Towards Better Multilingual Side-by-Side LLM Evaluation

Anonymous ACL submission

Abstract

With the rise of large language models, evaluating their outputs has become increasingly important. While supervised evaluation compares model responses to ground truths, dialogue models often use the Side-by-Side approach, where a judge compares the responses of baseline and candidate models using a predefined methodology. In this paper, we conduct an indepth analysis of the Side-by-Side approach for evaluating models in Russian, Arabic, as well as for code generation and investigate the circumstances under which LLM-evaluators can be considered an alternative to expert annotation. We propose and publicly release a methodology that can enhance the correlation between automatic evaluation and human annotation through careful prompt engineering and adding model reasoning. We demonstrate the problem of positional bias and propose metrics for measuring it, as well as ways to mitigate it.

1 Introduction

003

007

800

014

017

018

019

037

041

As large language models (LLMs) rapidly advance, evaluating them effectively has become a crucial task that can be approached from various angles. Evaluation methods for these models are typically divided into supervised and unsupervised approaches. The supervised method involves comparing the model's responses to ground-truth answers. Such methods imply a straightforward output format, where the model is required to classify, select, match options, or generate a short answer, after which automatic metrics like accuracy, exact match (EM), and F1-score are used. The model's abilities in tasks such as question answering, common sense, and reasoning are tested in this way.

However, for models intended for interacting with users, providing a perfect answer to every query is not always possible. In these situations, the Side-by-Side (SbS) approach is frequently employed, where an independent judge compares the responses of a candidate model with those of a baseline model. The comparative element in this method helps avoid bias in judges' evaluations while allowing for the assessment of the overall quality of the dialogue agent's response. 042

043

044

047

048

053

054

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

078

079

Due to the indeterminacy and inconsistency of the evaluation criteria for this method, we decided to explore its characteristics using the example of evaluating models in different languages. In this paper, we examine SbS evaluation by comparing its manual execution with execution using an LLMbased evaluator, and also present ways to improve this approach. We aim to answer two main research questions:

- 1. Is the issue of positional bias still relevant? How can we address it in SbS evaluation?
- 2. How does the formulation of the prompt for the LLM-as-judge help, and to what extent?

We propose a methodology that can increase the correlation between model-as-judge assessments and human assessments while exploring ways to significantly improve the performance of a judge with a relatively small number of parameters. We demonstrate that minor changes in the task formulation for the evaluator model can significantly enhance the quality of its evaluation. Our comparative analysis of various open and closed commercial models using our benchmark helps us assess the impact of prompt-engineering techniques on the quality of evaluation.

2 Related Works

Prior to the emergence of robust large language models, evaluations of natural language generation systems often relied on automated metrics such as BertScore (Zhang et al., 2020) and GPTScore (Fu et al., 2023). Although these metrics offer scalability, they do not fully capture the subtlety and context sensitivity that human judgments can

130

provide. Human evaluators, traditionally considered the "gold standard" for the assessment of NLG (Ouyang et al., 2022), remain critical for tasks that require deep linguistic and domain expertise. However, human evaluations introduce issues of subjectivity, potential biases, and reproducibility challenges (Clark et al., 2021; Belz et al., 2023). They are also time-consuming, resource-intensive, and limited by the slower processing speed of human annotators. As a result, leveraging powerful LLMs (e.g., GPT-4) to approximate or even replace human annotators has gained prominence (Zheng et al., 2023; Chiang et al., 2023; Liu et al., 2023; Zhu et al., 2023). These LLM-based evaluators can achieve high agreement with human preferences, yet they too exhibit specific biases, such as position bias, verbosity bias, and self-enhancement bias (Zheng et al., 2023). To address these shortcomings, researchers have introduced strategies such as generating chain of thought plans (Wei et al., 2023) for more transparent evaluations. However, this technique has limited effectiveness for tasks that do not involve mathematical or logical reasoning (Sprague et al., 2024).

081

094

100

103

104

105

106

107

110

111

112

113

114

115

116

117

118

The emergence of "thinking" models, which incorporate reasoning processes before delivering final outputs (Wu et al., 2024; Guo et al., 2025; OpenAI, 2024b), marks a significant advancement in the evaluation of other large language models (LLMs) (Hosseini et al., 2024; Saha et al., 2025). Our research builds upon these developments by focusing on the application of "thinking" models specifically designed to evaluate other LLMs. By examining the alignment of these models with human judgment, we seek to assess their accuracy, fairness, and transparency compared to traditional metrics.

3 Methodology

3.1 Data collection

We manually collect datasets for each of the SbS 119 setups: Russian, Arabic, and code generation. The first dataset contains 2396 instructions, each writ-121 ten manually, taking into account the specifics of 122 Russian culture and patterns characteristic of Rus-123 sian users. The instructions are classified by task 124 125 type, including question answering, creative writing, information extraction, summarization, clas-126 sification, and chat-based questions. For those in-127 structions requiring factual accuracy in responses, we also gathered factual references. Later, these ref-129

erences help judges make correct decisions based on the provided knowledge sources.

For the second setup, we asked native Arabic speakers to translate and adapt the instructions from this dataset to fit Arabic culture.

We also manually collect a dataset of 800 instructions with various coding tasks for comparing code model generations. The questions in the dataset are divided into five categories: writing docstrings, creating unit tests, text-to-code, refactoring and explaining a piece of code.

3.2 Side-by-Side method

We utilize a Side-by-Side method in which human experts and large language models serve as judges. We generate responses to instructions from the collected dataset for two models — the baseline and the candidate. After that the selected judge compares pairs of responses and delivers a verdict.

In many studies using a similar evaluation method, the judge performs a binary classification, indicating either that the response of the first model is better or the second one is. We decided to expand the set of options, implying that both responses could be equally good or equally bad. Thus in our case, the judge states that either **a**) whether the response from the candidate model is better than the baseline, **b**) vice versa, **c**) both responses are good or **d**) both responses are bad. The names of the models are concealed from the judges.

Naturally, this approach needs to be formalized to standardize the evaluations. With a well-defined task and properly specified criteria, we aim to align the model-based assessment results as closely as possible with human annotations. In the next section, we describe the design of our evaluation methodology.

3.2.1 Manual evaluation

We prepare instructions for the human experts to evaluate pairs of model responses. To determine which response is superior, each pair is assessed and compared against several criteria, listed in order of decreasing importance:

- 1. **[Safety]** The response should not contain information that could harm an individual.
- 2. **[Ethics]** The response must adhere to ethical standards: it should not be rude, offensive, biased, or judgmental.
- 3. [**Truthfulness**] The response should not contain inaccurate or questionable statements. 178

272

273

227

228

9	The expert refers to the attached factual refer-
D	ence to verify the truthfulness of the response.

17

18

181

182

186

187

191

192

193

194

196

198

199

201

202

207

209

210

212

213

214

215

216

217

218

219

221

225

226

- 4. [**Relevance**] The response should align with the request: it must follow the instructions, avoid answering unnecessary questions, and be in the required language.
 - 5. [**Completeness**] The response should be thorough and comprehensive.
 - 6. **[Style]** The response should be written with correct spelling, punctuation, and syntax, and should avoid informal language, unless explicitly stated otherwise in the instructions.

The order of the model responses in each pair is randomized. The manual evaluation is conducted with an overlap of three people. In cases where the verdicts for a particular pair of responses did not match, they are sent for reassessment with discussion among the experts.

Background information about the team of experts, including details about their age and education, can be found in the Appendix A.

3.2.2 Automatic evaluation

We develop several prompts for LLM evaluators that take into account the criteria described in the previous section.

Instead of randomizing the order of model responses, we perform two runs through the dataset. In the first run, the prompt places the candidate model's response first followed by the baseline model's response; in the second run, the order is reversed. The scores are averaged after the two runs are completed. We could shuffle the model responses within pairs for the LLM-judge input to save its runtime, as we do for human experts. However, conducting two separate runs allows us to analyze the presence of positional bias in the tested evaluators.

An important task in preparing the evaluator model is the preparation of the prompt. Prompt I is designed to succinctly describe the task of SbS evaluation. In Prompt II, we aim to address and describe all the criteria listed for the team of experts. We also attempt to add the following modifications to the prompt.

Reasoning

We ask the model to reflect before reaching a verdict, to analyze responses based on each criterion, and to aggregate scores when providing a comparison result. Some models have been specially trained to reason (Guo et al., 2025; OpenAI, 2024b), for which such an addition to the prompt presumably will not make any difference.

We find an issue with models trained on reasoning and those evaluated with reasoning prompts a significant portion of the answers (>10%) consists not of the expected symbols representing one of four classes, but a different response. Therefore, when using reasoning, we make the model strictly adhere to formatting.

Factology

We include factual information in the prompt and ask the model to rely on a knowledge source when evaluating responses where necessary. Factual references were collected from Wikipedia.

Multi-agent approach

The reasoning of the evaluator's language model really helps improve performance when evaluating responses from other models. However, despite the advantages of the Chain-of-thought (CoT) method, when the model reasons step by step, there is a problem called Degeneration-of-thought (Liang et al., 2023), when the LLM begins to be confident in its reasoning, even if it is not correct. Authors provide an example of a multi-agent approach that avoids this problem. To do this, agent-1 expresses his opinion on a task, agent-2 responds to this, and after the agents' dialogue, the agent-judge analyzes the agents' responses and issues a final verdict.

Based on this research, we propose the following two schemes of a multi-agent approach.

- 1. **Soft**. Agent-1 makes his assessment regarding a pair of proposals, and agent-2 either agree or disagree with agent-1. Next, the agent-judge makes his verdict based on the two previous verdicts.
- 2. Hard. Agent-1 makes his assessment regarding a couple of proposals, and agent-2 always disagrees with agent-1. After that the agentjudge makes his verdict based on the two previous verdicts.

All variations of the prompts can be found in Appendix B.

4 Experiments

For our experiments in Russian we select Qwen2.5-32B-Instruct (Yang et al., 2024) fine-tuned on the

T.,	daa	Davamatara	SBS 1	results				ADCC	MDCC	ՍԵՍոՒ	Citation
Ju	uge	rarameters	Α	В	С	D	Е	AFCC	MFCC	пг пир	Citation
	manual	21 experts	4.6	44.5	23.7	27.2	0.0				
s	llama3.1-405b	405B	36.6	30.9	31.6	0.8	0.0	-0.229	0.507	link	(Dubey et al., 2024)
del	llama3.3-70b	70B	34.1	39.0	26.0	1.0	0.0	0.041	0.595	link	(Dubey et al., 2024)
ŌШ	gpt-40	-	17.6	13.6	46.0	22.7	0.0	-0.144	0.495	-	(OpenAI, 2024a)
eq	o1-mini	-	32.7	40.3	18.1	9.0	0.0	0.165	-	-	(OpenAI, 2024b)
ğ	gpt4	-	23.8	24.0	51.3	0.8	0.1	-0.081	0.642	-	(Achiam et al., 2023)
-ali	deepseek-r1-dst.	70B	24.2	33.5	31.9	7.3	3.1	0.227	0.640	link	(Guo et al., 2025)
lish	deepseek-v3	671B (37B)	11.5	47.7	34.9	1.0	4.8	0.624	0.570	link	(Liu et al., 2024)
[gu	claude sonnet	175B	13.6	9.1	71.1	1.0	5.2	-0.122	0.427	-	(Anthropic, 2024)
щ	claude opus	137B	23.3	35.0	20.2	21.4	0.2	0.681	0.602	-	(Anthropic, 2024)
_	T-lite-it-1.0	7.6B	18.4	26.7	39.5	15.4	0.0	0.238	0.139	link	(T-bank, 2024)
siar	T-pro-it-1.0	32.8B	40.6	44.6	4.1	10.7	0.5	0.070	0.397	link	(T-bank, 2024)
sus	GigaChat-Max	70-100B	24.9	38.9	30.7	0.9	4.6	0.265	-	-	(Sber, 2024)
R	YandexGPT	-	29.8	43.2	7.4	0.2	19.5	0.231	0.572	-	(Yandex, 2024)

Table 1: **Comparative analysis of LLMs as judges for SbS Evaluation in Russian.** Various models of different sizes, aligned with both English and Russian languages, were selected as judges. Prompt I was used for obtaining verdicts. The percentage distribution of verdicts across the entire benchmark is presented by symbols: A) the candidate model's answer is better, B) the baseline model's answer is better, C) both models' answers are equally good, D) both models' answers are equally poor. For the LLM evaluators, the average value for each verdict across two benchmark runs is provided. Additionally, we include the APCC with expert assessments, where all verdict are aggregated by verdict class and MPCC, where we use a sliding window to go through the verdicts and calculate the median for each batch.

Russian language as the candidate model and GPT-40 (OpenAI, 2024a) as the baseline model. For the Arabic language we use Llama-3.0-70b (Dubey et al., 2024) fine-tuned on the Arabic language as the candidate model. The parameters for generating responses on the benchmark are the same for all models, their values can be found in the Appendix. The paired generations are shuffled and given to a team of experts for annotation (with the model names concealed) along with the evaluation methodology described in Section 3.2.2. These same generations are also evaluated by LLM judges.

274

276

277

281

282

289

294

297

298

301

4.1 Analysis of manual evaluation

We provide the expert evaluators with universal criteria for assessment through guideline; however, this does not guarantee full correlation among their responses. We believe it is expected and acceptable for annotators to have differing opinions when evaluating pairs of responses, which is precisely why our assessment involved an overlap.

The dataset is divided into parts consisting of 600 questions each, and each of them is evaluated independently by three different people. We calculate the Pearson correlation coefficient (PCC) between each pair of annotators for each of the four splits of the dataset. The PCC ranges from **0.667** to **0.937** depending on the dataset split. We conclude that if the correlation of any model evaluator falls within the specified range, it can likely be considered as a good judge option. The exact metrics can be found in Appendix A. 302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

4.2 Analysis of automatic evaluation

We select a range of models of different sizes for the evaluation, including both open-source and commercial models. The results of the manual and automatic evaluation in Russian and Arabic can be found in Table 1 and Table 5 respectively. We measure the correlation between the assessments of the models and the experts using the following metrics.

Aggregated Pearson Correlation Coefficient (APCC). We count how many verdicts fall into each class A/B/C/D and calculate the correlation between LLM-as-judges and experts assessments based on these four values. While calculating this metric, we lose a lot of information about individual verdicts, but we can estimate how close the model is to the experts in delivering a final verdict for the entire benchmark.

Median Pearson Correlation Coefficient (MPCC). We apply a sliding window with a size of 10 and a stride of 5 across all verdicts from the benchmark. For teach batch we calculate the median using formula: $Median = \frac{\sum A + \sum C}{\sum A + \sum B + 2 \cdot \sum C}$. We obtain a set of medians for the expert and model

Judge model	MPCC		PCon@AB	
Judge mouel	Cons.	Δ		
llama3.1-405b	0.526	0.058	0.336	
llama3.3-70b	0.599	0.249	0.476	
gpt-40	0.666	0.035	0.329	
gpt4	0.674	0.259	0.339	
deepseek-r1-dst.	0.846	0.012	0.500	
deepseek-v3	0.164	0.077	0.271	
claude sonnet	-0.275	0.241	0.106	
claude opus	0.598	0.125	0.431	
T-lite-it-1.0	-0.370	0.093	0.122	
T-pro-it-1.0	0.363	0.179	0.400	
YandexGPT	0.319	0.092	0.409	

Table 2: Comparative analysis of evaluator scores with and without swap of models' answers. Metrics MPCC-Consistency, MPCC- Δ and PCon@AB indicate the presence of positional bias among LLMevaluators. The closer the values of metrics MPCC-Consistency and PCon@AB are to one, the more consistent the model is when the positions of answers in prompt are changed; while lower MPCC- Δ indicates lower positional bias.

verdicts and calculate the PCC between them. This method retains almost all information for all verdicts but imposes a linear relationship between verdict classes, which may be somewhat incorrect to establish.

Both metrics are averaged over two runs: one with the direct order of responses in the prompt and the other with the reverse order. Overall we consider both metrics APCC and MPCC to assess the correlation between LLM and expert verdicts.

We suggest looking not only at the correlation coefficients but also at the proportions of verdict returned by the judges. In addition to high correlation with manual evaluation, it is important for the LLM to replicate significant statistical patterns. For example, in Table 1 according to expert judgment, we can see that the baseline model answers better significantly more often than the candidate model. For many evaluator models, however, the number of positive (A) statements is often close to the number of negative (B) statements. Judging by both APCC and BPCC we conclude that Claude Opus and Deepseek-r1-dst - distillation of Deepseek-r1 into Llama-70B - show the best correlation with manual assessments among all the tested LLM-as-judges for the Russian language.

Promnt	SBS 1	results			APCC	MPCC	PCon@AB
Tompt	A	В	С	D	in ee	in ee	reonend
experts	4.6	44.4	23.7	27.2			
I	24.2	33.5	31.9	7.3	0.226	0.589	0.500
П	14.5	26.0	44.1	15.4	0.277	0.623	0.361
II-fact	14.5	26.1	44.2	15.1	0.275	0.606	0.355
II-reason	17.9	29.4	43.9	8.4	0.226	0.639	0.412
II-fact+reason	17.9	28.9	43.9	8.9	0.218	0.636	0.384

Table 3: Analysis of deepseek-r1-distill-llama judge model with different prompts. The proportion of responses and the PCC with expert evaluation are provided for Prompt I, Prompt II, as well as variations of Prompt II with the additions of factual background and reasoning.

Promnt	SBS	results			APCC	MPCC	PCon@AB	
1 tompt	A	В	С	D			i concento	
experts	7.6	41.4	23.7	27.2				
I	13.1	20.3	63.8	2.8	0.041	0.573	0.573	
П	9.0	11.5	57.5	21.9	0.014	0.505	0.146	
II-fact	9.0	11.5	57.6	21.9	0.013	0.504	0.145	
II-reason	31.5	31.6	32.0	4.7	-0.089	0.653	0.499	
II-fact+reason	30.3	31.7	33.5	4.2	-0.051	0.639	0.497	

Table 4: Analysis of llama3.3-70b judge model with different prompts. The proportion of responses and the PCC with expert evaluation are provided for Prompt I, Prompt II, as well as variations of Prompt II with the additions of factual background and reasoning.

4.3 Impact of positional bias

Table 10 presents a study of LLM judges for positional bias. We perform two measurements for each evaluator model - without models' answers swap and with - and calculate the PCC of aggregated values with manual annotation. We introduce metric **PCon@AB** that indicate the presence of bias in the evaluator models.

PCon@AB =		36
$\sum_{BM} \mathbb{1}(J_{\text{swap}=0} = J_{\text{swap}=1} J = \mathbf{A} \vee \mathbf{B})$	(1)	36
$\sum_{BM} \mathbb{1} \left((J_{\text{swap}=0} = \mathbf{A} \lor \mathbf{B}) \lor (J_{\text{swap}=1} = \mathbf{A} \lor \mathbf{B}) \right)^{\cdot}$	(1)	50

356

357

359

360

361

362

363

366

367

368

369

370

371

372

373

374

376

377

This metric shows the consistency of the model's answers without swap and with - it indicates the proportion of matching answers among answers A and B given the different order of model responses.

The metric **MPCC-Consistency** is calculated as the Pearson correlation coefficient between two sets of medians obtained for the verdicts with and without swap, while the metric **MPCC-** Δ is the difference between the MPCC calculated separately for the verdicts obtained with and without swap.

PCon@AB, **MPCC-Consistency** and **MPCC-** Δ do not rely on manual annotation,

355

332

Judge model	verdi	cts			APCC	MPCC	PCon@AB
Judge model	A	В	С	D	in cc	in ee	rement
manual	27.6	14.8	46.4	11.2			
llama3.1-405b	21.6	13.6	37.8	26.8	0.735	0.390	0.464
llama3.3-70b	15.7	7.8	39.8	36.5	0.413	0.337	0.306
gpt-40	2.4	3.1	20.4	74.1	-0.379	-0.153	0.112
deepseek-r1-dst.	23.1	15.4	36.8	24.5	0.835	0.369	0.354
deepseek-v3	4.9	3.0	19.4	72.7	-0.390	0.280	N/A
claude sonnet	21.7	16.4	21.3	40.5	-0.423	0.259	0.448
claude opus	22.0	15.3	32.3	30.1	0.485	0.347	0.168

Table 5: Comparative analysis of LLMs as judges for SbS Evaluation in Arabic. Various models of different sizes were selected as judges. Prompt II was used for obtaining verdicts. The average value for each verdict across two benchmark runs and PCC with expert assessments is provided.

allowing us to determine how prone the model is to positional bias without expert involvement.

4.4 **Elevating LLM-as-judge performance**

In this section, we address two questions: a) how much can we increase the correlation with manual annotation by constructing prompts? b) can prompts help with the positional bias issue?

As suggested in Section 3.2.2, we create several prompt variations and measure two LLMs-asjudges with each: Deepseek-r1-distill-llama and Llama3.3-70b. Table 3 shows the comparison for the first model - after updating prompt I to II, the correlation with manual annotation significantly increased. However, modifying prompt II by adding factual information and reasoning cause the correlation to decrease. This is expected when requesting the model to reason - Deepseek-r1-distill-llama is already trained to do reasoning, therefore the additional step is redundant.

The patterns do not hold for the model Llama3.3-70b, as can be seen in 4. Adding factual information and reasoning for this model significantly increases the correlation with experts. It is noteworthy that adding a request to reason in the prompt not only slightly increases correlation with experts but also significantly enhances the model's robustness against positional bias.

We formulate several conclusions that we consider foundational for our methodology based on the results of these experiments. We recommend them as guidelines for performing similar evaluations.

• While the issue of positional bias remains significant for LLM-as-a-judge in the SbS task, it can be almost entirely avoided by using

Judge model	verdi	cts			APCC	MPCC	PCon@AB	
	А	В	С	D				
manual	28.4	17.8	23.1	30.8				
llama3.1-405b	29.0	4.4	9.0	57.6	0.819	0.427	0.364	
llama3.3-70b	41.3	9.2	17.3	32.2	0.898	0.432	0.439	
deepseek-r1-dst.	35.4	19.3	25.6	19.7	0.343	0.460	0.424	
deepseek-v3	24.4	6.8	4.1	64.8	0.818	0.089	0.241	
claude sonnet	17.4	0.8	6.4	75.4	0.797	0.312	0.098	
claude opus	25.4	11.3	27.1	36.3	0.903	0.392	0.089	

Table 6: Comparative analysis of LLMs as judges for SbS Evaluation for code. Various models of different sizes were selected as judges. Prompt II was used for obtaining verdicts. The average value for each verdict across two benchmark runs and PCC with expert assessments is provided.

models trained to reason. For other models, the effect can also be reduced by asking the model to reason beforehand.

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

· A well-crafted prompt can significantly increase correlation, but the prompt should be tailored individually for each model as it is not transferable between different LLMas-judges. From Table 3, we see that as the complexity of the prompt increases, the correlation of the Deepseek-r1-dst-llama model with human labeling rises, nearly reaching the quality of a model more parameters (Claude opus).

4.5 SbS evaluation for Arabic

We conduct similar studies on the Arabic version of our benchmark with Prompt II for a range of evaluator models.

From the Table 5, it follows that models Deepseek-r1-distill-llama and Llama3.1-405b show the best correlation with human judgment, while models Llama3.1-405b and Claude Sonnet are the least prone to positional bias.

4.6 SbS evaluation for code generation

In addition to SbS assessments in two languages, we also conduct similar experiments for models intended for code generation, using Prompt II for LLM evaluators. From the Table 6 we see that Llama models have generally the best performance in terms of the APCC and MPCC metrics, while Claude Opus and Deepseek-r1-dst can still be considered as strong options.

From Tables 1, 5 and 6, we conclude that regardless of the language in which the SbS task is performed, there is an issue with positional bias in

379

406

407

408

409

410

411

412

LLMs-as-judges. For each language (and for the programming languages), there are different families of models that perform best at evaluating texts in that language, and there isn't a single model that excels in everything.

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479 480

481

482

483

484

Judge model manual T-pro-it-1.0	method	verdi	PCC				
		A	В	С	D	Е	Mean
manual		7.6	41.4	23.7	27.2	0.0	
T-pro-it-1.0	base soft hard	40.6 37.9 36.1	44.6 45.6 36.7	4.1 9.5 14.82	10.7 6.0 3.6	0.5 1.0 8.78	0.070 0.118 -0.041

Table 7: **Multi-agent approach.** Measurement results for the T-pro-it model as a judge. We managed to increase the correlation with manual annotation using the soft approach.

4.7 Analysis of multi-agent approach

As can be seen from the Table 7, the hard method shows weak correlation with human markup. This is most likely due to the fact that T-pro-it already gives a response similar to the markup in the first iteration, and the other agents (the agent who disagrees or agrees with the statement and the agent judge) confuse the verdict of the first agent.

At the same time, the soft method increases correlation with experts, since it is most likely that the second agent does not necessarily contradict, but sometimes complements the reasoning of the first agent, and the agent judge re-evaluates all statements based on previous reasoning. This variation of Multi-Agent Debate is a strong method that develops the idea of CoT.

5 Conclusion

In this work, we present a methodology for conducting Side-by-Side evaluations using language model evaluators, which we apply to compare open and closed commercial large language models as judges. We highlight the significance of the positional bias issue and propose metrics for its evaluation during automatic SbS assessments, as well as suggest methods for mitigating its impact.

Additionally, we suggest ways to make language model evaluations align better with human ratings. This involves demonstrating the importance of prompts in conducting evaluations using our methodology, and emphasize the need for a tailored approach to crafting prompts for each evaluator model. We also assess the impact of adding factual information and reasoning on the judging model's capabilities and its influence on correlation485with manual annotations.486Limitations487

We acknowledge that there are other strong open 488 and proprietary models that we do not consider as 489 evaluators for the SbS task. We also do not research 490 on biases in expert assessments; there are patterns 491 in the candidate and baseline models' responses 492 inherent to these models that could lead experts to 493 guess which model produced a given response. Ad-494 ditionally, humans can also have positional biases. 495

References

496

497

498

499

502

503

504

505

510

511

512

513

514

515

516

517

518

521

523

524

525

527

530

531

532

534

535 536

537

538

539

541

542

544

546

547

548

549

552

553

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Anthropic. 2024. Introducing the next generation of claude. Available at: https://www.anthropic.com/news/claude-3-family.
- Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in nlp. Preprint, arXiv:2305.01633.
 - 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.
 - 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.
 - Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
 - Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. *Preprint*, arXiv:2302.04166.
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. arXiv preprint arXiv:2501.12948.
 - Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh Agarwal. 2024. V-star: Training verifiers for self-taught reasoners. *Preprint*, arXiv:2402.06457.

- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2023. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024.
 Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 2511–2522, Singapore. Association for Computational Linguistics.
- OpenAI. 2024a. Hello gpt-4o. Available at: https: //openai.com/index/hello-gpt-4o/.
- OpenAI. 2024b. Openai o1-mini. Available at: https://openai.com/index/ introducing-openai-o1-preview/.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.
- Swarnadeep Saha, Xian Li, Marjan Ghazvininejad, Jason Weston, and Tianlu Wang. 2025. Learning to plan reason for evaluation with thinking-llm-as-ajudge. *Preprint*, arXiv:2501.18099.
- Sber. 2024. Yandexgpt 4. Available at: https://giga. chat/.
- Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2024. To cot or not to cot? chain-ofthought helps mainly on math and symbolic reasoning. *Preprint*, arXiv:2409.12183.
- T-bank. 2024. T-lite and t-pro open russian-language open source models with 7 and 32 billion parameters. Available at: https://habr.com/ru/companies/ tbank/articles/865582/.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

589

590

592

595

596

598

599

600

601

602

603

604

605

606

554

555

- 60
- 6
- 611 612
- 613 614
- 615 616
- 617 618
- 619
- 6 6
- 623 624 625
- 626
- 627 628
- 629
- 63
- 631 632
- 633
- 635
- 636

637

638

640

641

642

645

- Tianhao Wu, Janice Lan, Weizhe Yuan, Jiantao Jiao, Jason Weston, and Sainbayar Sukhbaatar. 2024. Thinking llms: General instruction following with thought generation. *Preprint*, arXiv:2410.10630.
- Yandex. 2024. Yandexgpt 4. Available at: https://ya. ru/ai/gpt-4/.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. *Preprint*, arXiv:1904.09675.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Preprint*, arXiv:2306.05685.
- Lianghui Zhu, Xinggang Wang, and Xinlong Wang. 2023. Judgelm: Fine-tuned large language models are scalable judges. *Preprint*, arXiv:2310.17631.

A Appendix

The team of human experts consists of 21 Russianspeaking individuals, comprising 14 women and 7 men. The experts' ages range from 21 to 42 years, with a median age of 29. Eleven members have a higher education degree in linguistics, six have a degree in translation, and two have a degree in philology.

Benchmark split	Experts 1,2	Experts 2,3	Experts 1,3
0-599	0.989	0.809	0.844
600-1199	0.873	0.751	0.376
1200-1799	0.914	0.932	0.966
1800-2396	0.985	0.866	0.939

Table 8: Pearson correlation coefficient between each pair of annotators for each of the four splits of the benchmark.

B Appendix

B.1 Positional bias metrics

Table 10 presents a study of LLM judges for positional bias. We perform two measurements for each evaluator model- without models' answers swap and with - and calculate the PCC of aggregated values with manual annotation. We also calculate the number of matching verdicts (**accuracy**) and

Generation parameters	value
max_tokens	1024
temperature	0.0
top_p	0.1
frequency_penalty	1.0
vllm version	0.6.4

Table 9: **Parameters used for obtaining generations from LLMs-as-judges**. VLLM was used for inferencing open models from huggingface.

its difference between swaps (Δ), while also introducing metrics **PBias@AB**, **Con@ABCD** and **PCon@AB** that indicate the presence of bias in the evaluator models.

$$\mathbf{PBias} @\mathbf{AB} = 650$$

646

647

648

649

651

652

653

654

655

656

657

658

659

660

661

662

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

$$\sum_{\text{swap}=\{0,1\}} \sum_{BM} \mathbb{1}(J = \mathbf{A} | J = \mathbf{A} \vee \mathbf{B}) - 1, \quad (2)$$

where *BM* represents all samples from the benchmark, *J* is the judge's verdict, and swap = $\{0, 1\}$ refers to the order of the test model answers in the prompt ($\{C, B\}$ and $\{B, C\}$ respectively). **PBias@AB** is from the interval (-1, 1), where the absolute value indicates the magnitude of the positional bias, and the sign indicates whether the positional bias is direct or reverse. The closer the value is to zero, the more unbiased the model is.

$$\mathbf{Con@ABCD} = \sum_{BM} \mathbb{1}(J_{\mathrm{swap}=0} = J_{\mathrm{swap}=1}). \quad (3)$$

Con@ABCD shows the consistency of the model's answers without swap and with swap — it indicates the proportion of matching answers given the different order of model responses.

C Appendix

C.1 Prompt I

prompt:

Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question below. Select the assistant that adheres to the user's instructions and responds to the question with higher quality. Your evaluation must rigorously consider factors such as helpfulness, relevance, accuracy, depth,

Judge model	$swap \qquad \begin{array}{c c c c c c c c c c c c c c c c c c c $	verdicts					PCC		PBias@AB	Con@ABCD	PCon@AB
Judge mouel		Mean	I Dias CIID	concribed	leonend						
	{C, B}	36.7	14.8	43.5	3.3	1.8	0.008	0.267	0.455	0.412	0.222
mqu	$\{B, C\}$	15.5	46.0	34.9	3.6	0.0	0.726	0.307	0.433	0.413	0.235
llomo2 1 405h	$\{C, B\}$	43.9	19.9	35.5	0.7	0.0	0.016	0 272	0.274	0.411	0.226
nama5.1-4050	$\{B, C\}$	29.3	42.0	27.8	1.0	0.0	0.528	0.272	0.274	0.411	0.550
llomo2 2 70h	$\{C, B\}$	23.3	39.0	36.6	1.0	0.0	0.560	0.405	0.172	0.542	0 476
nania5.5-700	$\{B, C\}$	44.8	39.0	15.3	0.8	0.0	0.250	0.405	-0.175	0.342	0.470
ant da	$\{C, B\}$	22.3	7.0	48.4	22.3	0.0	0.176	0.276	0.250	0.606	0.220
gpt-40	$\{B, C\}$	12.9	20.3	43.7	23.1	0.0	0.576	0.570	0.359	0.606	0.329
ant 1	$\{C, B\}$	25.8	9.1	63.9	1.1	0.1	0.068	0.314	0.272	0.547	0.220
gpt4	$\{B, C\}$	21.9	38.9	38.7	0.5	0.0	0.560		0.373	0.347	0.339
deepseek-r1-distill-llama	$\{C, B\}$	25.1	30.4	35.0	7.3	2.3	0.503	0 564	0.074	0.574	0.500
	$\{B, C\}$	23.3	36.6	28.8	7.3	4.1	0.625	0.504	0.074	0.374	0.300
doonsoolt y?	$\{C, B\}$	0.7	64.3	32.2	0.8	2.2	0.784	0.558	0.477	0.276	0.271
deepseek-v5	$\{B, C\}$	22.3	31.2	37.7	1.4	7.4	0.392		-0.477	0.370	0.271
alauda connat	$\{C, B\}$	7.8	14.9	72.1	0.7	4.4	0.209	0.100	0.480	0.621	0.106
claude sonnet	$\{B, C\}$	19.5	3.3	70.0	1.3	5.9	-0.009	0.100	-0.480	0.031	0.106
alauda onus	$\{C, B\}$	19.3	41.7	23.1	15.6	0.3	0.866	0.766	0.170	0.508	0.431
claude opus	$\{B, C\}$	27.3	28.2	17.2	27.1	0.2	0.666	0.700	-0.170	0.508	0.431
T lite it 1.0	$\{C, B\}$	6.1	36.1	46.7	11.2	0.0	0.716	0.504	0.512	0.207	0.122
1-1110-11-110	$\{B, C\}$	30.8	17.2	32.3	19.6	0.0	0.292	0.304	-0.312	0.397	0.122
Three it 1.0	$\{C, B\}$	56.1	28.4	3.9	11.1	0.5	-0.01	0.226	0.296	0.454	0.400
1-p10-11-1.0	$\{B, C\}$	24.0	60.8	4.3	10.3	0.6	0.674	0.330	0.380	0.434	0.400
	$\{C, B\}$	24.9	38.9	30.7	0.9	4.6	0.25	0.25			
OrgaChat-Iviax	$\{B, C\}$							0.23			
VandayCPT	$\{C, B\}$	34.9	38.9	8.7	0.3	17.2	0.084	0 166	0.140	0.480	0.400
	$\{B, C\}$	24.6	47.4	6.0	0.2	21.8	0.247	0.100	0.140	0.400	0.409

Table 10: **Comparative analysis of evaluator scores with and without swap of models' answers.** Metrics **Accuracy**- Δ , **PBias@AB**, **Con@ABCD** and **PCon@AB** indicate the presence of positional bias among LLM-evaluators. The values of metrics **Accuracy**- Δ and **PBias@AB** being close to zero indicate an absence of positional bias, while the values of metrics **Con@ABCD** and **PCon@AB**, on the contrary — the closer they are to one, the more consistent the model is when the positions of prompt answers are changed.

680 681 682 683 684 685 686 687 688 689 690 691	creativity, and detail of the responses. Avoid any biases based on the position or order of responses to ensure an unbiased decision. The length of responses should not affect your evaluation. Maintain objectivity and neutrality towards assistant names. Output exactly one of the following symbols: A, B, C, or D. Use the following criteria strictly:
691 692 693 694 695 696 697 698 699 700	 Output 'A' if the first response is notably better. Output 'B' if the second response is notably better. Output 'C' if both responses are equally good. Output 'D' if both responses are equally inadequate.
701 702 703 704 705 706	<pre>IMPORTANT!! Be strict in your assessment : if there are errors in both responses, even minor ones, treat responses as inadequate and output ' D'.</pre>
707	instruction:
709	Question: {question}
710 711	First Response: {answer1}
712	Second Response: {answer2}
713	Answer:
714	C.2 Prompt II
714 715	C.2 Prompt II <pre>prompt :</pre>
714 715 716 717	C.2 Prompt II prompt: Please act as an objective and strict
714 715 716 717 718	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of</pre>
714 715 716 717 718 719	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's</pre>
714 715 716 717 718 719 720 721	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality</pre>
714 715 716 717 718 719 720 721 722	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response</pre>
714 715 716 717 718 719 720 721 722 723	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to</pre>
714 715 716 717 718 719 720 721 722 723 724	C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides
714 715 716 717 718 719 720 721 722 723 724 725 726	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727	C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step.
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step.</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 720	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria:</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly.</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria,</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 726 727 728 729 730 731 732 733 734 735 736	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate. 2. Ethics: Responses must be free from</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 730 731 732 733 734 735 736 737	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate. 2. Ethics: Responses must be free from rudeness, offensiveness, bias, or</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate. 2. Ethics: Responses must be free from rudeness, offensiveness, bias, or evaluative language. Any ethical violations automatically make the</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 731 732 733 734 735 736 737 738 739 740	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate. 2. Ethics: Responses must be free from rudeness, offensiveness, bias, or evaluative language. Any ethical violations automatically make the response unacceptable.</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 730 731 732 733 734 735 736 737 738 739 740 741	<pre>C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate. 2. Ethics: Responses must be free from rudeness, offensiveness, bias, or evaluative language. Any ethical violations automatically make the response unacceptable. 3. Truthfulness: The response must be</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 730 731 732 733 734 735 736 737 738 739 740 741 742	<pre>C.2 PromptII prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate. 2. Ethics: Responses must be free from rudeness, offensiveness, bias, or evaluative language. Any ethical violations automatically make the response unacceptable. 3. Truthfulness: The response must be factually accurate and aligned with the areavided factual information an area factual and aligned with the areavided factual information an and the factual and aligned with the areavided factual information and the factual and aligned with the areavided factual information and strict and the provided factual information and step-by-step and and</pre>
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744	 C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate. Ethics: Responses must be free from rudeness, offensiveness, bias, or evaluative language. Any ethical violations automatically make the response unacceptable. Truthfulness: The response must be factually accurate and aligned with the provided factual information. Any false or unverifiable claims
714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745	 C.2 Prompt II prompt: Please act as an objective and strict judge, evaluating the responses of two AI assistants to the user's question based on the provided factual information and strict quality standards. Assess each response against the following criteria to determine which assistant provides the best overall answer. Your evaluation must be thorough, reasoned, and step-by-step. Evaluation Criteria: 1. Safety: The response must not pose any harm, directly or indirectly. Even if it meets all other criteria, any unsafe content renders the response inadequate. 2. Ethics: Responses must be free from rudeness, offensiveness, bias, or evaluative language. Any ethical violations automatically make the response unacceptable. 3. Truthfulness: The response must be factually accurate and aligned with the provided factual information. Any false or unverifiable claims should be considered critical errors

4. Adequacy to the Request: The response must fully address the user's query without unnecessary deviations. It should adhere to specific instructions such as style, tone, and language. Failure to meet these requirements makes the response inadequate.

747

748

749

750

751

752 753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777 778

779

780

781 782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800 801

802

803

804

805 806

807

808

809

810

811

812

813

814

815

816

- 5. Completeness: The response should cover all relevant aspects of the query in a single reply, avoiding the need for follow-ups or additional clarifications.
- 6. Style: The response should be clear, well-structured, and professionally written. Poor readability, incoherence, or inappropriate style should result in a lower evaluation.
- Evaluation Method:

- Maintain objectivity, avoiding bias based on response position or length
- Penalize any response that fails to meet the standards outlined above.
- Compare each response to the factual information and the above criteria.
- Explicitly describe your train of thought for each criterion, explaining why one response is better than the other or if they are equal.

Decision Rules:

- Output '[[A]]' if the first response is clearly superior across all criteria.
- Output '[[B]]' if the second response is clearly superior across all criteria.
- Output '[[C]]' if both responses are equally good and fully meet the criteria.
- Output '[[D]]' if either response contains any factual inaccuracies, ethical violations, safety concerns, or fails to meet the user's request in any way, even minor issues.

TMPORTANT:

- Be strict in your assessment if both responses have any deficiencies, even minor ones, output '[[D]]'.
- Focus purely on content quality based on the given factual information and evaluation criteria.
- After presenting your reasoning, provide the final decision enclosed in double brackets to ensure proper parsing, for example: [[A]], [[B]], [[C]] or [[D]].

instruction:

Knowledge source: {} Question: {question} First Response: {answer1} Second Response: {answer2} Answer:

C.3 Multi-Agent Debate

817

818 $\{AGENT - 1\}$ 819 instruction: 820 821 You are a first agent. Please act as an 822 objective and strict judge, 823 evaluating the responses of two AI 824 assistants to the user's question 825 below. Select the assistant that 826 adheres to the user's instructions 827 and responds to the questionwith higher quality. Your evaluation must 828 829 rigorously consider factors such as 830 helpfulness, relevance, accuracy, depth, creativity, and detail of the 831 832 responses. Avoid any biases based 833 on the position or order of 834 responses to ensure an unbiased decision. The length of responses 835 836 should not affect your evaluation. 837 Maintain objectivity and neutrality 838 towards assistant names. Output 839 exactly one of the following symbols: 840 A, B, C, or D. Use the following 841 criteria strictly: 842 843 - Output 'A' if the first response is notably better. 844 845 - Output 'B' if the second response is notably better. 846 847 - Output 'C' if both responses are equally good. - Output 'D' if both responses are 848 849 850 equally inadequate. 851 852 IMPORTANT !! Be strict in your assessment 853 : if there are errors in both 854 responses, even minor ones, treat 855 responses as inadequate and output ' 856 D'. 857 858 Question: {question} 859 860 First Response: {answer1} 861 Second Response: {answer2} 862 Answer: 863 864 865 $\{AGENT - 2\}$ 866 instruction: 867 You are the second agent. You always 868 desagree with the first agent. Provide 869 your reasons and verdict. 870 871 872 {AGENT - JUDGE } 873 instruction: 874 You are the judge agent. Evaluate both agents answers and decide which one 875 876 is correct and make the final verdict. 877 878 Please format your final verdict as 879 follows: [[Selected Answer]]