Style Over Substance: Evaluation Biases for Large Language Models

Anonymous ACL submission

Abstract

As large language models (LLMs) continue to advance, accurately and comprehensively evaluating their performance becomes increasingly challenging. Ranking the relative performance of LLMs based on Elo ratings, according to human or LLM judgment, is gaining more popularity. However, the extent to which humans 800 and LLMs are capable evaluators remains uncertain. This study investigates the behavior of crowd-sourced and expert annotators, as well as LLMs, when comparing outputs from different models. To achieve this, we curate a dataset of intentionally flawed, machine-generated an-013 swers. Our findings reveal a concerning bias in the evaluation process, as answers with factual errors are rated more favorably than answers that are too short or contained grammatical er-017 rors. To address this issue, we propose independently evaluating machine-generated text across multiple dimensions, rather than merg-021 ing all the evaluation aspects into a single score. We instantiate this idea with the Elo rating system, resulting in the Multi-Elo Rating System (MERS). Empirical results from our study reveal that this proposed approach significantly enhances the quality of LLM-based evaluations, particularly in terms of factual accuracy. However, there is no significant improvement in crowd-sourced evaluations, indicating the need for further investigation.

1 Introduction

Recent advancements in the field of natural language processing have demonstrated that the utilization of supervised instruction fine-tuning and reinforcement learning from human feedback (RLHF) can yield substantial improvements in the performance of large language models (LLMs) with respect to their ability to comprehend and execute instructions (Ouyang et al., 2022; Wei et al., 2022; Sanh et al., 2022; Chung et al., 2022; OpenAI, 2023; Wu et al., 2023; Li et al., 2023a; Lyu et al., 2023). This progress signifies a significant stride in the domain of language model development. However, the assessment of these enhanced LLMs presents a notable challenge, particularly when confronted with more generalized instructions that entail open-ended responses. Such instructions often lack a definitive metric for evaluation within the confines of traditional natural language processing benchmarks. 042

043

044

047

048

051

053

057

059

061

062

063

064

065

067

068

069

070

071

073

074

075

076

077

In response to this challenge, recent studies commonly adopt the *Elo rating system* (Elo, 1967) for evaluation. This approach involves enlisting either human or LLM evaluators to compare and rate the outputs generated by two different models (Askell et al., 2021; Bai et al., 2022a; Srivastava et al., 2022; Dettmers et al., 2023; Zheng et al., 2023). This evaluation method enables the computation of an Elo-based leaderboard to rank the relative performance of LLMs (Chiang et al., 2023).¹ Nonetheless, an important question arises concerning the qualifications of human and LLM judges to serve as effective evaluators in this context. Evaluating model outputs encompasses a multifaceted decision-making process, and it remains an open question whether these judges possess the expertise to accurately determine the superior model output. Further research is needed to address this inquiry comprehensively and refine the evaluation procedures for enhanced LLMs.

In this study, we systematically generate a set of responses, considering factors such as language proficiency, factual accuracy, and response length, and conduct pairwise comparisons for compute Elo ratings.² We employ 40 general-purpose questions sourced from Chiang et al. (2023) that do not require specialized expertise to ensure the generaliza-

¹https://arena.lmsys.org/

²In this work, we refer to various settings used for answering questions as different 'models', even though all answers are generated by GPT-4. Furthermore, we describe the pairwise comparison as a game/battle played between two models.

	Answer Features				Elo Ratings						
	# of words	Language	# of Factual			Human	GPT-4	Claude-1			
		Errors	Errors	C	rowd	Expert					
Correct	≈ 100	N.A.	0	1091		1162	1482	1320			
+ Short	≈ 50	N.A.	0	970		1029	1096	1052			
One Minor Factual Error	ror ≈ 100 N.A. 1, minor 1074			1137	1415	1265					
+ Short	≈ 50	N.A.	1, minor	1002		964	988	997			
Several Minor Factual Errors	≈ 100	N.A.	\approx 3, minor	1032		1024	1206	1182			
+ Short	≈ 50	N.A.	\approx 3, minor	952		873	851	891			
Several Major Factual Errors	≈ 100	N.A.	\approx 3, major	1025		892	861	979			
+ Short	≈ 50	N.A.	\approx 3, major	937		832	710	782			
Advanced Learner	≈ 100	Spelling	0	1041		1138	1213	1126			
+ Short	≈ 50	Spelling	0	941		986	824	841			
Intermediate Learner	≈ 100	Grammatical	0	1015		1108	771	904			
+ Short	≈ 50	Grammatical	0	921		855	582	662			

Table 1: Elo ratings for answers in different settings based on the annotations given by crowd-sourced annotators, expert annotators, GPT-4, and Claude-1.

tion of our study and reduce annotation difficulty. The answers for these questions are generated by GPT-4 with specific instructions. To probe the potential impact of language proficiency towards human and LLM judgments, we instruct GPT-4 to emulate an advanced English learner, occasionally incorporating spelling errors, or an intermediate English learner, occasionally introducing grammatical mistakes during the response generation process. To probe factual accuracy, we direct GPT-4 to include varying degrees of factual errors in the responses. Lastly, we explore the influence of response length by instructing GPT-4 to generate answers of differing lengths. To ensure that the generated responses conformed to the desired criteria, we conduct manual reviews and carry out postediting as necessary. After obtaining the collection of responses, we conduct annotations with a diverse pool of annotators, including crowd-sourced annotators, expert annotators, GPT-4, and Claude-1. We calculate the Elo ratings for each "model" based on their judgments to understand their preferences when deciding on high-quality model outputs.

096

100

In our investigation (Table 1 and Section 3), we 101 uncover intriguing findings related to the assessment of answer quality. Firstly, we observe un-103 certainty in determining answer quality among hu-104 mans, including both crowd-sourced and expert 105 annotators, whereas LLMs exhibit greater certainty 106 107 in their evaluations. Furthermore, we notice a substantial difference in fact-checking capabilities 108 between humans and LLMs. Humans generally 109 do not thoroughly fact-check answers unless the 110 factual error is glaringly evident, whereas LLMs 111

demonstrate some degree of fact-checking ability, albeit with imperfections. Another significant finding is that both humans and LLMs favor longer answers. Interestingly, crowd-sourced annotators appear to be (almost) unbiased in their evaluation of answer ordering, while LLMs lean toward the first answer and experts prefer the second answer. 112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

We note that a unified measure is insufficient for evaluating LLMs. Furthermore, some aspects, such as factuality, are arguably more crucial than others. Yet, we observe judges preferring factually incorrect models over those that are grammatically incorrect or short. With these identified issues, we present the Multi-Elo Rating System (MERS) designed to assess machine-generated text from multiple dimensions. Our proposed approach requires human judges and LLMs to independently evaluate the machine-generated text from three aspects: "Accuracy", "Helpfulness", and "Language". This allows us to achieve a more comprehensive and transparent understanding of the quality of model outputs. Note that our approach is not limited to these three dimensions and can be easily extended to more dimensions. Our empirical findings demonstrate a significant improvement in the evaluation quality of GPT-4, particularly regarding factual accuracy. However, we observe that humans still exhibit uncertainty in their assessments.

Based on our findings, we highly recommend that practitioners evaluate machine-generated text from various perspectives rather than depending solely on a single unified measure. Additionally, we advise practitioners to exercise caution when using crowd-sourced annotators.

Model type	Output
Correct	Effective ways to deal with stress include regular exercise, maintaining a healthy diet, getting enough sleep, and taking time for relaxation activities like meditation. It's also important to have a support system in place, whether that's friends, family, or a mental health professional.
Minor Errors	Effective ways to deal with stress include physical activities like playing tennis, doing yoga, or running a marathon. Also, adopt healthy habits like sleeping for 10 hours, eating three full meals a day, and avoiding caffeine. Lastly, meditation, guided by a professional like John Zen, can be very soothing.
Major Errors	Some effective ways to deal with stress include drinking 10 cups of coffee a day, binge-watching Netflix for at least 5 hours, and attending Dr. Monty Farquhar's stress-relief seminars. Additionally, practicing the art of Floga, a combination of Flute-playing and Yoga, can significantly reduce stress levels.
Advanced Learner	Whenevr feeling stressed, always remember to take deep breathes. This can calm nerves. Regular physical exercise also reduce stress hormonal levels. A balance diet promoting overral well-being can lead to up healty stress response. Finally, finding hobby or activity you enjoy can help distract you from stressful situations.
Intermediate Learner	Deal with stress, you can do exercise regularly, practicing a meditation, get with plenty sleep, and eat healthy foods also. You can too to connect with others so express your feelings, and avoiding caffeine, alcohol, and nicotine, and take time to relax and have fun.

Table 2: Examples of different error types for the question "What are the most effective ways to deal with stress?". Factual errors are highlighted in red and language errors are highlighted in orange.

2 Evaluation Method

146

147

148

149

150

151

152

153

154

155

156

157

158

159

161

162

163

166

167

168

170

171

172

173

In this section, we cover the Elo rating system, methods for generating incorrect answers for models, and evaluations for humans and LLMs.

2.1 Elo Rating System

The Elo rating system is a method used to calculate the relative skill levels of players in two-player games, such as chess. Given two players \mathcal{A} and \mathcal{B} whose Elo ratings are \mathcal{R}_A and \mathcal{R}_B respectively, the expected score for these two players are:

$$\mathcal{E}_A = \frac{1}{1+10^{\frac{\mathcal{R}_B - \mathcal{R}_A}{400}}}, \text{ and } \mathcal{E}_B = \frac{1}{1+10^{\frac{\mathcal{R}_A - \mathcal{R}_B}{400}}}.$$
(1)

Suppose the player A is expect to obtain \mathcal{E}_A scores from this game but actually get \mathcal{S}_A scores, the updated Elo rating of player A is:

 $\mathcal{R}'_A = \mathcal{R}_A + \mathcal{K} \cdot (\mathcal{S}_A - \mathcal{E}_A), \qquad (2)$

where \mathcal{K} is adjustment parameter, called the \mathcal{K} -factor. Equation 2 is also used for obtaining \mathcal{R}'_B . Following Chiang et al. (2023), we set $\mathcal{K} = 32$. Additionally, if \mathcal{A} is better than \mathcal{B} , we set $\mathcal{S}_A = 1$ and $\mathcal{S}_B = 0$. If \mathcal{B} is better than \mathcal{A} , we set $\mathcal{S}_A = 0$ and $\mathcal{S}_B = 1$. We set both $\mathcal{S}_A = 0.5$ and $\mathcal{S}_B = 0.5$ if both players are equally good. Moreover, to minimize the influence of the game orders, the Elo rating calculation is performed 10K times with random orders as suggested by Dettmers et al. (2023).

2.2 Answer Generation

Chiang et al. (2023) release a set of 80 questions categorized into 8 categories to evaluate the gener-

ation capability of LLMs.³ This dataset is widely used by recent studies (Dettmers et al., 2023; Li et al., 2023a), due to its high quality. However, due to the requirement for specialized expertise to answer some of these questions, we exclude the "fermi", "coding", and "math" questions, as they typically demand extra efforts to evaluate the answers. This step aims to reduce the potential impact of varying human raters' capabilities on the evaluation process. Additionally, we also remove the "roleplay" and "writing" questions, as they involve creative writing and are prone to subjectivity in human assessment. As a result, our final question set consists of 40 questions, focusing on the "generic", "knowledge", "common sense", and "counterfactual" categories. We believe these retained questions can be easily understood and answered by the general public.

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

Once we have the set of questions, we require GPT-4 to generate answers with specific error types in addition to providing the correct answers. We provide some examples of these error types in Table 2. Regarding language quality errors, we ask GPT-4 to respond as either an advanced English learner or an intermediate English learner. The answers generated by an advanced English learner may occasionally contain spelling errors, while those from an intermediate English learner commonly include grammatical mistakes. In terms of factual accuracy, we expect GPT-4 to produce answers with either minor or major errors. Minor errors primarily involve fabricated names or incorrect numbers, while major errors contain incorrect facts and suggestions. Furthermore, we utilize GPT-4 to

³https://github.com/lm-sys/vicuna-blog-eval/ blob/main/eval/table/question.jsonl

generate both long (approximately 100 words) and 208 short (approximately 50 words) answers for each 209 question to investigate the preference of both hu-210 mans and LLMs regarding answer length. Hence, 211 there are 12 models (settings) in total in this study. 212 We present the prompts used for answer genera-213 tion in Appendix A. To ensure unbiased evaluation 214 results regarding answer ordering, all evaluators, 215 including crowd-sourced human annotators, expert 216 annotators, and LLMs, evaluate all answer pairs 217 from both forward and reversed directions. In the 218 end, we have 5280 unique pairwise comparisons 219 across all 12 models and 40 questions. During the generation process, we notice that GPT-4 may 221 not always fulfill our requirements. For instance, the output sometime may not contain any factual or language errors. Therefore, we manually review all the answers and make necessary edits to ensure they align with our requirements.

2.3 Crowd-Sourced Evaluation

227

232

237

240

241

242

243

244

245

246

247

248

250

257

Crowd-sourced human annotators are commonly used to assess the quality of machine-generated text. For this study, we utilize Amazon Mechanical Turk (AMT) to collect text evaluations from human raters in NLG evaluations. To minimize the potential impact of annotator demographics, we only recruit crowd-sourced human annotators from the United States. We also ensure the quality of annotations by exclusively offering tasks to annotators who have completed over 5,000 assignments and maintain a success rate of 99% or higher. To ensure the diversity of the annotators, we limit each annotator to provide no more than 20 annotations. As a result, we engage 377 unique annotators, each contributing an average of 14 annotations.

During annotation, we provide a simple interface that mainly follows the guidelines of Dettmers et al. (2023), as shown in Appendix C. Annotators are presented with a question and two model-generated responses placed side-by-side. Their task is to select the better output or indicate a tie between them. To ensure the annotators' attentiveness and thorough reading of the responses, we incorporate a mandatory 20-second delay before they can submit their answers. Furthermore, we anonymize the model name from our human annotators.

2.4 Expert Evaluation

To address concerns about the reliability of crowdsourced annotators, we have implemented a parallel system involving a team of 20 experts. Each mem-

[Question] What are the most effective ways to deal with stress?
[The Start of Assistant 1's Answer] Effective ways to deal with stress include [The End of Assistant 1's Answer]
[The Start of Assistant 2's Answer] Whenevr feeling stressed, [The End of Assistant 2's Answer]
[System] We would like to request your feedback on the performance of two AI \hookrightarrow assistants

Figure 1: A simplified example for the evaluation prompt employed by GPT-4 and Claude-1 is designed to assess the quality of responses. A complete example can be found in Figure 10.

258

259

260

261

262

264

265

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

284

285

286

287

ber of this dedicated team holds at least a master's degree from an English-speaking country, ensuring a comprehensive understanding of the language and proficiency in annotation nuances. It is important to note that employing expert annotators comes at a higher cost than using crowd-sourced alternatives. Therefore, we sample 200 games out of a total pool of 5280 games, which is statistically sufficient to yield meaningful insights into the model's performance.⁴ To ensure a fair comparison, we provide the same annotation instructions to both the expert team and the crowd-sourced participants. Each expert is assigned to annotate 10 games, which takes approximately 20 minutes. This approach enables us to accurately evaluate the effectiveness and reliability of each annotation method.

2.5 LLM Evaluation

As human evaluation can be costly and inefficient, there is an increase in the use of advanced LLMs, such as GPT-4, to evaluate model outputs. In our work, we also use LLMs as judges to assess answer quality. However, previous studies rely solely on GPT-4 as the LLM judge (Chiang et al., 2023; Li et al., 2023a; Zheng et al., 2023), which may not be appropriate for our work as our answers are refined by humans after being generated by GPT-4. This raises concerns about potential biases that GPT-4 may have towards its own outputs, which could skew the evaluation process. To ensure fair comparisons and mitigate any bias, we introduce Claude-1 from Anthropic (Bai et al., 2022b) as an additional

⁴The US Chess Federation believes that 25 games are enough to assess a player's ability, and in our sample of 200 games, each model participated in at least 28 games. Source: https://new.uschess.org/ frequently-asked-questions-member-services-area



Figure 2: The percentage distribution of decision choices made by humans and LLMs in terms of answer length (in words).

LLM judge, in addition to GPT-4.⁵ By incorporating multiple LLM judges, we can establish a more comprehensive and unbiased assessment of the generated answers in our study. We utilize the evaluation prompt from Dettmers et al. (2023), as presented in Figure 1. The prompt assesses the answers based on their helpfulness, relevance, accuracy, and level of detail, while also aiming to avoid bias related to answer ordering.

3 Analysis

289

290

291

294

301

303

310

311

313

315

317

319

322

323

In this study, we assess the performance of 12 models using Elo ratings, as evaluated by crowdsourced annotators, expert annotators, GPT-4, and Claude-1. The Elo ratings are presented in Table 1. In this section, we delve into a comprehensive analysis of these distinctions.

All the judges exhibit a bias toward longer texts. Text evaluation by both human and LLM judges often exhibits a bias towards longer responses, where GPT-4 demonstrates the most bias and the expert annotators demonstrate the least bias, as depicted in Figure 2. This inclination is expected since one of the scoring criteria requested is "the level of detail", which often correlates with the length of the output. However, such an inclination is arguably undesirable. For example, we observe that GPT-4 considers "Several Minor Factual Errors" (1206 Elo) to be better than "Correct + Short" (1096 Elo). When manually reviewing the justifications provided by GPT-4 for its evaluation decisions, we observe that GPT-4 sometimes considers the response of "Several Major Factual Errors" as superior to that of "Correct + Short", even when factual errors are detected. We present an example in Appendix D.



Figure 3: The distribution of decision choices made by humans and LLMs.

324

325

326

328

330

331

332

333

334

336

337

338

339

340

341

342

343

344

345

346

349

350

351

353

354

355

357

358

361

Humans are indecisive. Unlike LLM judges, both expert and crowd-sourced annotators demonstrate indecisiveness, leading to Elo scores that remain relatively close to the initial value of 1000. The Elo scores assigned by crowd-sourced annotators range from 926 to 1086, while those by expert annotators range from 832 to 1162. These ranges are significantly narrower than those observed for GPT-4 and Claude-1. Interestingly, human judges are more inclined to assign "tie" judgments compared to LLMs, as illustrated in Figure 3. These findings raise important questions about the level of reading comprehension and attention to detail demonstrated by humans, especially when they are crowd-sourced workers evaluating answers.

The order of answers affects the judges' decisions. Our analysis also reveals an interesting bias toward the order of answers, which is evident in the assessments made by both human judges and LLM judges. As illustrated in Figure 3, it can be observed that the crowd-sourced annotators, GPT-4, and Claude-1 exhibit a distinct and notable preference for Assistant 1. This intriguing finding is consistent with prior research (Wang et al., 2023a), further underscoring the significance of this bias. Interestingly, the expert annotators favor the answer given by Assistant 2 and are more likely to assign "tie" judgments. To ensure fairness and comprehensive evaluation, we strongly recommend considering both answer orders when analyzing the performance of the systems. We leave the further study of the behavior of humans for future work.

Crowd-sourced annotators lack fact-checking, while experts and LLMs can fact-check, albeit imperfectly. The close Elo scores of those models with factual errors from crowd-sourced annotators in Table 1 suggest inadequate fact-checking capability of humans, raising concerns about false

 $^{^5\}mbox{GPT-4}$ signature: gpt-4-0613 and Claude-1 signature: Claude-1.3

	Crowd	Expert	GPT-4	Claude-1
Crowd	_	0.08	0.11	0.10
Expert	0.08		0.09	0.14
GPT-4	0.11	0.09		0.51
Claude-1	0.10	0.14	0.51	—

Table 3: Pairwise inter-annotator agreement measured by Cohen's kappa coefficient κ . The pairs involving the experts only cover 200 games.

information and malicious intent in LLM-generated 362 363 outputs. People are vulnerable to believing such outputs, particularly when they appear convincing. 364 Relatedly, instances are found where humans overlooked basic verification tasks, like identifying fake URLs or citations (Nakov et al., 2021),⁶. Although 367 experts are more diligent in fact-checking, the general public's behavior towards LLM output tends to be closer to that of crowd-sourced annotators, 370 posing significant safety risks. On the other hand, 371 LLMs do notice errors, but not consistently. When 372 an LLM fails to detect inaccuracies, it often favors flawed outputs over shorter or grammatically 374 incorrect responses.

LLMs only reach a moderate consensus, while humans embrace more diversity in thought. We assess inter-annotator agreement using Cohen's kappa coefficient κ (Cohen, 1960) and present the results in Table 3. Our analysis, following the interpretation of κ by McHugh (2012), shows that only GPT-4 and Claude-1 achieve moderate agreement, while other comparisons demonstrate only slight agreement. This confirms that human annotators exhibit inconsistencies during annotation. Zheng et al. (2023) define the agreement between two types of judges as the probability of non-identical individuals of each type agreeing on a randomly selected question and shows an approximately 80% agreement between LLMs and crowd-sourced hu-390 mans, which seems to contradict our findings. However, this discrepancy arises because they remove ties and inconsistent annotations, while we consider all annotations. When ties and inconsistencies are retained, Zheng et al. (2023) report an 395 approximately 60% agreement between LLMs and crowd-sourced humans, which is slightly better than random guessing and aligns with our findings.

4 Multi-Elo Rating System

As discussed in Section 3, we identify the limitations in the evaluation of machine-generated answers based on humans and LLMs. We observe that the current widely used evaluation approach fails to yield satisfactory results and lacks a comprehensive understanding of the factors influencing the outcomes. Collapsing multiple decision components (e.g., accuracy, level of details, relevance, language, helpfulness, etc.) into a single score undermines the importance of individual components and introduces subjectivity regarding their relative significance in the final decision. 399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

To overcome these limitations, we propose a novel multi-dimensional evaluation approach for assessing the outputs of LLMs, which we refer to as the Multi-Elo Rating System (MERS). This approach is inspired by machine translation research, where evaluations are often divided into at least two parts: fluency and adequacy. Recent advancements in MT also advocate breaking down the evaluation into multiple aspects (Lommel et al., 2014). In this section, we present a detailed description of our approach and the results.

4.1 Methodology

The Multidimensional Quality Metrics (MQM) framework provides a comprehensive approach for evaluating and establishing standards for translation quality (Lommel et al., 2014). Drawing inspiration from this framework, we propose a similar approach to evaluate the outputs of LLMs from multiple dimensions. Our evaluation focuses on three dimensions of the generated text, as follows:

- Accuracy: The accuracy of text involves considering factual correctness and logical consistency.
- **Helpfulness**: The helpfulness of the text involves considering its relevance of the information and whether it addresses the question given, taking into account the depth of the response given.
- Language: The language of the text involves considering its clarity, coherence, grammar, syntax, and tone.

Note that our approach is not limited to these dimensions and can be easily extended.

The quality of an answer is dependent on its specific context. For instance, if a model gives a detailed but complicated explanation of black holes to an 8-year-old, the answer may be accurate but not useful. Conversely, if a model is asked to compose an email and produces a message with incorrect in-

⁶https://www.nytimes.com/2023/06/08/nyregion/ lawyer-chatgpt-sanctions.html

	Accuracy				Helpfulness		Language			
	Crowd	Expert	GPT-4	Crowd	Expert	GPT-4	Crowd	Expert	GPT-4	
C. +S.	1056 963	1180 1 040	1200 1 158	1045 983	1208 979	1384 1 009	1036 1007	1109 1068	1415 1199	
OMin.FE + S. SMin.F + S. SMaj.FE + S.	1026 978 1036 978 1030 955	1090 898 1044 931 963 787	1120 1016 993 857 794 746	1048 993 1051 956 1037 940	1153 941 1069 865 1015 766	1378 965 1248 845 926 726	1019 990 1029 996 1023 982	1114 1012 1096 935 1010 879	1334 1109 1200 988 995 871	
AL + S. IL + S.	1028 979 1015 956	1121 971 1076 898	1139 1051 1018 908	1032 969 1002 945	1146 891 1095 872	1196 804 908 612	1004 994 991 930	1039 863 992 884	1051 814 560 465	

Table 4: Elo ratings for different models with regard to "Accuracy", "Helpfulness", and "Language" given by crowd-sourced annotators, expert annotators, and GPT-4. C. stands for Correct. OMin.FE stands for One Minor Factual Error. SMin.F stands for Several Minor Factual Errors. SMaj.FE stands for Several Major Factual Errors. AL stands for Advanced Learner. IL stands for Intermediate Learner. S. stands for Short.

formation, the response may lack accuracy but still have some value. By taking this multi-dimensional approach, we can gain a clearer understanding of model performance and prioritize different aspects based on our individual requirements.

To facilitate the multi-dimensional evaluation by humans, we introduce a simple modification to the interface, asking them to rate the quality across the three different aspects, as shown in Appendix C. Additionally, we experiment with two approaches for GPT-4 evaluation: asking three independent queries versus a single query that requests judgments for all three aspects together. In this paper, we report the results given by asking three independent queries. More details are in Appendix E.

4.2 Stop Using Crowd-Sourced Annotators!

In this section, we compare the annotation outcomes provided by both crowd-sourced and experts, as presented in Table 4. Regarding "Accuracy", we find that experts are proficient in identifying factual errors in answers, although not entirely 469 flawless. However, crowd-sourced annotators ex-470 hibit indecisiveness in their evaluations. Notably, 471 the crowd-sourced annotators perceive Several 472 Major Factual Errors + Short (955 Elo) and 473 Correct + Short (963 Elo) as nearly equally 474 good. Regarding "Helpfulness", the experts display 475 a stronger preference for longer answers, while 476 477 the crowd-sourced annotators only slightly favor them, as evidenced by their Elo scores. In terms 478 of "Language", both expert and crowd-sourced an-479 notators face challenges in recognizing spelling or 480 grammatical errors, suggesting that humans may 481

be less sensitive to language errors. Overall, the experts outperform the crowd-sourced annotators in the evaluation, despite not being entirely errorfree themselves. These results serve as a warning against over-reliance on crowd-sourced judgments and highlight concerns regarding the general audience's ability to critically interpret LLM's output. 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

4.3 Experts versus GPT-4

In this section, we discuss the difference between the experts and GPT-4 in evaluation from multiple dimensions and present the results in Table 4. Regarding the "Accuracy" dimension, it is noteworthy that the Elo scores for factual accuracy closely align with the single Elo scores presented in Table 1, suggesting that experts indeed prioritize factual accuracy during evaluation. GPT-4 can also effectively rank models based on the severity of the errors. Regarding the "Helpfulness" dimension, both experts and GPT-4 consistently consider longer answers to be more helpful. Similar to the discussion in Section 4.2, we believe that this preference stems from the strong correlation between "helpfulness" and "the level of detail", as longer answers tend to convey more information, making them perceived as more helpful. Regarding the "Language" dimension, recognizing spelling or grammatical errors in text is challenging for experts, while GPT-4 effectively distinguishes between answers based on their language quality and appears to penalize grammatical errors more heavily during assessment. Overall, this comprehensive analysis sheds light on the evaluation process and reveals the differing perspectives of experts and GPT-4 in assessing various

449

450

451

452

515

516

517

518

519

520

521

524

526

529

531

532

533

534

535

536

539

541

542

543

544

545

550

553

555

556

557

563

aspects of model outputs.

5 Discussion

When assessing the quality of an LLM, it is not advisable to merge many factors into a single score, as there is no consensus on which factor should be prioritized, both for the model and human evaluators. Instead, a more detailed approach is needed. In this work, we suggest assessing three distinct factors individually for a thorough evaluation. Nonetheless, this categorization might not be exhaustive in capturing all necessary elements for an ideal response, indicating the need for further research in this area.

Another important consideration is the use of human judgments for evaluating LLM performance. While crowd feedback can provide a general indication of how the audience perceives the LMs' output, caution must be exercised. Crowd-sourced evaluators may not always involve rigorous fact-checking, thus giving high scores to factually incorrect answers. Expert evaluators are better in this aspect, with the caveat of more difficulty in scaling the process. Additionally, both human annotators demonstrate various biases, such as the length and order of the answers. Hence, it is crucial to supplement human opinions with other evaluation methods to gain a better understanding of LLMs.

6 Related Work

Large Language Models Large Language Models (LLMs) commonly refer to Transformer-based language models with billions of parameters (Vaswani et al., 2017). Examples of these models include GPT-3 (Brown et al., 2020), PaLM (Chowdhery et al., 2022), BLOOM (Scao et al., 2022), and LLaMA (Touvron et al., 2023). These models demonstrate impressive abilities in understanding natural language and handling complex tasks. It is found that supervised fine-tuning and reinforcement learning from human feedback can further enhance their performance following general language instructions (Mishra et al., 2022; Wang et al., 2022b; Wei et al., 2022; Sanh et al., 2022; Ouyang et al., 2022; Chung et al., 2022; Yin et al., 2022; Gupta et al., 2022; Muennighoff et al., 2022; Wu et al., 2023). Therefore, accurately and comprehensively assessing the performance of these LLMs remains an unsolved challenge.

Evaluation NLP models are commonly assessed using standardized benchmark test suites. Several recent studies suggest using a diverse set of NLP benchmarks for LLM evaluation (Hendrycks et al., 2021; Liang et al., 2022; Li et al., 2023b). As pointed out by Gudibande et al. (2023) and Zheng et al. (2023) that there is a gap between users' perception and standardized evaluation suites, recent LLM studies often incorporate human evaluation for a more nuanced understanding of model performance (Wang et al., 2022a; Chiang et al., 2023). As human evaluations can be costly, some recent works utilize state-of-the-art LLMs to evaluate model outputs (Wang et al., 2023b; Kim et al., 2023). More recently, several works employ the Elo rating system from chess games to gauge the LLMs' capabilities (Askell et al., 2021; Bai et al., 2022a; Srivastava et al., 2022; Dettmers et al., 2023; Zheng et al., 2023). However, these previous works operate under the assumption that human evaluations serve as the gold standard. In contrast, Clark et al. (2021) demonstrate that humans are no longer the gold standard and not able to distinguish between human-written and machine-generated text. We investigate the judgments of humans and LLMs in evaluating text across multiple aspects.

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

7 Conclusion

In this study, we investigate the limitations of humans and LLMs as evaluators by examining their behaviors in assessing machine-generated text. We deliberately introduce fabricated factual and grammatical errors into a set of machine-generated answers and analyze the responses of crowd-sourced annotators, experts, and LLMs. The primary goal is to gain insights into the limitations and biases exhibited by both humans and LLMs. We observe that both humans and LLMs demonstrate various biases. To address the observed issues, we propose independently assessing machine-generated text across multiple dimensions and demonstrating this idea with the Elo rating system, resulting in the Multi-Elo Rating System. Our empirical findings show that this approach effectively enhances GPT-4's evaluation quality, especially in terms of factual accuracy. However, crowd judges continue to exhibit indecisiveness. In light of these findings, we encourage practitioners to adopt a multi-dimensional evaluation approach for machine-generated text instead of relying solely on a single unified measure. We also advise caution in using crowd annotators to assess LLMs due to their indecisiveness, bias towards lengthy responses, and limited fact-checking capabilities.

671 672 673 674 675 676 677 678 679 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714

715

716

717

718

719

720

721

722

723

724

725

726

727

667

668

670

8 Limitations

614

623

631

633

638

639

640

641

644

645

647

648

651

653

654

655

657

658

662

615Question CoverageWe select only 40 questions616from Chiang et al. (2023), mainly due to explosive617annotation cost. We acknowledge that this limited618selection may not capture the full spectrum of ques-619tion types and variations. Consequently, there is620a potential risk that some aspects of the research621question may not receive sufficient representation622or exploration.

Evaluation Dimension Coverage In our proposed Multi-Elo Rating System, we only explore three crucial dimensions: "Accuracy", "Helpfulness", and "Language". We acknowledge that while these dimensions provide valuable insights, they may not encompass the entirety of the multifaceted nature of text evaluation. Furthermore, it is important to recognize that our definitions for the three dimensions we have chosen are not infallible. Different stakeholders may have diverse perspectives on these dimensions.

We leave the investigation of addressing these limitations to future work.

References

- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Benjamin Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam Mc-Candlish, Chris Olah, and Jared Kaplan. 2021. A general language assistant as a laboratory for alignment. *CoRR*, abs/2112.00861.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback. *CoRR*, abs/2204.05862.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosiute, Liane

Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemí Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022b. Constitutional AI: harmlessness from AI feedback. *CoRR*, abs/2212.08073.

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%* chatgpt quality.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways. CoRR, abs/2204.02311.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam

Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. *CoRR*, abs/2210.11416.

728

729

731

735

737

738

741

743

744

745

746

747

750

751

752

753

754

755

756

757

758

766

767

770

771

773

774

775

776

777

779

781

- Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. In *Proceedings* of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 7282–7296, Online. Association for Computational Linguistics.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.
 - Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *CoRR*, abs/2305.14314.
 - Arpad E Elo. 1967. The proposed uscf rating system, its development, theory, and applications. *Chess Life*, 22(8):242–247.
- Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. The false promise of imitating proprietary llms. *CoRR*, abs/2305.15717.
- Prakhar Gupta, Cathy Jiao, Yi-Ting Yeh, Shikib Mehri, Maxine Eskenazi, and Jeffrey Bigham. 2022. InstructDial: Improving zero and few-shot generalization in dialogue through instruction tuning. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 505– 525, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net.
- Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023. Prometheus: Inducing finegrained evaluation capability in language models. *CoRR*, abs/2310.08491.
- Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. 2023a. Bactrian-x : A multilingual replicable instruction-following model with low-rank adaptation. *CoRR*, abs/2305.15011.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023b. CMMLU: measuring massive multitask language understanding in chinese. *CoRR*, abs/2306.09212.

Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yüksekgönül, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2022. Holistic evaluation of language models. CoRR, abs/2211.09110.

782

783

784

785

786

789

790

791

792

793

794

795

796

797

799

800

801

802

803

804

805

806

807

808

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

- Arle Lommel, Hans Uszkoreit, and Aljoscha Burchardt. 2014. Multidimensional quality metrics (mqm): A framework for declaring and describing translation quality metrics. *Tradumàtica*, (12):0455–463.
- Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. 2023. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. *CoRR*, abs/2306.09093.
- Mary L McHugh. 2012. Interrater reliability: the kappa statistic. *Biochemia medica*, 22(3):276–282.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M. Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2022. Crosslingual generalization through multitask finetuning. *CoRR*, abs/2211.01786.
- Preslav Nakov, David P. A. Corney, Maram Hasanain, Firoj Alam, Tamer Elsayed, Alberto Barrón-Cedeño, Paolo Papotti, Shaden Shaar, and Giovanni Da San Martino. 2021. Automated fact-checking for assisting human fact-checkers. In Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021, pages 4551–4558. ijcai.org.
- OpenAI. 2023. GPT-4 technical report. CoRR, abs/2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang,

- 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949
- 950 951 952 953 954 955 956 957

Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In Advances in Neural Information Processing Systems.

839

860

863

864

870

871

872

873

874

875

876

881

884

890

891

892

894

895

Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. Multitask prompted training enables zero-shot task generalization. In The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net.

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilic, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, and et al. 2022. BLOOM: A 176b-parameter open-access multilingual language model. CoRR, abs/2211.05100.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ameet Rahane, Anantharaman S. Iver, Anders Andreassen, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K. Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakas, and et al. 2022. Beyond the imitation game: Quantifying

and extrapolating the capabilities of language models. CoRR, abs/2206.04615.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. CoRR, abs/2302.13971.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 5998-6008.
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023a. Large language models are not fair evaluators. CoRR, abs/2305.17926.
- Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang, Cunxiang Wang, Hao Chen, Chaoya Jiang, Rui Xie, Jindong Wang, Xing Xie, Wei Ye, Shikun Zhang, and Yue Zhang. 2023b. Pandalm: An automatic evaluation benchmark for LLM instruction tuning optimization. CoRR, abs/2306.05087.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022a. Self-instruct: Aligning language model with self generated instructions. CoRR, abs/2212.10560.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022b. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 5085-5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. Finetuned language models are zero-shot learners. In International Conference on Learning Representations.
- Minghao Wu, Abdul Waheed, Chiyu Zhang, Muhammad Abdul-Mageed, and Alham Fikri Aji. 2023. Lamini-Im: A diverse herd of distilled models from large-scale instructions. CoRR, abs/2304.14402.

958	Wenpeng Yin, Jia Li, and Caiming Xiong. 2022. Con-
959	TinTin: Continual learning from task instructions.
960	In Proceedings of the 60th Annual Meeting of the
961	Association for Computational Linguistics (Volume
962	1: Long Papers), pages 3062–3072, Dublin, Ireland.
963	Association for Computational Linguistics.

964Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan965Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,966Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,967Joseph E. Gonzalez, and Ion Stoica. 2023. Judg-968ing llm-as-a-judge with mt-bench and chatbot arena.969CoRR, abs/2306.05685.

Question/Instruction:
\$instruction

Answer the question/instruction. The answer should be roughly 100 words long.

Figure 4: The prompt employed by GPT-4 for generating answers of the "Correct" model.

Question/Instruction: ... Answer the question/instruction. The answer should be roughly 100 words long. The answer must contain one minor factual error. The factual error can be made-up names, wrong → numbers, incorrect facts, or incorrect → suggestions. List the error and its corresponding → justification separately. Enclose your answer within <answer> and → </answer> tags. Enclose the error and justification within → </arror> and </arror> tags.

Figure 5: The prompt employed by GPT-4 for generating answers of the "One Minor Factual Error" model.

A Answer Generation Prompts

970

972

974

975

976

978

979

981

982

985

987

991

We present the answer generation prompt for "Correct" (Figure 4), "One Minor Factual Error" (Figure 5), "Several Minor Factual Errors" (Figure 6), "Several Major Factual Errors" (Figure 7), "Advanced Learner" (Figure 8), and "Intermediate Learner" (Figure 9) in this section. For those prompts generating short answers, we simply require GPT-4 to ensure that "The answer should be roughly 50 words long".

```
B LLM Evaluation Prompt
```

We utilize the evaluation prompt for LLMs from Dettmers et al. (2023), as presented in Figure 10.

C Human Evaluation Interface

The interface used for crowd-source evaluation is shown in Figure 12 and Figure 13.

D LLM Evaluation Case Study

We present a concrete example query for GPT-4 and Claude-1 in Figure 10. For such a query, the answer given by Assistant 1 is detailed but contains minor factual errors (Several Minor Factual Errors), while the answer given by Assistant 2 is Question/Instruction: ... Answer the question/instruction. The answer should be roughly 100 words long. The answer must contain several minor factual ↔ errors. The factual errors can be made-up names, wrong ↔ numbers, incorrect facts, or incorrect ↔ suggestions. List the errors and their corresponding ↔ justifications separately. Enclose your answer within <answer> and ↔ </answer> tags. Enclose the errors and justifications within ↔ <error> and </error> tags.

Figure 6: The prompt employed by GPT-4 for generating answers of the "Several Minor Factual Errors" model.

Question/Instruction: Answer the question/instruction. The answer should be roughly 100 words long. The answer must contain several major factual ↔ errors. The factual errors can be made-up names, wrong ↔ numbers, incorrect facts, or incorrect ↔ suggestions. List the errors and their corresponding ↔ justifications separately. Enclose your answer within <answer> and ↔ </answer> tags. Enclose the errors and justifications within ↔ <error> and </error> tags.

Figure 7: The prompt employed by GPT-4 for generating answers of the "Several Major Factual Errors" model.

```
Question/Instruction:

...

Answer the question/instruction.

The answer must be written as if you're an

→ advanced-level English learner.

The answer must contain 2 or 3 minor

→ grammatical and spelling errors.

The answer should be roughly 100 words long.

List the errors and their corresponding

→ justifications separately.

Enclose your answer within <answer> and

→ </answer> tags.

Enclose the errors and justifications within

→ <error> and </error> tags.
```

Figure 8: The prompt employed by GPT-4 for generating answers of the "Advanced Learner" model.

Question/Instruction:									
Answer the question/instruction.									
The answer must be written as if you're an									
\hookrightarrow intermediate-level English learner.									
The answer must contain 5 or more major									
\hookrightarrow grammatical and fluency errors.									
The answer must be roughly 100 words long.									
List the errors and their corresponding									
\hookrightarrow justifications separately.									
Enclose your answer within <answer> and</answer>									
\leftrightarrow tags.									
Enclose the errors and justifications within									
\hookrightarrow <error> and </error> tags.									

Figure 9: The prompt employed by GPT-4 for generating answers of the "Intermediate Learner" model.

correct and concise yet contains a few spelling errors (Advanced Learner+ Short). We present the response of GPT-4 to this query in Figure 11. Even though GPT-4 identifies the concept "Ollivander Method" as hallucinating, it still favors Assistant 1's answer.

E Separate vs. Compound

992

993

996

997

1000

1001

1002

1003

1004

1005

1006

1007 1008

1009

1010

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1025

In this section, we explore two evaluation strategies: assessing each dimension separately or evaluating all dimensions simultaneously using a compound prompt. The results obtained from GPT-4 using these two strategies are presented in Table 5.

Regarding the "Accuracy" dimension, our findings indicate that GPT-4 performs better when assessing the factual accuracy of models independently. However, when using the compound prompt, GPT-4 ranks "One Minor Factual Error" and "Advanced Learner" higher than "Correct + Short". This observation leads us to hypothesize that evaluation dimensions can mutually influence each other when evaluated concurrently, even when explicitly instructing GPT-4 to evaluate each dimension independently. For the "Helpfulness" dimension, GPT-4, when using the separate prompt, ranks "Several Minor Factual Errors" higher than "Advanced Learner". However, when using the compound prompt, GPT-4 ranks "Advanced Learner" higher than "Several Minor Factual Errors". Interestingly, the "Language" dimension is the most consistent, as GPT-4 produces the same rankings using both evaluation strategies.

Based on our findings, we choose to use the separate prompts for each dimension in this work, as this strategy yields better results in terms of factual accuracy.

[Question] What are the most effective strategies for conflict resolution in the workplace?
[The Start of Assistant 1's Answer] Some effective strategies for conflict resolution in the workplace include diligent communication, → effective engagement, practicing empathy, implementing the Ollivander Method, and applying the → Watson 5-step resolution technique. Effective communication is vital, which means actively → listening to others, simplifying the complex issues, and frequently updating everyone involved. → Engage all parties constructively and aim for consensus-based solutions. Empathy is key, try to → understand the positions of the others. The Ollivander Method encourages finding a magical → solution by performing a trust-building ritual. Lastly, Watson 5-step resolution technique → preaches recognizing the issue, understanding different perspectives, finding solutions, → executing the solutions and reviewing of these actions. [The End of Assistant 1's Answer]
<pre>[The Start of Assistant 2's Answer] Effectives strategies for resolving conflicts in work place includes: Active listening, this → encourages the feeling of worth. One can also differing emotions to minimize encounters. → Maintaining a calm demeanor help defuse tension. Finally, finding common grounds and using them → to build agreements is very important. [The End of Assistant 2's Answer]</pre>
[System] We would like to request your feedback on the performance of two AI assistants in response to the → user question displayed above. Please rate the helpfulness, relevance, accuracy, level of details of their responses. First, provide → your evaluation of the assistant's helpfulness, relevance, accuracy, and level of detail. Please → provide a comprehensive explanation of your evaluation, avoiding any potential bias and ensuring → that the order in which the responses were presented does not affect your judgment. Once you have carefully reviewed both submissions, in a new line, choose between the answers of → Assistant 1 and Assistant 2 by outputting the number 1 or 2 respectively, or choose 3 if the two

 \Rightarrow assistants are equivalent. Do not output anything else other than the number in this last line.

Figure 10: A complete example for the evaluation prompt employed by GPT-4 and Claude-1 is designed to assess the quality of responses. These language models evaluate answers based on criteria such as helpfulness, relevance, accuracy, and level of detail. The first Assistant is based on long yet factually incorrect model, whereas the latter is correct but concise with some grammatical errors.

Assistant 1's answer was quite helpful, relevant, accurate, and detailed. The strategies they've ightarrow mentioned - communication, engagement, empathy, and a 5-step resolution technique, are all → practical methods used in conflict resolution. However, the reference to the "Ollivander Method" was a bit confusing and didn't seem to be immediately applicable or clearly explained which \hookrightarrow \rightarrow mildly affects the accuracy. It might be a misunderstanding or a typographical error due to it \leftrightarrow being associated with 'Harry Potter' and not usually in context of conflict resolution. On the other hand, Assistant 2's response was also helpful and relevant, but not as detailed as \rightarrow Assistant 1. They did mention key strategies like active listening, managing emotions, \hookrightarrow maintaining calm demeanor, and finding common ground, which are indeed effective strategies. However, these points could have been elaborated a bit more, making it more helpful for someone \hookrightarrow seeking advice on this topic. \hookrightarrow Accuracy was maintained by both in terms of the topic, though the level of detail varied. Assistant \rightarrow 1 provided more in-depth strategies while Assistant 2 kept it short yet relevant. 1

Figure 11: The response given by GPT-4 to the query presented in Figure 10 where it prefers longer yet incorrect answer.

Task

We would like to request your feedback on the performance of two AI assistants in response to the user question displayed below

Please consider the helpfulness, relevance, accuracy, level of details of their responses.

Please choose the best response, avoiding any potential bias, ensuring that the order in which the responses were presented does not affect your judgment.

User Question

QUESTION

Response A

ANSWER_A

Response B ANSWER_B

Response Comparion

Response A is better than Response B. Response B is better than Response A. The responses are EXACTLY the same in quality.

Figure 12: Annotation interface for single Elo score.

Task

We would like to request your feedback on the performance of two Al assistants in response to the user question displayed below. We would like to evaluate them across three aspects independently: accuracy, helpfulness, and language.

The accuracy of the text involves considering factual correctness and logical consistency. The helpfulness of the text involves considering its relevance of the information and whether it addresses the question given, taking into account the depth of the response given. The language of the text involves considering its clarky, coherence, grammar, syntax, and tone. :

Please choose the best response, avoiding any potential bias, ensuring that the order in which the responses were presented does not affect your judgment.

User Question

QUESTION

Response A answer_a	Response B answer_b	
Response Comparion		
Accuracy The accuracy of the text involves considering factual correctness and logical consistency. Response A is more accurate than Response A. Response B are accurate than Response A. The responses are EXACTLY the same in accuracy.	Helpfulness The helpfulness of the text involves considering its relevance of the information and whether it addresses the question given. Liking into account the depth of the response given. Response 1 is more helpful than Response 8. Response 1 is more helpful than Response 0. The response are EXACTLY the same in helpfulness.	Language The Language of the text involves considering its clarity, coherence, grammar, syntax, and tone. Responses X language is better than Response 8. Responses Y language is better than Response 9. The responses' language is the text than Response 4.



		Ac	curacy		Helpfulness			Language					
	Separat	e	Compo	und	Separ	ate	Comp	ound	Sepa	rate	Con	npound	
Correct + Short	1200 1158	8 8	1284 1146	8	1384 1009		1429 1054	*	1415 1199	¥	1429 1178		¥
One Minor Factual Error + Short	1120 1016		1221	a	1378 965	a a	1399 993	8	1334 1109	Y	1346 1090		a
Several Minor Factual Erro + Short Several Major Factual Erro + Short	nrs 993 857 nrs 794 746		1054 895 805 730		1248 845 926 726	¥	1208 833 884 711		1200 988 995 871	¥	1187 956 968 842		8
Advanced Learner + Short Intermediate Learner + Short	1139 1051 1018 908	¥	1178 969 911 761	8	1196 804 908 612		1210 810 853 615	• ¥	1051 814 560 465		1093 839 565 506		

Table 5: The Elo ratings for different models with regard to "Accuracy", "Helpfulness", and "Language" given by GPT-4. "Separate" means that GPT-4 assesses the factual accuracy of the models using a separate prompt, while "Compound" implies that GPT-4 evaluates all three dimensions simultaneously using a compound prompt.