

# Detecting Hallucination and Coverage Errors in Retrieval Augmented Generation for Controversial Topics

Tyler A. Chang<sup>\*1,2</sup>, Katrin Tomanek<sup>\*1</sup>, Jessica Hoffmann<sup>1</sup>, Nithum Thain<sup>1</sup>,  
Erin van Liemt<sup>1</sup>, Kathleen Meier-Hellstern<sup>1◇</sup>, Lucas Dixon<sup>1◇</sup>

<sup>1</sup>Google Research

<sup>2</sup>UC San Diego

tachang@ucsd.edu

{katrintomanek, jhoffmann, nthain, evanliemt, kathyhellstern, ldixon}@google.com

## Abstract

We explore a strategy to handle controversial topics in LLM-based chatbots based on Wikipedia’s Neutral Point of View (NPOV) principle: acknowledge the absence of a single true answer and surface multiple perspectives. We frame this as retrieval augmented generation, where perspectives are retrieved from a knowledge base and the LLM is tasked with generating a fluent and faithful response from the given perspectives. As a starting point, we use a deterministic retrieval system and then focus on common LLM failure modes that arise during this approach to text generation, namely hallucination and coverage errors. We propose and evaluate three methods to detect such errors based on (1) word-overlap, (2) salience, and (3) LLM-based classifiers. Our results demonstrate that LLM-based classifiers, even when trained only on synthetic errors, achieve high error detection performance, with ROC AUC scores of 95.3% for hallucination and 90.5% for coverage error detection on unambiguous error cases. We show that when no training data is available, our other methods still yield good results on hallucination (84.0%) and coverage error (85.2%) detection.

**Keywords:** conversational systems, natural language generation, evaluation methodologies

## 1. Introduction

Large Language Models (LLMs) have achieved state-of-the-art performance on a wide range of tasks, and a growing audience of users is engaging with LLM-driven chatbots.<sup>1</sup> While these chatbots are highly flexible and generalizable, they are known to struggle with factuality and bias (Sheng et al., 2019; Shuster et al., 2021; Chang and Bergen, 2024). In many real world scenarios, model developers require more precise control over LLM-based chatbot responses.

In this paper, we investigate how LLMs can be used with retrieval augmented generation for controversial topics, and we propose methods to detect errors in the tuned LLM responses. In retrieval augmented generation, factual information is retrieved and provided as additional context to an LLM (Lewis et al., 2020; Li et al., 2022; Azure, 2023; Iyer and Thallam, 2023). Through curated retrieval sources, retrieval augmented generation enables fine-grained control over LLM responses. However, in the case of controversial topics, users often seek information for which there are not agreed-upon factual answers. These topics range from the in-

consequential (e.g. “the superiority of the Yankees vs. the Red Sox”) to the fundamental (“What religious faith should I adhere to?”). Building useful LLMs requires the ability to ensure that LLM responses adhere to desired levels of neutrality and nuance in such cases.

Thus, we introduce the **NPOV Response Task**: given a query about a controversial topic, the model retrieves arguments for multiple perspectives and is tasked to generate a multi-perspective response, inspired by Wikipedia’s Neutral Point of View (NPOV) principle. We use a deterministic argument retrieval system, and we focus on the challenge of faithful response generation from provided arguments. We adapt a conversational LLM to this task and examine two common error types that violate faithfulness to inputs: (1) **hallucinations** (generating unprovided arguments), and (2) **coverage errors** (omitting provided arguments).

We build a dataset of model query-response pairs, conditioned on arguments from Britannica’s ProCon (ProCon.org, 2022). Using expert annotators, we identify instances of hallucination and coverage errors. We then propose methods for detecting such hallucination and coverage errors, both with and without access to human-labeled data.

Our main results demonstrate that with access to error-free examples and examples containing only synthetic errors, LLM-based classifiers can achieve ROC AUCs of 95.3% and 90.5% in de-

\*Joint first authorship.

◇ Research group leadership.

<sup>1</sup>Among others: <https://openai.com/blog/chatgpt>; <https://bard.google.com>; <https://www.anthropic.com/index/introducing-claude>.

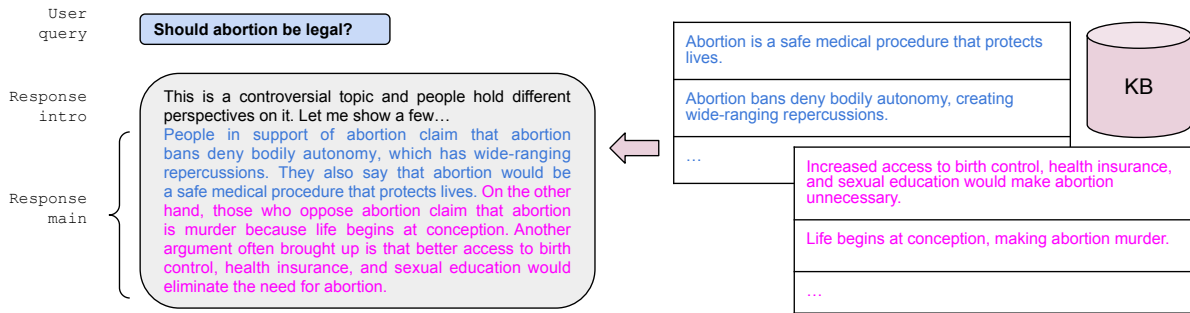


Figure 1: Example NPOV response to a user query on a sensitive topic (left) based on pro and con perspectives with two arguments each, as obtained from a knowledge base (right). Arguments taken from ProCon.org (2022). Our error detection methods focus on the NPOV main response.

detecting organic hallucinations and coverage errors respectively on our task. Even without access to annotated data, we can leverage salience and word overlap techniques to achieve ROC AUCs of 84.0% for hallucinations and 85.2% for coverage errors. While we focus on NPOV response generation, our approaches can be applied more generally to detect hallucination and coverage errors in retrieval augmented generation, facilitating finer-grained control over LLM responses.

## 2. Handling Controversial Topics

Our work is centered around how LLMs can be controlled to respond to queries about controversial topics for which there is no single correct answer. For example, in response to “Should abortion be legal?”, an LLM without direction might produce a highly opinionated or offensive response. To address such concerns, “guardrails” are often added to LLMs, either completely preventing the generation of responses to such topics or responding with canned answers (“I am just a language model and cannot answer this question...”). Such approaches can lead to erasure harm and reduce the usefulness of the system on potentially important topics. Another approach is to personalize responses to align with a user’s position; however, this can reinforce harmful biases and popular misconceptions, and act as a chatbot echo chamber.

As an alternative strategy, we propose to acknowledge the lack of agreement and surface main viewpoints instead. This approach is inspired by Wikipedia’s **Neutral Point of View (NPOV)** principle, which requires that content is written such that it represents “fairly, proportionately, and, as far as possible, without editorial bias, all the significant views that have been published by reliable sources on a topic.”<sup>2</sup> Figure 1 (left) gives an example of an NPOV response on a highly controversial topic. We

<sup>2</sup>From [https://en.wikipedia.org/wiki/Wikipedia:Neutral\\_point\\_of\\_view](https://en.wikipedia.org/wiki/Wikipedia:Neutral_point_of_view), last accessed 2023/10/20.

explore whether such responses can be generated by an LLM using retrieval augmented generation, and we detect common failure modes such as hallucination and coverage errors.

### 2.1. NPOV Response Generator

We separate *response* generation from *content* generation. For the scope of this paper, we assume that there is a content retrieval process and a knowledge base of curated arguments for different perspectives. The knowledge base we use in this paper consists of arguments from Britannica’s ProCon website (§2.2).

The NPOV Response Task is then: given the user query and retrieved perspectives (where perspectives consist of concatenated arguments), generate a response that consists of an introduction sentence, serving as a bridge from the user query, and a verbalization of the given perspectives. When generating the response, relevant aspects of the given arguments must not be dropped (ensure full coverage) and no other arguments should be added (avoid hallucinations). This task formulation gives model developers fine-grained control over LLM responses. An example is shown in Figure 1.

We use soft prompt-tuning (Lester et al., 2021) to adapt an LLM to generate NPOV responses given pro and con arguments. Our base LLM is a 64B decoder-only LaMDA model pre-trained on public dialog data and web text (Cohen et al., 2022). We use a soft prompt length of 5 tokens, and we train for 20K steps with batch size 16 and learning rate 0.1. We typically reach maximum dev set performance after 2-5K steps. Specific prompt format and detailed hyperparameters are in Appendix B.

Our training set consists of 80 query-response pairs covering 9 controversial topics from ProCon (§2.2). ProCon question headers (e.g. “Should abortion be legal?”) are used as user queries. For each topic, we randomly sample one, two, or three arguments from both the pro and con side in Pro-

Con<sup>3</sup> and then manually write several paraphrased responses capturing these arguments. We observe that after prompt-tuning, the NPOV Response Generator generalizes well beyond the topics and arguments seen during training.

## 2.2. ProCon as a Knowledge Base

Britannica’s ProCon (ProCon.org, 2022) is a website presenting pros and cons for commonly debated topics. Pros and cons are researched and compiled by ProCon research staff and editors, and they aim to be nonpartisan.<sup>4</sup> As of October 2022, ProCon contains 72 active (i.e. “non-archived”) topics. For both the pro and con perspective for each topic, several arguments are given, each consisting of a short argument phrase accompanied by a longer explanation. The median number of arguments per perspective per topic is 4, but some topics contain many more arguments (e.g. *Social Media* has 23 arguments per perspective). We randomly sample ProCon arguments as inputs to the NPOV Response Generator for each topic (§4.1). Each topic is associated with a leading question in ProCon (e.g. “Should abortion be legal?”), which we treat as the user query asked to the LLM.

## 3. Methods to Detect Hallucinations and Coverage Errors

We focus on hallucination and coverage error detection, adopting the following definitions:

- If the generated response contains at least one argument which was not provided, we call this a **hallucination**.
- If one or more of the given arguments is completely dropped from the response, we call this a **coverage error**.

We call these **full errors**, as they address the hallucination or coverage of a full argument. On top of these well-defined errors, we notice that the NPOV Response Generator sometimes produces other unfaithful changes to arguments, including: (1) partial hallucinations (slight meaning change, e.g. “consensus” becomes “unanimity”), (2) partial coverage errors (only a part of the argument is dropped), (3) repetitions (response contains the same given argument multiple times), and (4) perspective confusions (response inverts the perspectives, e.g. pro arguments are presented as cons). We call all of these **ambiguous errors**.

<sup>3</sup>We always ensure the same number of pro and con arguments.

<sup>4</sup>Of course, not all controversial topics can be framed as pro versus con debates, and such a binary framing of highly complex topics can omit important nuance (see Ethical Considerations).

We propose three methods for detecting hallucination and coverage errors in generated responses: ROUGE, salience, and LLM-based classifiers.

### 3.1. ROUGE

As a baseline, we use ROUGE-1 (word-matching) to compute hallucination and coverage error scores (Lin, 2004). For a given response from the NPOV Response Generator, ROUGE calculates the proportion of *response words* that are matched in the *input arguments* (ROUGE-1 precision) and the proportion of *input argument words* that are matched in the *response* (ROUGE-1 recall).<sup>5</sup> Low precision is indicative of hallucination, and low recall is indicative of a coverage error. Because the NPOV Response Task requires that both input perspectives be covered, we compute ROUGE-1 recall separately for each input perspective and then compute the minimum as our overall recall score. For ROUGE, words are defined using whitespace and punctuation separation, dropping stop words and using word stemming from NLTK (Bird et al., 2009).

### 3.2. Salience

Aside from word matching, previous work has proposed methods to attribute output subword tokens to input tokens in LLMs using model gradients (Denil et al., 2014; Li et al., 2016; Bastings and Filipova, 2020). These methods are computationally costly, but they can often capture more nuance (e.g. word synonyms and token interactions) than simple word-matching. One popular approach is to compute the logit (pre-softmax probability) gradient for each output token with respect to each input token embedding, producing a gradient vector for each input-output token pair. The attribution from each input to the output token is defined as the dot product between the corresponding gradient vector and the input token embedding (Denil et al., 2014).<sup>6</sup>

In the NPOV Response Generation scenario, there are attribution values from each input token (e.g. the given arguments per perspective and the user query) and each previously generated token to each output token. This produces a token-to-token salience map  $M_{\text{tokens}} \in \mathbb{R}^{(m+\ell) \times \ell}$ , where  $m$  is the number of input tokens and  $\ell$  is the number of model response tokens. Before any further processing, we square the salience map and normalize columns to sum to one (i.e. the attributions to each output token sum to one).

<sup>5</sup>We also implemented a hallucination and coverage error detection method that matched input and response arguments with BERTScore (Zhang et al., 2020), but we obtained similar results to ROUGE. We omit results due to space limitations.

<sup>6</sup>We obtain similar results using gradient L2 norms.

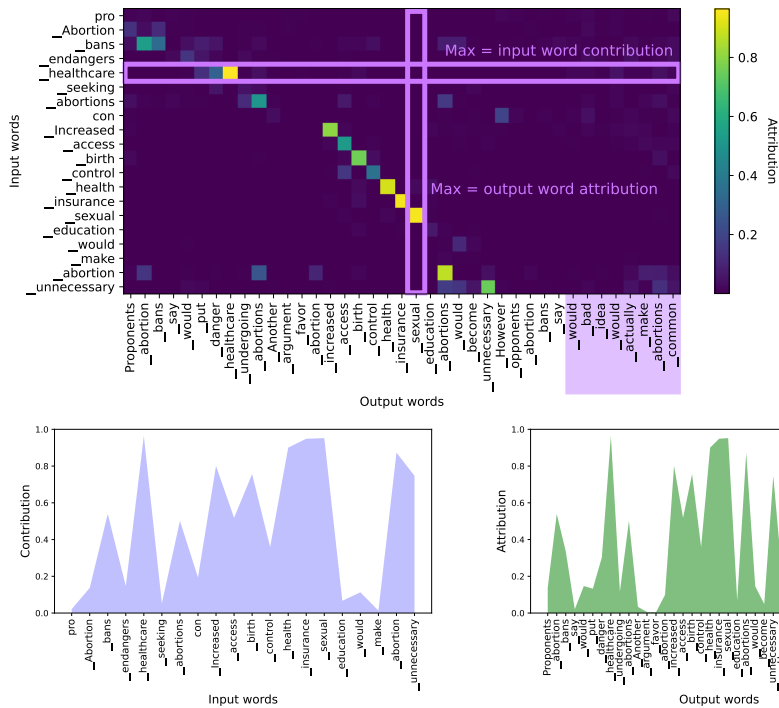


Figure 2: Top: saliency map from input argument content words (rows) to model response content words (columns). Bottom: individual word scores for contribution (input words; left) and attribution (response words; right). The purple highlighted words are hallucinated in the model response.

Because we are primarily concerned with hallucination and coverage errors for content words, we convert the subword token-to-token saliency map to a word-to-word saliency map  $M_{\text{words}}$ . We define words by concatenating consecutive LLM tokens that are not separated by punctuation or whitespace; we then drop stop words, as defined in NLTK (Bird et al., 2009). We define the attribution from an input word  $w_0$  to an output word  $w_1$  as the maximum attribution from any subword token in  $w_0$  to any subword token in  $w_1$ . We restrict our saliency maps to the input argument words (rows) and the output NPOV response words (columns). A sample word-to-word saliency map for a query-response pair is shown in Figure 2.

Qualitatively, we observe that covered words tend to have a high contribution to a single corresponding word in the response. Thus, we define the *contribution* score of an input argument word as its maximum contribution to any response word (i.e. maximum for each row of  $M_{\text{words}}$ ). We define the *attribution* score of a response word as its maximum attribution from any input argument word (i.e. maximum for each column of  $M_{\text{words}}$ ). Contribution and attribution scores for input words and response words respectively are shown in Figure 2.

To compute an example-level *contribution* score for a query-response pair, we compute the mean contribution score over words in each of the two input perspectives. As with ROUGE, we take the minimum of the two perspective contributions as a

final contribution score. To compute an example-level *attribution* score, we compute the mean attribution score over all response words. Finally, hallucination and coverage error scores in  $[0, 1]$  are computed by subtracting the attribution and contribution scores respectively from 1.0. Formal equations are in Appendix C.

### 3.3. LLM-Based Classifiers

The two previous methods for detecting hallucination and coverage errors are data-free, not requiring labeled model responses for training. For the non-data-free scenario, we explore how well LLM-based classifiers perform on these tasks, relying on a small set of human annotations of model responses ( $\sim 500$  examples; §5.2). Our classifiers are built on FLAN-PaLMChilla, a 62B decoder-only LLM (Chowdhery et al., 2023) which has been instruction-tuned on a large number of tasks (Chung et al., 2022). We use soft prompt-tuning to adapt this LLM into classifiers for hallucination and coverage error detection. The classifiers have as input: (1) the user query, (2) the generated NPOV response, and (3) the given arguments per perspective. We train the LLM to predict the label “NO” if there is a full error and “YES” otherwise. Prompt-tuning hyperparameters are the same as §2.1; specific prompt formats are in Appendix B. We tune the classifiers separately for the two error types.

For inference, we generate error classification scores in  $[0, 1]$  by obtaining the LLM’s log perplexity scores for the tokens corresponding to the two output class labels (“YES” and “NO”), apply softmax, and take the score of the negative class (“NO”).<sup>7</sup>

## 4. Dataset

To train and evaluate the hallucination and coverage error detection methods above on the NPOV Response Task, we construct datasets of organic (i.e. naturally occurring) and synthetic errors, with and without paraphrasing.

### 4.1. Annotation Procedure

For each of the 72 controversial topics from ProCon, we generate a unique query and up to 18 query-response pairs by first randomly sampling combinations of pro and con arguments, with either 1, 2, or 3 arguments per side, and then using the NPOV Response Generator to generate a response. We annotate these query-response pairs (also called *examples*) in three stages to (1) identify error-free examples, (2) identify examples with errors, and (3) generate paraphrased examples:

1. For the first three examples per topic, we sample two generator responses, with sampling temperatures 0.0 and 0.7. We annotate whether responses contain hallucinations or coverage errors, annotating examples with a mix of the two temperatures. We annotate the token spans in the response that cover each input argument, along with any hallucinated response spans and uncovered input argument spans.
2. Because examples with hallucination and coverage errors are less frequent than error-free examples even for high temperatures (20.0% errors in 0.7 temperature responses), we sample a single 0.7 temperature response for each of the remaining (up to) 15 examples per topic.<sup>8</sup> We annotate for hallucination and coverage errors, including full and ambiguous errors (§3).
3. Hallucination and coverage error detection methods should capture whether meaning is retained between input arguments and generated responses, even if the arguments are not copied verbatim. We therefore generate examples with enforced paraphrasing between the input arguments and the response. To do so,

<sup>7</sup>For single-token labels, this score equals the probability of “NO” conditioned on either “YES” or “NO” output. We obtain similar results training the models with flipped labels, i.e. “YES” for errors and “NO” otherwise.

<sup>8</sup>Preliminary experiments with the NPOV Response Generator suggest that temperatures above 1.0 tend to produce overly long and irrelevant responses.

we paraphrase the input arguments for all error-free examples generated in Step 1. For each argument, we use an off-the-shelf paraphrasing tool and manually verify that the paraphrasing does not induce substantial meaning change.<sup>9</sup>

In total, we identify 160 examples with no errors and 326 examples with at least one error, and we generate 152 paraphrased examples with no errors.

#### 4.1.1. Inter-Annotator Agreement

To validate the viability and coherence of our annotation task, we hired a team of 10 external annotators to re-identify both hallucination and coverage errors in our dataset. Our annotation provider was paid 49 USD per hour for a total of 25 hours of work (Appendix D). Annotators were presented with 188 of the query-response pairs annotated in annotation Step 1 (§4.1) and 86 pairs from Step 2. Given the user query, the provided arguments, and the response from the NPOV Response Generator, annotators were asked to mark whether each response had a hallucination or coverage error. Each query-response pair was annotated by 5 annotators. We compare the annotator majority vote to our annotated labels, finding 90% agreement for hallucinations and 94% for coverage errors. To measure inter-annotator agreement, we compute Krippendorff’s alpha for hallucinations ( $\alpha = 0.60$ ) and coverage errors ( $\alpha = 0.73$ ) across the 10 annotators. These values are in line with or above similar text classification tasks (Wulczyn et al., 2017).

### 4.2. Synthetic Errors Dataset

Due to the relative rarity of *organic* errors produced by the NPOV Response Generator, we synthetically generate examples with errors by modifying error-free query-response pairs. Specifically, we modify the list of given arguments while keeping the original response unchanged. For coverage errors, we add one randomly sampled unused argument for the given topic from ProCon and add it to the list of given arguments. This creates a full coverage error because the original response does not cover this argument. For hallucinations, we randomly remove one of the given arguments. This creates a hallucination because the original response still addresses the removed argument. We apply synthetic error generation to both paraphrased and unparaphrased examples that were annotated as error-free in §4.1 (312 examples), generating 667 new examples with synthetic hallucinations, synthetic coverage errors, or both.

<sup>9</sup>We use <https://quillbot.com/> for paraphrasing. We find it more efficient to paraphrase the input arguments than to paraphrase the whole response.

Test set error type	Hallucinations			Coverage Errors		
	ROUGE	Saliency	Classifier	ROUGE	Saliency	Classifier
Full organic	0.840	0.808	<b>0.953</b>	0.795	0.852	<b>0.905</b>
Unparaphrased synthetic	0.772	0.736	<b>0.998</b>	0.890	0.875	<b>0.986</b>
Paraphrased synthetic	0.680	0.708	<b>0.977</b>	0.746	0.831	<b>0.993</b>
Ambiguous organic	0.814	0.772	<b>0.851</b>	<b>0.834</b>	0.755	0.756

Table 1: ROC AUCs for example-level hallucination and coverage error detection on four test sets (§4.3).

### 4.3. Test Sets with Different Error Types

Taking the annotations and synthetic errors generated above, we split the 72 ProCon topics into a train set (9 topics), development set (28 topics), and test set (35 topics). We intentionally make our development and test sets substantially larger than our train set because our work focuses on evaluation (rather than training) of the NPOV Response Generator. Our dataset contains two types of query-response pairs (paraphrased and unparaphrased) and three types of errors (synthetic full, organic full, and organic ambiguous). We evaluate the performance of our error detection methods on different slices of the test set to better understand where different approaches have strengths or weaknesses. Hence, each table in the results section states the specific test set slices evaluated:

- **Full organic:** unparaphrased error-free examples vs. organic full errors.
- **Unparaphrased synthetic:** unparaphrased error-free examples vs. corresponding examples with synthetically-generated errors.
- **Paraphrased synthetic:** paraphrased error-free examples vs. corresponding examples with synthetically-generated errors.
- **Ambiguous organic:** unparaphrased error-free examples vs. ambiguous organic errors, including partial errors, repetition, and perspective confusion (§3).

## 5. Results

### 5.1. Example-Level Error Detection

First, we evaluate the three error detection methods (ROUGE, saliency, and classifiers) at the example-level, i.e. detecting whether a query-response pair contains an error. The classifiers shown here are trained only on query-response pairs which are either error-free or contain *synthetic* errors, including both paraphrased and unparaphrased versions (503 examples total); we explore the impact of training data on classifier performance in §5.2.

Table 1 shows ROC AUC scores on the different test sets (§4.3) for all three methods.<sup>10</sup> While

<sup>10</sup>The area under the receiver operating characteristic

the *full organic* set (organic error-free examples vs. organic full errors) is the most realistic, our synthetic sets allow for more controlled evaluations. For all four test sets and for both hallucination and coverage errors, the ROC AUC difference when comparing the best performing method to either other method is statistically significant ( $p < 0.001$ ), using the Wilcoxon statistic (Hanley and McNeil, 1983) and Bonferroni correction for multiple comparisons (Bonferroni, 1936; VanderWeele and Mathur, 2019).

Classifiers consistently outperform the other two methods by a large margin on all sets except ambiguous coverage errors (discussed below), with ROC AUCs above 90% for both hallucination and coverage error detection, for all full error types (organic and synthetic, paraphrased and unparaphrased). Comparing ROUGE and saliency, results are mixed. On the full organic errors, ROUGE performs better at detecting hallucinations (84.0% AUC), whereas saliency performs better at detecting coverage errors (85.2% AUC).

For copy-like tasks with few expected word changes, ROUGE outperforms saliency on both hallucination and coverage error detection (results on the unparaphrased synthetic errors set). However, on the paraphrased synthetic errors, saliency appears to capture the underlying semantics better than ROUGE, allowing it to more accurately detect both hallucination and coverage errors.

Finally, we evaluate our methods on ambiguous errors (including partial argument hallucination and coverage errors, argument repetition, and perspective confusion; see §3). ROUGE performs well here, likely due to minimal natural paraphrasing from the NPOV Response Generator. Classifier ROC AUC scores drop substantially on ambiguous errors, likely because classifiers are trained only on full errors. This discrepancy seems most problematic for coverage error detection, where classifiers perform even worse than ROUGE. Future work should establish clearer definitions of ambiguous errors, allowing larger sets of ambiguous errors to be annotated and used to train classifiers.

curve (ROC AUC) quantifies classification performance across classification thresholds by comparing the trade-off between true positive rate and false positive rate.

Test set error type	Hallucinations				Coverage Errors			
	Error-free +Synth	+Para	+Dev	+Org	Error-free + Synth	+Para	+Dev	+Org
Full organic	0.789	0.828	0.953	0.920	0.880	0.903	0.905	0.956
Ambiguous organic	0.807	0.820	0.851	0.862	0.702	0.529	0.756	0.640

Table 2: ROC AUC scores for classifiers trained on different amounts and types of data (§5.2), ordered from smallest to largest training set size. Table 1 results use the classifiers trained on +Dev.

Test set error type	Hallucinations		Coverage Errors	
	ROUGE	Saliency	ROUGE	Saliency
Full organic	0.673	<b>0.724</b>	0.669	<b>0.799</b>
Unparaphrased synthetic	0.697	<b>0.710</b>	0.693	<b>0.808</b>
Paraphrased synthetic	0.614	<b>0.673</b>	0.582	<b>0.742</b>
Ambiguous organic	<b>0.542</b>	<b>0.542</b>	0.738	<b>0.740</b>

Table 3: ROC AUC scores for word-level error detection using ROUGE and saliency.

## 5.2. Classifier Training Data Ablations

We analyze the impact of different types and amounts of training data on classifier performance, considering the following four scenarios:

- **Error-free +Synth:** all error-free query-response pairs, plus synthetic errors; training split only (70 examples).
- **+Para:** previous, plus equivalent paraphrased examples; training split only (138 examples).
- **+Dev:** previous, plus equivalent examples from the development split (503 examples).
- **+Org:** previous, plus examples with organic full errors; training and development splits (573 examples).

Table 2 shows classifier performance on the full organic and ambiguous organic test sets (§4.3). For coverage error detection, performance strictly improves on the full organic set as we add more training data. However, adding the organic error examples leads to a decline in performance on the ambiguous organic set. For hallucination detection as well, we see performance improvement when adding more training data. Adding the organic errors leads to a performance drop on the full organic set, but not the ambiguous organic set.

Overall, adding more data, even consisting of synthetic errors, leads to improvements on most test sets for both hallucination and coverage error detection. Surprisingly, adding organic errors on top leads to mixed results, showing that organic data is not necessarily always helpful or needed for good classifier performance. The **+Dev** scenario might already be large enough that the addition of organic errors does not provide benefit.

## 5.3. Word-Level Error Detection

In practice, it may also be useful to locate specific response words that are hallucinated, or specific input words that are uncovered. Of the methods in §3, ROUGE and saliency both produce hallucination and coverage error scores at the word level. Specifically, the ROUGE coverage error score would be 0 if an input word is matched in the response (and 1 otherwise), and the ROUGE hallucination score would be 0 if a response word is matched in the input arguments (and 1 otherwise). For saliency, before example-level aggregation, scores are already computed per word (§3.2). Sample word-level saliency scores for hallucination and coverage errors are shown in Figure 2.

We compare the word-level hallucination and coverage error scores from saliency and ROUGE with the ground truth annotations of hallucinated and uncovered words annotated in our test sets (§4.1). Results are computed over all non-stop words in each test set, defining words by merging LLM tokens (§3.2); we compute results over all response words for hallucination word detection, and over all input words for coverage error word detection. Results are reported in Table 3. Differences between saliency and ROUGE are statistically significant for all test sets and error types ( $p < 0.001$ ) except the ambiguous organic error set, using the Wilcoxon statistic corrected for multiple comparisons as in §5.1. Saliency performs equally to or better than ROUGE for detecting both hallucinated words in model responses and uncovered words in input arguments on all test sets. On the test set with paraphrased synthetic errors, saliency has the largest relative gains over ROUGE, likely due to its ability to capture semantics even in cases of word mismatch, similar to the trends for example-level error detection.

## 6. Discussion

Overall, LLM-based classifiers trained on relatively small amounts of data perform surprisingly well, outperforming all other methods detecting full errors and obtaining promising ROC AUC scores between 90% and 99%. This is especially notable given that the classifiers are trained only on *synthetic* hallucination and coverage errors and yet perform well on the organic test set.

While worse than the classifiers, the data-free methods presented here still achieve strong results. Our experiments show that ROUGE is a strong data-free baseline for hallucination and coverage error detection in tasks with minimal paraphrasing. When more paraphrasing is expected, salience provides stronger results, appearing to better capture semantics than simple word matching. Moreover, salience is effective for word-level hallucination and coverage error detection, allowing us to locate the parts of a generated response that are problematic.

Our experiments also show the value of different test set slices. While the synthetically constructed datasets might diverge from the true data distribution, they offer a way to analyze strengths and weaknesses of different methods in an isolated fashion, e.g. paraphrased examples demonstrating the shortcomings of ROUGE.

Finally, all methods struggle on ambiguous organic errors, although these results are inconclusive. Largely, this set is a “catch-all” for problematic and low-agreement errors, possibly explaining the poor performance of different error detection approaches. Training classifiers on this subset is important future work, but requires a larger dataset of more clearly-defined ambiguous errors.

## 7. Related Work

**Errors in controlled text generation.** Our approach to NPOV Response Generation using provided input perspectives is an example of *retrieval augmented generation*, where information (e.g. a document or paragraph) is retrieved from a knowledge source (e.g. a search engine) and used to condition a model response (Li et al., 2022). Like in our scenario, retrieval-augmented models sometimes exhibit hallucinations (Dziri et al., 2022) and coverage errors (Krishna et al., 2021) relative to the retrieved source. Our error detection methods may be applied to these scenarios more generally.

Specifically, our task is closely related to *table-to-text generation*, which aims to generate fluent and faithful natural language descriptions of tabular data. Table-to-text generation has been studied using a variety of datasets, including WikiBio (Lébreton et al., 2016), ToTTo (Parikh et al., 2020), DART (Nan et al., 2021), and WebNLG (Gardent et al.,

2017). Traditional metrics such as ROUGE, BLEU, and METEOR compare model responses to a reference output, but metrics developed specifically for table-to-text tasks (e.g. PARENT; Dhingra et al., 2019) often consider both the table source and reference output when scoring a model response, to better preserve faithfulness to the source (Liu et al., 2021; Thomson and Reiter, 2021). Our work similarly compares model responses to the input source; however, our input fields are perspectives composed of several full sentences (arguments) rather than short expressions (e.g. entities or numbers that allow minimal paraphrasing, as in most table-to-text tasks). For this reason, pure matching-based scoring approaches (e.g. ROUGE, BLEU, and PARENT) are less effective for our task.

More broadly, hallucinations are a common artifact in natural language generation (NLG). At a high level, they can be described as cases where generated output is “unfaithful” to provided or desired source content (Ji et al., 2023). Due to the fluency of modern NLG systems, hallucinations can remain undetected and mislead users. Tolerance to such errors is particularly low in summarization and table-to-text tasks, where a retrieved source is provided. In the NPOV Response Task, we focus on *full* errors, where a hallucinated or uncovered argument can be identified relatively unambiguously.

**Prompt-tuning.** Both the NPOV Response Generator (§2.1) and the classifiers (§3.3) use soft prompt-tuning, a method where only a small number of parameters are tuned and the base LLM is left unchanged (Lester et al., 2021). Mozes et al. (2023) show that LLMs can be prompt-tuned even on very small datasets to function as classifiers. Open-source code to train such classifiers is available through the Gemma Responsible Generative AI Toolkit (Google, 2024).

**Salience.** Previous work has identified hallucinations in machine translation using proportions of source contributions to output tokens (Dale et al., 2023; Voita et al., 2021), using aggregated layerwise token attribution (Ferrando et al., 2022). Our salience-based method for error detection is similar, but attributions are based on loss gradients (Bastings and Filippova, 2020). We focus on dot products between gradients and inputs, which are often used to roughly quantify model attributions from input tokens (Ding and Koehn, 2021; Boggust et al., 2023; Zhao et al., 2022).<sup>11</sup> Previous work has applied gradient-based salience methods to fine-tuned encoder-decoder and encoder-only classification models (Tenney et al., 2020). We extend this to decoder-only models, prompt-tuned on sequence-to-sequence tasks.

<sup>11</sup>We obtain similar results using gradient L2 norms.



## 8. Limitations

Our work has several limitations. The NPOV Response Generator is trained and evaluated only in English, and our NPOV Response Task does not address how to create the content in the perspectives and their arguments. The arguments used in our work are pulled from ProCon, which limits both our set of controversial topics and our sets of perspectives (i.e. only pro and con; see Ethical Considerations); future work might consider more nuanced methods of perspective identification, selection, and/or generation.

Our work also does not focus on biases in LLM hallucinated or omitted content. For example, the NPOV Response Generator may be more likely to hallucinate or omit arguments for specific topics or perspectives, e.g. based on the frequency of topics and perspectives in the LLM pre-training corpus (Durmus et al., 2023). Even when focusing just on error detection rather than error content, we focus primarily on errors that are easy to identify and have high levels of inter-annotator agreement. Based on our own annotations, inter-annotator agreement on ambiguous errors is much lower than for full errors. The majority of ambiguous errors that we observed can be classified as partial errors, repetition, or argument confusion (§3), but an important branch of future work is to establish more thorough taxonomies and annotation schemes for hallucination and coverage error types.

Finally, in future work we hope to evaluate whether our findings on LLM-based classifier performance generalize to other (ideally publicly available) LLMs. Many significant results involving LLMs have generalized to other LLMs (e.g. in-context learning, chain-of-thought reasoning, and parameter-efficient tuning methods; Brown et al., 2020; Wei et al., 2022; Lester et al., 2021), but our results should be verified for other LLMs.

## 9. Conclusion

In this paper, we introduce the NPOV Response Task as an approach to retrieval augmented generation for controversial topics. We focus on response generation, after pro and con arguments are provided to an LLM. We propose and evaluate methods for detecting hallucination and coverage errors in LLM-generated responses, and we demonstrate a synthetic error generation strategy that can be used to train and evaluate our proposed methods. We find that prompt-tuned LLM classifiers trained only on synthetic errors achieve high error detection performance on organic examples. Our other methods, while performing worse than our classifiers, still achieve strong results without the need for training data.

## Ethical Considerations

With the rise of LLM-based chatbots and broader societal concerns about echo chambers, filter bubbles, and polarization, the ability of LLMs to provide neutral, factual, and nuanced responses to controversial topics is an important avenue of work. However, having LLMs respond to queries about controversial topics is inherently challenging: who decides what is controversial, neutral, and factual, and how this is encoded in an LLM is a hard and nebulous problem. Moreover, as LLMs and chatbot technologies become increasingly easy to create, maliciously engineered and maliciously applied models are likely to become more prevalent. Retrieval augmented generation is a way to control LLM responses in a maximally transparent way.

In this paper, we assume the existence of a database with NPOV-expressed perspectives. However, such a database is not an easy artifact to create, and the contents will often be hotly contested. The dataset we use is derived from Britannica’s ProCon website (ProCon.org, 2022). However, this still reduces arguments to pro and con perspectives, which can reinforce a binary vision of the world. Our work also does not address how to best arrive at and reflect consensus on specific arguments. For example, when should the model express “many experts” vs. “a few experts” as a qualification for an argument? Failure here can serve to elevate fringe arguments. Even deciding whether a topic is controversial is already culturally charged. For instance, the subject of gun control might be a non-issue for some European countries yet remain polarizing in the United States. Similarly, omitting topics or arguments that are relevant for minorities or non-Western countries risks reinforcing systemic erasure and promoting socio-cultural biases. To address and mitigate these biases in a perspectives database, processes are necessary to ensure that the group of experts providing perspectives is diverse and multicultural.

The more basic question of when to apply an LLM in practical scenarios needs careful consideration. In some domains (e.g. medical information), even very low error rates may not be acceptable, while other domains (e.g. creative writing) have very different risk profiles. Proper evaluations, policies, and guardrails should be put in place before LLMs are applied in practice to new domains.

Finally, the computational footprints of the NPOV Response Generator and the LLM-based error classifiers are large, with each model built upon a 60B+ parameter LLM. Similarly, computing salience maps for error detection requires computing gradients from the NPOV Response Generator itself, thus inducing a large computational cost. Of the error detection methods evaluated in our work,

ROUGE is by far the most computationally efficient. Future work may consider more computationally efficient approaches, such as evaluating smaller models as error detection classifiers.

## Acknowledgements

We would like to thank Ian Tenney, Jasmijn Bastings, Vinodkumar Prabhakaran, and the anonymous reviewers for valuable feedback.

## 10. Bibliographical References

- Microsoft Azure. 2023. [Retrieval augmented generation \(RAG\) in Azure Cognitive Search](#). *Microsoft Azure Documentation*.
- Jasmijn Bastings and Katja Filippova. 2020. [The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?](#) In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 149–155, Online. Association for Computational Linguistics.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. [Natural language processing with Python: Analyzing text with the natural language toolkit](#). O’Reilly Media.
- Angie Boggust, Harini Suresh, Hendrik Strobelt, John V Gutttag, and Arvind Satyanarayan. 2023. [Beyond faithfulness: A framework to characterize and compare saliency methods](#). In *Proceedings of the ACM Conference on Fairness, Accountability, and Transparency (FAccT)*.
- C.E. Bonferroni. 1936. *Teoria statistica delle classi e calcolo delle probabilità*. Pubblicazioni del R. Istituto superiore di scienze economiche e commerciali di Firenze. Seeber.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- Tyler A. Chang and Benjamin K. Bergen. 2024. [Language model behavior: A comprehensive survey](#). *Computational Linguistics*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Aleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2023. [PaLM: Scaling language modeling with Pathways](#). *Journal of Machine Learning Research*, 24(240):1–113.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. [Scaling instruction-finetuned language models](#). *arXiv preprint*.
- Aaron Daniel Cohen, Adam Roberts, Alejandra Molina, Alena Butryna, Alicia Jin, Apoorv Kulkshreshtha, Ben Hutchinson, Ben Zevenbergen, Blaise Hilary Aguera-Arcas, Chung ching Chang, Claire Cui, Cosmo Du, Daniel De Freitas Adwardana, Dehao Chen, Dmitry (Dima) Lepikhin, Ed H. Chi, Erin Hoffman-John, Heng-Tze Cheng, Hongrae Lee, Igor Krivokon, James Qin, Jamie Hall, Joe Fenton, Johnny Soraker, Kathy Meier-Hellstern, Kristen Olson, Lora Mois Aroyo, Maarten Paul Bosma, Marc Joseph Pickett, Marcelo Amorim Menegali, Marian Croak, Mark Díaz, Matthew Lamm, Maxim Krikun, Meredith Ringel Morris, Noam Shazeer, Quoc V. Le, Rachel Bernstein, Ravi Rajakumar, Ray Kurzweil, Romal Thoppilan, Steven Zheng, Taylor Bos, Toju Duke, Tulsee Doshi, Vincent Y. Zhao, Vinodkumar Prabhakaran, Will Rusch, YaGuang Li, Yanping Huang, Yanqi Zhou, Yuanzhong Xu, and Zhifeng Chen. 2022. [LaMDA: Language models for dialog applications](#). *arXiv preprint*.
- David Dale, Elena Voita, Loic Barrault, and Marta R. Costa-jussà. 2023. [Detecting and mitigating hallucinations in machine translation: Model internal workings alone do well, sentence similarity even better](#). In *Proceedings of the 61st Annual Meeting*

- of the Association for Computational Linguistics (Volume 1: Long Papers), pages 36–50.
- Misha Denil, Alban Demiraj, and Nando De Freitas. 2014. [Extraction of salient sentences from labelled documents](#). *arXiv preprint*.
- Bhuwan Dhingra, Manaal Faruqui, Ankur Parikh, Ming-Wei Chang, Dipanjan Das, and William Cohen. 2019. [Handling divergent reference texts when evaluating table-to-text generation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4884–4895. Association for Computational Linguistics.
- Shuoyang Ding and Philipp Koehn. 2021. [Evaluating saliency methods for neural language models](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5034–5052, Online. Association for Computational Linguistics.
- Esin Durmus, Karina Nyugen, Thomas I. Liao, Nicholas Schiefer, Amanda Askeel, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, Liane Lovitt, Sam McCandlish, Orowa Sikder, Alex Tamkin, Janel Thamkul, Jared Kaplan, Jack Clark, and Deep Ganguli. 2023. [Towards measuring the representation of subjective global opinions in language models](#). *arXiv preprint*.
- Nouha Dziri, Sivan Milton, Mo Yu, Osmar Zaiane, and Siva Reddy. 2022. [On the origin of hallucinations in conversational models: Is it the datasets or the models?](#) In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5271–5285.
- Javier Ferrando, Gerard I. Gállego, Belen Alas truey, Carlos Escolano, and Marta R. Costajussà. 2022. [Towards opening the black box of neural machine translation: Source and target interpretations of the transformer](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8756–8769.
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. [The WebNLG challenge: Generating text from RDF data](#). In *Proceedings of the 10th International Conference on Natural Language Generation*.
- Google. 2024. [Gemma: Responsible generative AI toolkit](#). *Google AI for Developers*.
- James A Hanley and Barbara J McNeil. 1983. [A method of comparing the areas under receiver operating characteristic curves derived from the same cases](#). *Radiology*, 148(3):839–843.
- Anand Iyer and Rajesh Thallam. 2023. [Building generative AI applications made easy with Vertex AI PaLM API and LangChain](#). *Google Cloud Blog, AI and Machine Learning*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Comput. Surv.*, 55(12).
- Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. [Hurdles to progress in long-form question answering](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4940–4957.
- Rémi Lebret, David Grangier, and Michael Auli. 2016. [Neural text generation from structured data with application to the biography domain](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1203–1213. Association for Computational Linguistics.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059. Association for Computational Linguistics.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Proceedings of the 34th International Conference on Neural Information Processing Systems*.
- Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022. [A survey on retrieval-augmented text generation](#). *arXiv preprint*.
- Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016. [Visualizing and understanding neural models in NLP](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 681–691, San Diego, California. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81. Association for Computational Linguistics.

- Tianyu Liu, Xin Zheng, Baobao Chang, and Zhifang Sui. 2021. [Towards faithfulness in open domain table-to-text generation from an entity-centric view](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(15):13415–13423.
- Maximilian Mozes, Jessica Hoffmann, Katrin Tomanek, Muhamed Kouate, Nithum Thain, Ann Yuan, Tolga Bolukbasi, and Lucas Dixon. 2023. [Towards agile text classifiers for everyone](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 400–414, Singapore. Association for Computational Linguistics.
- Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xiangu Tang, Aadit Vyas, Neha Verma, Pranav Krishna, Yangxiaokang Liu, Nadia Irwanto, Jessica Pan, Faiaz Rahman, Ahmad Zaidi, Mutethia Mutuma, Yasin Tarabar, Ankit Gupta, Tao Yu, Yi Chern Tan, Xi Victoria Lin, Caiming Xiong, Richard Socher, and Nazneen Fatema Rajani. 2021. [DART: Open-domain structured data record to text generation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 432–447, Online. Association for Computational Linguistics.
- Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. [ToTTo: A controlled table-to-text generation dataset](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1173–1186, Online. Association for Computational Linguistics.
- ProCon.org. 2022. <https://www.procon.org/>. Accessed: 2022-10-12.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. [The woman worked as a babysitter: On biases in language generation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. [Retrieval augmentation reduces hallucination in conversation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3784–3803.
- Ian Tenney, James Wexler, Jasmijn Bastings, Tolga Bolukbasi, Andy Coenen, Sebastian Gehrmann, Ellen Jiang, Mahima Pushkarna, Carey Radebaugh, Emily Reif, and Ann Yuan. 2020. [The language interpretability tool: Extensible, interactive visualizations and analysis for NLP models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 107–118, Online. Association for Computational Linguistics.
- Craig Thomson and Ehud Reiter. 2021. [Generation challenges: Results of the accuracy evaluation shared task](#). In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 240–248. Association for Computational Linguistics.
- Tyler J. VanderWeele and Maya B. Mathur. 2019. [Some desirable properties of the bonferroni correction: Is the bonferroni correction really so bad?](#) *American Journal of Epidemiology*, 188(3):617–618.
- Elena Voita, Rico Sennrich, and Ivan Titov. 2021. [Analyzing the source and target contributions to predictions in neural machine translation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1126–1140.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Proceedings of the 36th International Conference on Neural Information Processing Systems*.
- Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2017. [Ex machina: Personal attacks seen at scale](#). In *Proceedings of the 26th International Conference on World Wide Web, WWW '17*, page 1391–1399, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [BERTScore: Evaluating text generation with BERT](#). In *International Conference on Learning Representations*.
- Yang Zhao, Zhang Yuanzhe, Jiang Zhongtao, Ju Yiming, Zhao Jun, and Liu Kang. 2022. [Can we really trust explanations? Evaluating the stability of feature attribution explanation methods via adversarial attack](#). In *Proceedings of the 21st Chinese National Conference on Computational Linguistics*, pages 932–944, Nanchang, China. Chinese Information Processing Society of China.

## Appendices

### A. ProCon Dataset Details

We use the perspectives and arguments for the different topics listed on Britannica’s ProCon website as of October 2022 (ProCon.org, 2022). We randomly split the 72 ProCon topics into train, dev, and test, as shown in Table 4, ensuring no overlap in topics across these splits. In line with ProCon’s usage guidelines, all arguments are used verbatim as stated on the specific topic website under the section “Pro & Con Arguments”. We scrape the subtitles of the pro and con columns as our arguments. The median number of arguments per pro and con perspective per topic is 4, with a maximum of 23 and a minimum of 2. The ProCon data is publicly available through their website, containing no personally-identifying information about individuals. We follow the guidelines specified by ProCon on “How to Use” their data (<https://www.procon.org/faqs/#II>).

### B. Prompt-Tuning Details

This section discusses implementation details of (1) the NPOV Response Generator and (2) the hallucination and coverage error classifiers, which are both based on prompt-tuning an LLM. We use the same prompt-tuning settings for both.

We deliberately refrain from resource-intense hyperparameter tuning and instead use configurations previously shown to work well (Mozes et al., 2023): we use soft prompt lengths of 5 tokens initialized with a random sample of the model’s 5K most frequent token vocabulary embeddings (Lester et al., 2021); we then train with a learning rate of 0.1 with 500 warm-up steps and linear decay, using small batch sizes of 16 for training and limiting training to 20K steps. In most cases, we reach the maximum development set performance after 2-5K steps. Prompt-tuning runs take a maximum of 4 hours per run on 64 TPUv4 chips.

For the task representations, we utilize a “curly braces format” to verbalize the task, consisting of several key-value pairs in the input and target sequence for the LLM. This format is easily picked up by modern LLMs, as they have typically been exposed to code during pre-training. Figure 3 shows how we format the task for the NPOV Response Task (§2.1). Figure 4 shows how we format the error classification tasks (§3.3).

#### B.1. Classifier Ablation: Annotation-Free Scenario

As an additional experiment, we analyze whether we can obtain good classifiers for hallucination and

coverage error detection by just re-utilizing the original training data from the NPOV Response Task, without the need to perform any of the manual annotations described in §4.1. We turn the data used to train the NPOV Response Generator into error classifier training data by (1) treating NPOV Response Task training examples as no error-examples, and (2) adding synthetic errors according to our procedure in §4.2. We call this approach “annotation-free” because we do not have to obtain any additional human annotations for classifier training. The resulting hallucination and coverage error classifiers are trained on 50 error-free examples and 131 examples with synthetic errors.

Table 5 shows results on the organic test sets for the “annotation-free” classifiers. Overall, these results are significantly worse than results with the non-annotation-free classifiers (compare to Table 2), and often worse than other data-free approaches (compare to ROUGE and salience in Table 1). This suggests that error classifier training may require organic model responses, even if the errors are synthetically generated.

### C. Salience Formulas

In §3.2, we describe how we compute a word-to-word salience map  $M_{\text{words}} \in \mathbb{R}^{m \times n}$ , where  $m$  is the number of non-stop words in the input arguments and  $n$  is the number of non-stop words in the generated NPOV response. Our salience maps are based on gradient times input attribution scores, but we obtain comparable results using gradient L2 norms. Here, we include formal equations defining our hallucination and coverage error detection metrics based on  $M_{\text{words}}$ .

Assume  $I_{\text{pro}}$  and  $I_{\text{con}}$  are the lists of non-stop words in the input pro and con arguments respectively. Assume  $O_{\text{resp}}$  is the list of non-stop words in the generated NPOV main response. For each input word  $w_i \in I_{\text{pro}} \cup I_{\text{con}}$ , we define its contribution score  $\alpha_i$  as its maximum contribution to any response word (i.e. the maximum across the corresponding row of  $M_{\text{words}}$ ):

$$\alpha_i = \max(M_{\text{words}}[i, :]) \quad (1)$$

Similarly, for each output word  $w_j \in O_{\text{resp}}$ , we define its attribution score  $\beta_j$  as its maximum attribution from any input argument word (i.e. the maximum across the corresponding column of  $M_{\text{words}}$ ):

$$\beta_j = \max(M_{\text{words}}[:, j]) \quad (2)$$

Sample contribution and attribution scores for input words and response words respectively are shown in Figure 2. For word-level error detection (§5.3), these word-level scores can be converted into coverage error scores  $1.0 - \alpha_i$  and hallucination scores  $1.0 - \beta_j$ .

Split	# of topics	Topics
Train	9	<i>Animal Dissection; Concealed Handguns; Cuba Embargo; Filibuster; Free College; GMOs (Genetically Modified Organisms); Net Neutrality; Obesity; Vaping E-Cigarettes</i>
Dev	28	<i>Binge-Watching; Cancel Culture; Churches and Taxes; College Education; Corporal Punishment; Daylight Saving Time; Dress Codes; Electoral College; Employer Vaccine Mandates; Fighting in Hockey; Golf; Homework; Kneeling during National Anthem; Marijuana (CBD) for Pets; Olympics; Penny; Pit Bull Bans; Pokémon; School Vouchers; Space Colonization; Standardized Tests; Student Loan Debt; Tablets vs. Textbooks; Teacher Tenure; Uber &amp; Lyft; US Supreme Court Packing; Video Games and Violence; Zoos</i>
Test	35	<i>Abortion; American Socialism; Animal Testing; Artificial Intelligence; Banned Books; Bottled Water Ban; Cell Phone Radiation; Climate Change; Corporate Tax Rate; DACA &amp; Dreamers; DC and Puerto Rico Statehood; Defund the Police; Drone Strikes Overseas; Fracking; Gold Standard; Gun Control; Historic Statue Removal; Mandatory National Service; Minimum Wage; OTC Birth Control; Paying College Athletes; Police Body Cameras; Prescription Drug Costs; Private Prisons; Recreational Marijuana Legalization; Reparations for Slavery; Right to Health Care; Sanctuary Cities; Saturday Halloween; School Uniforms; Social Media; Social Security Privatization; Universal Basic Income; Vaccines for Kids; Vegetarianism</i>

Table 4: ProCon topics assigned to the different dataset splits.

<p><b>Input Sequence:</b>  <b>User question:</b> {Should abortion be legal?}  <b>Topic:</b> {abortion}  <b>Perspective #1:</b> {pro: Abortion bans deny bodily autonomy, creating wide-ranging repercussions. pro: Abortion is a safe medical procedure that protects lives.}  <b>Perspective #2:</b> {con: Life begins at conception, making abortion murder. con: Increased access to birth control, health insurance, and sexual education would make abortion unnecessary.}  <b>Neutral response opening:</b> {</p> <p><b>Target Sequence:</b>  This is a controversial topic and people hold different perspectives on it. Let me show a few...}  <b>Neutral response core:</b> {People in support of abortion claim that abortion bans deny bodily autonomy, which has wide-ranging repercussions. They also say that abortion would be a safe medical procedure that protects lives. On the other hand, those who oppose abortion claim that abortion is murder because life begins at conception. Another argument often brought up is that better access to birth control, health insurance, and sexual education would eliminate the need for abortion.}</p>
--

Figure 3: Task format for the NPOV Response Task.

Test Set	Hallucination	Coverage
Full organic	0.739	0.896
Ambiguous org.	0.732	0.804

Table 5: Annotation-free classifier error detection ROC AUC scores.

For example-level error detection (§5.1), we compute an example-level coverage error score by (1) taking the geometric mean of word-level contribution scores for each input perspective, (2) taking the minimum of the two perspective scores (to reflect

the fact that both perspectives must contribute), and (3) subtracting from 1.0 (lower contributions are more likely to be coverage errors):

$$s_{\text{cov}} = 1.0 - \min \left( \text{gmean}_{w_i \in I_{\text{pro}}} (\alpha_i), \text{gmean}_{w_i \in I_{\text{con}}} (\alpha_i) \right)$$

We compute an example-level hallucination score by (1) taking the geometric mean of word-level attribution scores in the NPOV main response, and (2) subtracting from 1.0 (lower attributions are more

<p><b>Coverage Detection</b>  <b>Input Sequence</b>  <b>User question:</b> {Should abortion be legal?}.  <b>Neutral response:</b> {This is a controversial topic and people hold different perspectives on it. Let me show a few... People in support of abortion claim that abortion bans deny bodily autonomy, which has wide-ranging repercussions. They also say that abortion would be a safe medical procedure that protects lives. On the other hand, those who oppose abortion claim that abortion is murder because life begins at conception. Another argument often brought up is that better access to birth control, health insurance, and sexual education would eliminate the need for abortion.}.  <b>Given pro arguments:</b> {pro: Abortion is a safe medical procedure that protects lives. pro: Abortion bans deny bodily autonomy, creating wide-ranging repercussions.}.  <b>Given con arguments:</b> {con: Increased access to birth control, health insurance, and sexual education would make abortion unnecessary. con: Life begins at conception, making abortion murder.}.  <b>All the given arguments are covered by the neutral response:</b> {  <b>Target Sequence</b>  YES}</p>	<p><b>Hallucination Detection</b>  <b>Input Sequence</b>  <b>User question:</b> {Should abortion be legal?}.  <b>Neutral response:</b> {This is a controversial topic and people hold different perspectives on it. Let me show a few... People in support of abortion claim that abortion bans deny bodily autonomy, which has wide-ranging repercussions. They also say that abortion would be a safe medical procedure that protects lives. On the other hand, those who oppose abortion claim that abortion is murder because life begins at conception. Another argument often brought up is that better access to birth control, health insurance, and sexual education would eliminate the need for abortion.}.  <b>Given pro arguments:</b> {pro: Abortion is a safe medical procedure that protects lives. pro: Abortion bans deny bodily autonomy, creating wide-ranging repercussions.}.  <b>Given con arguments:</b> {con: Increased access to birth control, health insurance, and sexual education would make abortion unnecessary. con: Life begins at conception, making abortion murder.}.  <b>Only given arguments are contained in the neutral response:</b> {  <b>Target Sequence</b>  YES}</p>
---	---

Figure 4: Task format for LLM-based error classifiers.

likely to be hallucinations):

$$s_{\text{hall}} = 1.0 - \text{gmean}_{w_j \in O_{\text{resp}}}(\beta_j)$$

Note that  $s_{\text{cov}}, s_{\text{hall}} \in [0, 1]$  because entries of  $M_{\text{words}}$  are in  $[0, 1]$ . We evaluate these hallucination and coverage error scores for example-level error detection in §5.1.

### C.1. Alternative Salience Aggregation Methods

Above, we use the maximum function (in Equations 1 and 2) to aggregate a contribution score for each input word and an attribution score for each response word. This is based on the observation that covered input words tend to have a high contribution to at least one response word, and non-hallucinated response words tend to have a high attribution from at least one input argument word. In Table 6, we report ROC AUC results on the full organic test set using different methods to aggregate word contributions and attributions in Equations 1 and 2. Specifically, we consider (1) the sum (i.e. the sum over all response contributions for each input word to quantify coverage, and the sum over all input attributions for each response word to quantify non-hallucination), and (2) the (negative) entropy. Lower entropies indicate less distributed contributions/attribution, such as when most of the contribution/attribution is to/from a single word (a pattern which appears in the majority of covered and non-hallucinated words).

We find that entropies perform worse than the maximum and sum aggregation functions for both hallucination and coverage error detection. The sum performs best for hallucination detection (summing input attributions for each response word), but the maximum performs best for coverage error detection (taking the maximum response contribution for each input word). We use the maximum in the

Aggregation	Hallucination	Coverage
Max	0.808	<b>0.852</b>
Sum	<b>0.846</b>	0.809
Negative entropy	0.786	0.664

Table 6: Example-level error detection ROC AUC scores for salience using different methods to aggregate a contribution score for each input word (coverage) and an attribution score for each response word (hallucination).

main results for consistency and to avoid overfitting to the test set.

### D. Human Annotation Details

For the human annotations in §4.1.1, our annotation service provider was paid 49 USD per hour for a total of 25 hours of work; they state that they ensure fair payment to annotators. The 10 annotators were specialized workers in the United States contracted by our annotation provider. Our annotation provider reported self-disclosed genders and age brackets of annotators, but this information was not used in our analyses. Our annotations focused on attributes of our NPOV Response Generator query-response pairs, collecting annotation labels but no other data generated by the annotators. To reduce annotation bias, annotators were not told how the labeled examples would be used, and they were not told that the response was machine-generated.