 = .002$, indicating that older adults ($M = 7.14, SE = 0.41$) and younger adults ($M = 6.81, SE = 0.41$) elicited similar response times from GPT. A significant main effect of task type was found, $F(1, 88) = 347.12, p < .001, \eta_g^2 = .524$, indicating that GPT responded more slowly in the information request task ($M = 10.35, SE = 0.43$) than in the assisted creativity task ($M = 3.59, SE = 0.23$), averaged across age groups. The interaction between age group and task type was not significant, $F(1, 88) = 0.00, p = .95, \eta_g^2 = .001$, indicat-

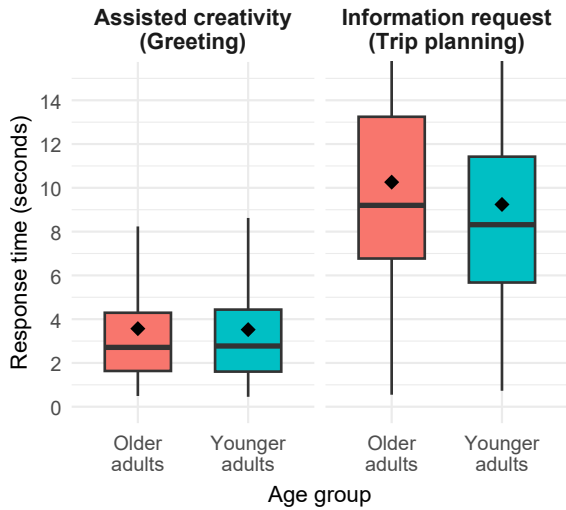


Figure 6: GPT response time (in seconds) by age group and task type. Boxes represent interquartile ranges, horizontal lines indicate medians, and diamonds denote means.

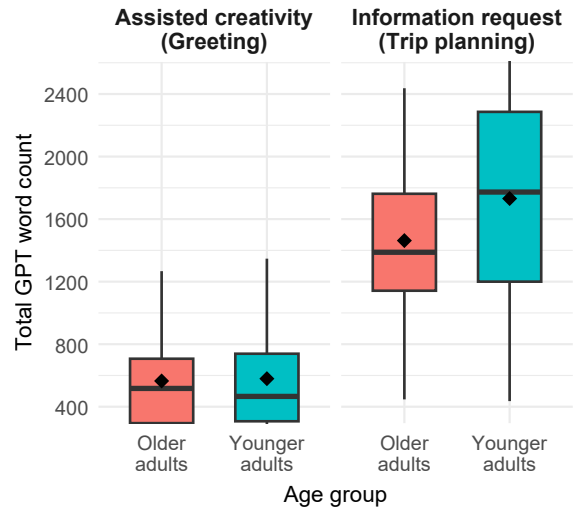


Figure 7: Total number of words generated by GPT by age group and task type. Boxes represent interquartile ranges, horizontal lines indicate medians, and diamonds denote means.

ing that the task-related increase in GPT response time was similar across both age groups and task types. These results indicate that GPT response times were similar across age groups but varied by task type. Specifically, responses were consistently slower in the information request task than in the assisted creativity task, irrespective of the user's age or interaction context.

Number of words generated by GPT A 2×2 mixed-design ANOVA was conducted to examine whether the total number of words generated by GPT differed by age group (younger adults vs. older adults; between-subjects) and task type (assisted creativity [greeting writing] vs. information request [trip planning]; within-subjects) (see Figure 7)

The main effect of age group was not significant, $F(1, 88) = 2.64, p = .11, \eta_g^2 = .020$, indicating that older adults ($M = 1,014, SE = 61.7$) and younger adults ($M = 1,155, SE = 61.7$) received similar amounts of text overall. A significant main effect of task type was found, $F(1, 88) = 287.43, p < .001, \eta_g^2 = .515$. GPT generated more words in the information request task ($M = 1,597, SE = 65.1$) than in the assisted creativity task ($M = 572, SE = 37.3$). The interaction between age group and task type was also significant, $F(1, 88) = 4.42, p = .038, \eta_g^2 = .016$. Follow-up comparisons indicated that the task-related increase in GPT output was large in both

age groups but more pronounced among younger adults (assisted creativity: $M = 579$; information request: $M = 1,731$) than among older adults (assisted creativity: $M = 565$; information request: $M = 1,463$). However, within each task, the two age groups did not differ significantly (all $ps > .17$). These results indicate that GPT's output length was primarily determined by task demands, with age exerting only a modest influence, reflected in the magnitude of the task effect rather than in absolute differences between age groups within each task.

Spelling Error Rate in Messages from GPT A 2×2 mixed-design ANOVA was conducted to examine whether spelling error rates in messages generated by GPT differed by age group (younger adults vs. older adults; between-subjects) and task type (assisted creativity [greeting writing] vs. information request [trip planning]; within-subjects) (see Figure 8)

The main effect of age group was not significant, $F(1, 88) = 1.45, p = .23, \eta_g^2 = .007$, indicating that older adults ($M = 0.012, SE = 0.001$) and younger adults ($M = 0.014, SE = 0.001$) had a similar spelling error rate in GPT-generated messages. A significant main effect of task type was found, $F(1, 88) = 235.83, p < .001, \eta_g^2 = .580$, indicating higher spelling error rates in GPT-generated messages in the information request task ($M = 0.025, SE = 0.001$) than in the assisted creativity task ($M = 0.001, SE = 0.000$). The in-

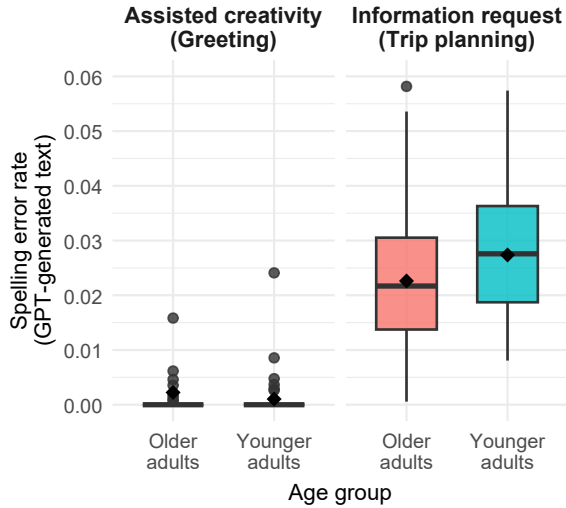


Figure 8: Spelling error rate in messages generated by GPT by age group and task type. Boxes represent interquartile ranges, horizontal lines indicate medians, and diamonds denote means.

1023 teraction between age group and task type was not
 1024 significant, $F(1, 88) = 3.83, p = .053, \eta_g^2 = .022$,
 1025 indicating that spelling errors in GPT-generated
 1026 text were similar by age groups and task types.
 1027 These results indicate that spelling errors in GPT-
 1028 generated text were primarily influenced by task de-
 1029 mands rather than participant age. Both age groups
 1030 exhibited similar increases in error rates in the in-
 1031 formation request task compared to the assisted
 1032 creativity task, regardless of interaction context.

1033 **Association Between Participant Input Length**
 1034 **and GPT Output Length** Following the identifi-
 1035 cation of differences in participant input and GPT
 1036 output across tasks and age groups, the relationship
 1037 between these two variables was systematically ex-
 1038 amined. Pearson correlations were calculated sepa-
 1039 rately for each age group to assess the relationship
 1040 between the number of words participants wrote to
 1041 GPT and the number of words GPT returned. Sig-
 1042 nificant positive associations were observed in both
 1043 groups: among younger adults, $r = .41, 95\% \text{ CI}$
 1044 $[.22, .57], p < .001$; among older adults, $r = .33,$
 1045 $95\% \text{ CI} [.13, .50], p = .002$. A Fisher's z test for
 1046 independent correlations indicated no significant
 1047 difference between the two coefficients, $z = 0.64,$
 1048 $p = .52$, and the confidence interval for the dif-
 1049 ference included zero ($-.17$ to $.34$). This finding
 1050 suggests that the association's strength was similar
 1051 across age groups.

1052 Correlations were also computed separately for

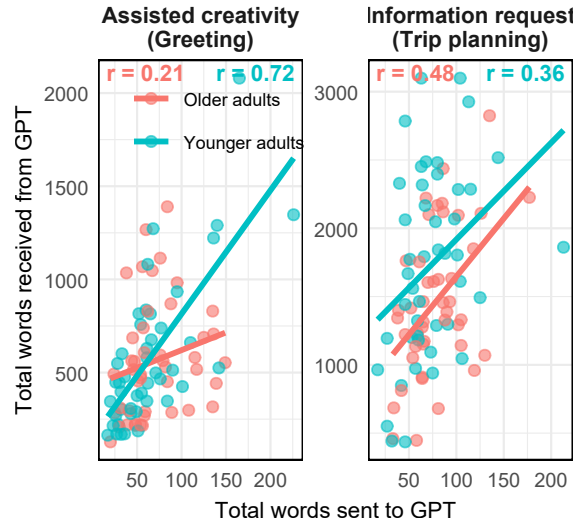


Figure 9: Relationship between the number of words participants sent to GPT and the number of words generated by GPT, shown separately by task type and age group. Lines represent linear fits within each group; correlation coefficients are reported for descriptive purposes.

1053 each task within each age group. For older adults,
 1054 the association between participant input length
 1055 and GPT output length was $r = .21$ in the as-
 1056 sisted creativity task [greeting writing] and $r = .48$
 1057 in the information request task [trip planning].
 1058 For younger adults, the corresponding correlations
 1059 were $r = .72$ in the assisted creativity task and
 1060 $r = .36$ in the information request task. As these
 1061 correlations were derived from repeated measure-
 1062 ments within individuals, no statistical compar-
 1063 isons were conducted among the four coefficients.
 1064 However, the descriptive pattern indicates that the
 1065 strength of the association between user input and
 1066 GPT output varies by task type.

1067 To further investigate whether GPT output length
 1068 varied as a function of participant input length, age
 1069 group, and task type, a linear mixed-effects regres-
 1070 sion model was fitted with participant included
 1071 as a random intercept. Participant input length
 1072 did not significantly predict GPT output length
 1073 overall, $\beta = 1.78, p = .36$. However, a signifi-
 1074 cant interaction between input length and task type
 1075 was observed, $\beta = 7.16, p = .008$, indicating
 1076 a stronger association between input and output
 1077 length in the information request task compared
 1078 to the assisted creativity task. No reliable interac-
 1079 tions involving age group and input length were
 1080 observed (all $ps > .07$). In contrast, a significant
 1081 interaction between age group and task type was

1162 as *one, much, sound, bit, travel, and like*. Additional distinctive words included *good, help, love,*
 1163 *need, card, and year*. Together, these words reflect a pattern of more reflective, evaluative, and relational language use in prompts produced by older participants.
 1164
 1165
 1166
 1167

1168 In contrast, younger adults were characterized by a different set of lexically distinctive words. The highest PPMI values were observed for action- and production-oriented terms such as *add, itinerari,*
 1169 *shorter, make, and write*. Other distinctive words included *messag, time, stay, day, best, place, funni,*
 1170 *visit, trip, and friend*, suggesting a stronger focus on task execution, content manipulation, and concrete planning in how younger participants formulated their prompts to the model.
 1171
 1172
 1173
 1174
 1175
 1176
 1177

1178 D.2 Lexical Distinctiveness in GPT Responses

1179 To examine whether the language model’s responses differed systematically by users’ age group, we conducted PPMI analyses on the model’s outputs. We first examined responses within each task type individually to identify context-specific adaptations, followed by an aggregated analysis across all interactions.
 1180
 1181
 1182
 1183
 1184
 1185

1186 **Information Request Task.** In the trip-planning task, the model’s responses exhibited clear age-related lexical distinctiveness (Figure 2).
 1187
 1188

1189 For older adults, the most lexically distinctive words in the model’s responses were dominated by transportation, infrastructure, and procedural travel information, including *London, Paris, train, airline, ticket, UK, station, platform, and journey*. Additional distinctive terms such as *current, save, communicate, and since*, further reflect an emphasis on logistical coordination and informational completeness. Together, these patterns suggest that when responding to older users’ trip-planning requests, the language model tended to generate more structured, transit-focused, and procedurally oriented information, emphasizing routes, transportation systems, and formal travel logistics.
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202

1203 In contrast, responses generated for younger adults were characterized by a different set of lexically distinctive words, reflecting experiential, social, and lifestyle-oriented content. The highest PPMI values were observed for terms such as *children, family-friendly, beach, boutique, café, snack, water, and dietary*, alongside leisure- and activity-related words such as *rest, late, and crowd*. Additional distinctive terms included *EUR and driver,*
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211

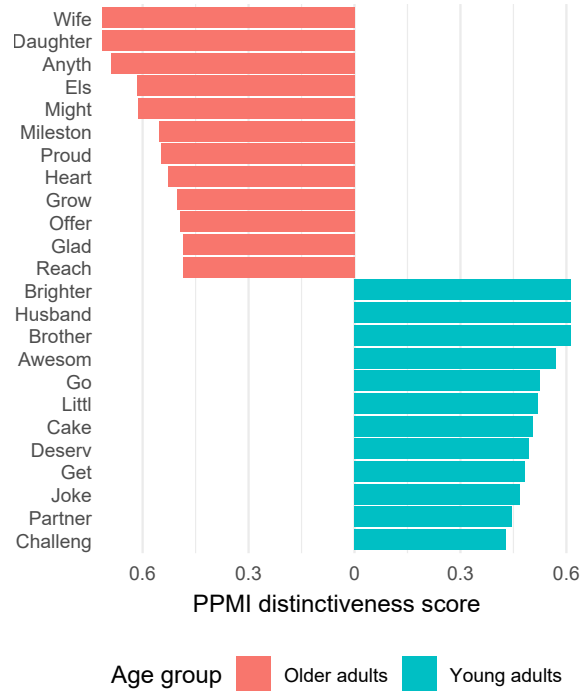


Figure 11: Top word stems per age group in the language model’s responses, ranked by PPMI, on the assisted creativity task (greeting writing).

1212 suggesting a focus on local experiences and practical on-the-ground considerations. These patterns indicate that, when addressing younger users, the language model more frequently emphasized experiential aspects of travel, amenities, and contextual lifestyle information, rather than formal transportation or infrastructural details.
 1213
 1214
 1215
 1216
 1217
 1218

1219 Taken together, these results suggest that age-related differences are reflected not only in how participants formulate prompts to the language model, but also in how the model structures its responses. Even when restricted to the same task type (information request) and controlling for overall word frequency, the model’s outputs show systematic lexical differences across age groups: responses to older adults emphasize logistical and transport-related information, while responses to younger adults emphasize experiential and lifestyle-oriented travel content.
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230

1231 **Assisted Creativity Task.** Age-related differences extended to the creative context (Figure 11).
 1232

1233 For older adults, the most lexically distinctive words in the model’s responses included *wife, daughter, proud, heart, milestone, and grow*, alongside evaluative and relational terms such as *glad, care, reach, and offer*. These words suggest that,
 1234
 1235
 1236
 1237

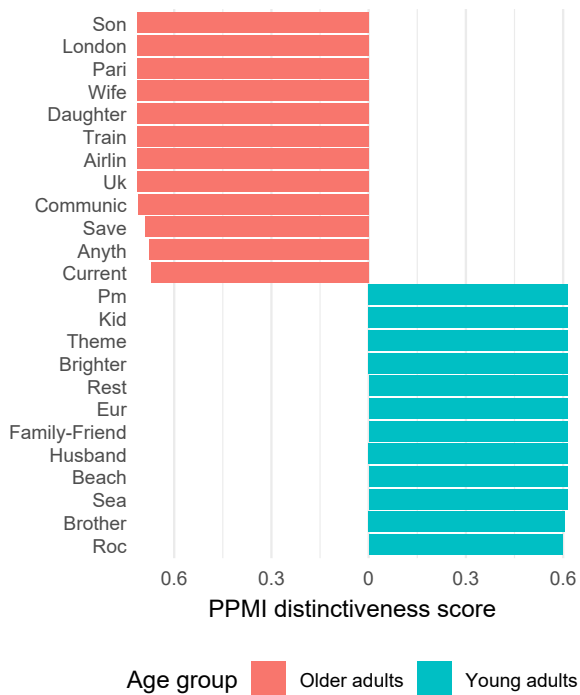


Figure 12: Top word stems per age group in the language model's responses, ranked by PPMI, pooled across task types.

when responding to older users in a creative context, the model tended to emphasize emotional significance, interpersonal relationships, and life-course framing.

In contrast, responses generated for younger adults were characterized by a different set of lexically distinctive words, including *brighter*, *awesome*, *joke*, *cake*, *little*, and *deserve*, as well as relational terms such as *husband*, *brother*, and *partner*. Additional distinctive words such as *go*, *get*, *now*, *make*, and *challenge* indicate a more action-oriented, upbeat, and immediate tone in responses to younger users. Together, these results indicate that age-related differences in the language model's outputs extend beyond informational content and persist in creative tasks. Specifically, the model appears to adapt not only the content but also the emotional framing and communicative tone of its responses to the user's age group.

Aggregated Analysis (Across Tasks). Finally, to assess stable response strategies independent of task content, we analyzed lexical distinctiveness across the entire corpus (Figure 12).

For responses addressed to older adults, GPT exhibited higher lexical distinctiveness for words related to formal planning, transportation infrastruc-

ture, and institutional contexts, including *son*, *London*, *Paris*, *train*, *airline*, *ticket*, *station*, and *platform*. In addition, several family-related terms (e.g., *wife*, *daughter*) and functionally oriented words (e.g., *communicate*, *current*, *save*) were more strongly associated with responses to older users. Together, these patterns suggest that GPT adopts a more structured, informational, and service-oriented style when generating responses for older adults.

In contrast, responses to younger adults were characterized by lexically distinctive words reflecting experiential, lifestyle-oriented, and informal content. These included terms related to leisure and environment (e.g., *beach*, *sea*, *rest*), social relations (*friend*, *brother*, *kid*), consumption and mobility (*euro*, *boutique*, *driver*), and affective or stylistic framing (e.g., *brighter*, *theme*). This pattern indicates that GPT's responses to younger users emphasize experiential framing and flexible, lifestyle-focused information.

Notably, these age-related lexical differences emerged despite collapsing across task types, indicating that they cannot be attributed solely to differences in task demands or topical content. Instead, the results suggest that the language model systematically adapts its lexical choices based on user age, reflecting age-sensitive response strategies in addition to differences observed in users' prompting behavior.

E Self-Evaluation of Model Responses

To examine how participants subjectively evaluated the quality of the model's responses, we analyzed self-reported ratings along four evaluation dimensions: relevance, usefulness, clarity, and sufficiency. Each dimension was analyzed using a 2×2 mixed-design ANOVA with age group (younger adults, older adults) as a between-subjects factor and task type (assisted creativity [greeting writing] vs. information request [trip planning]) as a within-subjects factor. This approach enabled us to identify potential interactions between task demands and age-related evaluation patterns, thereby allowing us to assess overall perceived response quality.

Perceived Relevance Relevance was first examined by assessing participants' self-reported extent to which the model's responses addressed their requests. The main effect of age group was not significant, $F(1, 88) = 0.19$, $p = .664$, $\eta_g^2 = .001$, indi-

1314 cating that older adults ($M = 4.38, SE = 0.09$)
1315 and younger adults ($M = 4.32, SE = 0.09$)
1316 provided similar relevance ratings overall. The
1317 main effect of task type was also not significant,
1318 $F(1, 88) = 3.04, p = .085, \eta_g^2 = .015$, with simi-
1319 lar relevance ratings for the assisted creativity task
1320 ($M = 4.44, SE = 0.08$) and the information re-
1321 quest task ($M = 4.26, SE = 0.09$).

1322 The interaction between age group and task
1323 type was significant, $F(1, 88) = 4.64, p = .034,$
1324 $\eta_g^2 = .022$. Follow-up comparisons showed that
1325 the effect of task type on relevance ratings differed
1326 by age group. Among younger adults, relevance
1327 ratings were significantly higher in the assisted cre-
1328 ativity task ($M = 4.53, SE = 0.11$) than in the
1329 information request task ($M = 4.11, SE = 0.12$),
1330 $t(88) = 2.76, p = .007$. In contrast, older adults
1331 showed no significant difference in relevance rat-
1332 ings between the assisted creativity ($M = 4.36,$
1333 $SE = 0.11$) and information request ($M = 4.40,$
1334 $SE = 0.12$) tasks ($p = .77$).

1335 These results indicate that relevance evaluations
1336 depended on the interaction between age and task
1337 demands. Younger adults were more sensitive to
1338 task type, with higher relevance ratings for assisted
1339 creativity. Older adults gave consistently high rele-
1340 vance ratings across tasks. As a result, older adults'
1341 relevance evaluations in the information request
1342 task were higher than those of younger adults.

1343 **Perceived Usefulness** Usefulness was then exam-
1344 ined by assessing participants' self-reported extent
1345 to which the model's responses were for their in-
1346 tended purpose. The main effect of age group was
1347 also not significant, $F(1, 88) = 1.10, p = .298,$
1348 $\eta_g^2 = .008$, as older adults reported usefulness rat-
1349 ings ($M = 4.27, SE = 0.11$) similar to those of
1350 younger adults ($M = 4.10, SE = 0.11$). The main
1351 effect of task type was not significant, $F(1, 88) =$
1352 $0.08, p = .781, \eta_g^2 = .000$, indicating that useful-
1353 ness ratings were highly similar for the assisted
1354 creativity task ($M = 4.20, SE = 0.10$) and the
1355 information request task ($M = 4.17, SE = 0.10$)
1356 by age groups. The interaction between age group
1357 and task type was not significant, $F(1, 88) = 1.04,$
1358 $p = .310, \eta_g^2 = .004$, indicating a similar pattern
1359 of task-related usefulness ratings across the two age
1360 groups. These results indicate that the usefulness
1361 evaluation of the model's responses was similar
1362 across age groups and task types, suggesting that
1363 both younger and older adults found the model's
1364 output equally useful regardless of the interaction

context.

1365
1366 **Perceived Clarity** Clarity was next examined
1367 by assessing participants' self-reported extent to
1368 which the model's responses were clear and un-
1369 derstandable. The analysis identified a signifi-
1370 cant main effect of age group, $F(1, 88) = 7.13,$
1371 $p = .009, \eta_g^2 = .050$, such that older adults as-
1372 signed higher clarity ratings overall ($M = 4.60,$
1373 $SE = 0.09$) compared to younger adults ($M =$
1374 $4.24, SE = 0.09$). A significant main effect of
1375 task type was found, $F(1, 88) = 18.15, p < .001,$
1376 $\eta_g^2 = .068$, with the assisted creativity task receiv-
1377 ing higher clarity ratings ($M = 4.63, SE = 0.06$)
1378 than the information request task ($M = 4.21,$
1379 $SE = 0.10$) by both age groups. The interaction
1380 between age group and task type was not signifi-
1381 cant, $F(1, 88) = 2.46, p = .12, \eta_g^2 = .009$, indi-
1382 cating that task type did not affect clarity ratings
1383 across age groups. These results indicate that clar-
1384 ity evaluations were influenced by both age and
1385 task demands: older adults rated the model's re-
1386 sponses as clearer than younger adults did. Across
1387 age groups, responses in the assisted creativity task
1388 were rated as clearer than those in the information
1389 request task, regardless of interaction context.

1390 **Perceived Sufficiency** Sufficiency was subse-
1391 quently examined by assessing participants' self-
1392 report of the adequacy of the model's responses,
1393 assuming no additional responses were provided.
1394 The main effect of age group was not signifi-
1395 cant, $F(1, 88) = 0.44, p = .508, \eta_g^2 = .003,$
1396 with similar overall sufficiency ratings for older
1397 adults ($M = 4.03, SE = 0.13$) and younger
1398 adults ($M = 3.91, SE = 0.13$). The main ef-
1399 fect of task type was significant, $F(1, 88) = 5.82,$
1400 $p = .018, \eta_g^2 = .024$, indicating higher sufficiency
1401 ratings for the assisted creativity task ($M = 4.14,$
1402 $SE = 0.11$) than for the information request task
1403 ($M = 3.80, SE = 0.12$) by age groups. The in-
1404 teraction between age group and task type was not
1405 significant, $F(1, 88) = 0.49, p = .485, \eta_g^2 = .002,$
1406 indicating that the difference in sufficiency ratings
1407 between tasks was consistent across age groups.
1408 These results indicate that sufficiency evaluations
1409 were influenced by task demands rather than age,
1410 with responses perceived as more sufficient in the
1411 assisted creativity task than in the information re-
1412 quest task across both age groups, regardless of the
1413 interaction context.

1414 Across the four evaluation dimensions, age-
1415 related differences appeared selectively. Older

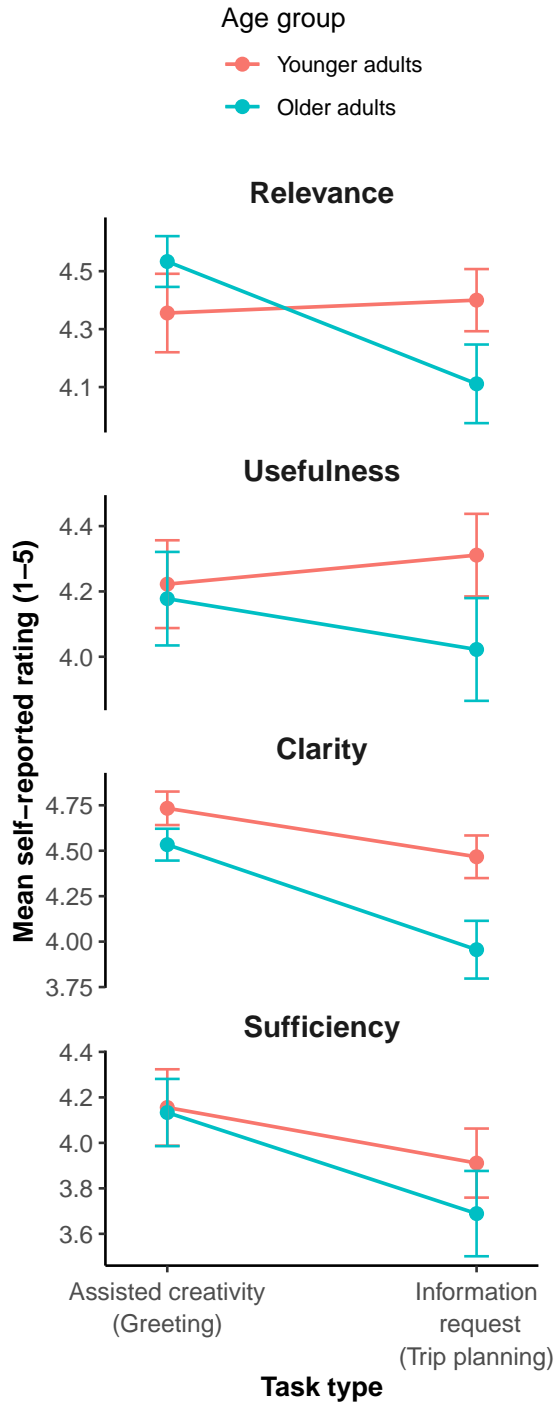


Figure 13: Age-related differences in self-reported evaluations of model responses across task types (Assisted creativity vs. Information request). Error bars represent standard errors.

adults rated the model’s responses as clearer than those of younger adults. Differences in relevance evaluations were also observed: older adults maintained high relevance ratings across tasks, whereas younger adults’ ratings varied by task type. No age-related differences were observed in usefulness or

sufficiency, indicating similar perceptions among younger and older adults. However, sufficiency ratings demonstrated a main effect of task type, with responses in the assisted creativity task perceived as more sufficient than those in the information request task for both age groups. These findings indicate that task demands influenced perceptions of response completeness, regardless of age.

F LLM-Judged Evaluation of Model Responses

Differences in AI performance across tasks and age groups were evaluated using an LLM-as-a-Judge framework, in which three independent language models rated each interaction across six evaluation dimensions: Relevance, Usefulness, Clarity, Sufficiency, Accuracy, and Conversational Flow. Ratings were averaged across models to yield a single score per participant, task, and dimension. Full methodological details are provided in Appendix B.2.

To determine whether these evaluations varied by task type and age group, separate 2×2 mixed-design ANOVAs were conducted for each evaluation dimension, with age group (younger adults vs. older adults) as the between-subjects factor and task type (assisted creativity [greeting writing] vs. information request [trip planning]) as the within-subjects factor.

LLM-judged Relevance Score Relevance Score, reflecting the extent to which the AI’s responses addressed users’ requests. The main effect of age group was not significant, $F(1, 88) = 0.31$, $p = .578$, $\eta_g^2 = .002$, indicating that conversations initiated by older adults ($M = 4.38$, $SE = 0.09$) and younger adults ($M = 4.32$, $SE = 0.09$) were judged as similarly relevant overall. A significant main effect of task type was found, $F(1, 88) = 9.41$, $p = .003$, $\eta_g^2 = .039$, indicating that AI responses were judged as more relevant in the assisted creativity task ($M = 4.72$, $SE = 0.05$) than in the information request task ($M = 4.48$, $SE = 0.07$) (see Figure 14). The interaction between age group and task type was not significant, $F(1, 88) = 0.64$, $p = .424$, $\eta_g^2 = .003$, indicating that the difference in Relevance ratings between tasks was similar by age groups. These results indicate that LLM judges evaluated the AI’s responses as more relevant as a function of task demands rather than age, with higher relevance ratings for the assisted creativity task across both age groups,

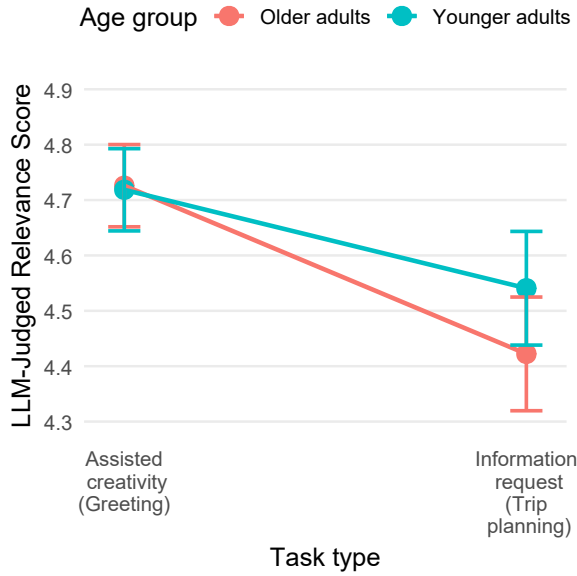


Figure 14: LLM-judged relevance scores by task type and age group. Points represent estimated marginal means; error bars indicate ± 1 SE.

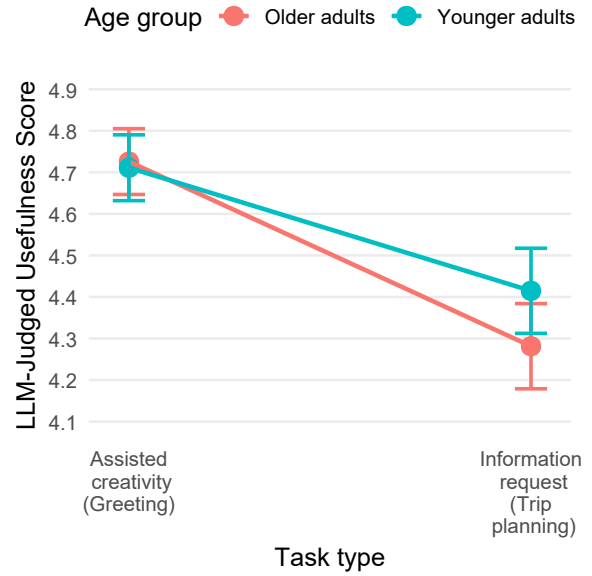


Figure 15: LLM-judged usefulness scores by task type and age group. Points represent estimated marginal means; error bars indicate ± 1 SE.

regardless of the interaction context.

LLM-judged Usefulness Score Usefulness Score, reflecting the extent to which the AI's responses were perceived as useful for the user's purpose. The main effect of age group was not significant, $F(1, 88) = 0.31, p = .581, \eta_g^2 = .002$, indicating that conversations initiated by older adults ($M = 4.50, SE = 0.08$) and younger adults ($M = 4.56, SE = 0.08$) were judged as similarly useful overall. A significant main effect of task type was found, $F(1, 88) = 25.84, p < .001, \eta_g^2 = .085$, indicating that AI responses were judged as more useful in the assisted creativity task ($M = 4.72, SE = 0.06$) than in the information request task ($M = 4.35, SE = 0.07$) (see Figure 15). The interaction between age group and task type was not significant, $F(1, 88) = 1.03, p = .312, \eta_g^2 = .004$, indicating that the difference in usefulness ratings between tasks was similar by age groups. These results indicate that LLM judges evaluated the AI's responses as more useful as a function of task demands rather than age, with higher usefulness ratings for the assisted creativity task than in the information request task across both age groups, regardless of the interaction context.

LLM-judged Clarity Score Clarity Score, reflecting how clear and understandable the AI's responses were judged to be. The main effect of

age group was not significant, $F(1, 88) = 1.49, p = .225, \eta_g^2 = .009$, indicating that conversations initiated by older adults ($M = 4.85, SE = 0.03$) and younger adults ($M = 4.90, SE = 0.03$) were judged as equally clear overall. The main effect of task type was also not significant, $F(1, 88) = 1.74, p = .191, \eta_g^2 = .009$, indicating that conversations initiated by the creativity task ($M = 4.72, SE = 0.06$) and information request task ($M = 4.35, SE = 0.07$) were judged as equally clear overall. The interaction between age group and task type was not significant, $F(1, 88) = 1.28, p = .261, \eta_g^2 = .007$, indicating that the difference in clarity ratings between tasks was similar by age groups. These findings indicate that LLM judges perceived the AI's responses as highly clear across both tasks and age groups, with no evidence for systematic differences attributable to task type or user age.

LLM-judged Sufficiency Score Sufficiency Score, whether the AI's responses would be sufficient for users to complete the task if no additional information were provided. The main effect of age group was not significant, $F(1, 88) = 1.23, p = .271, \eta_g^2 = .009$, indicating that conversations initiated by older adults ($M = 4.13, SE = 0.08$) and younger adults ($M = 4.25, SE = 0.08$) were judged as similarly sufficient overall. A significant main effect of task type was found, $F(1, 88) = 144.72, p < .001, \eta_g^2 = .377$, in-

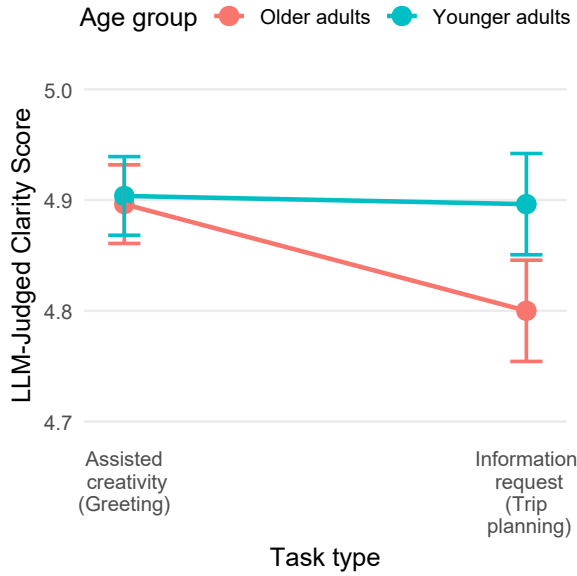


Figure 16: LLM-judged clarity scores by task type and age group. Points represent estimated marginal means; error bars indicate ± 1 SE

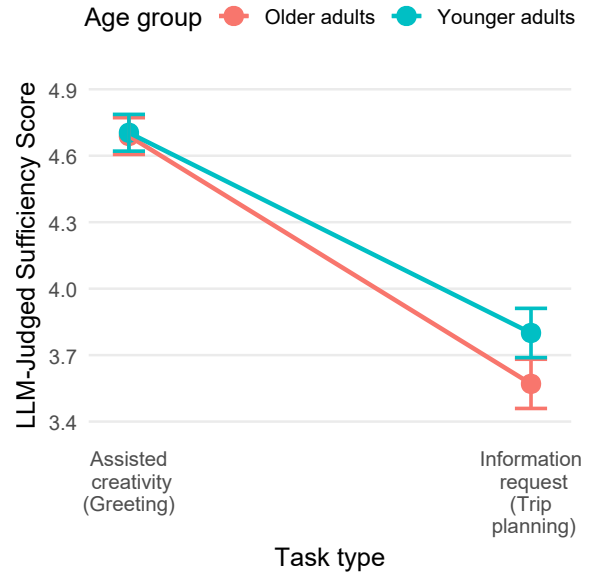


Figure 17: LLM-judged sufficiency scores by task type and age group. Points represent estimated marginal means; error bars indicate ± 1 SE.

1530 dicating that AI responses were judged as more
 1531 sufficient in the assisted creativity task ($M = 4.70$,
 1532 $SE = 0.06$) than in the information request task
 1533 ($M = 3.69$, $SE = 0.08$) (see Figure 17). The in-
 1534 teraction between age group and task type was not
 1535 significant, $F(1, 88) = 1.63$, $p = .205$, $\eta_g^2 = .007$,
 1536 indicating that the difference in sufficiency ratings
 1537 between tasks was similar by age groups. These
 1538 results indicate that LLM judges evaluated the AI's
 1539 responses as more sufficient as a function of task
 1540 demands rather than age, with higher sufficient
 1541 ratings for the assisted creativity task than in the
 1542 information request task across both age groups,
 1543 regardless of the interaction context.

1544 **LLM-judged Accuracy Score** Accuracy Score,
 1545 reflecting the extent to which the AI's responses
 1546 were judged to be factually correct and free from
 1547 errors or misleading information. The main effect
 1548 of age group was not significant, $F(1, 88) = 0.99$,
 1549 $p = .322$, $\eta_g^2 = .007$, indicating that conversations
 1550 initiated by older adults ($M = 4.72$, $SE = 0.06$)
 1551 and younger adults ($M = 4.64$, $SE = 0.06$) were
 1552 judged as similarly accurate overall. A significant
 1553 main effect of task type was found, $F(1, 88) =$
 1554 19.33 , $p < .001$, $\eta_g^2 = .070$, indicating that AI
 1555 responses were judged as more accurate in the as-
 1556 sisted creativity task ($M = 4.80$, $SE = 0.04$)
 1557 than in the information request task ($M = 4.56$,
 1558 $SE = 0.06$) (see Figure 18). The interaction be-

1559 tween age group and task type was not significant,
 1560 $F(1, 88) = 0.04$, $p = .844$, $\eta_g^2 < .001$, indicat-
 1561 ing that the difference in accuracy ratings between
 1562 tasks was similar by age groups. These results indi-
 1563 cate that LLM judges evaluated the AI's responses
 1564 as more accurate as a function of task demands
 1565 rather than age, with higher accurate ratings for
 1566 the assisted creativity task than in the informa-
 1567 tion request task across both age groups, regard-
 1568 less of the interaction context.

1569 **LLM-judged Conversational Flow Score** Con-
 1570 versational Flow Score, reflecting how smooth, co-
 1571 herent, and natural the interaction was across con-
 1572 versational turns. The main effect of age group
 1573 was not significant, $F(1, 88) = 0.32$, $p = .576$,
 1574 $\eta_g^2 = .002$, indicating that conversations initiated
 1575 by older adults ($M = 4.60$, $SE = 0.07$) and
 1576 younger adults ($M = 4.55$, $SE = 0.07$) were
 1577 judged as having similar conversational flow over-
 1578 all. The main effect of task type was also not signif-
 1579 icant, $F(1, 88) = 3.46$, $p = .066$, $\eta_g^2 = .017$, indi-
 1580 cating that conversations initiated by the creativity
 1581 task ($M = 4.66$, $SE = 0.07$) and informa-
 1582 tion request task ($M = 4.49$, $SE = 0.07$) were
 1583 judged as having similar conversational flow overall.
 1584 The interaction between age group and task type
 1585 was not significant, $F(1, 88) = 0.75$, $p = .387$,
 1586 $\eta_g^2 = .004$, indicating that the difference in con-
 1587 versational flow ratings between tasks was similar
 1588 by age groups.

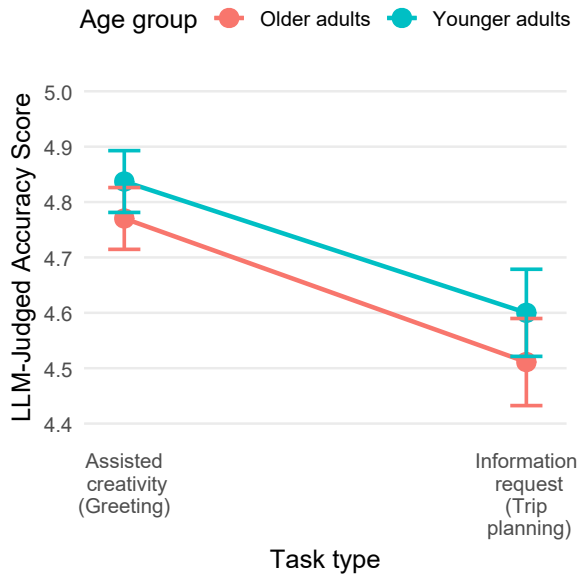


Figure 18: LLM-judged accuracy scores by task type and age group. Points represent estimated marginal means; error bars indicate ± 1 SE.

1588 These findings indicate that LLM judges perceived
 1589 the AI's responses as having high conversational
 1590 flow across both tasks and age groups, with no evi-
 1591 dence for systematic differences attributable to task
 1592 type or user age.

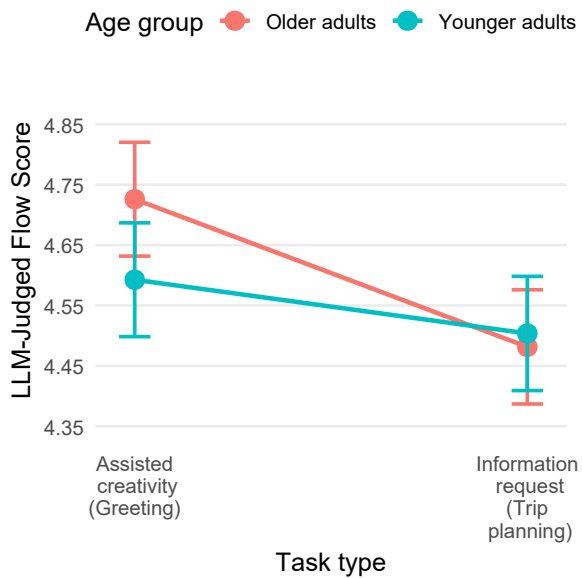


Figure 19: LLM-judged conversational flow scores by task type and age group. Points represent estimated marginal means; error bars indicate ± 1 SE.