

NeuS: Neutral Multi-News Summarization for Framing Bias Mitigation

Anonymous ACL submission

Abstract

Media framing bias can lead to increased political polarization, and thus, the need for automatic mitigation methods is growing. We propose a new task, a *neutral* summary generation from multiple news articles of the varying political spectrum, to facilitate balanced and unbiased news reading. In this paper, we first collect a new dataset, obtain some insights about framing bias through a case study, and propose a new effective metric and models for the task. Lastly, we conduct experimental analyses to provide insights about remaining challenges and future directions. One of the most interesting observations is that generation models can hallucinate not only factually inaccurate or unverifiable content, but also politically biased content.

1 Introduction

Media framing bias occurs when journalists make skewed decisions regarding which events or information to cover (informational bias), and how to cover them (lexical bias) (Entman, 2002; Groeling, 2013). Even if the reporting of the news is based on the same set of underlying issues or facts, the framing of that issue can convey a radically different impression of what had actually happened (Gentzkow and Shapiro, 2006). Since the news media plays a crucial role in shaping public opinion toward various important issues (De Vreese, 2004; McCombs and Reynolds, 2009; Perse and Lambe, 2016), bias in said media could reinforce the problem of political polarization.

Allsides.com (Sides, 2018) mitigates this problem by displaying articles from various media in a single interface along with an expert-written roundup of news headlines. This roundup is a neutral summary for readers to grasp the bias-free understanding of an issue before reading individual articles. Although Allsides fights framing bias, the scalability still remains a bottleneck due to a

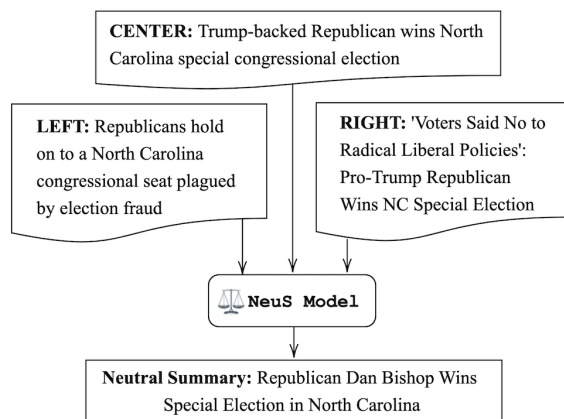


Figure 1: Illustration of the proposed task. We want to generate neutral summarization of news headlines from varying spectrum of political orientation.

lot of time-consuming human labor of composing the roundup. Multi-document summarization (MDS) models (Lebanoff et al., 2018; Liu and Lapata, 2019) could be one possible choice for automating the roundup generation as both multi-document summaries and roundups share the nature of extracting salient information out of multiple input articles. Yet, the ability of MDS models to provide *neutral* understanding of the issue has yet to be explored – a crucial aspect of the roundup.

In this work, we propose to fill in this research gap by proposing a task of Neutral multi-news Summarization (NEUS), which aims to generate a framing-bias-free summary out from news headlines with varying degrees and orientation of political bias (Fig. 1). To begin with, we construct a new dataset by crawling Allsides.com and investigate how framing bias manifests in the news to provide a deeper and comprehensive understanding of the problem. First, an important insight is the close association between framing bias and the polarity of the text. Grounded on this basis, we propose a polarity-based framing-bias metric that is simple yet effective in terms of alignment with human per-

065 ceptions. Then, based on the second insight that
066 titles serve as a good indicator of framing bias, we
067 propose NEUS models that leverage the news titles
068 as an additional signal to increase awareness of
069 framing bias.

070 Our experimental results provide rich insights
071 for understanding the problem of mitigating fram-
072 ing bias. Primarily, we explore whether existing
073 summarization models can already solve the prob-
074 lem and empirically demonstrate their shortcom-
075 ings in addressing the stylistic aspect of framing
076 bias. After that, we investigate and discover an
077 interesting relationship between framing bias and
078 hallucination, an important safety-related problem
079 in NLP. We empirically show that the hallucinatory
080 generation has the risk of being not only factually
081 inaccurate and/or unverifiable, but also politically
082 biased and controversial. To the best of our knowl-
083 edge, this aspect of hallucination has not been dis-
084 cussed. We want to encourage more attention to
085 hallucinatory framing bias to prevent a generation
086 from fueling political bias and polarization.

087 We conclude with a discussion about the remain-
088 ing challenges to provide insights for future work.
089 We hope our work with the proposed NEUS task
090 serves as a good starting point to promote the auto-
091 matic mitigation of media framing bias.

092 2 Related Works

093 Media Framing Bias Detection and Prediction

094 Media bias has been studied extensively in vari-
095 ous fields such as social science, economics, and
096 political science, and various methods have been
097 proposed to analyze the political preference and
098 framing bias of news outlets (Groseclose and Mi-
099 llyo, 2005; Miller and Riechert, 2001; Park et al.,
100 2011; Gentzkow and Shapiro, 2010; Haselmayer
101 and Jenny, 2017). Framing bias is selective re-
102 porting of an event to sway readers’ opinions with
103 different factors such as commission of extra in-
104 formation and word choices (Entman, 1993, 2007;
105 Gentzkow and Shapiro, 2006). In NLP, computa-
106 tional approaches for detecting media bias often
107 consider linguistic cues that induce bias in politi-
108 cal text (Recasens et al., 2013; Yano et al., 2010;
109 Lee et al., 2019; Hamborg et al., 2019b). For in-
110 stance, Gentzkow and Shapiro count the frequency
111 of slanted words within articles. These methods
112 mainly focus on the stylistic (“how to cover”) as-
113 pect of framing bias. There is relatively less effort
114 on the informational (“what to cover”) aspect of

115 framing bias (Park et al., 2011; Fan et al., 2019)
116 and they are constrained to detection tasks. In this
117 work, we attempt to tackle both by generating a
118 bias-free summary out of biased headlines.

119 **Media Bias Mitigation** News aggregation by
120 displaying articles from different news outlets
121 on a particular topic (e.g., Google News¹, Ya-
122 hoo News²), is the most common approach in
123 NLP to mitigate media bias, but it still has limi-
124 tations (Hamborg et al., 2019a). Other approaches
125 have been proposed to provide additional infor-
126 mation (Laban and Hearst, 2017), such as auto-
127 matic classification of multiple view points (Park
128 et al., 2009), multinational perspectives (Hamborg
129 et al., 2017), and detailed media profiles (Zhang
130 et al., 2019b). However, these methods focus on
131 providing a broader perspective from an enlarged
132 selection of articles to news readers, which still
133 puts burden on the readers. We propose instead to
134 automatically neutralize and summarize partisan
135 headlines to produce a neutral headline summary.

136 **Multi-document Summarization** As a chal-
137 lenging subtask of automatic text summarization,
138 multi-document summarization (MDS) aims to con-
139 dense a set of documents to a short and informative
140 summary (Lebanoff et al., 2018). Recently, re-
141 searchers apply deep neural models for MDS task
142 thanks to the introduction of large-scale datasets
143 (Liu et al., 2018; Fabbri et al., 2019). With the ad-
144 vent of large pre-trained language models (Lewis
145 et al., 2019; Raffel et al., 2019), researchers also
146 apply them to improve the MDS models perfor-
147 mance (Jin et al., 2020; Pasunuru et al., 2021).
148 In addition, many works have studied particular
149 subtopics of the MDS task, such as agreement-
150 oriented MDS (Pang et al., 2021), topic-guided
151 MDS (Cui and Hu, 2021) and MDS of medical stud-
152 ies (DeYoung et al., 2021). However, few works
153 have explored the field of generating framing bias-
154 free summaries from multiple news articles. In this
155 paper, we propose the NEUS task and create a new
156 benchmark.

157 3 Task and Dataset

158 3.1 Task Formulation

159 The main objective of NEUS is to generate a neu-
160 tral headline summary H_{neu} given multiple news

¹<https://news.google.com/>

²<https://news.yahoo.com/>

Issue A: Trump Put Hold On Military Aid To Ukraine Days Before Call To Ukrainian President	
Left:	Trump ordered hold on military aid days before calling Ukrainian president, officials say
Right:	Trump administration claims Ukraine aid was stalled over corruption concerns, decries media ‘frenzy’
Center:	Trump Put Hold on Military Aid Ahead of Phone Call With Ukraine’s President
Issue B: Michael Reinhoehl appeared to target right-wing demonstrator before fatal shooting in Portland, police say	
Left:	Suspect in killing of right-wing protester fatally shot during arrest
Right:	Portland’s Antifa-supporting gunman appeared to target victim, police say
Center:	Suspect in Patriot Prayer Shooting Killed by Police
Issue C: Trump Says the ‘Fake News Media’ Are ‘the true Enemy of the People’	
Left:	President Trump renews attacks on press as ‘true enemy of the people’ even as CNN receives another suspected bomb
Right:	‘Great Anger’ in America caused by ‘fake news’ — Trump rips media for biased reports’
Center:	Trump blames ‘fake news’ for country’s anger : ‘the true enemy of the people’

Table 1: Illustration of **difference in framing** from Left/Right/Center media with examples from ALL-SIDES dataset. We used titles for the analysis of bias, since they are simpler to compare and are representative of the framing bias that exist in the headline.

headlines $H_{0...N}$ with varying degrees and orientations of political bias. The neutral summary H_{neu} should (i) retain salient information and (ii) minimize as much framing bias as possible from the input headlines.

3.2 ALLSIDES Dataset

Allsides.com provides access to triplets of news, which report about the same event from left, right, and center American publishers, with an expert-written neutral summary of the headlines and its neutral title. The dataset language is English and mainly focuses on U.S. political topics that often result in media bias. The top-3 most frequent topics³ are ‘Elections’, ‘White House’, ‘Politics’.

We crawl the headline triplets to serve as the source inputs $\{H_L, H_R, H_C\}$, and the neutral headline summary to be the target output H_{neu} for our task. Note that “center” does not necessarily mean completely bias-free (all, 2021) as illustrated in Table 1. Although “center” medias are relatively less tied to particular political ideology, they may still contain framing bias because editorial judgement naturally leads to human-induced biases. In addition, we also crawl the title triplets $\{T_L, T_R, T_C\}$ and the neutral issue title T_{neu} that are later used in our modelling.

To make this dataset richer, we also crawled other meta-information such as date, topic-tags and media-name. In total, we crawled 3,564 triplets (10,692 headlines). We use 2/3 of the triplets, which is 2,276 triplets, to be our train and validation set (80 : 20 ratio), and the remaining 1,188 triple as our test set. We will publicly release this dataset for future research use.

³The full list is provided in appendix.

4 Analysis of Framing Bias

Based on literature of media framing bias from NLP community and political studies, we know the definition and types of framing bias (Goffman, 1974; Entman, 1993; Gentzkow et al., 2015; Fan et al., 2019) — *Informational framing bias* is the biased selection of information (tangential or speculative information) to sway the minds of readers; *Lexical framing bias* is the sensational writing style or linguistic attributes that may mislead readers. However, the definition is not enough to understand exactly how framing bias manifests in real examples that, in our case, is ALLSIDES dataset. We conduct case-study to obtain essential insights that can guide our design choices for defining the metric and methodology.

4.1 Case-Study Observations

First, we identify and share the examples of framing bias in accordance with the literature (Table 1).

Informational Bias This bias exists dominantly in form of “extra information” on top of the salient key-information about the issue that changes the overall impression of the issue. For example, in Table 1, when reporting about “Military Aid Hold To Ukraine” (Issue A), the right media reports the speculative claim that there was “corruption concerns” and tangential information “decries media ‘frenzy’” that amplifies the negative impression of the issue. Sometimes, media with different political leanings report additional information to convey completely different **focus** of the issue. For Issue C, left-media implies that Trump’s statement about fake news has led to “CNN receiving another suspected bomb”, whereas right-media implies that media is at fault by producing “biased reports”.

Lexical Bias This exists mainly as biased word choices that change the nuance of the information that is being delivered. For example, in Issue B, we can clearly observe that two media change the framing of the issue by using different terms “suspect” and “gunman” to refer to the shooter, and “protester” and “victim” to refer to the person shot. Also, in Issue A, when one media uses “(ordered) hold”, another media uses “stalled” that has a more negative connotation.

4.2 Main Insights from Case-Study

Next, we share important insights from case study observation that guide our metric and model design.

Relative Polarity Polarity is one of the commonly used attributes in identifying and analyzing framing bias (Fan et al., 2019; Recasens et al., 2013). Although informational and lexical bias are conceptually different, both are closely associated with polarity changes of concepts, i.e., positively or negatively, to induce strongly divergent emotional responses from the readers (Hamborg et al., 2019b). Thus, polarity can serve as a good indicator of framing bias. However, we observe that the polarity of text must be utilized with care in the context of framing bias. *It is the relative polarity that is meaningful to indicate the framing bias, not the absolute polarity.* To elaborate, if the news issue itself is about tragic events such as “Terror Attack in Pakistan” or “Drone Strike That Killed 10 people”, then the polarity of the neutral reporting will also be negative.

Indicator of Framing We discovered that news title is very representative of the framing bias that exists in the associated headline and article – this makes sense because title can be viewed as the succinct overview of the content that follows. For instance, in Table 3 source input example, right media’s title and headline are mildly mocking the “desperate” democrats’ failed attempts to take down President Trump. In contrast, left media’s title and headline show a completely different frame – implies that many investigations are happening and there’s “possible obstruction of justice, public corruption, and other abuses of power.”

5 Metric

We use three kinds of metrics to evaluate the neutral summaries to tackle the problem from different dimensions. For framing bias, we propose

polarity-based metric with a detailed articulation of our design choices (§5.1). For evaluating whether the summaries retain salient information, we adopt commonly used information recall-related metrics (§5.2). In addition, we use a hallucination metric to evaluate if the generations contain unfaithful hallucinatory information because the existence of such hallucinatory generations can make the summary fake news. (§5.3).

5.1 Framing Bias Metrics

5.1.1 Design Consideration

Our framing bias metric is developed upon the insight we obtained from our case-study in §4.

First of all, we propose to build our metric based on the fact that framing bias is closely associated with polarity. There are options of model-based and lexicon-based polarity detection approaches and we leverage the lexicon-based approach for the following reasons. 1) There is increasing demand for interpretability in the field of NLP (Belingov et al., 2020; Sarker et al., 2019), and the lexicon-based approach is more interpretable (provides token-level human interpretable annotation) compared to black-box neural models. 2) In the context of framing bias, distinguishing the subtle nuance of words between synonyms are crucial (e.g., dead vs murdered). Lexicon-resource provides such token-level fine-grain scores and annotations, making it useful for our usage.

Metric calibration is the second design consideration motivated by our insight about the relativity of framing bias. The absolute polarity of token itself does not necessarily indicate framing bias (i.e., word “riot” has negative sentiment but does not always indicate bias), so it is important to measure the relative degree of polarity. Therefore, calibration of the metric in reference to the neutral target is important. Any tokens existing in neutral target will be ignored in bias measurement for the generated neutral summary. For instance, if a word “riot” exists in neutral target, it will not be counted in bias measurement through calibration.

5.1.2 Framing Bias Metric Details

For our metric, we leverage Valence-Arousal-Dominance (VAD) (Mohammad, 2018) dataset which has a large list of lexicons annotated for valence, arousal and dominance scores. Valence, arousal and dominance represent the direction of polarity (positive, negative), the strength of the

327 polarity (active, passive) and the level of control
328 (powerful, weak) respectively.

329 Given the neutral summary generated from the
330 model \hat{H}_{neu} , our metrics are calculated using the
331 VAD lexicons in the following way:

- 332 1. Filter out all the tokens that appears in neutral
333 target H_{neu} to obtain set of tokens *unique* to
334 \hat{H}_{neu} . This ensures that we are measuring the
335 *relative* polarity of \hat{H}_{neu} in reference to the
336 neutral target H_{neu} – calibration effect.
- 337 2. We identify tokens with either positive or neg-
338 ative valence (v), which as result will further
339 filter out neutral words such as stopwords and
340 non-emotion provoking words.
- 341 3. Sum up the associated arousal scores for these
342 identified positive and negative tokens from
343 Step 2 – positive arousal score ($Arousal_+$)
344 and negative arousal score ($Arousal_-$).
345 We intentionally separate the positive and
346 negative scores for finer-grain interpreta-
347 tion. We also have the combined arousal
348 score ($Arousal_{sum}=Arousal_++Arousal_-$)
349 for coarse view.
- 350 4. Repeat for all $\{H_{neu}, \hat{H}_{neu}\}$ pairs in the test-
351 set, and calculate the average scores to use as
352 the final metric. We report these scores in our
353 experimental results section (§7).

354 In essence, our metric approximates the exist-
355 tence of framing bias by measuring the aroused
356 degree of the generated summary. The aroused
357 degree is a relative value between the generated
358 summary to the neutral target reference. We pro-
359 vide our code for reproducibility.

360 5.1.3 Human Evaluation

361 To ensure the quality of our metric, we evaluate the
362 correlation between our framing bias metrics with
363 the human judgement. We did A/B testing⁴ where
364 the annotators are given two generated headlines
365 about an issue, one with higher $Arousal_{sum}$ score
366 and another with lower score and are asked to se-
367 lect more biased headline summary. When asking
368 which is more “biased”, we adopt the question by
369 Spinde et al.. We also provide examples and defini-
370 tion of framing bias for better understanding of the
371 task. We obtained 3 annotations each for 50 sam-
372 ples and selected the ones with majority of voting.

373 One of the challenges of this evaluation is in
374 personal political bias of annotators. Although it is

⁴Please refer to appendix for more detail of the A/B testing

375 hard to eliminate such bias completely, we attempt
376 to avoid it by collecting annotations from those
377 who are less related to the issues of testset. Clearly
378 speaking, given that our testset covers mainly about
379 US politics, we restricted the nationality of anno-
380 tators to be non-US internationals who claim to be
381 bias-free from US political party.

382 After obtaining the human annotations from A/B
383 testing, we obtain another version of annotation
384 based on the metric score – i.e., the one with higher
385 $Arousal_{sum}$ is chosen to be more biased headline
386 generation. The Spearman correlation coefficient
387 between human-based and metric-based annota-
388 tions is 0.63615 with p-value < 0.001 and agree-
389 ment percentage is 80%. These indicate that the
390 association between the two annotations is statisti-
391 cally significant, suggesting that our metric is
392 providing good approximation of the framing bias
393 existence.

394 5.2 Salient Info

395 It is important for the generation to retain essen-
396 tial/important information while reducing the fram-
397 ing bias. Thus, we also report ROUGE (Lin, 2004)
398 and BLEU (Papineni et al., 2002) between gener-
399 ated neutral summary, \hat{H}_{neu} , and human written
400 summary, H_{neu} . Note that ROUGE measures the
401 recall (i.e., how much the n-grams in the human
402 reference text appeared in the machine generated
403 text) and BLEU measures the precision (i.e., how
404 much the n-grams in the machine generated text
405 appeared in the human reference text). The higher
406 the BLEU and ROUGE1-R score, the better the es-
407 sential information converges. In our results, we
408 only report Rouge-1, but Rouge-2 and Rouge-L
409 can be found in the appendix.

410 5.3 Hallucination Metric

411 Recent studies have shown that neural sequence
412 models can suffer from the hallucination of ad-
413 ditional content not supported by the input (Re-
414 iter, 2018; Wiseman et al., 2017; Nie et al., 2019;
415 Pagnoni et al., 2021; Maynez et al., 2020), conse-
416 quently adding factual inaccuracy to the generation
417 of NLG generations. Although not directly related
418 to the goal of NEUS, we evaluate the hallucination
419 level of the generations. We choose hallucination
420 metric called FeQA (Durmus et al., 2020) for our
421 work, because it is one of the publicly available
422 metric known to have high correlation with human
423 faithfulness scores. This is a QA-based metric that
424 is built on the assumption that same answers will

425 be given from hallucination-free generation and the
426 source document when asked same question.

427 6 Models and Experiments⁵

428 6.1 Baseline Models

429 Since one common form of framing bias is the
430 reporting of extra information (§4), summariza-
431 tion models—that extracts commonly shared salient
432 information—may already generate neutral sum-
433 mary to some extent. To answer this, we report
434 experimental results using the following baselines.

- 435 • LEXRANK (Erkan and Radev, 2004): an
436 extractive single-document summarization
437 (SDS) model that extracts representative sen-
438 tences that hold information common in both
439 left and right articles.
- 440 • BARTCNN: an abstractive SDS model
441 that fine-tunes BART-large (Lewis
442 et al., 2019) (406M parameters) using
443 CNN/DailyMail (Hermann et al., 2015)
444 dataset.
- 445 • BARTMULTI: a multi-document summariza-
446 tion (MDS) model that fine-tunes BART-large
447 using Multi-News (Fabbri et al., 2019) dataset.
- 448 • PEGASUSMULTI: a MDS model that fine-
449 tunes Pegasus-base (Zhang et al., 2019a)
450 (568M parameter) using Multi-News dataset.

451 Since the summarization models are not trained
452 with in-domain data, we provide another baseline
453 model trained with in-domain data for full picture.

- 454 • NEUSFT: a baseline that fine-tunes BART-
455 large model using ALLSIDES.

456 6.2 Our NEUS Models (NEUS-TITLE)

457 We designed our models based on one of the in-
458 sights from case-study (§4) — news title serves
459 as an indicator of the framing bias in the corre-
460 sponding headline. We hypothesize that it would
461 be helpful to divide-and-conquer by neutralizing
462 from title-level first, then leveraging the “neutral-
463 ized title” to guide the final neutral summary of
464 the longer headlines. Multi-task learning (MTL) is
465 a natural modelling choice because there are two
466 sub-tasks involved – title-level and headline-level
467 neutral summarization. However, we also have to
468 ensure a sequential relationship between the two
469 tasks in our MTL training, because headline-level

⁵Experimental details are in appendix for reproducibility.

470 neutral summarization leverages the generated neu-
471 tral title as the additional resource.

472 We propose a simple yet elegant trick to address
473 by adapting the idea of prompting, a method of
474 reformatting NLP tasks in the format of a natural
475 language response to natural language input (Sanh
476 et al., 2021). We train the BART’s autoregressive
477 decoder to generate the target text Y formatted as
478 follows:

$$479 \text{TITLE} \Rightarrow T_{neu}. \text{HEADLINE} \Rightarrow H_{neu}.$$

480 where T_{neu} and H_{neu} denote neutral title and neu-
481 tral headline summary.

482 The input X to our BART encoder is formatted
483 similarly to the target text Y :

$$484 \text{TITLE} \Rightarrow T_L. \text{HEADLINE} \Rightarrow H_L.[SEP]$$

$$485 \text{TITLE} \Rightarrow T_C. \text{HEADLINE} \Rightarrow H_C.[SEP]$$

$$486 \text{TITLE} \Rightarrow T_R. \text{HEADLINE} \Rightarrow H_R.$$

487 where $T_{L/C/R}$ and $H_{L/C/R}$ denote title and head-
488 line from left, center and right media, and [SEP]
489 denotes the special token that separates between
490 different inputs.

491 This trick allows us to easily optimize for both
492 title and headline neutral summarization tasks by
493 optimizing for the negative log likelihood of the
494 single target Y . The auto-regressive nature of the de-
495 coder ensures the sequential relationship between
496 title and headline as well.

497 7 Results and Analysis

498 In this section, we point out noteworthy observa-
499 tions from the quantitative results in Table 2 with
500 some insights obtained through qualitative analysis.
501 Table 3 shows some generation examples that are
502 most representative of the insight we share⁶.

503 7.1 Main Results

504 Firstly, summarization models can reduce the
505 framing bias to a certain extent (drop in
506 $Arousal_{sum}$ score from 10.40 to 4.76 and 3.32
507 for LEXRANK and BARTCNN). This is because in-
508 formational framing bias has been addressed when
509 summarization models extract the most salient sen-
510 tences that contain common information from the
511 inputs. However, summarization models, espe-
512 cially LEXRANK cannot handle the lexical framing
513 bias as shown in Table 3. Moreover, if we further
514 observe LEXRANK, it is one of the best performing

⁶More examples are provided in appendix.

Models	Avg. Framing Bias Metric			Salient Info		Hallucination
	Arousal ₊ ↓	Arousal ₋ ↓	Arousal _{sum} ↓	BLEU↑	ROUGE1-R↑	FeQA↑
All Source input	6.76	3.64	10.40	8.27	56.57%	-
LEXRANK	3.02	1.74	4.76	12.21	39.08%	53.44%
BARTCNN	2.09	1.23	3.32	10.49	35.63%	58.03%
PEGASUSMULTI	5.12	2.39	7.51	6.12	44.42%	22.24%
BARTMULTI	5.94	2.66	8.61	4.24	35.76%	21.06%
NEUSFT	1.86	1.00	2.85	11.67	35.11%	58.50%
NEUS-TITLE	1.69	0.83	2.53	12.05	36.07%	45.95%

Table 2: Experimental results for ALLSIDES testset. We provide the level of framing bias inherent in “source input” from ALLSIDES testset to serve as reference point for framing bias metrics. For framing bias metrics, the *lower* number is the better (i.e., ↓). For other scores, the *higher* number is the better (i.e., ↑).

SOURCE: <Left> **Title:** Here Are The 81 People And Entities Close To Trump Democrats Are Investigating. **Headline:** Democrats on the House Judiciary Committee on Monday sent document requests to 81 agencies, entities and individuals close to President Donald Trump as part of a broad investigation into possible obstruction of justice, public corruption and other abuses of power. The list includes Trump’s sons, Eric Trump and Donald Trump Jr., as well as his son-in-law, Jared Kushner.

<Center> **Title:** House Panel Requests Documents From Associates of Trump. **Headline:** House Democrats intensified their investigations into President Trump and his associates Monday, demanding records from more than 80 people and organizations related to his business dealings, interactions with the Justice Department and communications with Russian President Vladimir Putin.

<Right> **Title:** Dems Continue Their Assault on The Trump Administration By Launching Another Probe. **Headline:** Democrats are desperate to take down President Donald Trump. The Russia probe has proven to be ineffective and, quite frankly, a waste of time and taxpayer money. They didn’t find what they wanted so now they’re launching another probe.

TARGET: House Democrats launched a broad probe into President Trump on Monday, requesting documents from 81 agencies and individuals as they investigate his business dealings, interactions with Russia, and possible obstruction of justice.

Lexrank: Democrats are **desperate** to take **down** President Donald Trump. The Russia probe has proven to be **ineffective** and, quite frankly, a **waste** of time and taxpayer money.

NEUSFT: The Russia probe has proven to be **ineffective** and, quite frankly, a **waste** of time and taxpayer money.

NEUS-TITLE: TITLE=> House Panel Requests Documents. ARTICLE=> The House Select Committee on Intelligence has requested documents from 81 people and entities **close** to President Trump, including his sons Eric and Donald Trump Jr., as well as Jared Kushner.

Table 3: Generation examples for analysis purpose. Red highlights the tokens identified by VAD lexicons. Refer to appendix for more examples.

515 model in terms of ROUGE1-R (39.08%), standard
516 metric for summarization performance, but not in
517 framing bias metric. This suggests that having good
518 summarization performance (ROUGE1-R) does not
519 guarantee that the model also is neutral – i.e., the
520 requirement for the summary to be neutral adds
521 extra dimension to summarization task.

522 Second, one interesting pattern that requires at-
523 tention is that only the *single-document* summariza-
524 tion models (BARTCNN and LEXRANK) managed
525 to reduce framing bias well, not the *multi-document*
526 summarization models (PEGASUSMULTI and
527 BARTMULTI). This is rather surprising because
528 our task setup is more similar to MDS than SDS.
529 One potential contributor to high bias in MDS mod-
530 els could be the hallucination. MDS models appear

531 to be suffering drastically more from hallucination
532 than all other models (both MDS models PEGA-
533 SUSMULTI and BARTMULTI achieve 22.24% and
534 21.06% when most of the other models achieve
535 over 50%)⁷. This suggests that the framing bias of
536 MDS models may be related to the hallucination
537 of politically biased content. We investigate this
538 aspect separately in the next subsection.

539 Third, although summarization models helped
540 to reduce the framing bias scores, we observe the
541 bigger bias reduction when trained with in-domain
542 data as expected. NEUSFT shows further drop
543 across all framing bias metrics without sacrificing
544 the ability to keep salient information. However,

⁷Note that 22.24% and 21.06% are already high FeQA scores, however, comparatively low score in reference

SOURCE: ... President Trump on Saturday blasted what he called the “phony” BuzzFeed story and the mainstream media’s coverage of it....

MDS Hallucination: president trump on sunday slammed what he called called a “phony” story by the “dishonest” and “fake news” news outlet in a series of tweets. ... “the fake news media is working overtime to make this story look like it is true,” trump tweeted. “they are trying to make it look like the president is trying to hide something, but it is not true!”

Table 4: Illustration of hallucinatory framing bias from MDS models and the corresponding “most relevant source snippet” from the source input. Refer to the appendix for more examples with full context

we observe that NEUSFT often copies directly without any neutral re-writing – the NEUSFT example shown in Table 3 is also a direct copy of sentence from the input source.

Lastly, we can observe slightly more improvement with NEUS-TITLE across all metric except FeQA score. This model demonstrates a stronger tendency to paraphrase rather than direct copy, and comparatively has more neutral framing of the issue. As shown in Table 3, when LEXRANK and NEUSFT are focused on the “ineffectiveness of Russia probe”, the gold “target” and NEUS-TITLE focuses on the start of the investigation with the request for documents. It also generated a title that has a similar neutral frame as the target, suggesting this title generation guided the correctly framed generation.

7.2 Further Analysis and Discussion

Q: Is hallucination contributing to the high framing bias in MDS models? Through qualitative analysis, we discovered MDS generations hallucinating many politically controversial or sensational content that does not exist in the input sources. These are probably originating from the memorization of either the training data or LM-pretraining corpus. For instance, in Table 4, we can observe stylistic bias injected – i.e., “the ‘dishonest’ and ‘fake news’ news outlet”. Also, excessive elaboration of the president’s comment towards the news media, which does not appear in source nor target, can be considered informational bias – “they are trying to make it look like the president is trying to hide something, but it is not true!”. This analysis unveils the overlooked danger of hallucination, which is the risk of introducing political framing bias in summary generations. Note that this problem is not just confined to MDS models only. Other baseline models also have room for improvement in terms of FeQA hallucination score.

Q: What are the remaining challenges and future direction? The experimental result of

NEUS-TITLE suggests that there is room for improvement. We qualitatively checked some error cases and discovered that the title-generation is, unsurprisingly, not always accurate. The error propagating from the title-generation step has adversely affected the overall performance. Thus, one possible future direction will be to improve the neutral title generation, which will then improve the neutral summarization.

Another challenge is associated with the subtle lexical bias that involves nuanced word choices that manoeuvre readers to understand event from biased frames. For examples, “put on hold” and “stalled” both means the same with the latter having more negative connotation. Improving the model’s awareness towards such nuanced words, or devising ways to incorporate style-transfer-based bias mitigation approaches (Liu et al., 2021) could be another useful future direction.

We started the neutral summarization task from an assumption that framing bias originates from the source inputs. However, as shown from the results and discussed in the previous question, we found the hallucinatory content in generation is another contributor of framing bias. Thus, tackling hallucination is also an important future direction for NEUS task.

8 Conclusion

We introduce a new task of Neutral Multi-News Summarization (NEUS), to mitigate the media framing bias by providing neutral summary of headlines, along with dataset ALLSIDES and a set of metric. Throughout the work, we share insights to understand challenges and future direction in the task. We show the relationships among polarity, extra information and framing bias, which guides us into metric design. Also, the insight that title serves as an indicator of framing bias leads us to the model design. Our qualitative analysis reveals hallucinatory content generated by models may also be one of the contributors of framing bias.

Ethical Considerations

If we can automatically generate a neutralized version of media reporting, it would be one meaningful solution to framing bias. However, the idea of unbiased journalism has been challenged a number of times⁸, because different journalists and reporters have their own editorial judgments that cannot be guaranteed to be completely bias-free. Therefore, we aim to do *bias-aware/neutral headline summarization*, which provides comprehensive summary of headlines from different media, instead of trying to neutralize an article.

One of the concerns we need to take into consideration is the bias induced from the computational approach. The automatic approaches may replace a known source bias with another bias possibly caused from human-annotated data or the machine learning models. Understanding the risk of uncontrolled adoption of such automatic tools, a careful guidance should be provided. For instance, the automatically generated neutral summary should be provided with reference to original source instead of stand-alone use.

Throughout this paper we use news from English-language only, and largely American news outlets. Partisanship from this data refers to domestic American politics. We note that this work does not cover media bias in international-level or in other languages. It might be hard to directly apply this work in different cultures or languages as the bias may exist differently depending on cultures. However, we wish the paradigm of NEUS, providing multiple sides to neutralize the view of an issue, can encourage other future research in mitigating framing bias in other languages or cultures.

References

2021. [Center – what does a "center" media bias rating mean?](#)

Yonatan Belinkov, Sebastian Gehrmann, and Ellie Pavlick. 2020. Interpretability and analysis in neural nlp. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 1–5.

Peng Cui and Le Hu. 2021. Topic-guided abstractive multi-document summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1463–1472.

⁸<https://www.allsides.com/blog/does-unbiased-news-really-exist>

- Claes De Vreese. 2004. The effects of strategic news on political cynicism, issue evaluations, and policy support: A two-wave experiment. *Mass Communication & Society*, 7(2):191–214.
- Jay DeYoung, Iz Beltagy, Madeleine van Zuylen, Bailey Kuehl, and Lucy Lu Wang. 2021. Ms2: Multi-document summarization of medical studies. *arXiv preprint arXiv:2104.06486*.
- Esin Durmus, He He, and Mona Diab. 2020. FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5055–5070, Online. Association for Computational Linguistics.
- Robert M Entman. 1993. Framing: Toward clarification of a fractured paradigm. 689 690
- Robert M Entman. 2002. Framing: Towards clarification of a fractured paradigm. *McQuail's Reader in Mass Communication Theory*. London, California and New Delhi: Sage. 691 692 693 694
- Robert M Entman. 2007. Framing bias: Media in the distribution of power. *Journal of communication*, 57(1):163–173. 695 696 697
- Günes Erkan and Dragomir R Radev. 2004. Lexrank: Graph-based lexical centrality as salience in text summarization. *Journal of artificial intelligence research*, 22:457–479. 698 699 700 701
- Alexander R Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir R Radev. 2019. Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model. *arXiv preprint arXiv:1906.01749*. 702 703 704 705 706
- Lisa Fan, Marshall White, Eva Sharma, Ruisi Su, Prafulla Kumar Choubey, Ruihong Huang, and Lu Wang. 2019. In plain sight: Media bias through the lens of factual reporting. *arXiv preprint arXiv:1909.02670*. 707 708 709 710 711
- Matthew Gentzkow and Jesse M Shapiro. 2006. Media bias and reputation. *Journal of political Economy*, 114(2):280–316. 712 713 714
- Matthew Gentzkow and Jesse M Shapiro. 2010. What drives media slant? evidence from us daily newspapers. *Econometrica*, 78(1):35–71. 715 716 717
- Matthew Gentzkow, Jesse M Shapiro, and Daniel F Stone. 2015. Media bias in the marketplace: Theory. In *Handbook of media economics*, volume 1, pages 623–645. Elsevier. 718 719 720 721
- Erving Goffman. 1974. *Frame analysis: An essay on the organization of experience*. Harvard University Press. 722 723 724

725	Tim Groeling. 2013. Media bias by the numbers: Challenges and opportunities in the empirical study of partisan news. <i>Annual Review of Political Science</i> , 16:129–151.		
726			
727			
728			
729	Tim Groseclose and Jeffrey Milyo. 2005. A measure of media bias. <i>The Quarterly Journal of Economics</i> , 120(4):1191–1237.		
730			
731			
732	Felix Hamborg, Karsten Donnay, and Bela Gipp. 2019a. Automated identification of media bias in news articles: an interdisciplinary literature review. <i>International Journal on Digital Libraries</i> , 20(4):391–415.		
733			
734			
735			
736			
737	Felix Hamborg, Norman Meuschke, and Bela Gipp. 2017. Matrix-based news aggregation: exploring different news perspectives. In <i>2017 ACM/IEEE Joint Conference on Digital Libraries (JCDL)</i> , pages 1–10. IEEE.		
738			
739			
740			
741			
742	Felix Hamborg, Anastasia Zhukova, and Bela Gipp. 2019b. Illegal aliens or undocumented immigrants? towards the automated identification of bias by word choice and labeling. In <i>International Conference on Information</i> , pages 179–187. Springer.		
743			
744			
745			
746			
747	Martin Haselmayer and Marcelo Jenny. 2017. Sentiment analysis of political communication: combining a dictionary approach with crowdcoding. <i>Quality & quantity</i> , 51(6):2623–2646.		
748			
749			
750			
751	Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. <i>Advances in neural information processing systems</i> , 28:1693–1701.		
752			
753			
754			
755			
756	Hanqi Jin, Tianming Wang, and Xiaojun Wan. 2020. Multi-granularity interaction network for extractive and abstractive multi-document summarization. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6244–6254.		
757			
758			
759			
760			
761			
762	Philippe Laban and Marti A Hearst. 2017. newslens: building and visualizing long-ranging news stories. In <i>Proceedings of the Events and Stories in the News Workshop</i> , pages 1–9.		
763			
764			
765			
766	Logan Lebanoff, Kaiqiang Song, and Fei Liu. 2018. Adapting the neural encoder-decoder framework from single to multi-document summarization. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4131–4141.		
767			
768			
769			
770			
771			
772	Nayeon Lee, Zihan Liu, and Pascale Fung. 2019. Team yeon-zi at semeval-2019 task 4: Hyperpartisan news detection by de-noising weakly-labeled data. In <i>Proceedings of the 13th International Workshop on Semantic Evaluation</i> , pages 1052–1056.		
773			
774			
775			
776			
777	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019.		
778			
779			
		Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. <i>arXiv preprint arXiv:1910.13461</i> .	780 781 782
		Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	783 784 785
		Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. In <i>International Conference on Learning Representations</i> .	786 787 788 789 790
		Ruibo Liu, Chenyan Jia, Jason Wei, Guangxuan Xu, Lili Wang, and Soroush Vosoughi. 2021. Mitigating political bias in language models through reinforced calibration .	791 792 793 794
		Yang Liu and Mirella Lapata. 2019. Hierarchical transformers for multi-document summarization. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5070–5081.	795 796 797 798 799
		Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919, Online. Association for Computational Linguistics.	800 801 802 803 804 805
		Maxwell McCombs and Amy Reynolds. 2009. How the news shapes our civic agenda. In <i>Media effects</i> , pages 17–32. Routledge.	806 807 808
		M Mark Miller and Bonnie Parnell Riechert. 2001. The spiral of opportunity and frame resonance: Mapping the issue cycle in news and public discourse. <i>Framing public life: Perspectives on media and our understanding of the social world</i> , pages 107–121.	809 810 811 812 813
		Saif Mohammad. 2018. Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 english words. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 174–184.	814 815 816 817 818
		Feng Nie, Jin-Ge Yao, Jinpeng Wang, Rong Pan, and Chin-Yew Lin. 2019. A simple recipe towards reducing hallucination in neural surface realisation. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2673–2679.	819 820 821 822 823 824
		Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. 2021. Understanding factuality in abstractive summarization with frank: A benchmark for factuality metrics. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4812–4829.	825 826 827 828 829 830 831
		Richard Yuanzhe Pang, Adam D Lelkes, Vinh Q Tran, and Cong Yu. 2021. Agreesum: Agreement-oriented multi-document summarization. <i>arXiv preprint arXiv:2106.02278</i> .	832 833 834 835

836	Kishore Papineni, Salim Roukos, Todd Ward, and Weijing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	
837		
838		
839		
840		
841	Souneil Park, Seungwoo Kang, Sangyoung Chung, and Junehwa Song. 2009. Newscube: delivering multiple aspects of news to mitigate media bias. In <i>Proceedings of the SIGCHI conference on human factors in computing systems</i> , pages 443–452.	
842		
843		
844		
845		
846	Souneil Park, Kyung-Soon Lee, and Junehwa Song. 2011. Contrasting opposing views of news articles on contentious issues. In <i>Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies</i> , pages 340–349.	
847		
848		
849		
850		
851		
852	Ramakanth Pasunuru, Mengwen Liu, Mohit Bansal, Sujith Ravi, and Markus Dreyer. 2021. Efficiently summarizing text and graph encodings of multi-document clusters. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4768–4779.	
853		
854		
855		
856		
857		
858		
859	Elizabeth M Perse and Jennifer Lambe. 2016. <i>Media effects and society</i> . Routledge.	
860		
861	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>arXiv preprint arXiv:1910.10683</i> .	
862		
863		
864		
865		
866	Marta Recasens, Cristian Danescu-Niculescu-Mizil, and Dan Jurafsky. 2013. Linguistic models for analyzing and detecting biased language. In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1650–1659.	
867		
868		
869		
870		
871		
872	Ehud Reiter. 2018. A structured review of the validity of bleu. <i>Computational Linguistics</i> , 44(3):393–401.	
873		
874	Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. <i>arXiv preprint arXiv:2110.08207</i> .	
875		
876		
877		
878		
879		
880	Abeed Sarker, Ari Z Klein, Janet Mee, Polina Harik, and Graciela Gonzalez-Hernandez. 2019. An interpretable natural language processing system for written medical examination assessment. <i>Journal of biomedical informatics</i> , 98:103268.	
881		
882		
883		
884		
885	All Sides. 2018. Media bias ratings. <i>Allsides.com</i> .	
886	Timo Spinde, Christina Kreuter, Wolfgang Gaissmaier, Felix Hamborg, Bela Gipp, and Helge Giese. 2021. Do you think it’s biased? how to ask for the perception of media bias. In <i>2021 ACM/IEEE Joint Conference on Digital Libraries (JCDL)</i> , pages 61–69. IEEE.	
887		
888		
889		
890		
891		
	Sam Wiseman, Stuart Shieber, and Alexander Rush. 2017. <i>Challenges in data-to-document generation</i> . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2253–2263, Copenhagen, Denmark. Association for Computational Linguistics.	892 893 894 895 896 897
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. <i>Transformers: State-of-the-art natural language processing</i> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	898 899 900 901 902 903 904 905 906 907 908 909
	Tae Yano, Philip Resnik, and Noah A Smith. 2010. Shedding (a thousand points of) light on biased language. In <i>Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk</i> , pages 152–158.	910 911 912 913 914
	Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2019a. <i>Pegasus: Pre-training with extracted gap-sentences for abstractive summarization</i> .	915 916 917
	Yifan Zhang, Giovanni Da San Martino, Alberto Barrón-Cedeno, Salvatore Romeo, Jisun An, Hae-woon Kwak, Todor Staykovski, Israa Jaradat, Georgi Karadzhov, Ramy Baly, et al. 2019b. Tanbih: Get to know what you are reading. <i>EMNLP-IJCNLP 2019</i> , page 223.	918 919 920 921 922 923
	Appendix	924
	A Topics covered in dataset	925
	The dataset language is English and mainly focuses on U.S. political topics that often result in media bias. The top-5 most frequent topics are ‘Elections’, ‘White House’, ‘Politics’, ‘Coronavirus’, ‘Immigration’.	926 927 928 929 930
	The full list is as follow (in a descending order of frequency): [‘Elections’, ‘White House’, ‘Politics’, ‘Coronavirus’, ‘Immigration’, ‘Violence in America’, ‘Economy and Jobs’, ‘Supreme Court’, ‘Middle East’, ‘US House’, ‘Healthcare’, ‘World’, ‘US Senate’, ‘National Security’, ‘Gun Control and Gun Rights’, ‘Media Bias’, ‘Federal Budget’, ‘Terrorism’, ‘US Congress’, ‘Foreign Policy’, ‘Criminal Justice’, ‘Justice Department’, ‘Trade’, ‘Impeachment’, ‘Donald Trump’, ‘North Korea’, ‘Russia’, ‘Education’, ‘Environment’, ‘Free Speech’, ‘FBI’, nan, ‘Abortion’, ‘General News’, ‘Disaster’, ‘US Military’, ‘Technology’, ‘LGBT Rights’, ‘Sexual Misconduct’, ‘Voting Rights and Voter Fraud’,	931 932 933 934 935 936 937 938 939 940 941 942 943 944

945 ‘Joe Biden’, ‘Race and Racism’, ‘Economic Pol-
 946 icy’, ‘Justice’, ‘Holidays’, ‘Taxes’, ‘China’, ‘Polar-
 947 ization’, ‘Democratic Party’, ‘Religion and Faith’,
 948 ‘Sports’, ‘Homeland Security’, ‘Culture’, ‘Cyber-
 949 security’, ‘National Defense’, ‘Public Health’,
 950 ‘Civil Rights’, ‘Europe’, ‘Great Britain’, ‘Banking
 951 and Finance’, ‘Republican Party’, ‘NSA’, ‘Busi-
 952 ness’, ‘State Department’, ‘Facts and Fact Check-
 953 ing’, ‘Media Industry’, ‘Labor’, ‘Veterans Affairs’,
 954 ‘Campaign Finance’, ‘Life During COVID-19’,
 955 ‘Transportation’, ‘Marijuana Legalization’, ‘Agri-
 956 culture’, ‘Arts and Entertainment’, ‘Fake News’,
 957 ‘Campaign Rhetoric’, ‘Nuclear Weapons’, ‘Israel’,
 958 ‘Asia’, ‘CIA’, ‘Role of Government’, ‘George Floyd
 959 Protests’, “Women’s Issues”, ‘Safety and Sanity
 960 During COVID-19’, ‘Animal Welfare’, ‘Treasury’,
 961 ‘Science’, ‘Climate Change’, ‘Domestic Policy’,
 962 ‘Energy’, ‘Housing and Homelessness’, ‘Bridging
 963 Divides’, ‘Mexico’, ‘Inequality’, ‘COVID-19 Mis-
 964 information’, ‘ISIS’, ‘Palestine’, ‘Bernie Sanders’,
 965 ‘Tulsi Gabbard’, ‘Sustainability’, ‘Family and Mar-
 966 riage’, ‘Pete Buttigieg’, ‘Welfare’, ‘Opioid Cri-
 967 sis’, ‘Amy Klobuchar’, ‘Food’, ‘EPA’, ‘South Ko-
 968 ree’, ‘Alaska: US Senate 2014’, ‘Social Security’,
 969 ‘US Constitution’, ‘Tom Steyer’, ‘Andrew Yang’,
 970 ‘Africa’]

971 B Additional Salient Information Score 972 Results

973 We report additional Salient information F1 (Ta-
 974 ble 5) and Recall (Table 6) scores for ROUGE1,
 975 ROUGE2 and ROUGEL.

	ROUGE1 F1	ROUGE2 F1	ROUGEL F1
LXRANK	33.60%	13.60%	29.77%
BARTCNN	33.76%	13.67%	30.57%
PEGASUSMULTI	30.03%	10.28%	26.70%
BARTMULTI	23.01%	6.84%	20.55%
NEUSFT	36.76%	16.27%	32.86%
NEUS-TITLE	35.49%	15.69%	32.05%

Table 5: Additional Salient Info Scores. F1 scores for ROUGE1, ROUGE2 and ROUGEL for ALLSIDES test-set. For the scores, the *higher* number is the better.

976 C Details for Human Evaluation (A/B 977 testing)

978 We first presented the participants with the defi-
 979 nition of framing bias from our paper, and also

	ROUGE1 RECALL	ROUGE2 RECALL	ROUGEL RECALL
LXRANK	39.08%	17.66%	34.69%
BARTCNN	35.63%	15.32%	32.22%
PEGASUSMULTI	44.42%	16.99%	39.45%
BARTMULTI	35.76%	12.48%	32.08%
NEUSFT	35.11%	15.74%	31.43%
NEUS-TITLE	36.07%	16.47%	32.63%

Table 6: Additional Salient Info Scores. Recall scores for ROUGE1, ROUGE2 and ROUGEL for ALLSIDES testset. For the scores, the *higher* number is the better.

980 showed examples in Table 1 to ensure they under-
 981 stand what framing bias is. Then we asked the
 982 following question: “Which one of the articles do
 983 you believe to be more biased toward one side or
 984 the other side in the reporting of news?” This is
 985 modified to serve as a question for AB testing based
 986 on “To what extent do you believe that the article
 987 is biased toward one side or the other side in the
 988 reporting of news?” The original question is one of
 989 the 21 questions which are suitable and reliable for
 990 measuring the perception of media bias, designed
 991 by Spinde et al. (2021).

992 The participants (research graduate students)
 993 have different nationalities including Canada,
 994 China, Indonesia, Iran, Italy, Japan, Poland and
 995 South Korea (ordered in an alphabetical order). All
 996 of participants answered to be not having political
 997 leaning towards U.S. politics. All participants are
 998 fully explained on the usage of collected data in
 999 this particular work and agreed on it.

1000 D Experimental Setup Details

1001 All our experimental codes are based on the Hug-
 1002 gingFace (Wolf et al., 2020). We used the following
 1003 hyperparameters during training and across models:
 1004 10 epoch size, $3e - 5$ learning rate, and a batch size
 1005 of 16. We did not do hyper-parameters tuning since
 1006 our objective is to provide various baselines and
 1007 analysis. Training run-time for all of our experi-
 1008 ments are fast (< 6 hr). We ran all experiments with
 1009 one NVIDIA 2080Ti GPU with 16 GB of memory.
 1010 The experiment was a single-run.

1011 E Generation Examples from Different 1012 Models

1013 To help better understand performances of each
 1014 models, we provide more examples of generation

1015 from all baseline models and our proposed models
1016 along with the target neutral headline summary.
1017 The examples can be found in Table 7, 8, 9.

1018 **F Illustration of hallucinatory framing** 1019 **bias**

1020 In Table 10, we provide more examples of halluci-
1021 natory framing bias that was discussed in analysis.

Target: *The House Ways and Means Committee has requested that the IRS turn over the last six years of tax returns from President Trump.*

LEXRANK: Representative Richard Neal invoked law that allows House Ways and Means Committee chair to obtain any person's tax information. The request by House Democrats for six years of Donald Trump's tax returns opens a serious new political, legal, constitutional and personal front in the party's battle to use its new majority to check a defiant President.

BARTCNN: Rep. Richard Neal invoked law that allows House Ways and Means Committee chair to obtain any person's tax information. The request by House Democrats for six years of Donald Trump's tax returns opens a serious new political, legal, constitutional and personal front in the party's battle to use its new majority to check a defiant President.

PEGASUSMULTI: The battle over President Trump's tax returns is officially on. Rep. Richard Neal, chairman of the House Ways and Means Committee, on Wednesday formally asked the IRS for six years of Trump's personal and business tax returns, as well as the returns of some of his companies, reports the Washington Post. The panel has the power to subpoena someone's tax information, and Neal has signaled that he plans to do just that. "The American people have a right to know who the president of the United States is, and we have a responsibility to do our part to ensure that the tax laws are administered in an efficient and thorough manner," he said in a statement, per the Hill. The White House has signaled that it will fight the request, with press secretary Sarah Huckabee Sanders calling it "an abuse of the tax-writing committee's authority." The letter to the IRS was signed by Neal, a Massachusetts Democrat, and was obtained by the Wall Street Journal. The Journal notes that Neal has the authority under a law passed in 1969 that allows the House Ways and Means Committee to obtain tax information from anyone.

BARTMULTI: the head of the powerful house tax-writing committee has asked the government for six years of president trump's tax returns, a move that is expected to lead to a long court battle with the white house, the washington post reports. in a letter to the irs on wednesday, rep. richard h. Neal, chairman of the house ways and means committee, invoked a law that allows any person to obtain any person's personal tax information, the wall street journal reports. according to the hill, the move is a " serious new political, legal, constitutional, and personal front in the party's battle to use its new majority to check a defiant president. " it is a bid to solve one of the most tantalizing and enduring mysteries of the trump campaign and his presidency — what is it in the president's unseen financial and business life that he doesn't want americans to see?

NEUSFT: House Ways and Means Committee Chair Richard Neal (D-MA) has formally requested President Trump's tax returns from the IRS.

NEUS-TITLE: TITLE=> House Committee Seeks Trump Tax Returns. ARTICLE=> The House Ways and Means Committee formally requested President Trump's tax returns on Wednesday, kicking off what could be a protracted legal fight between Congress and the Trump administration.

Table 7: Neutral Summary Generation Examples from baseline models and NEUS-TITLE.

Target: *The State Department has blocked Gordon Sondland, the U.S. Ambassador to the European Union, from testifying at a scheduled deposition on the Ukraine controversy. Sondland is a key figure as he had texted with Bill Taylor, a top U.S. diplomat in Ukraine, arguing over allegations of a quid pro quo. This move by the State Department has triggered a new clash with House Democrats as part of their impeachment inquiry, and President Trump has dismissed the deposition as a “a kangaroo court” as he defended the State Department decision.*

LEXRANK: The Trump administration on Tuesday blocked a planned deposition from Gordon Sondland, the U.S. ambassador to the European Union and a central figure in the Democrats’ impeachment inquiry, according to a statement by his lawyer. Democrats swiftly called the administration’s move an act of obstruction and accused President Trump of hiding information about his efforts to pressure the Ukrainian government to investigate his political rivals.

BARTCNN: The State Department instructed Ambassador to the European Union Gordon Sondland not to appear for the deposition, according to his counsel. Democrats swiftly called the administration’s move an act of obstruction and accused President Trump of hiding information about his efforts to pressure the Ukrainian government to investigate his political rivals.

PEGASUSMULTI: The Trump administration has blocked US ambassador to the European Union Gordon Sondland from testifying before a House panel investigating possible obstruction of justice by the president, the New York Times reports. Sondland is a key witness in the probe into whether Trump obstructed justice when he fired FBI chief James Comey. Sondland had been scheduled to testify in front of the House Judiciary Committee on Tuesday, but the State Department ordered him not to appear, according to the Washington Post. Sondland is a former British ambassador to the US. In a Tuesday tweet, President Trump said he would love to send Sondland to testify, "but unfortunately he would be testifying before a totally compromised kangaroo court, where Republican’s rights have been taken away, and true facts are not allowed out for the public to see." He added: "I believe you are incorrect about President Trump’s intentions."

BARTMULTI: a key witness in the house impeachment inquiry was ordered not to appear at a scheduled deposition on tuesday, ramping up tensions between congressional democrats and the trump administration. the state department instructed ambassador to the eu and key figure in the impeachment inquiry, Gordon sondland, not to show up for the deposition, according to his counsel. democrats called the move an act of obstruction and accused president trump of hiding information about his efforts to pressure the ukraine government to investigate his political rivals, the new york times reports. " this is a clear attempt to block testimony from key witnesses in the ongoing impeachment investigation, " said house minority leader nancy pelosi, per the washington post. " it’s an impeachable act. " in a tweet tuesday night, trump said he " would love to send ambassador [sundland], a really good man and great american, to testify, but unfortunately he would be testifying before a totally compromised kangaroo court, where republican ’ s rights have been taken away, and true facts are not allowed out for the public to see. "

NEUSFT: The Trump administration blocked a planned deposition from Gordon Sondland, the U.S. ambassador to the European Union and a central figure in the Democrats’ impeachment inquiry, according to a statement from his lawyer. President Trump responded via Twitter, saying that he “would love to send Ambassador Sonderland, a really good man and great American, to testify, but unfortunately he would be testifying before a totally compromised kangaroo court, where Republican’s rights have been taken away, and true facts are not allowed out for the public to see.”

NEUS-TITLE: TITLE=> State Dept Orders Ambassador Sondland Not to Testify Before House. ARTICLE=> The State Department has ordered Ambassador Gordon Stondland not to testify before the House of Representatives in the impeachment inquiry of President Donald Trump. Democrats have accused the Trump administration of blocking statements from witnesses.

Table 8: Continued from Previous Page: Neutral Summary Generation Examples from baseline models and NEUS-TITLE.

Target: *Ukrainian police have opened an investigation into whether or not U.S. Ambassador Marie Yovanovitch came under surveillance before she was recalled from her post in Ukraine last May. Democrats have released documents that show Lev Parnas, an associate of Rudy Giuliani, communicating about Yovanovitch's removal.*

LEXRANK: Ukraine's government announced Thursday that police are investigating whether ousted U.S. ambassador Marie Yovanovitch was subject to illegal surveillance, in response to new documents released ahead of President Trump's impeachment trial. Those documents, released by Democratic lawmakers, showed Lev Parnas – an associate of Trump lawyer Rudy Giuliani – communicating about the removal of Marie Yovanovitch as the ambassador to Ukraine.

BARTCNN: Police in Ukraine have opened a criminal investigation into whether allies of President Trump had the U.S. ambassador to the country under surveillance. Newly released documents suggest that Ms. Yovanovitch was being watched in Kyiv. The Ministry of Internal Affairs of Ukraine says it's asking the FBI to help its inquiry.

PEGASUSMULTI: The Ukraine has opened a criminal investigation into allegations that President Trump's allies had the US ambassador to the country under surveillance while she was stationed in the country, the New York Times reports. According to the Wall Street Journal, the allegations surfaced in newly released documents released by Democrats ahead of Trump's impeachment trial. The documents show that Rudy Giuliani's associate, Lev Parnas, communicated about the removal of Marie Yovanovitch as ambassador to Ukraine with Yovanovitch's former boss at the State Department, Sergey Kislyak. Yovanovitch was removed from her post in May of last year. Ukraine's Ministry of Internal Affairs says it's asking the FBI to help with the investigation. "Ukraine cannot ignore such illegal activities on the territory of its own state," the ministry says in a statement

BARTMULTI: new documents released ahead of president trump's impeachment trial suggest that the us ambassador to the country was under surveillance while she was stationed in the country, the new york times reports. according to the wall street journal, the allegations were made in a text message between a us politician and an associate of trump lawyer rudy giuliani. in the messages, the politician says he'd like to see former ambassador to ukraine and current ambassador to russia, former deputy foreign minister oleksandr turchynov, removed from his post. " i would like to remove him from his position, " the politician wrote in the message, which was sent in april 2017. the messages were sent in response to newly released documents from the house intelligence committee. those documents show that former trump lawyer paul ryan had been trying to get rid of yovanovitch, who was ambassador at the time, for more than a year, the journal reports. the messages also show that ryan's chief of staff, george w. bush, had been in touch with ryan about removing her from her post.

NEUSFT: Ukraine's national police are investigating allegations that U.S. Ambassador Marie Yovanovitch was subjected to illegal surveillance while she was stationed in Kyiv.

NEUS-TITLE: TITLE=> Ukraine Opens Investigation Into Claims U.S. Ambassador Yovanovitch Was Under Surveillance. ARTICLE=> Ukraine's national police are investigating allegations that former U. S. ambassador Marie Yovanovitch was subjected to surveillance while she was stationed in Kyiv, following the release of documents by Democratic lawmakers.

Table 9: Continued from Previous Page: Neutral Summary Generation Examples from baseline models and NEUS-TITLE.

[1] Source: Newly released documents suggest that Ms. Yovanovitch was being watched in Kyiv while she was the U.S. ambassador there. The police in Ukraine have opened a criminal investigation into whether allies of President Trump had the United States ambassador to the country under surveillance while she was stationed in Kyiv, the Ukrainian government said on Thursday. [SEP] Ukraine's government announced Thursday that police are investigating whether ousted U.S. ambassador Marie Yovanovitch was subject to illegal surveillance, in response to new documents released ahead of President Trump's impeachment trial. Those documents, released by Democratic lawmakers, showed Lev Parnas – an associate of Trump lawyer Rudy Giuliani – communicating about the removal of Marie Yovanovitch as the ambassador to Ukraine. [SEP] Ukraine's national police are investigating allegations that U.S. Ambassador Marie Yovanovitch was under surveillance in Kyiv last spring — an allegation made in WhatsApp messages between a U.S. politician and Lev Parnas. "Ukraine cannot ignore such illegal activities on the territory of its own state," the Ministry of Internal Affairs of Ukraine said in a statement Thursday. The ministry says it's asking the FBI to help its inquiry.

[1] MDS Hallucination: new documents released ahead of president trump's impeachment trial suggest that the us ambassador to the country was under surveillance while she was stationed in the country, the new york times reports. according to the wall street journal, the allegations were made in a text message between a us politician and an associate of trump lawyer rudy giuliani. in the messages, the politician says he'd like to see former ambassador to ukraine and current ambassador to russia, former deputy foreign minister oleksandr turchynov, removed from his post. "i would like to remove him from his position," the politician wrote in the message, which was sent in april 2017. the messages were sent in response to newly released documents from the house intelligence committee. those documents show that former trump lawyer paul ryan had been trying to get rid of yovanovitch, who was ambassador at the time, for more than a year, the journal reports. the messages also show that ryan's chief of staff, george w. bush, had been in touch with ryan about removing her from her post.

[2] Source: "If I threaten someone, you'll know it," White House counselor Kellyanne Conway told reporters Friday morning in a bizarre hours-long effort to discredit an audio recording that seemingly showed her threatening to use West Wing resources to investigate the personal life of a reporter a day earlier. "I never threatened anyone," Conway told a gaggle of journalists at the White House. "Don't use those words. I didn't threaten. It's not a threat. I never threatened her." [SEP] A tense phone conversation between a reporter for the Washington Examiner and White House senior counselor Kellyanne Conway was published by the newspaper on Thursday. In the conversation, Conway objected that a story written by the reporter, Caitlin Yilek, mentioned that her husband George Conway is a fierce critic of President Trump on Twitter. Yilek was writing a story on Conway possibly becoming President Trump's next White House chief of staff if Trump decides to move on from the official now in the position, Mick Mulvaney. [SEP] Kellyanne Conway defended her phone call with a Washington Examiner reporter during an appearance on Fox News. Fox News host Sandra Smith described the phone call Friday between reporter Caitlin Yilek and the counselor to the president as a "heated clash" during which things got "very personal." Conway claimed, "Everything I said in that phone call, I've said publicly before. Every single thing. I said it again." Smith asked Conway if there was anything she regretted about how she handled her phone call, which she denied.

[2] MDS Hallucination: a recording of a tense phone call between senior white house adviser and reporter catherine yilek was published by the washington examiner on friday, and it's not going over well. in the recording, yalek says she was told by the white house that she would be fired if she didn't back off from her reporting on president trump's controversial son-in-law, and that she'd have to go to court if she did not. "i'm not going to do that," she says in the recorded conversation, per the new york times. "you're going to get fired if you don't shut the f — up." in the call, she also says she'll use the office of management and budget to investigate the personal life of the reporter. "if i threaten someone, you'll know it," the caller can be heard saying in the audio recording, per politico. "don ' t use those words. it ' s not a threat. i never threatened anyone." but on monday, white house counselor to the president katie holmes told fox news that she had never threatened the reporter.

Table 10: Examples of hallucinatory framing bias from MDS models and the corresponding the source input.