# Predicting Text Preference Via Structured Comparative Reasoning

Anonymous ACL submission

#### Abstract

Comparative reasoning plays a crucial role in 002 predicting text preferences; however, large language models (LLMs) often demonstrate incon-003 sistencies in their reasoning, leading to incorrect preference predictions. While approaches 006 like Chain-of-Thought improve accuracy in 007 many settings, they struggle to consistently distinguish the similarities and differences of complex texts. We introduce  $SC^2$ , a model that prompts LLMs to predict text preferences by generating structured intermediate comparisons.  $SC^2$  begins by proposing aspects for comparison, followed by generating textual comparisons under each aspect. We select con-014 015 sistent comparisons with a pairwise comparator that ensures each comparison of a given as-017 pect clearly distinguishes differences between texts, significantly reducing hallucination and improving consistency. Our empirical studies across various NLP tasks, including summa-021 rization, retrieval, and automatic rating, demonstrate that SC<sup>2</sup>'s enhanced performance in text preference prediction is significant.

#### 1 Introduction

024

034

040

Comparative reasoning is crucial for predicting text preferences, as deciding the best out of a set of texts requires careful examination of the similarities and differences across the documents. Hence, comparative reasoning has been especially useful in NLP tasks such as text summarization (Yang et al., 2023; Lee et al., 2023), search ranking (Qin et al., 2023), and automatic evaluation (Adlakha et al., 2023), where text preference prediction is a key step.

However, as corpora grow more dense and complex across domains, accurate comparative reasoning becomes increasingly challenging. Existing approaches rely on pretraining or fine-tuning models (Yu et al., 2023a; Iso et al., 2022) at the cost of massive human annotation and computation. With the emergence of large language models (LLMs) (OpenAI, 2023; Touvron et al., 2023a; Anil et al., 2023; Jiang et al., 2023a), prompting approaches like Chain-of-Thought (CoT) (Wei et al., 2022) offer a promising solution for enhancing comparative reasoning. These approaches leverage LLMs' exceptional language generation capabilities without incurring significant overhead.

042

043

044

045

046

047

048

050

051

054

059

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

078

079

081

Nonetheless, LLMs exhibit arbitrary and erroneous outputs when prompted for comparative reasoning (Adlakha et al., 2023). Specifically, LLMs demonstrate inconsistency in their comparative analyses of texts. Figure 1 (bottom left) provides an example of logically inconsistent LLM reasoning using zero-shot CoT prompting. The LLM's generated explanation initially describes a property as common to the text pair (highlighted in green), but later implies that the same property is a strength of just one of the documents (highlighted in yellow). This inconsistency in the LLM's comparative analysis leads to an incorrect prediction.

To address these challenges, we present  $SC^2$ , a StruCtured Comparative reasoning model that constructs an intermediate structured representation contrasting two text corpora for more accurate text preference prediction, as illustrated in Figure 1. First,  $SC^2$  proposes a set of aspects from text pairs to guide the comparison step. Second,  $SC^2$  generates textual comparisons for every aspect. We refer to aspects and comparisons together as intermediate structured representations. To improve the consistency of reasoning (e.g., a contrastive comparison of a aspect should not overlap with a common comparison),  $SC^2$  adopts approximate inference: SC<sup>2</sup> samples multiple responses in generative process and uses a pairwise comparator to select the most consistent one as the final output.

We demonstrate the effectiveness of  $SC^2$  in improving text preference prediction across various tasks including text summarization (Stiennon et al., 2020), document retrieval (Soboroff et al., 2018), and helpfulness and harmlessness detection (Bai et al., 2022) with average 2.5 and 7.0 points gain



Figure 1: An example illustrating the differences in predicting the text preference between CoT-Zero-Shot prompting and  $SC^2$ . The top portion shows a query between a human and a chatbot, along with two candidate responses  $\leftarrow$  and  $\rightarrow$  to it. The table in the middle of the figure presents a structured intermediate representation produced by  $SC^2$ . Small phrases in the first row are *aspects*. *Comparisons* are entries that are not in the first column or row in the table. The Venn diagram visualizes the atomic comparisons with respect to  $\leftarrow$  and  $\rightarrow$ .

over the top and bottom baselines, respectively.

Our analysis further confirms the effectiveness of  $SC^2$  without incurring expensive LLM usage, and ablation studies emphasize the importance of the pairwise comparator. Our extensive human evaluations also indicate that  $SC^2$  aids in interpretation and assists users in making decisions.

#### 2 Related Work

**Prompting Large Language Models** LLMs have recently advanced the state-of-the-art performance across many NLP tasks (Anil et al., 2023; OpenAI, 2023; Chowdhery et al., 2022; Touvron et al., 2023a,b; Jiang et al., 2023a). These LLMs have demonstrated the capability to provide chain-of-thought explanations that elucidate their reasoning processes (Wei et al., 2022; Ko-jima et al., 2022). However, the chain-of-thoughts generated by LLMs are often arbitrary or contradictory (Turpin et al., 2023; Chen et al., 2023; Dhuliawala et al., 2023), unfaithful to the facts(Lyu et al., 2023, 2024) or lacking robustness to rephrased questions. To mitigate these issues, several works

aim to leverage consistency-based (Wang et al., 2023; Zhou et al., 2022), or verification-based approach (Ling et al., 2023; Lyu et al., 2024) to improve the reasoning capacity of LLMs, yet the benefit of such additional techniques are still ambivalent (Huang et al., 2023). Furthermore, all these advanced techniques still concentrate on processing raw-text inputs, thereby overlooking the integration of structural information. Moreover, they lack the implementation of explicit consistency constraints, which is crucial for maintaining logical coherence in the generated outputs.

**Comparative Reasoning and Summarization** Comparative reasoning involves comparing and contrasting different documents (Yu et al., 2023a), which has applications for a broad range of NLP tasks including text ranking (Jiang et al., 2023b; Qin et al., 2023), reward modeling (Ouyang et al., 2022; Lee et al., 2023) and automatic text generation evaluation (Liu et al., 2023a). Initial explorations focused on mining comparative content from text corpora (Jindal and Liu, 2006; Li et al., 2010). More recent studies have developed models

177

178

179

189 190

187

188

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

for generating comparative text, including generating arguments for answering comparative questions (Chekalina et al., 2021; Amplayo et al., 2021) and summarizing comparative opinions (Iso et al., 2022). Additionally, Zhong et al. (2022, 2023) prompts LLM to describe the differences between two text distributions in natural language and Dunlap et al. (2023) further extends to discover differences given a set of images from ImageNet.

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

171

One challenge of directly prompting LLMs for comparative reasoning is that the input text often contains a mixture of diverse patterns. As such, it is crucial to incorporate fine-grained aspects to guide LLMs for generating more comprehensive summarizations (Sun et al., 2023; Xu et al., 2023; Wu et al., 2023; Yu et al., 2023b). Early works (Lin and Hovy, 2000; Titov and McDonald, 2008) used clustering or topic modeling to identify aspects in documents. Lekhtman et al. (2021) fine-tune a pretrained language model for aspect extraction, which relies on manual labeling of comparative data. On the other hand, Goyal et al. (2022); Yang et al. (2023) leverage LLMs to perform summarization with the fixed aspects provided by humans. Differently, we leverage LLMs to automatically discover aspects to guide comparative reasoning, which provides a flexible way to incorporate fine-grained task-relevant signals while requiring minimal labeling efforts.

#### 3 Methods

Our model,  $SC^2$ , produces comparative reasoning for text preference prediction that applies to densely written texts, generalizes to multiple domains, and ensures consistency. In this section, we give the generative process and inference procedure for  $SC^2$ . Our primary focus is ensuring the comparisons consistently distinguish similarities and differences between texts.

#### 3.1 **Generative Process**

The generative process has three steps. First, given a text pair,  $SC^2$  simplifies the task by delineating a set of aspects, as depicted in Figure 1. These 168 aspects, consisting of concise phrases, enable the 169 structured comparison between the texts. Second, 170  $SC^2$  produces comparisons, which are concise, aspect-focused comparative statements that clearly 172 express how the texts are similar and different. We 173 require the comparisons to be consistent: simi-174 larities identified as shared between the text pair 175 should not overlap with what's unique to each of 176

them. Given the aspects and comparisons, the final step predicts which text is preferred.

Formally, for a text pair problem, we denote the text pair as  $\leftarrow$  and  $\rightarrow$ , along with a query. SC<sup>2</sup> has three components: Aspects  $a = \{a_1, a_2, \dots, a_n\},\$ comparisons  $c = \{c_1, c_2, \dots, c_n\}$ , and text preferences  $y \in \{\leftarrow, \rightarrow\}$ . The comparison c has three columns:  $\{c^{\leftarrow}, c_i^{\rightarrow}, c_i^{\cap}\}, c_i^{\rightarrow}$  and  $c_i^{\leftarrow}$  refers to properties exclusive to  $\rightarrow$  and  $\leftarrow$  respectively, and  $c_i^{\cap}$ to properties shared by both texts.

 $SC^2$  follows the following generative process: First, it generates the aspects conditioned on the text using an *aspect model*, P(a). Second, comparisons for each aspect are generated from the comparison model

$$P(c|a) \propto \prod_{i} l(c_i) \times P(c_i|c_{\langle i}, a)$$

where the function  $l : C \to \mathbb{R}^+$  evaluates the consistency of  $c_i$ . A higher value of  $l(c_i)$  indicates a greater degree of consistency. Finally, preference model P(y|c, a) produces the preference label y. Parameterization We use LLMs with specific prompts to parameterize each model. With LLMs generating reliable scalar values of consistency is unreliable (Imani et al., 2023; Liu et al., 2023b). Instead of directly regressing a consistency score, we rely on pairwise comparisons, which have been observed to be more reliable (Qin et al., 2023). We define a pairwise comparator  $l'(c, c') = \mathbb{1}(l(c) >$ l(c')<sup>1</sup>, which takes a pair of comparisons (c, c')and determines the more consistent one.

To facilitate this, we recruit experts to develop few-shot prompts that demonstrate a direct comparison of two structured representations based on consistency within itself. We guide our annotators to assess pairs (c, c') against consistency criteria, emphasizing that elements of the comparison should ideally exhibit no overlap. Detailed instructions are attached in the Appendix.

#### **Tournament-based Inference** 3.2

Given the generative model, the goal of inference is to produce aspects and comparisons that are high probability under the model and consistent. We take a step-wise approach, choosing aspects, comparisons, and then finally predicting preferences.

When choosing aspects, we follow prior work by employing a variety of sampling strategies to obtain near-optimal aspects  $a^*$  from P(a) (Wang

<sup>&</sup>lt;sup>1</sup>We break the tie randomly.



Figure 2: Illustration of tournament inference. Given a set of samples,  $C = \{c^1, c^2, c^3, c^4, c^5, c^6\}$ , the tournament approach randomly partitions them into three groups in the first round, and each two is paired as input to l' and output from l' will be entering the next round. In this way, we only need to use l' 5 times.

et al., 2023; Amplayo et al., 2021). We provide more details on these strategies in Section 4.

Given aspects  $a^*$ , our next goal is to find comparisons that are likely under the comparison model  $\arg \max_c P(c|a^*) = \arg \max_c l(c) \cdot P(c|a^*)$ . There are two challenges with this objective: First, the set of possible comparisons is intractably large. Second, the consistency function l(c) is unreliable. We approach the first challenge by sampling a set C of high probability comparisons from P(c|a), and the second challenge by selecting the most consistent comparison by applying the pairwise consistency comparator l'(c, c') in a binary reduction. Formally, we select the most consistent comparison by optimizing

$$c^* = \underset{c \in C}{\operatorname{arg\,max}} \sum_{c' \in C \setminus \{c\}} l'(c, c')$$

Naively, this optimization problem above requires  $O(|C|^2)$  pairwise comparisons to optimize exactly.

To reduce the number of pairwise comparisons, we utilize a tournament approach that performs O(|C|) comparisons. The tournament approach utilizes a binary reduction: Each step of the binary reduction takes a pair of comparisons and eliminates the less logically consistent one into the successive rounds. We illustrate the tournament approach in Figure 2. The naive and tournament approaches are equivalent if transitivity holds in the consistency comparator l'(c, c'). In practice, transitivity does not always hold with LLM parameterizations, resulting in the tournament approach trading off accuracy for efficiency.

Finally, with structured intermediate representation  $(a^*, c^*)$ , SC<sup>2</sup> decides between  $\leftarrow$ ,  $\rightarrow$  which one is preferred by taking arg max  $P(y|a^*, c^*)$ .

#### 4 Experimental Setup

**Aspect Model** We experiment with two models for generating aspects: the online aspect model and the offline aspect model. Both models use PaLM-2-L to obtain aspects.

The online aspect model dynamically generates aspects using the CoT paradigm to deduce aspects based on text inputs and applies self-consistency to select the most agreeable aspect for each text pair. However, this model is costly due to the extensive use of LLM API calls for every pair of tasks.

The offline aspect model extracts aspects from a collection of text corpora, adapting the concept from but employing LLMs. Specifically, this model prompts an LLM to extract aspects from each text within the collected corpora (50 pairs of texts for each task in this paper). It then prompts an LLM to refine and consolidate all generated aspects. Finally, we identify five fixed aspects as to use directly for any text pair of one task. This aspect model significantly reduces costs by allowing offline refinement of aspects. Refined aspects are fixed, thus they can be directly used without any additional expense.

In our experimental studies, we report only the best results for both baselines and in this section. To understand the impact of the aspect model, we detail its effects in our analysis section.

**Comparison Model** We use PaLM-2-L as the major LLM backbone of comparison model of  $SC^2$  to produce intermediate structured representations.

**Preference Model** For the final text preference prediction model, we experiment with two other LLM backbones differing in their model capacity. We aim to prove that the intermediate structured representations produced by  $SC^2$  with PaLM-2-L can help any backbone LLMs to predict text preference more accurately, regardless of their capacity. Specifically, we have used OpenAI's GPT-3.5, and GPT-4 (OpenAI, 2023) in our experiments.

**Prompting Templates** Prompts used in different models can be found in our Appendix. Note that we do not tailor the preference model's prompts; instead, we adapted the templates from Rafailov et al. (2023) for a fair comparison across baselines<sup>2</sup>.

247

248

249

254

258

259

260

261 262 263

264

266

267

268

270

271

272

273

274

275

276

277

278

279

281

282

283

287

288

290

291

292

293

294

295

297

298

299

300

<sup>&</sup>lt;sup>2</sup>In the original DPO paper (Rafailov et al., 2023), the authors did not use the Anthropic-Harmless dataset, we adapted their templated for Harmless datasets.

388

389

390

391

392

393

394

395

396

397

399

400

353

**Hyperparameters** As SC<sup>2</sup> searches for the best comparisons during the inference stage, as a result, we have a hyperparameter |C|, referring to the number of samples generated by the comparison model. |C| is an important parameter that might affect the quality of the intermediate structured representation produced by . For the reported results in this section, we set |C| = 8. We study the influence of this hyperparameter in Sec 6.

Baselines For evaluation, we consider several
baselines, primarily focused on the LLM-based
prompting approaches. Below is a detailed
overview of these baselines:

(1) *Direct Prompting (DP):* This method directly
prompts LLMs to predict text preference.

(2) *DP w/Aspects:* This approach is a variation of
DP. The difference is that DP w/Aspects incorporates aspects generated by the aspect models.

(3) *CoT-0-shot:* This baseline utilizes a standard
CoT-0-shot template for task preference prediction
(with "let's think step by step"). More details of the
prompt template are available in the appendix.

(4) *CoT-1-shot*: In addition to zero-shot prompting,
we also carry out experiments using a 1-shot example within the CoT paradigm. For that purpose, we
craft our 1-shot examples across different datasets.
(5) *CoT-SelfCon*: This baseline integrates selfconsistency to CoT-0-shot baseline aiming to remove the arbitrariness.

331

333

335

339

340

Specifically, CoT-SelfCon first samples multiple responses from the LLM using the same prompt and text pair input. Subsequently, CoT-SelfCon aggregates all responses to identify the most frequently generated answer. In our experimental studies, we set the number of sample responses to 8 and use a majority vote to determine the desired response, randomly selecting a response in the event of a tie. A mathematical explanation of Self-Con can be found in our appendix.

341Datasets (1) TL;DR (Stiennon et al., 2020): We342use OpenAI's filtered Reddit and CNN/Daily Mail343TL;DR dataset. OpenAI also created a preference344dataset from this, where labelers rated two gener-345ated summaries per post. For the CNN/Daily Mail346part, for a given news, we extracted two graded347summaries and used the overall score to create the348label. More details are in the original paper.

(2) **RLAIF-HH (Bai et al., 2022)**: The RLAIFHH from Anthropic dataset comprises dialogues
from interactions between crowdworkers and large
language models. In these exchanges, workers ei-

ther seek assistance or provoke potentially harmful responses from the AI. The responses are then labeled based on their helpfulness or harmfulness. (3) **TREC News (Soboroff et al., 2018):** The

TREC News (soboron et al., 2013): The TREC News dataset contains query-document pairs focused on ad-hoc ranking and filtering tasks from the late 1980s to early 2000s. We modify the dataset as follows for preference prediction: for a given query, we extract two document answers to construct the triplet and use the relevance score provided by the original dataset to decide which document is more preferred.

**Dataset Sampling** As datasets that have been used in the past are in large volumes, we only sampled a small ratio of them due to the cost of running all experiments. We sample roughly 250-300 data points from each dataset uniformly. More details of the sampled dataset can be found in the Appendix.

**Metrics** We report the accuracy of all approaches in our experiment  $\left(\frac{\text{Correctly Predicted Instances}}{\text{All Instances}}\right)$  to measure the performance.

#### 5 Results

Experimental results in Table 1 demonstrate  $SC^{2}$ 's strong performance across all evaluation domains, with average gains of  $\sim 2.5$  and  $\sim 7.0$  points over the top and bottom baselines respectively. This confirms the benefits of structured comparative reasoning for enhanced text preference prediction. Using structured intermediate representations produced by  $SC^{2}$ , the preference prediction model better handles these comparative reasoning difficulties.

Moreover, we observe the input length as an additional factor impacting performance. For instance, the TREC News dataset comprises considerably longer texts than other corpora. Here, the DP method lags behind  $SC^2$  by over 9 points, compared to the average 7 point deficit across baselines. Though input length serves as an imperfect proxy for complexity, the results also signaled the potential benefit of using our method for longer inputs.

We also want to point out that  $SC^2$  could be further improved by coupling with some of the existing general prompting techniques, for example, self-consistency (Wang et al., 2023) and selfverification (Madaan et al., 2023).

#### 6 Analysis

To further understand the benefit of using  $SC^2$  to produce an intermediate structured representation,

Preference Model	Comparison Model	TLDR		RLAIF			Document Ranking	
	1110001	Reddit	CNN/DM	AVG	Helpful	Harmless	AVG	TREC News
	DP	62.89	61.39	62.14	58.40	58.15	58.27	44.36
	DP w/Aspects	62.50	62.55	62.52	59.20	53.72	56.46	46.18
GPT-3.5	CoT-0-shot	63.67	64.48	64.08	59.00	56.94	57.97	47.64
	CoT-1-shot	64.06	63.71	63.88	59.20	58.55	58.88	50.18
	CoT-SelfCon	64.92	63.32	64.12	60.60	58.75	59.68	50.55
	$\mathbf{SC}^2$	68.36	68.34	68.55	63.20	59.76	61.49	52.95
	DP	66.41	64.86	65.63	62.60	58.85	60.58	52.00
	DP w/Aspects	65.63	65.25	65.44	60.60	60.97	60.78	55.64
GPT-4	CoT-0-shot	68.75	68.34	68.54	63.00	60.56	61.78	59.64
	CoT-1-shot	69.92	69.50	69.71	63.80	60.16	61.98	61.09
	CoT-SelfCon	71.67	69.12	69.90	64.00	60.76	62.38	61.82
	$\mathbf{SC}^2$	73.83	70.65	72.25	66.60	62.98	64.79	64.73

Table 1: Experimental results of  $SC^2$  across different datasets in three different domains. DP refers to direct prompting. We use accuracy to measure the performance and report averaged the results from 5 rounds.

	Random Selection	Tournament Scheme	Exact Search
# LLM calls	1	7	56
Decoded Len	372	2,651	13,272
Accuracy	0.63	0.71	0.73

Table 2: Cost and accuracy analysis of different inference approach of  $SC^2$ .

in this section, we conduct ablation studies and indepth analysis. We also implement a user study to explore the potential of using  $SC^2$  to inform human beings' decisions.

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

#### 6.1 Effectiveness of Pairwise Comparator

To calibrate the effective gain arising from the pairwise comparator l', we first compare variants of  $SC^2$  with the comparators and those with different hyperparameter configurations of  $SC^2$ . We use different intermediate structured representations produced by variants of  $SC^2$  to predict the text preference. Results are shown in Figure 3.

With |C| = 1, where there is effectively no pairwise comparator l', the performance of the preference model was found to be comparable to baseline results shown in Table 1. This suggests that inconsistent structured representations could potentially degrade the performance of the preference model. An increase in accuracy was observed with larger values of |C|, indicating the benefits of pairwise comparator. However, this improvement plateaued when |C| exceeded 8, hinting at a potential ceil-

ing effect for our approach, irrespective of further increases in |C|.

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

449

450

451

#### 6.2 Impact of Different Aspect Models

To understand the effect of different aspect models, we conduct ablation studies comparing the baseline that used aspects and  $SC^2$  with aspect models proposed in our experimental study.

Table 3 presents the results. It shows that  $SC^2$  with the offline aspect model consistently outperforms or performs as well as  $SC^2$  with the online aspect model. However, for the DP w/Aspects baseline, neither the online nor the offline aspect model demonstrates superiority over the other. This indicates that  $SC^2$  does not require online LLM calls which dynamically generate aspects and can effectively utilize the offline aspect model to obtain aspects for the given task at pretty low cost.

## 6.3 Cost Analysis of SC<sup>2</sup> and Few-shot CoT-SelfCon

As discussed in our experimental study, CoT-SelfCon has no pairwise comparator components, resulting in lower LLM usage. On the other hand, in our primary experimental studies, we utilize PaLM-2-L to create intermediate structured representations and other LLM for preference prediction to avoid potential overfitting. In contrast, the CoT-SelfCon baseline consistently employs the same LLM (GPT-4) all the way.

To ensure a fair comparison and eliminate biases

Preference	Aspect	Model	TLDR		RLAIF		Document
Model	Model						Ranking
			Reddit	CNN/DM	Helpful	Harmless	TREC News
	Online	DP w/Aspects	62.89	62.16	59.20	53.72	46.18
	Online	$\mathbf{SC}^2$	67.97	67.95	63.00	59.15	53.09
GPT3.5	Offline	DP w/Aspects	62.50	62.55	58.80	53.32	44.96
	Offline	$\mathbf{SC}^2$	68.36	68.34	63.20	59.76	53.09
	Online	DP w/Aspects	65.63	63.32	60.40	60.97	54.81
	Online	$\mathbf{SC}^2$	73.05	70.27	66.00	62.37	63.81
GPT4	Offline	DP w/Aspects	64.84	65.25	60.60	59.76	55.64
	Offline	$\mathbf{SC}^2$	73.83	70.65	66.60	62.98	64.73

Table 3: Calibration of different aspect models. Online refers to Online Aspect Model, and Offline refers to Offline Aspect Model. DP w/Aspects refers to Direct Prompting with Aspects.

Total LLM calls	8	15	24
$\mathbf{SC}^2$	0.682	0.738	0.750
CoT-SelfCon	0.678	0.728	0.730

Table 4: Accury of text preference prediction of  $SC^2$ against CoT-SelfCon with the same # of LLM calls.

that might arise from using different LLMs and # total LLM calls, we only use GPT-4 for both  $SC^2$ and CoT-SelfCon in this analysis. We use a fixed number of total LLM calls, including the generation of intermediate structured representations and the final preference prediction. We limit our experiments to a single dataset with 100 samples and average the results over 5 rounds for the cost consideration. The results are shown in Table 4. Our analysis indicates that with the same # total LLM calls and the same LLM backbone,  $SC^2$  predicts preference consistently more accurately.

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

477

Furthermore, we evaluate against few-shot CoT-SelfCon, widely regarded as a strong baseline. Given that  $SC^2$  operates in a zero-shot setting in experimental studies, for a fair comparison, we compare few-shot SC<sup>2</sup> with few-shot CoT-SelfCon, varying # LLM calls and # few-shot examples. We limit # few-shot examples to 5. This makes sure the context length is within the LLM's length window for both approaches.

Results are shown in Table 5. When the total LLM calls are low, CoT-SelfCon maintains a slight advantage over  $SC^2$ . However, as the number of 475 LLM calls increases, SC<sup>2</sup> consistently outperforms 476 few-shot CoT-SelfCon with the margin widening. This trend is attributed to the necessity for pair-478 wise comparators to produce logically consistent 479 intermediate-structured representations, leading to 480 more accurate predictions. 481

Total	# Few-shot	CoT-	$SC^2$
LLM calls	Examples	SelfCon	
8	3	0.678	0.672
	5	0.694	0.685
15	3	0.718	0.733
	5	0.742	0.756
24	3	0.778	0.797
	5	0.796	0.812

Table 5: Accury of text preference prediction of fewshot  $SC^2$  against few-shot CoT-SelfCon with the same number of LLM calls.

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

#### 6.4 Efficiency of Tournament Approach

We study the efficiency and effectiveness of the tournament approach w.r.t. other inference methods. Random Selection refers to the process of randomly selecting one sample from C during the inference stage, while Exact Search involves running all possible comparisons, which takes  $O(n^2)$ . We measure the cost using the total input length and the number of LLM calls, as this is common practice for the actual cost calculation in commercial Large Language Models (LLMs). We used the same dataset from the previous subsection.

We find a significant gap between the Random Selection approach and the other two approaches as shown in Figure 2. Although Exact Search yields the best results, it requires 4 times the token length and 49 more LLM calls, potentially leading to a substantial increase in cost.

#### **Human Evaluation** 6.5

We conduct additional human evaluations to see how the intermediate structured representations produced by SC<sup>2</sup> inform human decision-making.



Figure 3: Impact of # samples |C| in SC<sup>2</sup>.

**Annoators** We recruit our annotators from an internal pool. Demographic and geographic characteristics of the annotator population are not accessible to our researchers. Information can be used to identify annotators that are fully anonymized. Consent forms have to be signed by annotators to take part in this study.

505

506

507

508

510

Study Design In consideration of ethical stan-511 dards and the requirement to avoid directly test-512 ing annotators, we structure our human evalua-513 tion as follows: Annotators are presented with a 514 query alongside a pair of text options, denoted as 515  $(\leftarrow, \rightarrow)$ . They determine which text, either  $\leftarrow$  or 516  $\rightarrow$ , is preferable. They have three options:  $\leftarrow$  is 517 better,  $\rightarrow$  is better, and tie. Following their initial 518 decision, annotators are then shown the intermediate structured representations generated by different variants of  $SC^2$ . They decide if this additional 521 information leads them to reconsider their initial choice and provide reasons for any change in their 523 decision. This evaluation process uses two variants of SC<sup>2</sup>: |C| = 1 and |C| = 8 respectively. For 525 ethical considerations, we only experiment with RLAIF-helpfulness and TL;DR-Reddit, ensuring the content is not harmful or violent manually. We 529 instantiate 100 data points for each dataset and assign each question to three annotators. We collect 530 96 and 98 questions with useful responses from all three annotators for RLAIF-helpful and TL;DR-Reddit respectively. 533

**Metrics** We use the ground truth to gather the scores: we assign 1 for any correct answer, 0 for a tie, and -1 for any other incorrect answers.

**Findings** As shown in Figure 4, with the aid of more consistent intermediate structured representations (|C| = 8), annotators are inclined to revise and flip their previous wrong answers to correct ones. This suggests that the intermediate structured representation may facilitate better decision-making among human evaluators. However, we also observe that intermediate structured



Figure 4: Human evaluation on structured representation produced by different settings of  $SC^2$ .

545

546

547

548

549

550

551

552

553

555

556

557

558

559

560

561

562

563

564

565

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

584

representations without using a pairwise comparator (|C| = 1) could mislead annotators, deterring them from selecting the correct preference. This amplifies the importance of the pairwise comparator to ensure consistency.

We also look into quantitative justifications provided by annotators. Most annotators stated that intermediate structured reprethem better understand sentations helped One mentioned, "the table two texts. gives the concise comparison", while another pointed out, "this [table] helped me to understand better the implications of the two answers, and I changed my mind after reading [the table]". Besides, we also observe complaints about the intermediate structured representations being hallucinatory and not factual. The issue is more noticeable in cases where the structured representation is produced by  $SC^2$  without a pairwise comparator. This suggests that enforcing a pairwise comparator might mitigate the arbitrariness of LLM's output for better consistency, but still poses the risk of presenting hallucinated results to annoators.

#### 7 Conclusion

This paper presents  $SC^2$ , a structured comparative reasoning model for improving text preference prediction.  $SC^2$  constructs intermediate structured representations to explicitly contrast text pairs, incorporating a consistency comparator to enhance accuracy and coherence. Comprehensive experiments across text summarization, retrieval, and response rating tasks demonstrated that  $SC^2$  significantly improves consistency and achieves state-of-the-art performance. Analyses confirm the effectiveness of  $SC^2$ 's structured reasoning approach and consistency enforcement. Our human evaluations show that  $SC^2$  interpretations can assist users in making informed decisions.

# 585 586

588

589

590

595

596

611

613

614

615

617

618

619

621

623

627

631

## 8 Limitations

This work has several limitations that provide opportunities for future investigation. First, the evaluation was conducted on a sample set of datasets that, while spanning diverse domains, might not fully characterize the breadth of real-world textual comparison needs. Expanding SC<sup>2</sup>'s testing to larger, multilingual corpora is essential to assess its full potential and limitations beyond English. Furthermore, there are likely upper bounds on  $SC^2$ 's effectiveness imposed by the reasoning capacity of the underlying language model backbone. As more advanced LLMs emerge, exploring their integration could help quantify this ceiling effect. On a technical level, in this paper, measuring consistency relies on approximate metrics, so developing more rigorous evaluation schemes could better highlight  $SC^{2}$ 's benefits. We also do not include other prompting techniques that have been well-studied in the community, which we leave for future work.

#### 9 Ethical Considerations

This research paper might risk potential biases that could arise from textual comparisons, particularly around sensitive attributes.  $SC^2$  is trained on established corpora like Wikipedia and books that may inherently contain societal biases. While a full analysis of these biases is beyond the scope here, we acknowledge the risk that  $SC^2$  may inherit problematic biases from its training data. Applying recent advancements in language bias detection to  $SC^2$  could help quantify and mitigate these risks. We are interested in exploring this as part of future work. Furthermore, this research focused solely on English; extending to other languages is an important direction that would require non-trivial adaptation. Overall, while showing promise,  $SC^2$  has significant scope for improvement as limitations around evaluation, multilingual capabilities, consistency measurement, bias, and applied usage are addressed through future work.

#### References

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2023. Evaluating correctness and faithfulness of instructionfollowing models for question answering. *ArXiv preprint*, abs/2307.16877.
- 632 Reinald Kim Amplayo, Stefanos Angelidis, and Mirella

Lapata. 2021. Aspect-controllable opinion summarization. In *Proc. of EMNLP*, pages 6578–6593, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *ArXiv preprint*, abs/2305.10403.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *ArXiv preprint*, abs/2204.05862.
- Viktoriia Chekalina, Alexander Bondarenko, Chris Biemann, Meriem Beloucif, Varvara Logacheva, and Alexander Panchenko. 2021. Which is better for deep learning: Python or MATLAB? answering comparative questions in natural language. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 302–311, Online. Association for Computational Linguistics.
- Yanda Chen, Ruiqi Zhong, Narutatsu Ri, Chen Zhao, He He, Jacob Steinhardt, Zhou Yu, and Kathleen McKeown. 2023. Do models explain themselves? counterfactual simulatability of natural language explanations. *ArXiv preprint*, abs/2307.08678.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *ArXiv preprint*, abs/2204.02311.
- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. Chain-of-verification reduces hallucination in large language models. *ArXiv preprint*, abs/2309.11495.
- Lisa Dunlap, Yuhui Zhang, Xiaohan Wang, Ruiqi Zhong, Trevor Darrell, Jacob Steinhardt, Joseph E Gonzalez, and Serena Yeung-Levy. 2023. Describing differences in image sets with natural language. *arXiv preprint arXiv:2312.02974*.
- Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of gpt-3. *ArXiv preprint*, abs/2209.12356.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2023. Large language models cannot self-correct reasoning yet. *ArXiv preprint*, abs/2310.01798.
- Shima Imani, Liang Du, and Harsh Shrivastava. 2023. Mathprompter: Mathematical reasoning using large language models. *ArXiv preprint*, abs/2303.05398.

## 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658

633

634

635

683 684 685

675

676

677

678

679

680

681

682

798

- Hayate Iso, Xiaolan Wang, Stefanos Angelidis, and Yoshihiko Suhara. 2022. Comparative opinion summarization via collaborative decoding. In *Findings of the Association for Computational Linguistics: ACL* 2022, pages 3307–3324, Dublin, Ireland. Association for Computational Linguistics.
  - Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023a. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023b. LLM-blender: Ensembling large language models with pairwise ranking and generative fusion. In *Proc.* of ACL, pages 14165–14178, Toronto, Canada. Association for Computational Linguistics.

705

706

711

713

714

715

717

718

728

731

734

736

741

- Nitin Jindal and Bing Liu. 2006. Identifying comparative sentences in text documents. In *Proceedings* of the 29th annual international ACM SIGIR conference on Research and development in information retrieval, pages 244–251.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199– 22213.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. ArXiv preprint, abs/2309.00267.
- Entony Lekhtman, Yftah Ziser, and Roi Reichart. 2021. DILBERT: Customized pre-training for domain adaptation with category shift, with an application to aspect extraction. In *Proc. of EMNLP*, pages 219–230, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Shasha Li, Chin-Yew Lin, Young-In Song, and Zhoujun Li. 2010. Comparable entity mining from comparative questions. In *Proc. of ACL*, pages 650–658, Uppsala, Sweden. Association for Computational Linguistics.
- Chin-Yew Lin and Eduard Hovy. 2000. The automated acquisition of topic signatures for text summarization.
  In COLING 2000 Volume 1: The 18th International Conference on Computational Linguistics.
- Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su. 2023. Deductive verification of chain-of-thought reasoning. *ArXiv preprint*, abs/2306.03872.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023a. Gpteval: Nlg evaluation using gpt-4 with better human alignment. *ArXiv preprint*, abs/2303.16634.

- Yixin Liu, Avi Singh, C Daniel Freeman, John D Co-Reyes, and Peter J Liu. 2023b. Improving large language model fine-tuning for solving math problems. *ArXiv preprint*, abs/2310.10047.
- Qing Lyu, Marianna Apidianaki, and Chris Callison-Burch. 2024. Towards faithful model explanation in nlp: A survey. *Computational Linguistics*, pages 1–70.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-ofthought reasoning. *arXiv preprint arXiv:2301.13379*.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. *ArXiv preprint*, abs/2303.17651.
- OpenAI. 2023. GPT-4 technical report. ArXiv preprint, abs/2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, et al. 2023. Large language models are effective text rankers with pairwise ranking prompting. *ArXiv preprint*, abs/2306.17563.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *ArXiv preprint*, abs/2305.18290.
- Ian Soboroff, Shudong Huang, and Donna Harman. 2018. Trec 2018 news track overview.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F. Christiano. 2020. Learning to summarize with human feedback. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2023. Principle-driven selfalignment of language models from scratch with minimal human supervision. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Ivan Titov and Ryan T. McDonald. 2008. Modeling online reviews with multi-grain topic models. In Proceedings of the 17th International Conference on World Wide Web, WWW 2008, Beijing, China, April 21-25, 2008, pages 111–120. ACM.

- 799
- 802

Azhar, et al. 2023a.

abs/2302.13971.

abs/2307.09288.

Representations.

abs/2306.01693.

tems.

cient foundation language models. ArXiv preprint,

Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-

bert, Amjad Almahairi, Yasmine Babaei, Nikolay

Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti

Bhosale, et al. 2023b. Llama 2: Open founda-

tion and fine-tuned chat models. ArXiv preprint,

Miles Turpin, Julian Michael, Ethan Perez, and

Samuel R. Bowman. 2023. Language models don't

always say what they think: Unfaithful explanations

in chain-of-thought prompting. In Thirty-seventh

Conference on Neural Information Processing Sys-

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,

Ed H. Chi, Sharan Narang, Aakanksha Chowdhery,

and Denny Zhou. 2023. Self-consistency improves

chain of thought reasoning in language models. In

The Eleventh International Conference on Learning

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten

Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le,

and Denny Zhou. 2022. Chain of thought prompt-

ing elicits reasoning in large language models. In Advances in Neural Information Processing Systems. Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, and Hannaneh Hajishirzi. 2023. Fine-grained human feedback gives better rewards

Ran Xu, Hejie Cui, Yue Yu, Xuan Kan, Wenqi Shi,

Yuchen Zhuang, Wei Jin, Joyce Ho, and Carl Yang.

2023. Knowledge-infused prompting: Assessing and advancing clinical text data generation with large

language models. ArXiv preprint, abs/2311.00287. Xianjun Yang, Yan Li, Xinlu Zhang, Haifeng Chen, and Wei Cheng. 2023. Exploring the limits of chatgpt

for query or aspect-based text summarization. ArXiv

Mengxia Yu, Zhihan Zhang, Wenhao Yu, and Meng Jiang. 2023a. Pre-training language models for comparative reasoning. ArXiv preprint, abs/2305.14457.

Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander Ratner, Ranjay Krishna, Jiaming Shen,

and bias. ArXiv preprint, abs/2306.15895.

and Chao Zhang. 2023b. Large language model as attributed training data generator: A tale of diversity

Ruiqi Zhong, Charlie Snell, Dan Klein, and Jacob Steinhardt. 2022. Describing differences between text

distributions with natural language. In International Conference on Machine Learning, ICML 2022, 17-23

ArXiv preprint,

for language model training.

preprint, abs/2302.08081.

- 810
- 811 812 813
- 814 815
- 816
- 817 818
- 819
- 820 821
- 822 823 824

825

- 834 835

836 837

- 841
- 842

- 847

850

853

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier July 2022, Baltimore, Maryland, USA, volume 162 of Martinet, Marie-Anne Lachaux, Timothée Lacroix, Proceedings of Machine Learning Research, pages Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal 27099-27116. PMLR. Llama: Open and effi-

Ruiqi Zhong, Peter Zhang, Steve Li, Jinwoo Ahn, Dan Klein, and Jacob Steinhardt. 2023. Goal driven discovery of distributional differences via language descriptions. Neural Information Processing Systems.

855

856

857

858

859

860

861

862

863

864

865

866

867

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. arXiv preprint arXiv:2205.10625.



Figure 5: Entailment and coverage score of  $SC^2$ .

### **Appendix: Additional Experiments and Dataset Information**

Datasets Dataset sampling details can be found in Table 6 871

Coverage and Entailment We also analyze the 872 entailment and coverage scores of the intermedi-873 ate structured representations produced by  $SC^2$ . 875 We use internal coverage and entailment prediction models, leveraging an internal pre-trained language 876 model expansion of T5, designed to predict entail-877 ment and coverage score. The score is in [0, 1], and a higher score indicates better coverage or entailment. On one hand, we aim to check the coverage and entailment depth of the structured representations, and on the other, we seek to determine 882 whether existing metrics could serve as good indicators: whether computing those metrics alone can select better responses from LLMs. Specifically, we consider the entire structured representation from  $SC^2$  and calculate the coverage and entailment score given the concatenation of  $\leftarrow$  and  $\rightarrow$ . We compare different representations from variants of  $SC^2$ . Results are shown in Figure 5.

> As the comparisons are distilled into comparative statements, a relatively lower score is expected. However, we find that the representations of various  $SC^2$  configurations are not drastic in variance.

891

Dataset	# Samples	Avg. Length
TL;DR-CNN/DM	256	572
TL;DR-Reddit	259	362
Antropic-Helpful	250	102
Antropic-Harmless	249	93
TREC News	291	947
AVG	278	433

Table 6: Statistics of Datasets in Experimental Studies

This indicates that LLMs might be hallucinating entities, names, and locations mentioned in  $\leftarrow$  and  $\rightarrow$ , which existing metrics such as coverage and entailment find challenging to identify. This also suggests that incorporating a pairwise comparator aids in obtaining reliable comparative reasoning.

This predictive model assesses the entailment and coverage scores for the tables generated by the comparison model concerning the text pair  $\leftarrow$ ,  $\rightarrow$ . Scores range from 0 to 1, with 1 indicating complete identity. We will add this background in the next revision.

Our goal of this analysis was to demonstrate that relying solely on coverage and entailment scores is insufficient for eliciting consistent comparisons produced by LLMs. As depicted in Figure 4, there is a similar trend in entailment and coverage scores regardless of the number of ICI (as discussed in our previous analysis, |C| is a proxy for final accuracy, with larger |C| normally yielding a better score). This supports our argument for the importance of pair-wise consistency comparison. If the reviewer finds this analysis not necessary, we can remove this analysis as well.

CoT-SelfCon We illustrate the details of CoT-SelfCon in the below.

- Specifically, CoT-SelfCon, in our experiment, has used the same input text  $\rightarrow$  and text  $\leftarrow$ , typically N times and is asked to generate a response, a choice between two options  $\rightarrow$  or ←.
- After each iteration, the model's response is recorded. Upon completing all iterations, the responses are aggregated to identify the most frequently given answer.
- · This method leverages repetition to find a consistent model's outputs, thus leading to a more consistent and potentially more accurate conclusion.
- Mathematically, CoT-SelfCon can be written 934 as given we received the answer from LLM to 935 be  $r = \{r_1, \ldots, r_n\}, a_i = ans(y_i),$

Most frequent answer = 
$$\arg \max_{x} \sum_{i=1}^{n} 1_{r_i=x}$$
 937

n

#### **Appendix: Prompts**

938

895

896

897

898

899

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

# Which of the following summaries does a better job of summarizing the most important points in the given forum article, without including unimportant or irrelevant details?

A good summary is both precise and concise.

**Original Article:** {article}

**Summary A:** {contextA}

**Summary B:** {contextB}

Take a deep breath and think about this question step by step! FIRST, think step by step to have a comparison of the two summaries, explaining which you prefer and why. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 6: Preference model prompt for CoT Zero-shot Prompting for TL;DR

Which of the following summaries does a better job of summarizing the most important points in the given forum article, without including unimportant or irrelevant details?

A good summary is both precise and concise.

**Original Article:** {article}

**Summary A:** {contextA}

**Summary B:** {contextB}

FIRST, explaining which you prefer and why. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 7: Preference model prompt for Direct Prompting for TL;DR

## Which of the following summaries does a better job of summarizing the most important points in the given forum article, without including unimportant or irrelevant details? You are also given some aspects to help you make the decision

A good summary is both precise and concise.

**Original Article:** {article}

Summary A: {contextA}

**Summary B:** {contextB}

Aspects: {aspects}

FIRST, explaining which you prefer and why. In your evaluation, you need to consider aspects that are given above. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 8: Preference model prompt for Direct Prompting with Aspects for TL;DR

Which of the following summaries does a better job of summarizing the most important points in the given forum article, without including unimportant or irrelevant details? You are also given a comparative reasoning table that analyzes the differences and similarities between the two summaries.

A good summary is both precise and concise.

**Original Article:** {article}

**Summary A:** {contextA}

**Summary B:** {contextB}

**Comparative Reasoning Table:** {table}

FIRST, explain which you prefer and why. In your evaluation, you can use the comparative reasoning table above to help you make the justification and the decision. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 9: Preference model prompt for  $SC^2$  for TL;DR

#### Which of the following documents aligns better with the query given?

A good retrieved document should be relevant to the query.

**Query:** {query}

**Document A:** {contextA}

**Document B:** {contextB}

Take a deep breath and think about this question step by step! FIRST, think step by step to have a comparison of the two retrieved documents, explaining which you prefer and why. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 10: Preference model prompt for Zero-shot CoT Prompting for TREC News

### Which of the following documents aligns better with the query given?

A good retrieved document should be relevant to the query.

**Query:** {query}

**Document A:** {contextA}

**Document B:** {contextB}

FIRST, have a comparison of the two retrieved documents, explaining which you prefer and why. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 11: Preference model prompt for Direct Prompting for TREC News

Which of the following documents aligns better with the query given? You are also given some aspects to help you make the decision

A good retrieved document should be relevant to the query.

**Query:** {query}

**Document A:** {contextA}

**Document B:** {contextB}

Aspects: {aspects}

FIRST, have a comparison of two retrieved documents, explaining which you prefer and why. In your evaluation, you need to consider aspects that are given above. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 12: Preference model prompt for Direct Prompting with Aspects for TREC News

Which of the following documents aligns better with the query given? You are also given a comparative reasoning table that analyzes the differences and similarities between the two documents.

A good retrieved document should be relevant to the query.

**Query:** {query}

**Document A:** {contextA}

**Document B:** {contextB}

**Comparative Reasoning Table:** {table}

FIRST, explaining which you prefer and why. In your evaluation, you can use the comparative reasoning table above to help you make the justifications and decision. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 13: Preference model prompt for SC<sup>2</sup> for TREC News

For the following query to a chatbot, which response is more helpful?

**Query to a Chatbot:** {article}

**Response A:** {contextA}

**Response B:** {contextB}

Take a deep breath and think about this question step by step! FIRST, think step by step to have a comparison of the two responses generated, explaining which you prefer and why. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 14: Preference model prompt for Zero-shot CoT Prompting for RLAIF-HH

#### For the following query to a chatbot, which response is more helpful?

**Query to a Chatbot:** {article}

**Response A:** {contextA}

# **Response B:** {contextB }

FIRST, have a comparison of the two generated responses, explaining which you prefer and why. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 15: Preference model prompt for Direct Prompting for RLAIF-HH

For the following query to a chatbot, which response is more helpful? You are also given some aspects to help you make the decision

**Query to a Chatbot:** {article}

**Response A:** {contextA}

**Response B:** {contextB}

Aspects: {aspect}

FIRST, have a comparison of the two generated responses, explaining which you prefer and why. In your evaluation, you need to consider aspects that are given above. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 16: Preference model prompt for Direct Prompting with Aspects for RLAIF-HH

For the following query to a chatbot, which response is more helpful? You are also given a comparative reasoning table that analyzes the differences and similarities between the two generated responses.

**Query to a Chatbot:** {article}

**Response A:** {contextA}

**Response B:** {contextB}

**Comparative Reasoning Table:** {table}

FIRST, explain which you prefer and why. In your evaluation, you can use the comparative reasoning table above to help you make the justifications and decisions. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 17: Preference model prompt for  $SC^2$  for RLAIF-HH

Which of the following summaries does a better job of summarizing the most important points in the given forum article, without including unimportant or irrelevant details? You are also given some aspects to help you make the decision

A good summary is both precise and concise. **Example Article:** {article}

**Example Summary A:** {contextA}

**Example Summary B:** {contextB}

Example Aspects: {aspects}

FIRST, explain which you prefer and why. In your evaluation, you need to consider aspects that are given above. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">. Example Answer: {example answer}

Now, Based on the example above, take a deep breath and think about this question step by step to answer the following question.

Which of the following summaries does a better job of summarizing the most important points in the given forum article, without including unimportant or irrelevant details? You are also given some aspects to help you make the decision

A good summary is both precise and concise.

**Original Article:** {article}

**Summary A:** {contextA}

**Summary B:** {contextB}

Aspects: {aspects}

FIRST, explain which you prefer and why. In your evaluation, you need to consider aspects that are given above. SECOND, on a new line, state only "A" or "B" to indicate your choice.

Your response should use the format: Comparison: <step by step comparison> Preferred: <"A" or "B">.

Figure 18: Preference Model Prompts for CoT-1.

**Instructions:** Your task is to conduct a consistency analysis of two generated comparative table responses. Your evaluation should focus solely on the consistency of the responses. Each comparative table is constructed to delineate similarities and differences about a given query, juxtaposing candidate Summary 1 against candidate Summary 2. Consistency in this context refers to the logical coherence within each table. Specifically, for each row corresponding to an aspect-level comparison, the entries of three columns that denote similarities should be distinct and non-overlapping with the entries that denote differences. A consistent response will differentiate between the commonalities and disparities, ensuring that the information under the 'similarities' column does not overlap what is presented under the 'differences' column. This clear segregation is crucial in assessing the quality of the responses and their effectiveness in summarizing and contrasting the key points from the summaries.

**Query to a Chatbot:** {article}

**Summary 1:** {contextA}

Summary 2: {contextB}

**Comparative Table Response A:** {contextA}

**Comparative Table Response B:** {contextB}

More consistent: <"A" or "B">. Justifications: <Justifications>.

Figure 19: Instructions to Craft prompts for Pairwise Comparator.