

Coverage-based Fairness in Multi-document Summarization

Anonymous ACL submission

Abstract

Fairness in multi-document summarization (MDS) measures whether a system can generate a summary fairly representing information from documents with different social attribute values. Fairness in MDS is crucial since a fair summary can offer readers a comprehensive view. Previous works focus on quantifying summary-level fairness using Proportional Representation, a fairness measure based on Statistical Parity. However, Proportional Representation does not consider redundancy in input documents and overlooks corpus-level unfairness. In this work, we propose a new summary-level fairness measure, **Equal Coverage**, which is based on coverage of documents with different social attribute values and considers the redundancy within documents. To detect the corpus-level unfairness, we propose a new corpus-level measure, **Coverage Parity**. Our human evaluations show that our measures align with the human perception of fairness. Using our measures, we evaluate the fairness of ten different LLMs. We find that Llama2 is the fairest among all evaluated LLMs. We also find that almost all LLMs overrepresent different social attribute values.

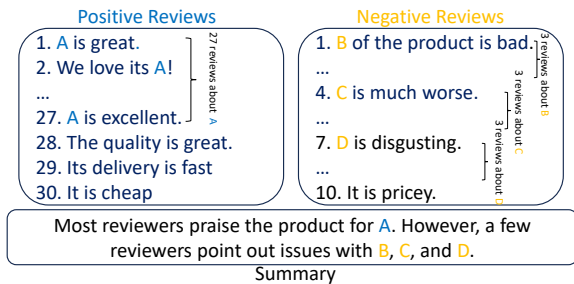
1 Introduction

Multi-document summarization (MDS) systems summarize the salient information from multiple documents about an entity, such as news articles about an event or reviews of a product. Typically, such documents are associated with *social attributes* e.g. political ideology in news and sentiments in reviews. Documents with different social attributes tend to have diverse information or conflicting opinions. A summary for them should fairly represent differing opinions across documents.

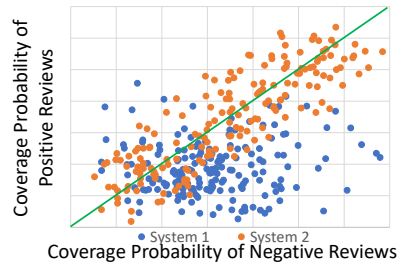
Fairness in MDS measures whether a system can generate summaries fairly representing information from documents with different social attribute values. It can be measured at a *summary-level*– quanti-

fying how fairly an individual summary represents documents with different social attribute values or at a *corpus-level*–quantifying how fairly a corpus of summaries as a whole represents documents with different social attribute values. Previous works in this area measured fairness in extractive or abstractive settings (Shandilya et al., 2018; Olabisi et al., 2022; Zhang et al., 2023; Huang et al., 2024). These works generally evaluate the fairness of a system as the aggregated summary-level fairness of its generated summaries. It is measured using Proportional Representation—a fairness measure based on Statistical Parity (Verma and Rubin, 2018). It requires that a document sentence being selected as a summary sentence is independent of its originating document’s social attribute value.

The definition of Proportional Representation suffers from two key problems. The first problem is that Proportional Representation does not consider the redundancy in input documents, common in multi-document settings. For example, in review summarization, suppose 75% of reviews are positive and most of them discuss topic A, while 25% of reviews are negative and evenly discuss topics B, C, and D (see Fig. 1a). According to the definition of Proportional Representation, if the system selects information from input independent of the social attribute value, a fair summary should ideally have 75% of information from positive reviews and 25% percent of information from negative reviews. The summary shown in the figure will be considered unfair according to this definition because, unlike the input, it contains more negative information than positive information. However, considering both redundancy and input sentiment distribution, this summary is fair since it covers equal proportions of reviews for both sentiments while avoiding redundant information. To address this, we propose a new summary-level coverage-based fairness measure, **Equal Coverage**. Unlike Proportional Representation, Equal Coverage requires a docu-



(a) Existing summary-level fairness measure, Proportional Representation, does not consider redundancy common in MDS. It incorrectly considers the illustrated summary unfair since, unlike the input, it contains more negative (but diverse) information than positive (but redundant) information. Our proposed measure, Equal Coverage, correctly considers this summary fair because it covers equal proportions of negative and positive reviews.



(b) Existing summary-level fairness measures can overlook corpus-level unfairness. Each point in this figure represents a summary sample. System 2 is fairer than System 1 since it has equal chances of overrepresenting negative (below the green line) and positive (above the green line) reviews while System 1 tends to overrepresent negative reviews. Our proposed measure, Coverage Parity, can correctly identify System 2 as fairer than System 1.

Figure 1: Issues with existing fairness measures for summary-level (a) and corpus-level (b) fairness.

ment being *covered* by a summary sentence to be independent of its social attribute value. Since a summary sentence can cover multiple documents with similar contents, Equal Coverage can address redundancy common for MDS.

The second problem is evaluating the fairness of an MDS system only using summary-level fairness measures can overlook corpus-level unfairness. Consider System 1 in Fig. 1b. Most of its summaries (blue dots) have a higher coverage probability for negative reviews than positive reviews. We observe that System 1 is unfair because its summaries *overrepresent* negative reviews. System 2 is fairer because its summaries (orange dots) have an equal chance of overrepresenting negative or positive reviews. Since individual summaries from both systems overrepresent negative or positive reviews, their summary-level fairness scores may be comparable. Hence, aggregated summary-level fairness scores cannot identify that System 2 is collectively fairer than System 1. To address this problem, we propose a new corpus-level fairness metric, **Coverage Parity**. Motivated by Best-Worst Scaling (Louviere et al., 2015), Coverage Parity is based on the principle that different social attribute values should have equal chances of being most overrepresented or underrepresented in summaries. Therefore, it can check whether the systems are equally (un)fair on different social attributes and identify which social attribute is overrepresented or underrepresented.

Our human evaluation shows that our measures align more with the human perception of fairness than Proportional Representation. Using our mea-

asures, we then evaluate the fairness of ten different Large Language Models (LLMs) (Ouyang et al., 2022; Touvron et al., 2023) in diverse domains: news, tweets, and reviews. For these domains, we consider social attributes having significant real-world impacts: ideologies and stances for news and tweets and sentiment for reviews. Our experiments find that for Llama2 and Claude3, larger models are fairer than smaller models, but this trend is inconsistent for GPTs and Mixtral. Our experiments also find that almost all LLMs tend to overrepresent certain social attribute values in each domain. It is an important finding that users can use to calibrate their perception before using LLM-generated summaries. It can also be used by developers to build fairer LLM-based summarization systems.

To conclude, our contributions are three-fold:

- We propose a new summary-level fairness measure, Equal Coverage, which incorporates redundancy of input information, common in MDS;
- We propose a new corpus-level fairness measure, Coverage Parity to detect corpus-level unfairness;
- We evaluate the fairness of LLMs using these two measures in various domains.

2 Related Work

Shandilya et al. (2018, 2020); Dash et al. (2019) propose to measure the summary-level fairness in MDS under the extractive setting using Proportional Representation. They propose an in-processing method to improve the fairness of extractive summaries. Similarly, Keswani and Celis (2021) uses Proportional Representation to measure the fairness of summaries under different dis-

tributions of social attributes in input documents and proposes a post-processing method to improve fairness. Olabisi et al. (2022) uses the same measure to measure fairness and proposes a clustering-based pre-processing method to improve fairness.

Recently, Zhang et al. (2023); Huang et al. (2024) extend Proportional Representation to measure fairness under the abstractive setting. To estimate the distribution of social attributes in a summary, Zhang et al. (2023) maps the summary back to the originating documents, while Huang et al. (2024) uses a finetuned model to estimate the social attribute of each summary sentence. Lei et al. (2024) measures the fairness similarly as Huang et al. (2024) and proposes to improve the fairness of abstractive summaries using reinforcement learning with a polarity distance reward. However, these works generally measure fairness using Proportional Representation, which has the limitations discussed in the introduction. In our experiments, we use Proportional Representation as a baseline.

3 Notation

We use G to denote all samples for evaluating the fairness of a system on a social attribute. Each sample $(D, S) \in G$ contains a document set $D = \{d_1, \dots, d_n\}$ and a summary S generated by the system for these documents, where d_i denotes the i -th document. Each input document d_i is labeled with an social attribute value $a_i \in \{1, \dots, K\}$.

4 Fairness Measures

In this section, we describe our proposed measures Equal Coverage (Sec. 4.1) and Coverage Parity (Sec. 4.2).

4.1 Equal Coverage

Equal Coverage is a summary-level fairness measure for measuring the fairness of a summary, S , for document sets, D . Equal Coverage is based on the principle that the documents with different social attribute values should have equal probabilities for being covered by a summary sentence. We denote the probability that a random document $d \in D$ whose social attribute value a is k is covered by a random summary sentence $s \in S$ as $p(d, s|a = k)$. This is referred to as the coverage probability for the social attribute value k . Similarly, we denote the probability that a random document, d , is covered by a random summary sentence $s \in S$ as $p(d, s)$, which is referred to as the coverage prob-

ability for a document. For a fair summary S according to Equal Coverage, coverage of a random document d should be independent of its social attribute value:

$$p(d, s|a = k) = p(d, s), \forall k. \quad (1)$$

However, two issues arise with summaries' complex sentence structures. First, a summary sentence can combine information from several documents, making it difficult to attribute the sentence to any single document. Second, sentence lengths vary significantly based on social attribute values. For example, summary sentences about positive sentiment can be long, such as 'the delivery is fast and it is worth the price.'. In contrast, those about negative sentiment can be much shorter, such as 'delivery is too slow.'. The length difference can skew the coverage probability for different social attribute values $p(d, s|a = i)$. To address these issues, Equal Coverage splits the summary sentences $s \in S$ into multiple simpler sentences by prompting GPT-3.5 motivated by Bhaskar et al. (2023); Min et al. (2023). For example, compound sentences are split into simple sentences that describe a single fact. We denote the j -th summary sentence after split as s_j . Further details are in App. A.1.

Equal Coverage estimates the coverage probability for different social attribute values $p(d, s|a = k)$ and for a document $p(d, s)$ based on the probability $p(d_i, s_j)$ that a document d_i is covered by a summary sentence s_j . The probability $p(d_j, s_k)$ is estimated as the entailment probability that the document d_j entails the summary sentence s_k by a textual entailment model (Williams et al., 2018). Motivated by Laban et al. (2022), Equal Coverage divides documents into chunks of M words. The l -th chunk of the document d_j is denoted as $d_{j,l}$. The entailment model estimates the probability $p(d_j, s_k)$ as the maximum entailment probability $p(d_{j,l}, s_k)$ between the chunk $d_{j,l}$ and the summary sentence s_k :

$$p(d_i, s_j) = \max\{p(d_{i,l}, s_j)|d_{i,l} \in d_i\}, \quad (2)$$

where $p(d_{i,l}, s_j)$ is the probability that the document chunk $d_{i,l}$ entails the summary sentence s_j as per the entailment model. Based on the probability $p(d_i, s_j)$ that the document d_i is covered by the summary sentence s_j , the coverage probability for social attribute value i , $p(d, s|a = k)$ is empirically estimated as:

$$p(d, s|a = k) = \frac{1}{|D_k||S|} \sum_{d_i \in D_k} \sum_{s_j \in S} p(d_i, s_j), \quad (3)$$

where D_k is the set of documents d_i whose social attribute value a_i is k . Similarly, the coverage probability for a document, $p(d, s)$ is estimated as:

$$p(d, s) = \frac{1}{|D||S|} \sum_{d_j \in D} \sum_{s_k \in S} p(d_j, s_k). \quad (4)$$

Recall that for a fair summary S according to Equal Coverage, coverage probabilities for different social attribute values $p(d, s|a = k)$ should equal the coverage probability for a document $p(d, s)$. Therefore, Equal Coverage measures the summary-level fairness $EC(D, S)$ as:

$$EC(D, S) = \frac{1}{K} \sum_{k=1}^K |p(d, s) - p(d, s|a = k)|. \quad (5)$$

A high Equal Coverage value $EC(D, S)$ indicates less fairness because there are big differences between coverage probabilities for different social attribute values $p(d, s|a = k)$.

To determine if the differences between coverage probabilities for different social attribute values $p(d, s|a = k)$ are statistically significant, we perform a permutation test (Pitman, 1937). We randomly permute social attribute values a_i of documents and use them to calculate the permuted Equal Coverage value $\hat{EC}(D, S)$ 5000 times. The statistical significance of the difference is estimated as the proportion of permuted Equal Coverage values $\hat{EC}(D, S)$ greater than itself. We denote the proportion of samples (D, S) where the difference between coverage probabilities for different social attribute values is statistically significant ($p < 0.05$), $R_{EC}(G)$, as the proportion of unfair summaries, which is used to evaluate the fairness of the system.

4.2 Coverage Parity

Coverage Parity is a corpus-level fairness measure designed to measure the fairness of a system of all samples G . Motivated by Best-Worst Scaling (Louviere et al., 2015), Coverage Parity is based on the principle that different social attribute values should have equal chances of being the most overrepresented or underrepresented in summaries. It only considers the most overrepresented and underrepresented social attribute values because different samples' document sets D have different combinations of social attribute values. For instance, a document set can contain only positive and neutral reviews, while another can contain reviews with all sentiments. Therefore, Coverage Parity ensures that each sample contributes equally to it.

For a fair MDS system, Coverage Parity requires that the average coverage probability difference $c_k(d, s) = p(d, s) - p(d, s|a = k)$ is close to zero when the social attribute value k is the most overrepresented or underrepresented among all samples. In a sample (D, S) , we define the social attribute value i as the most overrepresented if its coverage probability difference $c_k(d, s)$ is the maximum $\max_k \{c_k(d, s)\}$ among all social attribute values $k \in \{1, \dots, K\}$, and the most underrepresented if its coverage probability difference $c_k(d, s)$ is the minimum $\min_k \{c_k(d, s)\}$. We collect these coverage probabilities differences $c_k(d, s)$ from all samples where social attribute value k is the most overrepresented or underrepresented in a set C_k .

Based on the average of the set of coverage probability differences C_i , Coverage Parity measures the fairness of the MDS system:

$$CP(G) = \frac{1}{K} \sum_{i=1}^K |\mathbb{E}(C_k)|. \quad (6)$$

A high Coverage Parity value $CP(G)$ indicates less fairness since it suggests that the chances of being overrepresented or underrepresented are very different for some social attribute values. Based on whether the average coverage probability differences, $\mathbb{E}(C_k)$, is greater or less than zero, we can also identify which social attribute value tends to be overrepresented or underrepresented.

5 Experimental Setup

We now describe experiments to evaluate the fairness of LLMs using our measures.

5.1 Datasets

We conduct experiments on five different datasets from the three domains: reviews, tweets, and news. These datasets are the Amazon (Ni et al., 2019), MITweet (Liu et al., 2023), SemEval (Mohammad et al., 2016), the Article Bias (Baly et al., 2020), and the News Stance (Ferreira and Vlachos, 2016; Pomerleau and Rao, 2017; Hanselowski et al., 2019) datasets. Tab. 1 shows the statistics of these datasets along with their social attribute values.

We observe that for some datasets, the fairness of summaries depends on the distributions of social attribute values in the input documents (Sec. 5.7). To balance social attribute values' impacts on fairness, we perform stratified sampling to collect 300 input document sets, D , for each dataset. Among these sampled sets, input document sets D dominated by

	Domain	Social Attribute	Social Attribute Values	Input Doc. Set Size	Avg. Doc. Len.
Amazon	Review	Sentiment	{negative, neutral, positive}	8	40
MITweet	Tweet	Political Ideology	{left, center, right}	20	34
Article Bias	News	Political Ideology	{left, center, right}	4-8	436
SemEval	Tweet	Stance toward Target	{support, against}	30	17
News Stance	News	Stance toward Target	{support, against}	4-8	240

Table 1: Dataset statistics for fairness evaluation in MDS

different social attribute values have equal proportions. The stratified sampling does not affect the calculation of our fairness measures. More details of preprocessing are in App. A.3.

5.2 Implementation Details

We evaluate the fairness of five families of LLMs: GPTs (Ouyang et al., 2022; Achiam et al., 2023) (GPT-3.5-0124, GPT-4-turbo-2024-04-09), Llama2 (Touvron et al., 2023) (Llama2-7b, Llama2-13b, Llama2-70b), Mixtral (Jiang et al., 2023, 2024) (Mixtral-7B-Instruct-v0.1, Mixtral-8x7B-Instruct-v0.1), Gemma (Team et al., 2024) (gemma-1.1-7b-it), and Claude3 (Bai et al., 2022) (claude-3-haiku-20240307, claude-3-sonnet-20240229). For comparison, we also evaluate the fairness of COOP (Iso et al., 2021) on the Amazon dataset and the fairness of PEGASUS (Zhang et al., 2020) and PRIMERA (Xiao et al., 2022) finetuned on the Multi-News (Fabbri et al., 2019) on the other datasets. The summary length is 100 words for the Article Bias and News Stance datasets and 50 for the other datasets. We describe summarization prompts in App. A.5.

To estimate the probability that a document is covered by a summary sentence (Eqn. 2), we use RoBERTa-large finetuned on the MNLI dataset (Williams et al., 2018). However, we show that our measures are independent of this choice (Sec. 5.3). To estimate the entailment probability using RoBERTa-large, we divide documents into chunks of $W = 100$ words. The chunk size W is tuned to maximize the proportion of summary sentences whose originating documents can be identified by RoBERTa-large (App. A.2) since LLMs are less prone to factual errors (Goyal et al., 2022).

5.3 Choice of Entailment Models

The implementations of our measures are independent of the choice of entailment model. To demonstrate this, we calculate our measures using three different textual entailment models: RoBERTa finetuned on the MNLI dataset; DeBERTa-large finetuned on multiple entailment datasets (Laurer et al.,

	Ro. vs De.	Ro. vs Al.	De. vs Al.
	Equal Coverage		
Amazon	<u>0.544</u>	<u>0.439</u>	<u>0.358</u>
MITweet	<u>0.430</u>	<u>0.442</u>	<u>0.357</u>
Article Bias	<u>0.470</u>	<u>0.709</u>	<u>0.429</u>
SemEval	<u>0.426</u>	<u>0.370</u>	<u>0.306</u>
News Stance	<u>0.712</u>	<u>0.780</u>	<u>0.703</u>
	Coverage Parity		
Amazon	<u>0.867</u>	<u>0.939</u>	<u>0.733</u>
MITweet	-0.127	0.006	<u>0.673</u>
Article Bias	0.624	<u>0.915</u>	0.612
SemEval	0.612	<u>0.891</u>	0.539
News Stance	<u>0.903</u>	<u>0.915</u>	<u>0.939</u>

Table 2: Spearman correlations between Equal Coverage (top) values and Coverage Parity (bottom) values computed using different entailment models: RoBERTa (Ro.), DeBERTa (De.), and ALBERT (Al.). Correlations that are statistically significant are underlined. We observe strong correlations, indicating that our measures are independent of the choice of the entailment model.

2024); and ALBERT-xl (Lan et al., 2019) finetuned on the MNLI and VitaminC (Schuster et al., 2021) datasets. We report the Spearman correlations between the Equal Coverage values, $EC(D, S)$, of each summary generated by all LLMs, obtained using these entailment models. We also report the Spearman correlations between the Coverage Parity value, $CP(G)$, of each LLM. The results are in Table 2. From the table, we can observe strong correlations between measures obtained using different textual entailment models on most datasets. It shows that these measures are independent of the choice of entailment models.

5.4 Human Perception of Fairness

We perform a human evaluation to determine which measure, Equal Coverage or Proportional Representation, aligns more with human perception of fairness. For Proportional Representation, we use the BARTScore (Yuan et al., 2021) implementation proposed by Zhang et al. (2023), as it shows the highest correlation with human perception.

Evaluating human perception of fairness is challenging due to the need to review entire input document sets. Therefore, we perform experiments on the Amazon dataset which only contains eight short reviews per input document set (Tab. 1). To

simplify the evaluation, we only consider input document sets containing only negative and positive but not neutral reviews, displayed in two randomized columns. For each input document set, we consider the summary generated by GPT-3.5 since it shows medium-level fairness (Tab. 3, 4). To further simplify the evaluation, we focus on summaries where Equal Coverage and Proportional Representation disagree on their fairness. We randomly select 25 samples containing input document sets and corresponding summaries that meet these criteria. Each sample is annotated by three annotators recruited from Amazon Mechanical Turk. The annotators should be from English-speaking countries and have HIT Approval Rates greater than 98%. More details are in App. A.4.

For each sample, annotators are asked to identify all unique negative and positive opinions in the input document set. They then evaluate whether the summary reflects these opinions and classify the summary as leaning negative, fair, or leaning positive. Out of these 25 samples, human perception of fairness aligns more with Equal Coverage in 17 samples, while aligns more with Proportional Representation in 8 samples. The difference is statistically significant ($p < 0.05$) using paired bootstrap resampling (Koehn, 2004). It shows that Equal Coverage aligns more with human perception of fairness than Proportional Representation.

We compare Coverage Parity with second-order fairness based on Proportional Representation to evaluate their effectiveness in identifying overrepresented sentiments on the group level. For this, we perform bootstrapping on the 25 samples to generate 5000 groups of 25 bootstrap samples each. For each groups of bootstrap samples, we define the sentiment that most bootstrap samples’ summary leaning toward based on human annotation as the sentiment overrepresented on the group level. We then compare the overrepresented sentiment based on human annotation with those identified by the Coverage Parity and second-order fairness. We find that Coverage Parity aligns with human annotation in 96% of the groups, while second-order fairness aligns in 4% of the groups. It shows that Coverage Parity aligns more with human perception of fairness than the second-order fairness.

5.5 Summary-level Fairness Evaluation

To evaluate the summary-level fairness of different LLMs, we report the proportion of unfair summaries according to Equal Coverage, $R_{EC}(G)$. We

	Amazon	MITweet	Article Bias	SemEval	News Stance	Overall
GPT-3.5	0.127	0.050	0.210	0.103	0.433	0.517
GPT-4	0.103	0.067	0.187	0.147	0.380	0.469
Llama2-7b	0.127	0.090	0.236	0.083	0.427	0.686
Llama2-13b	0.107	0.080	0.239	0.087	0.315	0.441
Llama2-70b	0.090	0.057	0.177	0.137	0.308	0.224
Mistral-7b	0.127	0.082	0.193	0.115	0.334	0.526
Mixtral-8x7b	0.117	0.097	0.172	0.128	0.372	0.549
Gemma	0.101	0.047	0.200	0.090	0.407	0.264
Claude3-haiku	0.110	0.063	0.174	0.093	0.500	0.411
Claude3-sonnet	0.107	0.070	0.174	0.120	0.401	0.401
COOP	0.204	-	-	-	-	-
PEGASUS	-	0.104	0.158	0.199	0.390	-
PRIMERA	-	0.118	0.134	0.118	0.377	-

Table 3: Proportion of unfair summaries and overall scores on different datasets according to Equal Coverage. A lower value indicates a fairer system. **Bold** indicates the fairest system.

also report an *Overall* score which is the average of normalized $R_{EC}(G)$ ($[0, 1]$) on all datasets. The results are in Tab. 3.

From the table, we observe Llama2-70b is the fairest. Among smaller LLMs (with 7-billion parameters), Gemma is the fairest. We also observe that almost all evaluated LLMs are fairer than COOP on the Amazon dataset, and PEGASUS and PRIMERA on the MITweet and Article Bias datasets. For the comparison within families of GPTs, Llama2, Claude3, we observe that larger models are generally fairer. However, Mixtral-8x7b is less fair than Mistral-7b although Mixtral-8x7b is larger, which might be because it scales its size using the mixture of smaller models.

Previous works by Zhang et al. (2023) and Lei et al. (2024) evaluate the fairness of LLMs using Proportional Representation. When evaluating the fairness on sentiments in the review domain, our results using Equal Coverage are consistent with these works in the finding that GPT-4 is fairer than GPT-3.5, and both are fairer than COOP. However, they find that Llama2-13b is the fairest, while our results show that larger models are fairer for Llama2. When evaluating the fairness on political ideologies in the tweet domain, our results are consistent with these works on that GPT-4 is fairer than GPT-3.5. However, they find that smaller models are fairer for Llama2, and GPTs are less fair than Llama2 models. Contrarily, we show that larger models are fairer for Llama2, and GPTs are fairer than Llama2-7b. As shown in Sec. 5.4, Equal Coverage aligns more with human perception, suggesting our results better reflect human judgments.

5.6 Corpus-level Fairness Evaluation

To evaluate the corpus-level fairness of different LLMs, we report the Coverage Parity, $CP(G)$, on

	Amazon			MITweet			Article Bias			SemEval			News Stance			Overall
	$CP(G)$	over	under	$CP(G)$	over	under	$CP(G)$	over	under	$CP(G)$	over	under	$CP(G)$	over	under	
GPT-3.5	0.032	<u>neg.</u>	pos.	0.012	<u>left</u>	<u>center</u>	0.015	center	<u>left</u>	0.012	<u>sup.</u>	<u>aga.</u>	0.018	aga.	<u>sup.</u>	0.475
GPT-4	0.029	<u>neg.</u>	pos.	0.011	<u>left</u>	<u>center</u>	0.024	<u>right</u>	<u>left</u>	0.011	<u>sup.</u>	<u>aga.</u>	0.038	<u>aga.</u>	<u>sup.</u>	0.647
Llama2-7b	0.037	<u>neg.</u>	pos.	0.013	<u>left</u>	<u>center</u>	0.026	<u>right</u>	<u>left</u>	0.010	<u>sup.</u>	<u>aga.</u>	0.018	<u>aga.</u>	<u>sup.</u>	0.597
Llama2-13b	0.017	<u>neg.</u>	pos.	0.009	<u>left</u>	<u>center</u>	0.025	<u>right</u>	<u>left</u>	0.014	<u>sup.</u>	<u>aga.</u>	0.022	<u>aga.</u>	<u>sup.</u>	0.460
Llama2-70b	0.008	<u>neg.</u>	neu.	0.008	<u>left</u>	<u>center</u>	0.025	<u>right</u>	<u>left</u>	0.007	<u>sup.</u>	<u>aga.</u>	0.012	<u>aga.</u>	<u>sup.</u>	0.177
Mistral-7b	0.024	<u>neg.</u>	pos.	0.009	<u>left</u>	<u>center</u>	0.029	<u>right</u>	<u>left</u>	0.014	<u>sup.</u>	<u>aga.</u>	0.023	<u>aga.</u>	<u>sup.</u>	0.551
Mixtral-8x7b	0.017	<u>neg.</u>	pos.	0.010	<u>left</u>	<u>center</u>	0.036	<u>right</u>	<u>left</u>	0.003	<u>sup.</u>	<u>aga.</u>	0.020	<u>aga.</u>	<u>sup.</u>	0.394
Gemma	0.025	<u>neg.</u>	pos.	0.015	<u>left</u>	<u>center</u>	0.020	<u>right</u>	<u>left</u>	0.014	<u>sup.</u>	<u>aga.</u>	0.032	<u>aga.</u>	<u>sup.</u>	0.718
Claude3-haiku	0.028	<u>neg.</u>	pos.	0.013	<u>left</u>	<u>center</u>	0.016	<u>center</u>	<u>left</u>	0.012	<u>sup.</u>	<u>aga.</u>	0.042	<u>aga.</u>	<u>sup.</u>	0.651
Claude3-sonnet	0.025	<u>neg.</u>	pos.	0.008	<u>left</u>	<u>center</u>	0.018	<u>right</u>	<u>left</u>	0.002	<u>sup.</u>	<u>aga.</u>	0.030	<u>aga.</u>	<u>sup.</u>	0.264
COOP	0.068	<u>pos.</u>	<u>neg.</u>	-	-	-	-	-	-	-	-	-	-	-	-	-
PEGASUS	-	-	-	0.009	<u>right</u>	<u>left</u>	0.035	<u>right</u>	<u>left</u>	0.001	<u>sup.</u>	<u>aga.</u>	0.022	<u>sup.</u>	<u>aga.</u>	-
PRIMERA	-	-	-	0.014	<u>right</u>	<u>left</u>	0.032	<u>right</u>	<u>left</u>	0.002	<u>aga.</u>	<u>sup.</u>	0.023	<u>sup.</u>	<u>aga.</u>	-

Table 4: Coverage Parity, $CP(G)$, and the most overrepresented (over) and underrepresented (under) social attribute values, and overall scores on different datasets. A lower value of $CP(G)$ indicates a fairer system. **Bold** indicates the fairest system. The social attribute values whose average coverage probability differences are statistically significantly ($p < 0.05$) different from zero based on Bootstrapping are underlined.

different datasets. For each dataset, we report the most overrepresented and underrepresented social attribute value i whose average coverage probability difference, $\mathbb{E}(C_i)$, is the maximum and minimum respectively. We also report an *Overall* score which is the average of normalized $CP(G)$ ($[0, 1]$) on all datasets. The results are in Tab. 4.

From the table, we observe that Llama2-70b is the fairest. Among smaller LLMs (with 7-billion parameters), Mistral-7b is the fairest. We also observe that evaluated LLMs are less fair than PEGASUS and PRIMERA on the SemEval dataset. While comparing within each family of LLMs, we observe that larger models are fairer for the families of Llama2, Mixtral, and Claude3. However, for the family of GPTs, GPT-4 is less fair than GPT-3.5, suggesting that larger models are not necessarily more fair on the system level. Besides, we observe that the fairness measured by Coverage Parity and Equal Coverage are different on some datasets. The difference indicates that we should consider both summary-level fairness and corpus-level fairness for comprehensively measuring fairness in MDS.

We can also observe that most LLMs overrepresent and underrepresent certain social attribute values on different datasets. For the Amazon dataset, most LLMs overrepresent negative reviews and underrepresent positive reviews. Contrarily, COOP overrepresents positive reviews and underrepresents negative reviews. For the MITweet and Article Bias datasets, all LLMs overrepresent left-leaning documents in the tweet domain, while most LLMs overrepresent right-leaning documents in the news domain. Contrarily, PEGASUS and PRIMERA overrepresent right-leaning documents for all domains. We can observe the same pat-

tern for the SemEval and News Stance datasets. All LLMs overrepresent documents supporting the target in the tweet domain but overrepresent documents against the target in the news domain. These results indicate that summaries generated by LLMs might overrepresent documents with certain social attribute values. Users should know this before they make judgments based on these summaries. For example, users should know that a review summary generated by LLMs for a product might unfairly overrepresent negative reviews so they can calibrate their perception of the product accordingly. Developers can also build fairer LLMs for MDS based on these results.

5.7 Fairness under Different Distributions of Social Attributes

We perform experiments to evaluate whether the fairness of LLMs changes under different distributions of social attribute values in input document sets. For this, we divide all samples G into K non-overlapping sets: $\{G_1, \dots, G_K\}$ based on distributions of social attribute values. Each set G_i includes samples where most documents $d \in D$ have a social attribute value of i . We denote the set G_i as the set dominated by social attribute value i . To measure differences of the summary-level fairness measured by Equal Coverage under different distributions, we use maximum differences of proportions of unfair summaries $R_{EC}(G_i)$ (Sec. 4.1) on sets dominated by different social attribute values G_i . To measure differences of the corpus-level fairness measured by Coverage Parity under different distributions, we use maximum differences of average coverage probability differences for the same social attribute value $\mathbb{E}(C_j)$ (Sec. 4.2) on

	Amazon	MITweet	Article Bias	SemEval	News Stance
Equal Coverage					
GPT-3.5	0.080	0.050	0.070	0.067	0.061
GPT-4	0.070	0.040	0.040	0.067	0.044
Llama2-7b	<u>0.110</u>	<u>0.040</u>	<u>0.030</u>	<u>0.020</u>	<u>0.032</u>
Llama2-13b	0.010	0.020	0.060	0.053	0.022
Llama2-70b	<u>0.100</u>	<u>0.060</u>	<u>0.122</u>	<u>0.052</u>	<u>0.028</u>
Mistral-7b	0.070	<u>0.071</u>	<u>0.051</u>	<u>0.103</u>	<u>0.031</u>
Mixtral-8x7b	0.120	0.039	0.051	0.053	0.025
Gemma	<u>0.069</u>	<u>0.050</u>	<u>0.070</u>	<u>0.067</u>	<u>0.038</u>
Claude3-haiku	0.060	0.040	0.042	0.023	0.053
Claude3-sonnet	<u>0.110</u>	0.050	0.060	0.028	<u>0.107</u>
Coverage Parity					
GPT-3.5	<u>0.085</u>	<u>0.030</u>	0.029	<u>0.011</u>	<u>0.093</u>
GPT-4	0.086	0.030	0.075	0.009	0.078
Llama2-7b	<u>0.128</u>	<u>0.027</u>	<u>0.038</u>	<u>0.011</u>	<u>0.091</u>
Llama2-13b	0.088	0.016	0.074	0.014	0.116
Llama2-70b	0.099	0.023	0.044	0.008	0.068
Mistral-7b	<u>0.072</u>	<u>0.012</u>	<u>0.046</u>	<u>0.017</u>	<u>0.120</u>
Mixtral-8x7b	0.069	<u>0.022</u>	0.027	0.008	0.061
Gemma	0.089	<u>0.029</u>	0.036	<u>0.017</u>	<u>0.063</u>
Claude3-haiku	<u>0.103</u>	0.021	<u>0.061</u>	<u>0.013</u>	<u>0.064</u>
Claude3-sonnet	<u>0.083</u>	0.022	<u>0.072</u>	<u>0.010</u>	<u>0.093</u>

Table 5: Differences of Equal Coverage and Coverage Parity under different distributions of social attribute values in input document sets. Statistically significant ($p < 0.05$) differences based on Bootstrapping are underlined. We observe significant differences in Coverage Parity for most LLMs. It indicates that the social attribute values overrepresented or underrepresented change significantly under different distributions.

different sets G_i . The results are in Table 5.

From the table, we observe that Equal Coverage remains stable under different distributions. However, for most LLMs, the differences of Coverage Parity are significant. It suggests that the social attribute values overrepresented or underrepresented change significantly under different distributions of social attribute values in input document sets. We further analyze the data from the Amazon and News Stance datasets, where the differences are significant for all LLMs. On these two datasets, almost all LLMs tend to overrepresent the social attribute values that dominate the input document sets. Therefore, corpus-level fairness measures like Coverage Parity might mistakenly judge which social attribute values are overrepresented based on the size of set dominated by different social attribute values $|G_i|$, showing the need of balanced datasets for evaluating the fairness. We address this issue through stratified sampling (Sec. 5.1).

5.8 LLM Perception of Fairness

We perform experiments to evaluate which measure, Equal Coverage or Proportional Representation, aligns more with LLMs’ perception of fairness. This is an exploratory experiment, and we do not assume the LLMs’ perception of fairness as ground truth. We prompt an LLM to generate a

	Amazon	MITweet	Article Bias	SemEval	News Stance
GPT-3.5	0.088	-0.035	0.070	0.051	-0.018
Llama2-70b	0.062	0.021	0.016	0.059	-0.042
Mixtral-8x7b	-0.025	-0.087	0.008	0.038	0.018
Gemma	-0.039	-0.035	-0.055	-0.020	0.061
Claude3-haiku	-0.094	-0.101	-0.095	0.041	-0.008

Table 6: Comparison of relative changes of Equal Coverage vs. Proportional Representation when LLMs are prompted to generate fair summaries. A positive value indicates the summary-level fairness measured by Equal Coverage decreases more compared to Proportional Representation, which suggests that Equal Coverage aligns more with LLMs’ perception of fairness.

summary for an input document set, then prompt it again to generate a fair summary for the same set. The second prompt requires that the summary fairly represent documents with different social attributes (App. A.6). However, it does not provide any other details about fairness, allowing the LLM to decide. The prompt also includes the social attribute value for each document. We compute the Equal Coverage and Proportional Representation for both summaries and consider the relative change in values before and after the LLM is prompted to generate a fair summary. If a measure aligns more with LLM’s perception of fairness, the score for the ‘fair’ summary should be lower. The differences between average relative changes of Equal Coverage and Proportional Representation are in Tab. 6.

From the table, we observe positive differences for most LLMs, suggesting that Equal Coverage decreases more compared to Proportional Representation. It means that Equal Coverage aligns more with LLM’s perception of fairness. Specifically, Proportional Representation aligns more with the perception of fairness of Gemma and Claude3-haiku, while Equal Coverage aligns more with the remaining LLMs.

6 Conclusion

We propose two coverage-based fairness measures for MDS, Equal Coverage for measuring summary-level fairness and Coverage Parity for measuring corpus-level fairness. Using these measures, we find that Llama2-70b is the fairest among all LLMs. We also find that almost all LLMs overrepresent certain social attribute values in each domain.

Future works can explore the effect of training data, especially instruction tuning and preference tuning data, on the fairness of LLMs. Future works can also finetune LLMs based on our measures to develop fairer models.

7 Limitations

The effectiveness of two proposed measures, Equal Coverage and Coverage Parity, relies on whether the probability that a document entails a summary sentence estimated by the entailment model is accurate. To evaluate the performance of the entailment model for such a task, previous works generally use the accuracy of the entailment model on the fact verification dataset or the correlation between the factuality scores of summaries annotated by humans with the factuality scores estimated by the entailment model on the summarization evaluation benchmark. Although there are several fact verification datasets and summarization evaluation benchmarks in the news domain, there are no such datasets in the reviews and tweets domain to our best knowledge. Therefore, we cannot evaluate the accuracy or perform calibration for the entailment models in these two domains. However, as shown in Sec. 5.3, Equal Coverage and Coverage Parity based on different commonly used entailment models are mostly correlated. These entailment models are also widely used for measuring factuality in summarization tasks (Maynez et al., 2020; Laban et al., 2022).

8 Ethical Consideration

The datasets we use are all publicly available. We do not annotate any data on our own. All the models used in this paper are publicly accessible. We do not do any training in this paper. For the inference of Llama2-7b, Llama2-13b, Mistral-7b, and Gemma, we use on Nvidia A6000 GPU. For the inference of Llama2-70b and Mixtral-8x7b, we use 4 Nvidia A6000 GPUs. For all other experiments, we use one Nvidia V100 GPU.

We perform human evaluation experiments on Amazon Mechanical Turk. The annotators were compensated at a rate of \$15 per hour. During the evaluation, human annotators were not exposed to any sensitive or explicit content.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini,

Cameron McKinnon, et al. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Ramy Baly, Giovanni Da San Martino, James Glass, and Preslav Nakov. 2020. We can detect your bias: Predicting the political ideology of news articles. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4982–4991.

Adithya Bhaskar, Alex Fabbri, and Greg Durrett. 2023. Prompted opinion summarization with gpt-3.5. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9282–9300.

Iman Munire Bilal, Bo Wang, Adam Tsakalidis, Dong Nguyen, Rob Procter, and Maria Liakata. 2022. [Template-based abstractive microblog opinion summarization](#). *Transactions of the Association for Computational Linguistics*, 10:1229–1248.

Arthur Bražinskas, Mirella Lapata, and Ivan Titov. 2019. Unsupervised opinion summarization as copycat-review generation. *arXiv preprint arXiv:1911.02247*.

Abhisek Dash, Anurag Shandilya, Arindam Biswas, Kripabandhu Ghosh, Saptarshi Ghosh, and Abhijnan Chakraborty. 2019. Summarizing user-generated textual content: Motivation and methods for fairness in algorithmic summaries. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–28.

Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. 2019. [Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1074–1084, Florence, Italy. Association for Computational Linguistics.

William Ferreira and Andreas Vlachos. 2016. Emergent: a novel data-set for stance classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies. ACL*.

Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of gpt-3. *arXiv preprint arXiv:2209.12356*.

Andreas Hanselowski, Christian Stab, Claudia Schulz, Zile Li, and Iryna Gurevych. 2019. A richly annotated corpus for different tasks in automated fact-checking. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 493–503.

Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022. [TRUE: Re-evaluating factual consistency evaluation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3905–3920, Seattle,

736	United States. Association for Computational Linguistics.	
737		
738	Nannan Huang, Haytham Fayek, and Xiuzhen Zhang.	
739	2024. Bias in opinion summarisation from pre-	
740	training to adaptation: A case study in political bias.	
741	<i>arXiv preprint arXiv:2402.00322</i> .	
742	Hayate Iso, Xiaolan Wang, Yoshihiko Suhara, Stefanos	
743	Angelidis, and Wang-Chiew Tan. 2021. Convex ag-	
744	gregation for opinion summarization. In <i>Findings</i>	
745	<i>of the Association for Computational Linguistics:</i>	
746	<i>EMNLP 2021</i> , pages 3885–3903.	
747	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	
748	sch, Chris Bamford, Devendra Singh Chaplot, Diego	
749	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	
750	laume Lample, Lucile Saulnier, et al. 2023. Mistral	
751	7b. <i>arXiv preprint arXiv:2310.06825</i> .	
752	Albert Q Jiang, Alexandre Sablayrolles, Antoine	
753	Roux, Arthur Mensch, Blanche Savary, Chris Bam-	
754	ford, Devendra Singh Chaplot, Diego de las Casas,	
755	Emma Bou Hanna, Florian Bressand, et al. 2024.	
756	Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	
757	Vijay Keswani and L Elisa Celis. 2021. Dialect diversity	
758	in text summarization on twitter. In <i>Proceedings of</i>	
759	<i>the Web Conference 2021</i> , pages 3802–3814.	
760	Philipp Koehn. 2004. Statistical significance tests for	
761	machine translation evaluation . In <i>Proceedings of the</i>	
762	<i>2004 Conference on Empirical Methods in Natural</i>	
763	<i>Language Processing</i> , pages 388–395, Barcelona,	
764	Spain. Association for Computational Linguistics.	
765	Philippe Laban, Tobias Schnabel, Paul N Bennett, and	
766	Marti A Hearst. 2022. Summac: Re-visiting nli-	
767	based models for inconsistency detection in summa-	
768	rization. <i>Transactions of the Association for Compu-</i>	
769	<i>tational Linguistics</i> , 10:163–177.	
770	Zhenzhong Lan, Mingda Chen, Sebastian Goodman,	
771	Kevin Gimpel, Piyush Sharma, and Radu Soricut.	
772	2019. Albert: A lite bert for self-supervised learning	
773	of language representations. In <i>International Confer-</i>	
774	<i>ence on Learning Representations</i> .	
775	Moritz Laurer, Wouter Van Atteveldt, Andreu Casas,	
776	and Kasper Welbers. 2024. Less annotating, more	
777	classifying: Addressing the data scarcity issue of su-	
778	pervised machine learning with deep transfer learning	
779	and bert-nli. <i>Political Analysis</i> , 32(1):84–100.	
780	Nayeon Lee, Yejin Bang, Tiezheng Yu, Andrea Madotto,	
781	and Pascale Fung. 2022. NeuS: Neutral multi-news	
782	summarization for mitigating framing bias. In <i>Pro-</i>	
783	<i>ceedings of the 2022 Conference of the North Amer-</i>	
784	<i>ican Chapter of the Association for Computational</i>	
785	<i>Linguistics: Human Language Technologies</i> , pages	
786	3131–3148, Seattle, United States. Association for	
787	Computational Linguistics.	
788	Yuanyuan Lei, Kaiqiang Song, Sangwoo Cho, Xiaoyang	
789	Wang, Ruihong Huang, and Dong Yu. 2024. Po-	
790	larity calibration for opinion summarization. <i>arXiv</i>	
791	<i>preprint arXiv:2404.01706</i> .	
	Songtao Liu, Ziling Luo, Minghua Xu, LiXiao Wei,	792
	Ziyao Wei, Han Yu, Wei Xiang, and Bang Wang.	793
	2023. Ideology takes multiple looks: A high-quality	794
	dataset for multifaceted ideology detection. In <i>The</i>	795
	<i>2023 Conference on Empirical Methods in Natural</i>	796
	<i>Language Processing</i> .	797
	Y Liu, X Zhang, D Wegsman, N Beauchamp, and	798
	L Wang. 2022. Politics: Pretraining with same-story	799
	article comparison for ideology prediction and stance	800
	detection. <i>Findings of the Association for Computa-</i>	801
	<i>tional Linguistics: NAACL 2022</i> .	802
	Jordan J Louviere, Terry N Flynn, and Anthony Al-	803
	fred John Marley. 2015. <i>Best-worst scaling: Theory,</i>	804
	<i>methods and applications</i> . Cambridge University	805
	Press.	806
	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and	807
	Ryan McDonald. 2020. On faithfulness and factu-	808
	ality in abstractive summarization. <i>arXiv preprint</i>	809
	<i>arXiv:2005.00661</i> .	810
	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis,	811
	Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettle-	812
	moyer, and Hannaneh Hajishirzi. 2023. Factscore:	813
	Fine-grained atomic evaluation of factual precision	814
	in long form text generation. In <i>Proceedings of the</i>	815
	<i>2023 Conference on Empirical Methods in Natural</i>	816
	<i>Language Processing</i> , pages 12076–12100.	817
	Saif Mohammad, Svetlana Kiritchenko, Parinaz Sob-	818
	hani, Xiaodan Zhu, and Colin Cherry. 2016.	819
	SemEval-2016 task 6: Detecting stance in tweets .	820
	In <i>Proceedings of the 10th International Workshop</i>	821
	<i>on Semantic Evaluation (SemEval-2016)</i> , pages 31–	822
	41, San Diego, California. Association for Computa-	823
	tional Linguistics.	824
	Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019.	825
	Justifying recommendations using distantly-labeled	826
	reviews and fine-grained aspects . In <i>Proceedings</i>	827
	<i>of the 2019 Conference on Empirical Methods in</i>	828
	<i>Natural Language Processing and the 9th Interna-</i>	829
	<i>tional Joint Conference on Natural Language Pro-</i>	830
	<i>cessing (EMNLP-IJCNLP)</i> , pages 188–197, Hong	831
	Kong, China. Association for Computational Lin-	832
	guistics.	833
	Olubusayo Olabisi, Aaron Hudson, Antonie Jetter, and	834
	Ameeta Agrawal. 2022. Analyzing the dialect diver-	835
	sity in multi-document summaries . In <i>Proceedings of</i>	836
	<i>the 29th International Conference on Computational</i>	837
	<i>Linguistics</i> , pages 6208–6221, Gyeongju, Republic	838
	of Korea. International Committee on Computational	839
	Linguistics.	840
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	841
	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	842
	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	843
	2022. Training language models to follow instruc-	844
	tions with human feedback. <i>Advances in Neural</i>	845
	<i>Information Processing Systems</i> , 35:27730–27744.	846

847 Edwin JG Pitman. 1937. Significance tests which may
848 be applied to samples from any populations. *Supple-*
849 *ment to the Journal of the Royal Statistical Society*,
850 4(1):119–130.

851 Dean Pomerleau and Delip Rao. 2017. The fake news
852 challenge: Exploring how artificial intelligence tech-
853 nologies could be leveraged to combat fake news.
854 *Fake news challenge*.

855 Tal Schuster, Adam Fisch, and Regina Barzilay. 2021.
856 [Get your vitamin C! robust fact verification with](#)
857 [contrastive evidence](#). In *Proceedings of the 2021*
858 *Conference of the North American Chapter of the*
859 *Association for Computational Linguistics: Human*
860 *Language Technologies*, pages 624–643, Online. As-
861 sociation for Computational Linguistics.

862 Anurag Shandilya, Abhisek Dash, Abhijnan
863 Chakraborty, Kripabandhu Ghosh, and Saptarshi
864 Ghosh. 2020. Fairness for whom? understanding the
865 reader’s perception of fairness in text summarization.
866 In *2020 IEEE International Conference on Big Data*
867 *(Big Data)*, pages 3692–3701. IEEE.

868 Anurag Shandilya, Kripabandhu Ghosh, and Saptarshi
869 Ghosh. 2018. Fairness of extractive text summariza-
870 tion. In *Companion Proceedings of the The Web*
871 *Conference 2018*, pages 97–98.

872 Gemma Team, Thomas Mesnard, Cassidy Hardin,
873 Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,
874 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale,
875 Juliette Love, et al. 2024. Gemma: Open models
876 based on gemini research and technology. *arXiv*
877 *preprint arXiv:2403.08295*.

878 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-
879 bert, Amjad Almahairi, Yasmine Babaei, Nikolay
880 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti
881 Bhosale, et al. 2023. Llama 2: Open founda-
882 tion and fine-tuned chat models. *arXiv preprint*
883 *arXiv:2307.09288*.

884 Sahil Verma and Julia Rubin. 2018. Fairness defini-
885 tions explained. In *Proceedings of the international*
886 *workshop on software fairness*, pages 1–7.

887 Adina Williams, Nikita Nangia, and Samuel Bowman.
888 2018. [A broad-coverage challenge corpus for sen-](#)
889 [tence understanding through inference](#). In *Proceed-*
890 *ings of the 2018 Conference of the North American*
891 *Chapter of the Association for Computational Lin-*
892 *guistics: Human Language Technologies, Volume 1*
893 *(Long Papers)*, pages 1112–1122. Association for
894 Computational Linguistics.

895 Wen Xiao, Iz Beltagy, Giuseppe Carenini, and Arman
896 Cohan. 2022. [PRIMERA: Pyramid-based masked](#)
897 [sentence pre-training for multi-document summariza-](#)
898 [tion](#). In *Proceedings of the 60th Annual Meeting of*
899 *the Association for Computational Linguistics (Vol-*
900 *ume 1: Long Papers)*, pages 5245–5263, Dublin,
901 Ireland. Association for Computational Linguistics.

	Sent.	50	100	200	400
GPT-3.5	0.800	0.858	0.876	0.874	0.716
Llama2-7b	0.725	0.773	0.792	0.788	0.642

Table 7: Proportion of summary sentences whose origi-
nating documents are identified by the entailment model
when using a document sentence (Sent.) or chunks with
different sizes as the premises. **Bold** indicates the opti-
mal chunk size for identifying originating documents of
summary sentences.

Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. 902
[Bartscore: Evaluating generated text as text genera-](#) 903
[tion](#). In *Advances in Neural Information Processing* 904
Systems, volume 34, pages 27263–27277. Curran As- 905
sociates, Inc. 906

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Pe- 907
ter Liu. 2020. Pegasus: Pre-training with extracted 908
gap-sentences for abstractive summarization. In *In-* 909
ternational conference on machine learning, pages 910
11328–11339. PMLR. 911

Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, 912
Kathleen McKeown, and Tatsunori B. Hashimoto. 913
2024. [Benchmarking Large Language Models for](#) 914
[News Summarization](#). *Transactions of the Associa-* 915
tion for Computational Linguistics, 12:39–57. 916

Yusen Zhang, Nan Zhang, Yixin Liu, Alexander Fab- 917
bri, Junru Liu, Ryo Kamoi, Xiaoxin Lu, Caiming 918
Xiong, Jieyu Zhao, Dragomir Radev, et al. 2023. Fair 919
abstractive summarization of diverse perspectives. 920
arXiv preprint arXiv:2311.07884. 921

A Appendix 922

A.1 Split and Rephrase Summary Sentences 923

Motivated by [Bhaskar et al. \(2023\)](#); [Min et al. \(2023\)](#), we split all summary sentences into simple sentences. The goal of this step is to ensure that each summary sentence after the split only discusses a single fact. Specifically, compound sentences are split into simple sentences, while sentences with compound subjects or objects are split into sentences with simple subjects or objects. After the split, summary sentences are then rephrased to remove the reported speech, like ‘documents say what’. For splitting and rephrasing, we prompt GPT-3.5 with demonstrations. We show the example prompt and the example summary sentences after splitting and rephrasing in Fig. 2. 924
925
926
927
928
929
930
931
932
933
934
935
936
937

A.2 Document Chunking 938

To estimate the probability that a document is covered by a summary sentence (Eqn. 2), we divide the document into chunks of no more than W words. Each chunk contains one or several neighboring 939
940
941
942

Prompt
<p>Each atomic content unit contain an atomic fact and does not need further split for the purpose of reducing ambiguity in human evaluation. If a sentence contain compound objects or compound subjects, extract an atomic content unit for each object or subject. Therefore, different atomic content units can share some similar contents. Do not use reported speech for any atomic content units. For each atomic content unit, replace 'the articles', 'the news articles', 'the media', 'the narrative' or their synonyms at the begining with 'it'. Below are some examples of extracted atomic content units of corresponding summaries.</p> <p>Summary: In a series of articles discussing the vice-presidential debate between Mike Pence and Tim Kaine, Pence emerged as the clear winner with his calm demeanor and strong debating skills. Pence strategically avoided defending Trump's controversial statements and focused on attacking Hillary Clinton instead. His performance provided a blueprint for other GOP candidates on how to navigate challenging situations. Despite Pence's success, the articles highlighted the ongoing challenges faced by the Republican Party in managing the fallout from Trump's actions and statements. Ultimately, Pence's debate victory may not have a significant impact on the overall election dynamics.</p> <p>Extracted Atomic Content Units:</p> <ol style="list-style-type: none"> (1) It discussed the vice-presidential debate between Mike Pence and Tim Kaine. (2) Pence emerged as the clear winner. (3) Pence emerged as the clear winner with his calm demeanor. (4) Pence emerged as the clear winner with his strong debating skills. (5) Pence strategically avoided defending Trump's controversial statements. (6) Pence focused on attacking Hillary Clinton instead. (7) His performance provided a blueprint for other GOP candidates on how to navigate challenging situations. (8) Pence is successful. (9) It highlighted the ongoing challenges faced by the Republican Party in managing the fallout from Trump's actions. (10) It highlighted the ongoing challenges faced by the Republican Party in managing the fallout from Trump's statements. (11) Ultimately, Pence's debate victory may not have a significant impact on the overall election dynamics. <p>Extract atomic content units for the following summary following the form of the above examples.</p> <p>Summary: The news highlights the recent wave of restrictive abortion laws passed in several Republican-controlled states, including Missouri and Alabama, which aim to challenge the landmark Roe v. Wade Supreme Court decision. The bills ban abortions after as early as 6-8 weeks, with limited exceptions, and are seen as part of a broader effort to overturn or erode legal abortion rights in the United States. These measures have sparked outcry from pro-choice advocates who argue they infringe on women's rights and could push women towards dangerous, illegal abortions. The issue has also become a political battleground, with the 2020 presidential election looming and the potential to reshape the Supreme Court's abortion-related jurisprudence.</p> <p>Extracted Atomic Content Units:</p>
Output
<ol style="list-style-type: none"> (1) It highlights the recent wave of restrictive abortion laws passed in several Republican-controlled states. (2) The bills ban abortions after as early as 6-8 weeks, with limited exceptions. (3) The bills are seen as part of a broader effort to overturn or erode legal abortion rights in the United States. (4) These measures have sparked outcry from pro-choice advocates. (5) Pro-choice advocates argue they infringe on women's rights. (6) Pro-choice advocates argue they could push women towards dangerous, illegal abortions. (7) The issue has become a political battleground. (8) The issue has become a political battleground with the 2020 presidential election looming. (9) The issue has the potential to reshape the Supreme Court's abortion-related jurisprudence.

Figure 2: Example prompt for splitting and rephrasing summary sentences (top) and summary sentences after splitting and rephrasing (bottom).

943	sentences of the document. Since LLMs are less	or against the claim. We also filter out duplicated	992
944	prone to factual errors (Goyal et al., 2022; Zhang	news or news longer than 600 words or shorter than	993
945	et al., 2024), the chunk size W is tuned to maximize	75 words. Each input document set contains 4 to 8	994
946	the proportion of summary sentences whose origi-	news supporting or against the same claim.	995
947	inating documents are identified by the Roberta-		
948	large is the highest. We tune the chunk size W	MITweet (Liu et al., 2023) consists of tweets	996
949	based on the average proportions of summary sen-	with labels of political ideologies on different facets	997
950	tences generated by GPT-3.5 and Llama2-7b on all	about different topics. We cluster tweets about the	998
951	datasets. The results are shown in Tab. 7. We can	same topic based on their TFIDF similarity into	999
952	observe that when the chunk size is 100, the pro-	clusters. We then divide these clusters into input	1000
953	portion of summary sentences whose originating	document sets of 20 tweets about the same topic.	1001
954	documents are identified by the Roberta-large is	The social attribute of a tweet will be left if it is left	1002
955	the highest. The result is consistent with the finds	on most facets, right if it is right on most facets, oth-	1003
956	of Honovich et al. (2022). Therefore, we set the	erwise neutral. Compared with the Election dataset	1004
957	chunk size W as 100.	(Shandilya et al., 2018), the MITweet dataset con-	1005
		tains tweets about more diverse topics other than	1006
958	A.3 Datasets	election, such as ‘Abortion’ and ‘Energy Crisis’.	1007
959	In this section, we describe the reason for choos-	Compared with the MOS dataset (Bilal et al., 2022)	1008
960	ing these datasets and how we preprocess these	used by Huang et al. (2024), the MITweet dataset	1009
961	datasets.	covers more diverse topics and has the ground-truth	1010
		label of social attribute value.	1011
962	Amazon (Ni et al., 2019) consists of reviews with	Tweet Stance (Mohammad et al., 2016) consists	1012
963	labels of their ratings of different products. We	of tweets with labels of stance toward a short phrase	1013
964	filter out reviews that are non-English or without	such as Climate Change or Hillary Clinton. We	1014
965	ratings. The input document set of this dataset con-	cluster tweets about the same short phrase based	1015
966	tains 8 reviews of the same product. We obtain the	on their TFIDF similarity into clusters. We then	1016
967	social attribute of each review based on its rating	divide these clusters into input document sets of 20	1017
968	provided in the dataset. The social attribute of a	tweets about the same short phrase.	1018
969	review will be positive if its rating is 4 or 5, neutral		
970	if its rating is 3, and negative if its rating is 1 or 2.	A.4 Human Evaluation	1019
971	Article Bias (Bražinskas et al., 2019) consists	For each sample of an input document set and its	1020
972	of news with labels of their political ideologies.	corresponding summary, annotators are asked to	1021
973	We run the clustering algorithm on this news to	identify all unique negative and positive opinions in	1022
974	generate a cluster of news about the same event	the input document set. They then evaluate whether	1023
975	following Liu et al. (2022). We then divide these	the summary reflects these opinions and classify	1024
976	clusters into input document sets of 4 to 8 news of	the summary as leaning negative, fair, or leaning	1025
977	the same event. For each news, we also perform	positive. To simplify the annotation, we provide	1026
978	truncation from the beginning to fit the context	annotators with unique opinions extracted by GPT-	1027
979	length restriction of Llama2. Compared with the	3.5. The interface for human evaluation is shown	1028
980	NeuS dataset (Lee et al., 2022) used by Lei et al.	in Fig. 3. A sample will be annotated as leaning	1029
981	(2024), the input document sets of the Article Bias	negative if more annotators annotate it as leaning	1030
982	dataset contain more input document per set and the	negative, leaning positive if more annotators an-	1031
983	distribution of social attributes in input documents	notate it as leaning positive, otherwise fair. For	1032
984	are more diverse.	a sample leaning negative or positive, we say the	1033
985	News Stance consists of news with labels of their	human perception of fairness aligns more with a	1034
986	stances toward claims, such as ‘Meteorite strike	fairness measure if the overrepresented sentiment	1035
987	in Nicaragua puzzles experts’. The dataset com-	identified by the measure is the same as the senti-	1036
988	combines news from three news stance datasets (Fer-	ment that the sample leans toward. For a sample	1037
989	reira and Vlachos, 2016; Pomerleau and Rao, 2017;	annotated as fair, we say the human perception of	1038
990	Hanselowski et al., 2019). For each claim, we only	fairness aligns more with a fairness measure if its	1039
991	keep news whose stances are directly supporting	normalized absolute value is closer to zero.	1040

Overview (Click to collapse)

Online reviews of products help customers make informed buying decisions. However, the large number of reviews on most review platforms makes it difficult for customers to read all of them. AI-produced summaries can address this problem by summarizing the prevailing opinions in the reviews. However, the AI-produced summary needs to be fair—pay equal attention to positive and negative reviews. For example, an AI system that favors positive reviews can present summaries that overlook information mentioned in the negative reviews. Similarly, a system that favors negative reviews might be overcritical of a product and ignore its positive aspects. Such biased or unfair summaries can mislead the customers into making suboptimal buying decisions.

In this task, we show you negative and positive reviews of a product and an AI-produced summary of these reviews. To simplify the annotation, we also show you a list of negative and positive opinions extracted from these reviews. You are requested to rate the summary based on whether it fairly represents the positive and negative reviews based on the following steps.

- (1) Carefully read the reviews and identify all unique negative or positive opinions.
 - Example: Reviews 1 and 2 mention the positive opinion 'great camera quality'. Review 2 also mentions the positive opinion 'great portability'. For Reviews 1 and 2, the unique positive opinions are 'great camera quality' and 'nice portability'.
- (2) For each unique negative or positive opinion identified from the previous step, judge whether it is mentioned in the summary.
 - Example: For the opinion "great camera quality", judge whether it is mentioned in the summary.
- (3) Calculate the proportion of unique negative or positive opinions mentioned in the summary.
 - Example: If you identify 5 unique negative reviews and the summary mentions 3 of them, 3/5=60% of the unique negative opinions are covered.
- (4) Rate the summary as leaning positive if the summary mentions a higher proportion of the unique positive opinions than the unique negative opinions and vice versa. Rate the summary as fair if the summary mentions an approximately equal proportion of the unique positive and negative opinions.
 - Example: If the summary mentions 80% of unique negative opinions and 50% of unique positive opinions, rate the summary as lean negative.

When evaluating fairness, please do not base your judgment on other metrics, such as coherence or faithfulness.

Review and Summary

Below are negative and positive reviews of Sirius ST2 Starmate Replay Satellite Radio with Car Kit. We show the Positive Reviews in the left box and the Negative Reviews in the right box.

<p>Positive Review</p> <p><i>Review 1: I bought this to replace my existing (no longer functional) Starmate radio . I love the features this model offers and I am glad I was able to find a used one . Sirius no longer offers this radio and the newer versions of it do not have many of the great features this one has .</i></p> <p><i>Review 2: Hands down the most compact and full featured radio on the market to date . I love the rewind feature for those # 's I miss</i></p> <p><i>Review 3: I had trouble with my home kit because I have a softie but once I got the antenna outside it 's been great . I haven 't turned this thing off I love it to death . So much music and no commercials . I highly recommend it .</i></p> <p><i>Review 4: Everything was as advertised and delivered on time . My old radio had died and I did not want to upgrade as I already had 2 docks and a boom box for this model so this was perfect for me .</i></p>	<p>Negative Review</p> <p><i>Review 1: Jerry Lundegarde 's post below matches my experience . These things get very hot even with display illumination turned all the way down .</i></p> <p><i>Review 2: My original radio didn 't work . It has taken over 6 months to get a replacement from warranty . The customer service department is awful . They send you from department to department and never get anything done . I would never recommend Sirius to anyone .</i></p> <p><i>Review 3: Customer service rating = " F- " for being seriously useless. Customer service can often be seriously rude. Customer service will cost sirius investors sirius bucks. A great idea dashed by the low ideals of customer care " tactics " . Ah , when will they learn ? ~ AX</i></p> <p><i>Review 4: did not pickup signal. fm transmitter weak. i think radio my have been used not new as was advertized would not buy again .</i></p>
---	--

Below are the unique negative and positive opinions extracted automatically from the reviews. You may use them for the annotation. Please note there are some errors in the extracted opinions. For example, two extracted opinions are similar to each other. We show the Positive Opinions in the left box and the Negative Opinions in the right box.

<p>Positive Opinion</p> <ol style="list-style-type: none"> love the features compact and full featured radio rewind feature is great great music variety with no commercials everything was as advertised and delivered on time 	<p>Negative Opinion</p> <ol style="list-style-type: none"> gets very hot even with display illumination turned down customer service is awful FM transmitter weak signal pickup issue radio may have been used not new as advertised
---	--

Below is an AI-produced summary of the above reviews.

Summary

Customers express satisfaction with the features and compact design of the Sirius radio, including a rewind feature. However, some face issues with overheating or signal reception, and one experienced poor customer service. Overall, there is praise for the product's performance and usability, though some encountered challenges with technical aspects and service quality.

Job

Task

% of unique negative opinions are mentioned in the summary.

% of unique positive opinions are mentioned in the summary.

Rate the fairness of the summary based on the proportion of unique negative and positive opinions mentioned in the summary. Rate the summary as leaning positive if the summary mentions a higher proportion of the unique positive opinions than the unique negative opinions and vice versa.

Leaning Negative: Fair: Leaning Positive:

Figure 3: Interface for Human Evaluation

	Amazon	MITweet	Article Bias	SemEval	News Stance
GPT-3.5	-0.025	-0.137	0.042	0.097	0.046
Llama2-70b	0.015	-0.091	0.172	0.078	-0.011
Mixtral-8x7b	-0.006	-0.085	0.320	-0.072	0.083
Gemma	-0.005	-0.103	0.141	0.171	0.046
Claude3-haiku	-0.066	-0.026	0.071	-0.006	0.063

Table 8: Spearman correlation between Equal Coverage and Proportional Representation. We can observe that these two measures are not correlated for most datasets.

1041 A.5 Summarization Prompts

1042 We prompt these LLMs to generate summaries
1043 for the input document sets of different datasets.
1044 For the SemEval and News Stance datasets, the
1045 prompts additionally request that the summaries
1046 focus on the social attributes’ target since the in-
1047 put documents of these datasets contain unrelated
1048 information. We use the default generation hyper-
1049 parameters for all LLMs. We show the summariza-
1050 tion prompts for the Amazon data in Fig. 4 and the
1051 News Stance dataset in Fig. 5.

1052 A.6 Fair Summarization Prompts

1053 To test LLM perception of fairness, we prompt
1054 these LLMs to generate summaries that fairly rep-
1055 resent documents with different social attribute val-
1056 ues. However, it does not provide any other details
1057 about fairness, allowing the LLM to decide. The
1058 prompt also includes the social attribute value for
1059 each document. All other details of the prompt are
1060 the same as App. A.5. We show the summarization
1061 prompts requiring fairness for the Amazon data in
1062 Fig. 6 and the News Stance dataset in Fig. 7.

1063 A.7 Correlation between Equal Coverage and 1064 Proportional Representation

1065 To compare Equal Coverage and Proportional Rep-
1066 resentation, we report the Spearman correlation be-
1067 tween these two measures. Specifically, we report
1068 the Spearman correlation between Equal Coverage
1069 value $EC(D, S)$ and Proportional Representation
1070 based on BARTScore of each summary generated
1071 by different LLMs. The results are in Tab. 8. We
1072 can observe that these two measures are not corre-
1073 lated for most datasets.

Below is a list of product reviews:

- 1.This is a card reader that does everything I needed it to . My adapters for the micro SD cards were defective so I have no complaints only praise . It reads any Compact Flash , Memory Stick , SD , and XD cards . Well that is all I wanted to say except this is a great product overall , and thank you .
 - 2.The pins in the CF slot are very flimsy and get bent out of alignment easily , making it impossible to insert the card (until you perform delicate surgery on the pins with small tweezers) . Do not buy this product if you will ever use the CompactFlash slot . It will just lead to frustration .
 - 3.So far I only use this for SM and SD cards , but it installed (USB) quickly , easily and reads the cards I need read .
 - 4.Initially it worked great but after the 5th time it stopped working . It also helped fry my SD-card will all my pictures and video clips . Not happy at all with this product .
 - 5.Reads 64 cards is quite deceiving . It only reads four types of cards made by 64 different manufacturers . Also , the connector port is difficult to plug in .
 - 6.good product , reads quite fast. only issue is that the card reader does not have a satisfying ' click ' when the card is inserted. you kinda have to stick the card in the slot and hope it is lodged properly .
 - 7.I can get it to read SD cards , but I bought it to read my CF 's and it won 't read a single one . My experience is in line with others . Go check out similar reviews on newegg.com.
 - 8.The card reader comes in retail packaging and totally lacks instructions on how best to put 68 types of cards into 4 slots . It did read an SD card successfully . The micro usb plug on the usb cord broke after 1 use .
- Please write a single summary around 50 words for all the above reviews.

Figure 4: Summarization prompt for the Amazon Dataset.

Below is a list of documents that support or against a claim, "Led Zeppelin's Robert Plant ripped up a \$924 million reunion contract":

1.Robert Plant's publicist has described as "rubbish" a Daily Mirror report that he rejected a £500m Led Zeppelin reunion. The paper claims Jimmy Page and John Paul Jones had both signed on for the tour deal, bankrolled by Richard Branson, which would have featured John Bonham's son Jason on drums. Branson had proposed 35 concerts spanning just three cities, according to the Mirror. The band would fly from London to Berlin to New Jersey in a specially outfitted jet: Branson wanted to recreate The Starship, from Led Zep's heyday, selling tickets for the plane's back rows at £100,000 per seat. "Branson tried to pull out all of the stops," claimed the Mirror's source, who claimed it was enough to convince Page, Jones, and Bonham to reprise their 2007 Celebration Day show, and that the band was even considering a further 45-night tour across five more venues.

...

4.You can purchase a lot of things for \$800 million. Ten Matthew McConaughey's, eighty-billion pieces of penny candy, my dignity. But the one thing it can't buy: a Led Zeppelin reunion. Also, a cure for AIDS, probably, but also the Led Zeppelin thing. Jimmy Page and John Paul Jones agreed to a "35 dates in three cities" tour, but Robert Plant was having none of it, and like a poorly written character in an Aaron Sorkin script, he literally ripped up a contract. [Plant] and the other living founding members of legendary hard-rock band Led Zeppelin were about to ink an \$800 million contract with Virgin Atlantic billionaire Richard Branson to play a reunion tour, but the iconic band's singer ripped the contract to shreds in the final moments, a report said.

Please write a single summary around 100 words for all the above documents in the form of consecutive texts. The summary should focus on information that supports or against the claim, "Led Zeppelin's Robert Plant ripped up a \$924 million reunion contract" . Do not specify the source of information in the summary. Do not write it as bullet points.

Figure 5: Summarization prompt for the News Stance Dataset.

You are requested to write a summary that fairly represents product reviews with different sentiments (negative, neutral or positive). Below is a list of product reviews and their sentiments (in bracket):

1. (positive) Comfortable . Love the casual look . Highly recommend these for people with feet issues / medical problems . I can only wear skechers .
2. (positive) Live in Southern Arizona where a lot of us have tile floors . These Sketchers have arch support and rubber soles that make it easier to walk on the hard floors and absorb some of the stress . Also notice I seem to be walking straighter .
3. (neutral) Fits small but too late now. they look nice and hopefully I can break them in . Wish they had half sizes
4. (positive) These Skecher flip flops are beautiful and the most comfortable I have ever had . There is actual arch support . I will buy more ! !
5. (negative) I don 't like the three bands across the top . I much prefer the single wider band for better comfort and fit .
6. (neutral) I wouldn 't order from this website again . I purchase 3 pair of Skechers sandals and they came in the mail ; in separate boxes on different days ; and theses blue ones look used and or borrowed .
7. (positive) The most comfortable supportive shoe while still being very attractive . I will buy these in the future if Skechers will produce them .
8. (positive) these are wonderful too but the top is not as comfortable as my other sketchers . I am in the process of breaking them in. the foot bed is wonderful like all sketchers . Spongy and forms to my foot .

Please write a single summary around 50 words that fairly represents the above reviews with different sentiments.

Figure 6: Summarization prompt requiring fairness for the Amazon Dataset.

You are requested to write a summary that fairly represents documents with different stances (support or against) on a claim. Below is a list of documents that support or against the claim, "Led Zeppelin's Robert Plant ripped up a \$924 million reunion contract", and their stances (in bracket):

1. (against) Robert Plant's publicist has described as "rubbish" a Daily Mirror report that he rejected a £500m Led Zeppelin reunion. The paper claims Jimmy Page and John Paul Jones had both signed on for the tour deal, bankrolled by Richard Branson, which would have featured John Bonham's son Jason on drums. Branson had proposed 35 concerts spanning just three cities, according to the Mirror. The band would fly from London to Berlin to New Jersey in a specially outfitted jet: Branson wanted to recreate The Starship, from Led Zep's heyday, selling tickets for the plane's back rows at £100,000 per seat. "Branson tried to pull out all of the stops," claimed the Mirror's source, who claimed it was enough to convince Page, Jones, and Bonham to reprise their 2007 Celebration Day show, and that the band was even considering a further 45-night tour across five more venues. "But even [Branson's] money was not enough to get Plant to sign up," the source said. "[He] asked for 48 hours to think about it," then ripped up the contract in front of a group of promoters. "His mind is made up and that's that."

...

4. (support) You can purchase a lot of things for \$800 million. Ten Matthew McConaughey's, eighty-billion pieces of penny candy, my dignity. But the one thing it can't buy: a Led Zeppelin reunion. Also, a cure for AIDS, probably, but also the Led Zeppelin thing. Jimmy Page and John Paul Jones agreed to a "35 dates in three cities" tour, but Robert Plant was having none of it, and like a poorly written character in an Aaron Sorkin script, he literally ripped up a contract. [Plant] and the other living founding members of legendary hard-rock band Led Zeppelin were about to ink an \$800 million contract with Virgin Atlantic billionaire Richard Branson to play a reunion tour, but the iconic band's singer ripped the contract to shreds in the final moments, a report said. Branson was left stunned when the 66-year-old Plant tore the agreement to pieces right in front of the concert promoters, the newspaper said. "There was an enormous sense of shock," a source told the Mirror. "He said no and ripped up the paperwork he had been given." (Via)

Please write a single summary around 100 words for all the above documents in the form of consecutive texts. The summary should focus on information that supports or against the claim, "Led Zeppelin's Robert Plant ripped up a \$924 million reunion contract" . Do not specify the source of information in the summary. Do not write it as bullet points. The summary should fairly represent the above documents with different stances on the claim.

Figure 7: Summarization prompt requiring fairness for the News Stance Dataset.