



Human vs. AI: Exploring Trust in Predictions with Rationales

PHILIPPE DE BEKKER, BO VAN DEN BERG, MAYA ELASMAR, SAM HESLENFELD, and ILIAS PAPADIMITRIOU

Delft University of Technology, The Netherlands

Abstract. Artificial Intelligence (AI) has been evolving rapidly over the last few years, and it is being applied to a wide range of applications and domains. Trust plays a critical role in decision-making processes. This paper explores and compares the extent to which people trust humans versus AI when considering predictions with rationales. This study adopts a two-stage approach. In the first stage, crowd workers and AI are tasked with predicting the price of a house. In the second stage, crowd workers are asked to indicate their amount of trust in the price advice provided by humans and AI in addition to a final verdict, i.e. human, AI or no preference. A diverse sample of individuals from various demographic backgrounds and professional fields is recruited to ensure a representative dataset. Ultimately, the results of the last stage indicate no preference without rationales, preference for the rationale if only one is provided and most interestingly, AI when both provide rationales. A hypothesis (requiring future work) is the possibility of AI being trusted more due to having more digestible and convincing rationales compared to the human rationales, even though the accuracy was actually lower.

CCS Concepts: • Human-centered computing; • Computing methodologies → Artificial Intelligence;

Additional Key Words and Phrases: house price prediction, rationales, crowd computing, artificial intelligence, image context, trust

ACM Reference Format:

Philippe de Bekker, Bo van den Berg, Maya Elasmar, Sam Heslenfeld, and Ilias Papadimitriou. 2023. Human vs. AI: Exploring Trust in Predictions with Rationales. In *CS4145-Q4-2023: Crowd Computing, TU Delft*. 13 pages.

1 INTRODUCTION

A lot of significant advancements have occurred in the last years in the field of Artificial Intelligence (AI) – AI systems generate more accurate and coherent responses, making them viable sources of advice and information [1–4]. They have been widely adopted in various aspects of our lives and process vast amounts of data in very little time. On the other hand, humans can easily capture contextual information and are exceptionally wired in this regard [5–9].

Seeking advice is closely related to trust. Advice from people can be trusted for various reasons, such as the ability to empathize or to easily perceive human reasoning. On the other hand, AI can perform complex calculations, taking into account patterns and factors that humans cannot analyze, but are still widely considered as black boxes, which cannot provide accurate reasoning that justifies their prediction. Therefore, it is extremely valuable to compare the trust of people in AI when predictions include rationales.

The task for which trust will be evaluated is also crucial. House price prediction is a difficult and non-trivial question for both humans and AI. It depends on many factors, some of which are interdependent and non-linear. This task captures the human's ability to perceive situational context, as most people can differentiate between poor and wealthy homes and they can recognize luxurious interior design.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Association for Computing Machinery.

Manuscript submitted to ACM

53 This problem is addressed by combining crowd computing with AI. The central research question in this paper is
 54 formulated as follows:
 55

56 *"Do people have more trust in humans or AI when considering predictions with rationales?"*

57 Resulting from this research question, three hypotheses have been defined:
 58

- 59 H 1. A rationale increases trust towards its source. Advice including a rationale will be preferred over one without it.
- 60 H 2. If both human and AI do **not** provide rationales, then the human will be trusted more.
- 61 H 3. If both human and AI do provide rationales, then the human will be trusted more.
 62

63 2 RELATED WORK

64 Various research already exists in the domain of exploring trust in AI technology. These works provide some valuable
 65 insights and inspiration for this research. The following categories of related work are identified:
 66

67 2.1 Trust in Artificial Intelligence

68 Trust in artificial intelligence is a complex issue that depends on many factors, such as user education, past experiences,
 69 user biases, and perception towards AI [10]. Previous research has shown that trust in artificial intelligence does not
 70 only depend on the reliability of the results it provides, but also tangibility and transparency [11]. In particular, people
 71 resist the black-box nature of AI [12–14]. Hence, AI is a highly intriguing topic that requires exploration across various
 72 genres of tasks. It is important to investigate the generalizability of each task, i.e. whether or not the task is suitable to
 73 extract information from to assess the trust in AI advice. Characteristics of the decision task, such as task risk, required
 74 expertise, and decision objectivity, can determine the task's generalizability [15]. Some sources report that humans
 75 tend to utilize the advice from AI more compared to advice from peers, though other findings claim a contrasting
 76 view - it boils down to the stakes involved in the operation. House price prediction is a task that belongs to the sector
 77 of business and finance, and therefore is considered a high-stakes task. In this task, there exists a ground truth and
 78 therefore quantitative measures of performance are suitable, which enables the objective evaluation of the human and
 79 AI predictions. Since house price prediction is considered an objective prediction task of high stakes, it makes for a
 80 highly generalizable task. Subsequently, the results of comparing human and AI predictions with rationales can provide
 81 a valuable contribution to understanding trust in this domain.
 82

83 2.2 Providing Rationales

84 The rationale was first introduced as a way to increase the performance of machine learning for text classification [16].
 85 It has been utilized in various annotation tasks, such as information retrieval relevance, sentiment analysis, or text
 86 classification. The rationale has proven to encourage more thoughtful decision-making and enable crowd verification [17].
 87 It has also been proved to discourage cheating and provide additional domain value [18]. In this research, the additional
 88 domain value is the information that mentions which criteria were important for the crowd worker in order to estimate
 89 the price of a property and subsequently serves for great context and verification of the decision.
 90

91 2.3 Framework of Similar Experiments

92 Several prior studies have explored the trustworthiness of advice from both human and AI sources [19–21]. A two-stage
 93 process is notably the most prevalent framework to compare trust. The first stage includes predictions made by both
 94 human and AI, and the second stage involves a verdict given by people to determine which advice is more trusted.
 95

105 Unlike other studies, this study also explores experimenting with randomly in- and excluding rationales for the human
106 and AI in order to gather more statistical evidence and evaluate whether any prior belief without rationales is skewing
107 the distribution of verdicts with all rationales included.
108

109 3 SYSTEM DESIGN

111 One of the main objectives of our task was to ensure an appropriate comparison between a human and an AI. For
112 example, the same input needs to be processed by both of them. To summarize all these requirements, a list of functional
113 and non-functional requirements is provided below. Recall that the system is a twofold of prediction and comparison as
114 explored in Section 2.3 and elaborated on more in-depth when discussing the experimental setup in Section 4. Ultimately,
115 these requirements form a blueprint for the system that will be deployed during the experiment.
116

117 3.1 Functional Requirements

- 120 FR 1. The model representing the AI must be able to process images and answer questions about them.
- 121 FR 2. When providing crowd workers with input to form a prediction, additional context of the area from the property
122 should be presented by providing a map with the surrounding area highlighted.
- 123 FR 3. When comparing human and AI predictions, the Graphical User Interface (GUI) must show the crowd workers
124 the complete input first and then both human and AI outputs in randomized order to prevent any bias.
- 125 FR 4. Respondents should be able to provide feedback on the helpfulness of the human and AI predictions.
- 126 FR 5. The comparison task should allow respondents to pick a preference between a human and an AI prediction.

127 3.2 Non-Functional Requirements

- 131 NFR 1. During prediction, both human and AI need to receive the same attributes and require to answer in a similar
132 format or select a fixed range of options.
- 133 NFR 2. The human and AI advice (prediction + rationale) should be checked for accuracy and reliability.
- 134 NFR 3. Each rationale should be free from errors and misleading information.
- 135 NFR 4. The crowdsourcing platform should allow worker selection based on selected qualities.
- 136 NFR 5. The crowd-working task must contain an attention check and require certain properties such as completion
137 time and the number of characters in open responses to be above a certain limit.
- 138 NFR 6. The crowd workers must be paid a fair hour salary.
- 139 NFR 7. The crowd workers should remain anonymous to ensure their privacy and keep results unbiased.
- 140 NFR 8. The crowd workers must be able to understand English proficiently.
- 141 NFR 9. The GUI should be user-friendly, intuitive, and visually appealing.

142 4 EXPERIMENTAL SETUP AND EVALUATION

143 The experiment was deployed using a two-stage approach in order to explore trust in predictions with rationales. Both
144 stages were performed by crowd workers using Toloka surveys [22]. In addition, the second stage required AI input
145 which was prompted via MiniGPT-4 [23] (see Appendix B). This tool is free, however, using OpenAI's GPT-4 [1] for
146 prompting might be more promising. The used dataset, obtained via Ahmed and Moustafa [24], contains 535 individual
147 houses. Each house is described by four images (frontal view, kitchen, bathroom, bedroom), as depicted in Figure 1,
148 and five textual attributes (number of bedrooms, number of bathrooms, surface area in square feet, zip code, price).
149 Generating the AI advice with the free tool (Section 4.3) was time-consuming, so the choice was eventually made to
150

157 continue with only 21 houses. This reduction in house quantity was compensated by adding more overlap in both
 158 stages. Moreover, the prices for each task suite were calculated using the US minimum salary of \$6 per hour (as advised
 159 by Toloka), resulting in a total experiment cost of \$141.
 160

161 The remaining part of this section discusses the data preparation (Section 4.1) needed to enable the first stage, divided
 162 by the human prediction (Section 4.2) and AI prediction (Section 4.3), and the second stage (Section 4.4). As quality
 163 control is universal throughout these sections, it is discussed separately in Section 4.5. An example of each normal and
 164 control task can be found in Appendix A.
 165



177 Fig. 1. Illustrative example of images provided per house in the dataset [24]. Displays a frontal view, kitchen, bathroom and bedroom.
 178

180 4.1 Data Preparation

182 To ensure the data used in this experiment was suitable for task deployment, analysis and comparison, the data
 183 underwent several modifications via the use of scripts written in Python.

184 As the dataset originates from 2016, the prices need to be aligned to the respective economic landscape of the
 185 respondent in order to properly evaluate the predictions against the ground truth and enable fair comparisons between
 186 human and AI predictions. Human knowledge translates to current time, however, the underlying AI training data from
 187 LAION is not up to date and has knowledge up until 2021 [25]. Utilising a web tool built around statistics from the U.S.
 188 Bureau of Labor Statistic allows for adjustment using inflation rates and a house price index [26]. The data indicates an
 189 adjustment factor for the human perspective of 1.30 from 2016 to 2023. For the AI perspective, factor 1.26 from 2016 to
 190 2022 resulting in approximately 1.057 from 2022 to 2023.
 191

192 In order to enhance the contextual understanding of the dataset for humans (especially crowd workers from all
 193 across the world), several measurements were taken. A Python zip-code package *uszipcode* based on US governmental
 194 data allowed for highlighting the respective areas on an interactive map provided in the Toloka survey (see Figure 6),
 195 facilitating easy exploration and visualization of the geographical distribution of the houses [27–29]. For example,
 196 users can check whether the house is located in an expensive part of the city or in a distant neighbourhood with no
 197 supermarkets nearby. In addition, some training data is also provided before doing any predictions to form a better
 198 mental model of the data, which also serves as quality control as mentioned in Section 4.5.
 199

200 Furthermore, the surface area is only provided in square feet. Catering to diverse user preferences and understanding,
 201 the surface area attribute is also converted from square feet to square meters - a widely recognized format - in order to
 202 show both.
 203

204 Lastly, as covered in the upcoming sections, the data needed to be translated into Toloka survey tasks built on JSON
 205 data, input prompts for MiniGPT-4 and results needed to be post-processed.
 206

209 4.2 Stage 1a: House Price Estimation by Humans

210 In the first stage, crowd workers were asked to provide an estimation of the price of the house based on pictures of
211 the frontal view, the bathroom, the bedroom, and the kitchen. They were also provided with numerical information
212 about the number of bedrooms and bathrooms, and the location of the zip code in a map. From this map, someone can
213 deduct information about whether the house is in a rural or urban area and if any significant landmarks are nearby. The
214 estimations of the price are presented as predefined ranges, such as \$200K - \$300K. The crowd workers are also asked
215 to provide a short rationale of 2-3 sentences, in which they justify the reason for which they provide this price range.
216 Each stage 1 task had an overlap of 3 reviewers. Before the task begins, two examples are given where the price of the
217 apartment is shown, in order to familiarize the crowd worker with both the user interface and some example prices.
218

222 4.3 Stage 1b: House Price Estimation by Artificial Intelligence

223 In order to calculate the AI estimations, we used the MiniGPT-4 API [23]. MiniGPT-4 is an AI tool for image processing,
224 which can then be used to then extract information from it. After processing the images of a house, this API was used
225 for two purposes, the first being the prediction of the house price. This prediction was asked to be made, based on the
226 images and properties of the house, and in the same price ranges that we included in the first stage of the experiment.
227 In addition, the API was used to generate a rationale for the predicted price range. This rationale should make the
228 AI reasoning more clear and more understandable, as this reasoning is often considered to be a black box. The exact
229 prompt that has been used to collect the MiniGPT-4 results, can be found in Appendix B. The generated house price
230 estimates already account for inflation. The tool also considers prices of similar houses in the zip code.
231

235 4.4 Stage 2: Advice Evaluation

236 For the second stage of the experiment, crowd workers were asked to compare the estimations provided by humans
237 and AI. They were provided with the information that is provided in the first stage, as well as the human estimation
238 and rationale and the AI estimation. In order to measure the rationale's impact, each individual explanation randomly
239 included or excluded the rationale. Each stage 2 task had an overlap of 7 reviewers.

244 4.5 Quality Control

245 This subsection outlines the quality control measures taken during the experiment. These resulted in approval rates of
246 81% for stage 1 and 48% for stage 2. The relatively lower acceptance rate for stage 2 might indicate that multiple-choice
247 forms are more hastily completed and that workers are more focused on images than text.

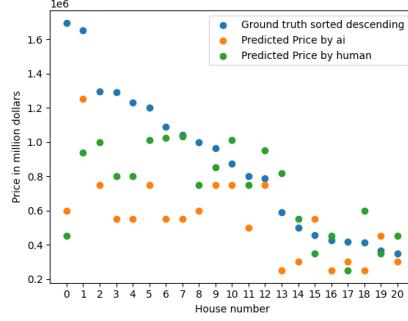
- 248 QC 1. Each task suite contains one control task, where the desired answer is clearly provided. This question functions
249 as an attention check.
QC 2. The first-stage rationales have an enforced minimum length of 50 characters to avoid very small generic answers.
QC 3. Submissions were manually reviewed and were rejected if the submission time was unreasonably low, the
control task was answered incorrectly, or if the rationales were insufficiently poor or not unique.
QC 4. Users were automatically banned after multiple fast responses, rejections, or failed control tasks.
QC 5. Only the top 20% of users quality-wise were allowed to participate in the tasks.

261 5 RESULTS AND DISCUSSION

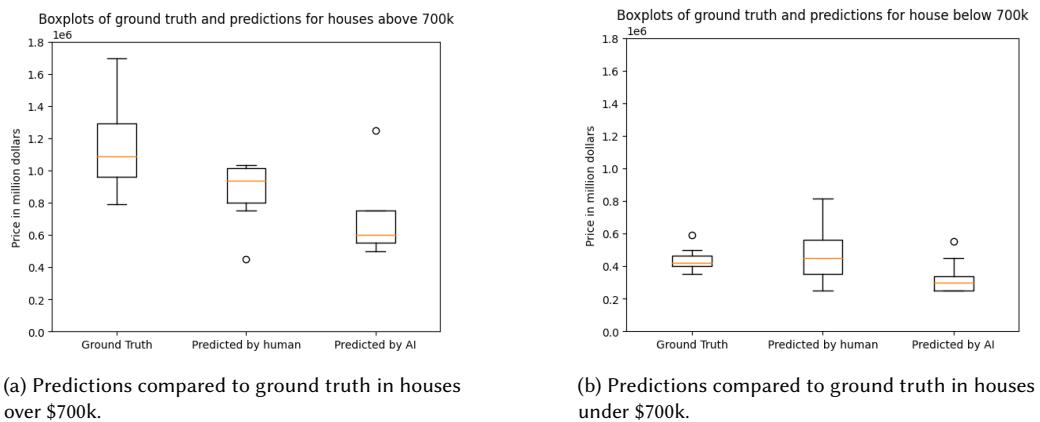
262 In this section, the findings of our experiments and analysis are presented. The results of both stages will be visualized
 263 and analyzed. In the first stage section, the accuracy of the human and AI predictions is compared and the dataset is
 264 divided in expensive and more affordable properties, in order to provide an intuition about what properties tend to be
 265 underestimated or overestimated by humans or AI. In the second stage section, statistical tests take place, aiming to
 266 investigate the hypotheses that were earlier introduced, i.e. whether a rationale increases trust towards its source, and
 267 whether the human will be trusted more, in case both humans and AI do or do not provide rationales.
 268

271 5.1 First Stage Results

272 In the first stage of the experiment, the crowd workers achieved higher accuracy in their predictions than AI. As we
 273 observe in Figure 2, both of them provide relatively accurate predictions. The Mean Absolute Percentage Error (MAPE)
 274 was calculated by comparing the prices to the mean of each range of predictions and was found to be 24.2% for the
 275 crowd workers and 36.3% for the AI. Particularly, as depicted in Figure 3a, the prices of the houses that cost over 700
 276 thousand dollars got significantly underestimated by the AI and slightly underestimated by humans.
 277



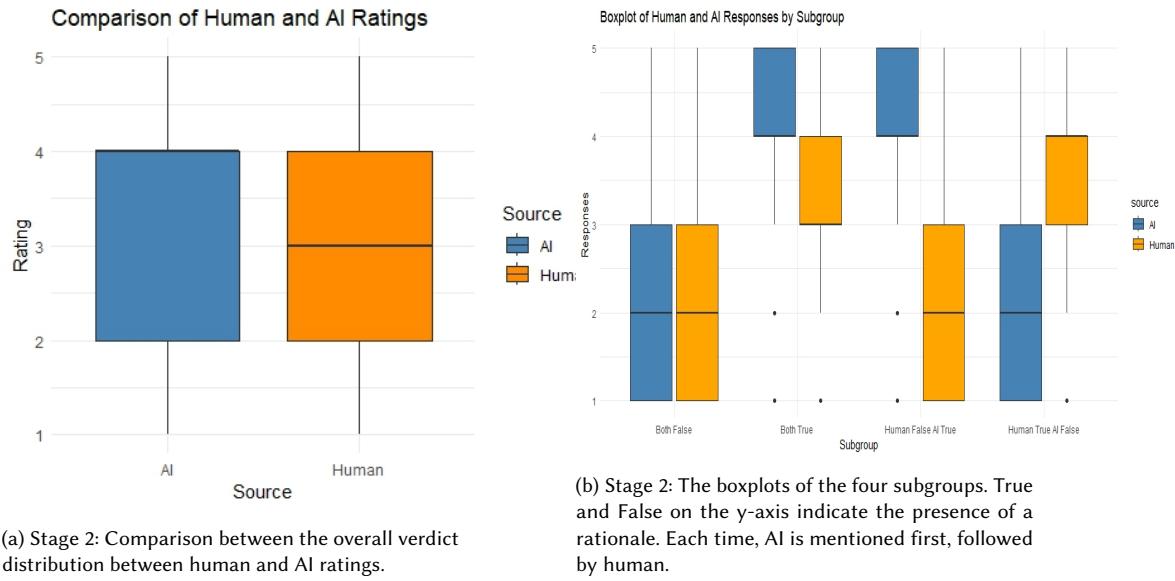
291 Fig. 2. Stage 1 results of ground truth price sorted descending compared to the predicted prices.
 292



310 Fig. 3. Predictions compared to ground truth boxplots.
 311

313 5.2 Second Stage Results

314 For the second stage, the crowd workers were asked to rate how helpful the human and AI were on a scale from 1 to 5.
 315 In R studio, a t-test was performed on these scales to compare humans with AI. The t-test assumes the data has a
 316 Gaussian distribution [30]. The t-test tests the null hypothesis: there is no difference between the mean of the human
 317 scales and AI scales. This hypothesis is rejected because the p-value < 0.05, which means there is a significant difference
 318 between Human and AI scales. The mean difference is 0.3819672 in favour of the AI. [Figure 4a](#) shows that AI ratings
 319 have indeed a higher mean.
 320



343 Fig. 4. The boxplots of the ratings from various perspectives. Ratings are from 1 to 5, indicating the helpfulness of the provided advice.
 344

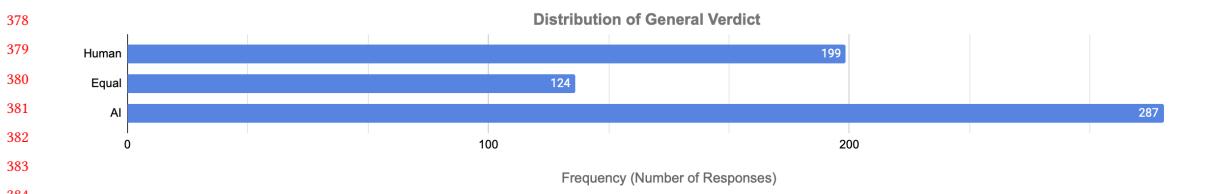
345 Furthermore, a Wilcoxon test was performed on the ratings of humans and AI. This test does not make any assumption
 346 about the data distribution. This test also rejected the null hypothesis and showed that there is a significant difference
 347 between Human and AI ratings with a p-value < 0.05. Additionally, the Wilcoxon test was performed on each task of
 348 the second stage. Results of these tests are depicted in [Table 1](#).
 349

352 Table 1. Wilcoxon signed rank test with continuity correction

	Human Rationale	AI Rationale	P-Value
	✓	✓	9.817e-12
	✓	✗	2.2e-16
	✗	✓	2.2e-16
	✗	✗	0.4255

361 The Wilcoxon test, when providing no rationales for both humans and AI, shows a p-value > 0.05 (no significant
 362 difference). For all other possibilities, i.e. only one rationale for both or both rationales, the test shows a p-value < 0.05,
 363

365 which indicates a significant difference between the two rating groups. The Wilcoxon test does not show which source
 366 (Human or AI) has a higher mean for these subgroups. However, this can be read from the boxplots in [Figure 4b](#). The
 367 figure shows that when the human rationales are provided, the mean ratings for humans are higher, with a mean
 368 difference of 2. The same holds for AI, when its rationales are provided, it has a higher mean. The mean difference for
 369 that case is 3. When both rationales are provided, the mean rating for AI is higher by 2. Based on these results, our
 370 first hypothesis: a rationale increases trust towards its source. Advice including a rationale will be preferred over one
 371 without it, is accepted. The second hypothesis: if both humans and AI do **not** provide rationales, then the human advice
 372 will be trusted more, is rejected. The final hypothesis: if both humans and AI do provide rationales, then the human
 373 advice will be trusted more, is rejected. Finally, it is shown in [Figure 5](#) that people generally preferred AI over humans
 374 when considering (house price) predictions with rationales.
 375
 376



377
 378 Fig. 5. Summarized preferences of people regarding (house price) predictions. Verdict can either be human, AI or equal.
 379
 380
 381
 382
 383
 384

385 6 CONCLUSION & FUTURE WORK

386 The aim of this study was to investigate whether people have more trust in humans or AI when considering (house
 387 price) predictions with rationales. After conducting a two-stage experiment, results indicate that people generally trust
 388 AI more than humans, even though this is not fully backed up by the accuracy compared to the ground truth. When
 389 crowd workers were not provided with any rationales, they rated human and AI advice evenly. When only one rationale
 390 was shown, that advice was rated substantially more useful. Essentially, this follows direct intuition. However, the
 391 mean difference was higher when only AI rationales were provided compared to only human rationales. This already
 392 shows that, even when only human rationales were provided, people thought AI was relatively promising. Eventually,
 393 when both rationales were provided, our study shows that people rate AI higher than human advice instead of equal as
 394 was the case with both not providing rationales. Summarizing all the groups, AI is cumulatively rated the highest and
 395 preferred over human advice.

396 Based on certain limitations, we deem the following as valuable for future work. The experiment could be extended
 397 by using residents nearby or from the region itself as crowd workers/respondents (as opposed to distant ones in the
 398 performed research). Enhanced contextual understanding of factors that influence the price of a property, such as the
 399 neighbourhood's reputation, could affect trust due to personal experience in this region. It would be interesting to make
 400 a controlled experiment that measures the effect of personal experience to the results. Another similar approach could
 401 be adding another stage where information is collected about critical indicators of the region, such as the reputation
 402 or the proximity to amenities (schools, hospitals, shopping centers), allowing for more precisely composed responses
 403 and thus results. Furthermore, other AI models should be compared as well for a broader perspective on the study. For
 404 example, using GPT-4 instead of MiniGPT-4 should increase the quality of the AI predictions and rationales. Lastly, it
 405 would be valuable to research if AI writing more digestible advice compared to humans influences the trust verdict
 406 positively, no matter the underlying accuracy.

417 REFERENCES

- 418 [1] OpenAI. Gpt-4 technical report, 2023.
- 419 [2] Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*, 2023.
- 420 [3] Hammeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. Evaluating the logical reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*, 2023.
- 421 [4] Namkee Oh, Gyu-Seong Choi, and Woo Yong Lee. ChatGPT goes to the operating room: evaluating GPT-4 performance and its potential in surgical education and training in the era of large language models. *Annals of Surgical Treatment and Research*, 104(5):269, 2023. doi: 10.4174/astr.2023.104.5.269. URL <https://doi.org/10.4174/astr.2023.104.5.269>.
- 422 [5] JE Korteling, GC van de Boer-Visschedijk, RA Boswinkel, RC Boonekamp, Rubricering rapport Ongerubriceerd, Managementuittreksel Ongerubriceerd, Samenvatting Ongerubriceerd, Rapporttekst Ongerubriceerd, and Bijlagen Ongerubriceerd. Effecten van de inzet van non-human intelligent collaborators op opleiding en training [v1719]. *Report TNO 2018 R11654. Soesterberg: TNO defense safety and security*, 2018.
- 423 [6] J.E. (Hans) Korteling, Gillian van de Boer-Visschedijk, Romy Blankendaal, Rudy Boonekamp, and A. Eikelboom. Human- versus artificial intelligence. *Frontiers in Artificial Intelligence*, 4:622364, 03 2021. doi: 10.3389/frai.2021.622364.
- 424 [7] Ben Shneiderman. Design lessons from ai's two grand goals: human emulation and useful applications. *IEEE Transactions on Technology and Society*, 1(2):73–82, 2020.
- 425 [8] Ben Shneiderman. Human-centered artificial intelligence: Reliable, safe & trustworthy. *International Journal of Human–Computer Interaction*, 36(6): 495–504, 2020.
- 426 [9] Yanyan Dong, Jie Hou, Ning Zhang, and Maocong Zhang. Research on how human intelligence, consciousness, and cognitive computing affect the development of artificial intelligence. *Complexity*, 2020:1–10, October 2020. doi: 10.1155/2020/1680845. URL <https://doi.org/10.1155/2020/1680845>.
- 427 [10] Onur Asan, Alparslan Emrah Bayrak, and Avishek Choudhury. Artificial intelligence and human trust in healthcare: Focus on clinicians. *J Med Internet Res*, 22(6):e15154, Jun 2020. ISSN 1438-8871. doi: 10.2196/15154. URL <https://doi.org/10.2196/15154>.
- 428 [11] Ella Glikson and Anita Williams Woolley. Human trust in artificial intelligence: Review of empirical research. *Academy of Management Annals*, 14(2):627–660, 2020.
- 429 [12] Stanford Law School. Artificial Intelligence in the Health Care Space: How We Can Trust What We Cannot Know, 4 2019. URL <https://law.stanford.edu/publications/artificial-intelligence-in-the-health-care-space-how-we-can-trust-what-we-cannot-know/>.
- 430 [13] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier, 2016. URL <https://arxiv.org/abs/1602.04938>.
- 431 [14] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267:1–38, February 2019. doi: 10.1016/j.artint.2018.07.007. URL <https://doi.org/10.1016/j.artint.2018.07.007>.
- 432 [15] Vivian Lai, Chacha Chen, Q. Vera Liao, Alison Smith-Renner, and Chenhao Tan. Towards a science of human-ai decision making: A survey of empirical studies, 2021.
- 433 [16] Omar Zaidan, Jason Eisner, and Christine Piatko. Using “annotator rationales” to improve machine learning for text categorization. In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pages 260–267, Rochester, New York, April 2007. Association for Computational Linguistics. URL <https://aclanthology.org/N07-1033>.
- 434 [17] Tyler McDonnell, Matthew Lease, Mucahid Kutlu, and Tamer Elsayed. Why is that relevant? collecting annotator rationales for relevance judgments. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, 4(1):139–148, Sep. 2016. doi: 10.1609/hcomp.v4i1.13287. URL <https://ojs.aaai.org/index.php/HCOMP/article/view/13287>.
- 435 [18] Mucahid Kutlu, Tyler McDonnell, Tamer Elsayed, and Matthew Lease. Annotator rationales for labeling tasks in crowdsourcing. *Journal of Artificial Intelligence Research*, 69:143–189, 2020.
- 436 [19] Gaole He, Lucie Kuiper, and Ujwal Gadira. Knowing about knowing: An illusion of human competence can hinder appropriate reliance on ai systems, 01 2023.
- 437 [20] Kailas Vodrahalli, Roxana Daneshjou, Tobias Gerstenberg, and James Zou. Do humans trust advice more if it comes from ai? an analysis of human-ai interactions, 2021. URL <https://arxiv.org/abs/2107.07015>.
- 438 [21] Xinru Wang and Ming Yin. Effects of explanations in ai-assisted decision making: Principles and comparisons. *ACM Trans. Interact. Intell. Syst.*, 12(4), nov 2022. ISSN 2160-6455. doi: 10.1145/3519266. URL <https://doi.org/10.1145/3519266>.
- 439 [22] Toloka AI. Powering ai with human insight - toloka ai, 2023. URL <https://toloka.ai/>.
- 440 [23] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023. URL <https://arxiv.org/abs/2304.10592>.
- 441 [24] Eman Ahmed and Mohamed Moustafa. House price estimation from visual and textual features. *CoRR*, abs/1609.08399:1–7, 2016. URL <http://arxiv.org/abs/1609.08399>.
- 442 [25] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs, 2021. URL <https://arxiv.org/abs/2111.02114>.
- 443 [26] Official Data Foundation / Alioth LLC. Housing price inflation calculator. <https://www.in2013dollars.com/Housing/price-inflation/2016-to-2023?amount=1>, 2023. [Accessed 7-Jun-2023].

469 [27] S Hu. Documentation uszipcode, 2022. URL <https://uszipcode.readthedocs.io/>.
470 [28] U.S. Census Bureau. Data python package uszipcode. URL <https://data.census.gov/cedsci/table?q=94103>.
471 [29] Toloka. Toloka template-builder map documentation. URL <https://toloka.ai/docs/template-builder/reference/field.map/>.
472 [30] Muhammad Usman. Power efficiency of sign test and wilcoxon signed rank test relative to t-test. 2015. Retrieved from <https://core.ac.uk/download/pdf/234680288.pdf>.
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520

521 A TASK EXAMPLES

522 For stage 1, an example standard task looks like Figure 6.

523

524 **House price estimating**

525 Consider a house described by the following
526 attributes and images...

527 Number of bedrooms:
528 4
529 Number of bathrooms:
530 4
531 Surface area (in sq. ft):
532 4053
533 Surface area (in m²):
534 377

535 Frontal view:
536 

537 Kitchen:
538 

539 Bedroom:
540 

541 Bathroom:
542 

543 Map of zip code area:
544 

545 1. How much would you estimate this
546 house to cost?
547 \$0 - \$100K
 \$100K - \$200K
 \$200K - \$300K
 \$300K - \$400K
 \$400K - \$500K
 \$500K - \$600K
 \$600K - \$700K
 \$700K - \$800K
 \$800K - \$900K
 \$900K - \$1M
 \$1M - \$1.25M
 \$1.25M - \$1.5M
 \$1.5M - \$1.75M
 \$1.75M - \$2M
 \$2M+

548 2. Provide a 2-3 sentence explanation
549 on why you chose this price.
550

551 Fig. 6. Example of a normal task in stage 1.

552 For stage 1, an example control task looks like Figure 7. As validation of the control task, we check if the worker entered
553 the price range \$800K - \$900K, as shown in the picture. The rationale provided in the control task is not validated due
554 to its complexity.

555 **Attention check**

556 This is an attention check. No need to look at the
557 house details. Just select the correct answers.

558 Number of bedrooms:
559 4
560 Number of bathrooms:
561 4
562 Surface area (in sq. ft):
563 4053
564 Surface area (in m²):
565 377

566 Frontal view:
567 

568 Kitchen:
569 

570 Bedroom:
571 

572 Bathroom:
573 

574 Map of zip code area:
575 

576 1. How much would you estimate this
577 house to cost?
578 \$0 - \$100K
 \$100K - \$200K
 \$200K - \$300K
 \$300K - \$400K
 \$400K - \$500K
 \$500K - \$600K
 \$600K - \$700K
 \$700K - \$800K
 \$800K - \$900K
 \$900K - \$1M
 \$1M - \$1.25M
 \$1.25M - \$1.5M
 \$1.5M - \$1.75M
 \$1.75M - \$2M
 \$2M+

579 2. Provide a 2-3 sentence explanation
580 on why you chose this price.
581

582 Fig. 7. Example of a control task in stage 1. The clear answer here is the range \$800K - \$900K. The rationale is not checked.

583 For stage 2, an example standard task looks like Figure 8, and an example control task looks like Figure 9. As validation
584 of the control task, we check if the worker sequentially selected the answers '3', '1', and 'no preference', as indicated in
585 the questions and advice.

573 **House price estimating**

574 Consider a house described by the following attributes and images...

575 Number of bedrooms:
4

576 Number of bathrooms:
4

577 Surface area (in sq. ft):
4053

578 Surface area (in m²):
377

579 Frontal view:

580 

581 Kitchen:

582 

583 Bedroom:

584 

585 Bathroom:

586 

587 Map of zip code area:

588 

589 Please now consider this human advice on the house price...

590 **Human advice:**
Price:
\$1,250,000 - \$1,500,000

591 Please now consider this AI-generated advice on the house price...

592 **AI-generated advice:**
Price:
\$500,000 - \$600,000

593 Explanation:
The property has 5 bedrooms and 5.5 bathrooms. It is a good size for a family. The size of 3912 square feet is also reasonable for the number of bedrooms and bathrooms. The property is located in a desirable area with a high demand for housing. The build quality appears to be of good quality, with a well-maintained exterior and a spacious interior. The property also has a large backyard and a garage, which is a plus. However, the price range of \$500,000 to \$600,000 is slightly lower than the average price for similar properties in the area. Therefore, the price prediction for this property is \$500,000 to \$600,000.

594 **1. How helpful is the human advice?**

595 1 -- Not helpful
 2
 3
 4
 5 -- Very helpful

596 **2. How helpful is the AI advice?**

597 1 -- Not helpful
 2
 3
 4
 5 -- Very helpful

598 **3. Whose advice do you prefer?**

599 Human
 AI
 No preference

Fig. 8. Example of a normal task in stage 2.

595 **Attention check**

596 This is an attention check. No need to look at the house details. Just select the correct answers.

597 Number of bedrooms:
4

598 Number of bathrooms:
4

599 Surface area (in sq. ft):
4053

600 Surface area (in m²):
377

601 Frontal view:

602 

603 Kitchen:

604 

605 Bedroom:

606 

607 Bathroom:

608 

609 Map of zip code area:

610 

611 Please now consider this human advice on the house price...

612 **Human advice:**
Price:
\$100K - \$200K

613 Please now consider this AI-generated advice on the house price...

614 **AI-generated advice:**
Price:
\$1M - \$1.25M

615 Explanation:
This is an attention check. Please select the option described in the question.

616 **1. Please select the 3rd option.**

617 1 -- Not helpful
 2
 3
 4
 5 -- Very helpful

618 **2. Please select the 1st option.**

619 1 -- Not helpful
 2
 3
 4
 5 -- Very helpful

620 **3. Please select 'no preference'.**

621 Human
 AI
 No preference

Fig. 9. Example of a control task in stage 2. The required answers are '3', '1', and 'no preference', as indicated in the questions.

B MINIGPT-4 PROMPTS

For the prediction of house prices and the generation of its corresponding rationale, the following prompt has been used (on the respective images of a certain house which were put in a 2 x 2 grid consisting of the frontal view, kitchen, bedroom and bathroom image), in which the data from the dataset was entered where necessary:

Please give a price prediction for the property in the provided image and textual attributes below. Take into account the number of bedrooms, bathrooms, and areas are the most important features to focus on.

The property has the following attributes:
- Number of Bedrooms: [number of bedrooms]
- Number of bathrooms: [number of bathrooms]
- Area: [area in sqft]
- Zipcode: [zipcode]

Furthermore, provide your rationale for the answer. Please refer in your rationale to the image, e.g. by detecting build quality, garden, garage (if present), etc., and the importance of the attributes.

Format your answer in the form of a JSON:

```
{  
  "zipcode": number,  
  "price_prediction_lower_bound": number,  
  "price_prediction_upper_bound": number,  
  "rationale": string  
}
```

Note that `price_prediction_upper_bound - price_prediction_lower_bound < 100000`.

Received 23 June 2023