

Plot-guided Adversarial Example Construction for Evaluating Open-domain Story Generation

Sarik Ghazarian,¹ Zixi Liu,¹ Akash SM,²

Ralph Weischedel,¹ Aram Galstyan,¹ Nanyun Peng^{1,3}

¹University of Southern California / Information Sciences Institute

²Indian Institute of Technology Roorkee

³Computer Science Department of University of California, Los Angeles

{sarik, zixiliu, weisched, galstyan}@isi.edu, akashsm@ce.iitr.ac.in, violetpeng@cs.ucla.edu

Abstract

With the recent advances of open-domain story generation, the lack of reliable automatic evaluation metrics becomes an increasingly imperative issue that hinders the fast development of story generation. According to conducted researches in this regard, learnable evaluation metrics have promised more accurate assessments by having higher correlations with human judgments. A critical bottleneck of obtaining a reliable learnable evaluation metric is the lack of high-quality training data for classifiers to efficiently distinguish plausible and implausible machine-generated stories. Previous works relied on *heuristically manipulated* plausible examples to mimic possible system drawbacks such as repetition, contradiction, or irrelevant content in the text level, which can be *unnatural* and *oversimplify* the characteristics of implausible machine-generated stories. We propose to tackle these issues by generating a more comprehensive set of implausible stories using *plots*, which are structured representations of controllable factors used to generate stories. Since these plots are compact and structured, it is easier to manipulate them to generate text with targeted undesirable properties, while at the same time maintain the grammatical correctness and naturalness of the generated sentences. To improve the quality of generated implausible stories, we further apply the adversarial filtering procedure presented by Zellers et al. (2018) to select a more nuanced set of implausible texts. Experiments show that the evaluation metrics trained on our generated data result in more reliable automatic assessments that correlate remarkably better with human judgments compared to the baselines.

1 Introduction

The surge of downstream applications for open-domain natural language generation (NLG), such as dialog systems (Zhang et al., 2020) and story

Human Written Story: jenny liked fresh fish. she decided to go fishing to catch her own. she brought her worms and pole and a chair. she sat there all day but didn't catch anything. she packed it up and went home disappointed.	
Sentence Manipulation: jenny liked fresh fish. she decided to go fishing to catch her own. she wrote songs every single day. she sat there all day but didn't catch anything. she packed it up and went home disappointed.	
Keyword Manipulation: jenny liked fresh fish. she decided to go fishing to catch her own. she brought her worms and pole and a chair. she sat there all day but didn't catch anything. she unpacked it up and went home disappointed.	
UNION: jenny liked fresh fish. jim has a very structured workout program to help him achieve goals. she brought her worms and pole and a relaxer. she sat there all day but didn't catch anything. she unpack it up and went home disappointed.	
Plot: jenny fresh fish -> decided fishing catch -> brought worms chair -> sat -> packed home disappointed	Manipulated Plot: jenny fresh fish -> tasha offered woman store -> brought worms chair -> sat -> got wet packed home disappointed
Manipulated Plot Guided Generation (Ours): jenny was out of fresh fish. tasha offered to buy her some from the woman at the store. she brought her worms and a chair and decided to play with them. jenny sat down and laid down on the chair. when she got wet, she packed up and went home disappointed.	

Figure 1: Heuristically generated implausible stories (the second block) for a given human-written story (the first block) using sentence, keyword and UNION manipulations versus injecting implausible sources into the story plot (the third block, from the left plot to the right one) and generating a more natural implausible story (the last story). Blue highlights show the implausible sections.

generators (Rashkin et al., 2020a) necessitates *automatic* evaluation metrics for quality assessment. The existence of accurate automatic evaluation metrics can accelerate the development cycle by facilitating the process of model comparison and hyperparameter search. Many existing reference-based approaches such as BLEU (Papineni et al., 2002) or ROUGE (Lin, 2004) fail to correlate well with human judgment in open-domain settings due to the fact that there can be potentially many plausible generations that do not have significant overlap with the limited set of given references. This failure invites research on more sophisticated and reliable evaluation metrics.

Recently, learning-based approaches have been proposed to overcome this limitation by training classifiers to distinguish between plausible and implausible texts (Li and Jurafsky, 2016; Holtzman et al., 2018). The choice of training data for learning such classifiers is a key determinant of the metric effectiveness. Existing works take human-written texts as plausible (positive) examples, while

the negative samples are heuristically generated by randomly substituting keywords or sentences (See Figure 1) (Li and Jurafsky, 2016; Guan and Huang, 2020). Guan and Huang (2020) further improved the quality of evaluators by applying heuristic rules such as adding repetition, reordering and negation (See the UNION story in Figure 1).

In this work, we hypothesize that heuristically generated data cannot adequately reflect the characteristics of the implausible texts generated by language models, thus result in suboptimal trained evaluation metrics. This deficiency can be mitigated by generating high-quality implausible examples that are closer to the test data. Toward this goal, we propose an approach based on the manipulation of *plots*, which are high-level structured representations of generated texts originally used as a content-planning tool for better text generation (Fan et al., 2019; Goldfarb-Tarrant et al., 2020). Specifically, we propose to manipulate plots by injecting incoherence sources into them. The generation models conditioned on such manipulated plots lead to implausible texts that have pertinent similarities with implausible machine-generated texts and thus can serve as good negative examples for training evaluation metrics.

We further improve the quality of training data by incorporating the adversarial filtering technique proposed by Zellers et al. (2018) to select more challenging negative samples generated from the manipulated plots (See Figure 1). Eventually, these samples result in more reliable evaluation metrics. The contributions of this work are four-fold:

- We study the importance of training data for learnable automatic evaluation metrics in open-domain story generation task and show the inadequacy of heuristically generated negative examples in this setting.
- We propose a novel technique to generate negative samples by introducing plot-level incoherence sources that guide generation models to produce implausible texts.
- We show the affirmative role of adversarial filtering techniques in constructing training data for learnable open-domain story generation evaluation metrics.
- We demonstrate that the evaluation metrics trained on our generated data have a significantly higher correlation with human judgments compared to strong baselines.

2 Related Work

Existing work on automatic evaluation of generation models can be classified into two subgroups, non-learning-based and learning-based methods, which we briefly summarize below.

Non-learning-based Metrics. Some metrics in this group consider the centrality of a text around a specific topic as a proxy for measuring its quality. The transitions of entities in neighbor sentences and their distribution across text have been served as a measurement for quality assessment (Mitsakaki and Kukich, 2004; Lapata and Barzilay, 2005). Perplexity is another commonly used metric to evaluate the quality of text and story generation models (Fan et al., 2018; Peng et al., 2018).

Learning-based Metrics. This group of metrics is based on neural-based classifiers trained on a set of positive (plausible) and negative (implausible) texts. The common point between these metrics is using random sentence substitution to construct training examples, while the architectures are slightly different. Li and Jurafsky (2016) trained a neural network with a sigmoid function on top of sentence embeddings extracted from LSTM. Lai and Tetreault (2018) designed SENTAVG that gets the sentence vectors from LSTM, takes the average of these vectors to represent the whole text, and then passes it through a hidden layer.

Recently, Guan and Huang (2020) proposed a more accurate automatic evaluation metric called UNION. This metric achieved better performance by using BERT (Devlin et al., 2019) as a more effective classification model and have a broader set of negative samples coming from different heuristics. For all learning-based metrics, the simplicity of heuristically generated data samples makes them inadequate for an accurate evaluation of plausibility in open-domain generated texts.

3 Implausible Text Construction

We formulate the evaluation of open-domain story generation as a binary classification task where the goal is to distinguish plausible and implausible generated stories, also referred to as positive and negative examples. Clearly, the availability of high-quality positive and negative examples is essential for training reliable and generalizable metrics. While human-generated stories can be considered as positive examples, what constitutes good negative examples is a non-trivial question. Specifically, consider a hypothetical decision boundary

that separates positive and negative stories. While any point on one side of the boundary will be a negative example, intuitively we want examples that are not too far away from that boundary. To achieve this, we will start from positive examples, and modify them in a controllable manner to generate corresponding negative samples.

3.1 Heuristic Negative Samples

There are some widely-used approaches to heuristically manipulate positive examples and change their structure to generate negative examples.

Sentence Substitution. Sentence substitution (briefly HEUR_SENT_SUB) replaces a fraction of sentences in the plausible text with random ones (See Figure 1). This breaks the discourse-level coherence, making a story not interpretable (Li and Jurafsky, 2016; Holtzman et al., 2018).

Keyword Substitution. Guan and Huang (2020) proposed to apply random substitutions at the keyword-level (briefly HEUR_KEY_SUB), where a fraction of keywords are randomly substituted with their corresponding antonyms from a commonsense knowledge base such as ConceptNet (Speer and Havasi, 2012) to corrupt the plausibility in the text. ConceptNet consists of (*object, relation, subject*) triplets. For each selected keyword that exists as an object or subject in the ConceptNet, its counterpart is extracted from one of the contradiction-type relations; *Antonym, NotDesires, NotCapableOf, or NotHasProperty*. For instance, *packed* word in the second example of the implausible text in Figure 1 is substituted by its antonym *unpacked*.

UNION Manipulations. Alongside the keyword and sentence substitutions, Guan and Huang (2020) proposed to use repetition, reordering, and negation techniques to generate a more complete and nuanced set of implausible examples. The sentences and keywords are repeated throughout the text to reflect the repetition issue of language models. The order of sentences is changed and negation words are added to make texts implausible due to wrong causal dependencies and conflicted logic. They simultaneously apply some of these techniques to human-written texts to construct negative examples (See third negative story in Figure 1). We refer to this data as UNION_DATA. Despite the demonstrated effectiveness of UNION_DATA in open-domain story evaluation, heuristically constructed negative samples are quite far from machine-

generated texts, and thus inadequate to represent a broad set of machine-generated implausible texts.

3.2 Proposed Approach

As we stated above, applying heuristic rules at the *utterance* level result in negative examples that are usually unnatural and do not reflect the complex characteristics of machine-generated texts. Instead, we propose to introduce perturbations at a more abstract plot level. Namely, we seek to improve the quality of negative samples using plot-controlled generation with adversarial filtering techniques.

3.2.1 Plot Manipulations

Studies have shown that high-quality fluent stories can be generated by planning in advance and leveraging lucrative plots (Yao et al., 2019; Fan et al., 2019; Goldfarb-Tarrant et al., 2019, 2020; Rashkin et al., 2020b; Brahman et al., 2020). Yao et al. (2019) leverage a sequence of keywords as the plot representation (also called storyline). Fan et al. (2019) use semantic role labeling tool to extract plots as abstract presentation of stories over actions and entities. Their experiments affirm that plots have positive effects on generating high-quality stories.

Here we leverage this idea for generating implausible texts, by controllable injection of implausibility sources, or perturbations, into the ground-truth plots. The resulting plot-level manipulations will force the model to reflect applied implausibility in the generated text and will negatively impact the text’s plausibility. In contrast to Guan and Huang (2020), our proposed plot-level manipulations (MANPLTS) do not directly change the text at the token level instead, we inject incoherence into language at the concept level. The plot-guided generation guarantees the naturalness of generations since it leverages a well-trained conditional language model. The generated samples are also anticipated to be closer and congruous to the machine-generated texts that will be assessed during the inference time. Concept-level incoherence creates implausible factors that guide models to include that implausible sources. Figure 2 demonstrates various proposed plot-level manipulations in dotted boxes.¹ All proposed manipulations are described in the following sections. We refer this data as **ManPlts**.

¹Our proposed data, trained models and code is released at <https://github.com/PlusLabNLP/Plot-guided-Coherence-Evaluation>

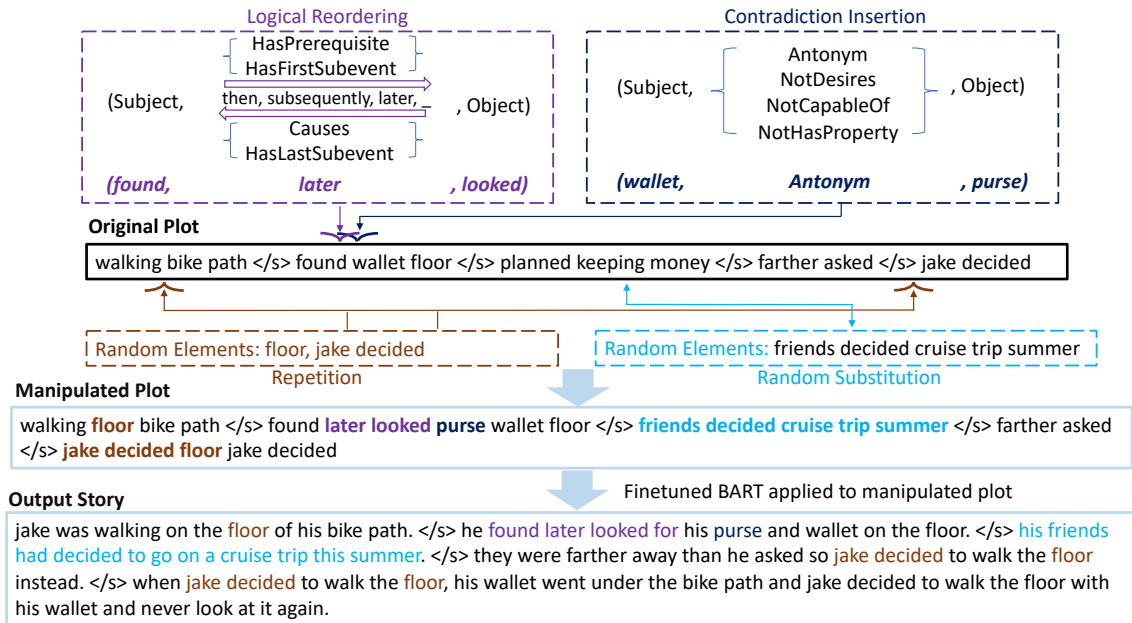


Figure 2: The plot-level manipulations applied to a human-written story’s plot (presented in the black box) to create an implausible story. Each dotted box with a specific color shows a distinct manipulation technique. The manipulated plots are passed through a generation model to generate implausible samples for the evaluation task.

Non-logically Ordered Plots. Logical conflict is one of the sources for implausibility that results from not-logically ordered concepts in the text. While Guan and Huang (2020) covered this type of implausibility by changing the order of sentences, we hypothesize that disrupting the logical order at the concept-level is more efficient. To accomplish concept reordering, we first randomly choose verbs from the plot and leverage the COMET (Bosselut et al., 2019) model to predict their subsequent events. Then we dislocate the resulted concept pairs. COMET, which is trained on tuples of the form (*subject*, *relation*, *object*), can be used to predict an *object* given a pair of *subject* and *relation*. As an example, given the pair (*work*, *causes*) COMET will predict *get pay* to show that work causes to get paid. We focus on COMET relations *HasPrerequisite*, *HasFirstSubevent*, *Causes* and *HasLastSubevent* that imply ordering. In the first two relations, object should appear before subject, while in the other two the order is reversed. Therefore, subject *work* comes before *get pay* due to the *causes* relation that holds between them. We flip the correct order of concepts and attach them with or without randomly selected connection words such as *then*, *later*, *subsequently* to generate implausible texts (the purple box in Figure 2).

Contradiction Insertion. The plausibility of a text is supported in part by the logical and correct

relationship between its words. It can be harmed by accompanying words with their antonyms or other conflicting concepts that add contradiction to the text and make it hard to grasp. In order to add such kind of implausibility, we propose to insert contradictory sources of randomly selected plots in consecutive positions. For each selected plot, we use ConceptNet (Speer and Havasi, 2012) to extract concepts that hold negation relations such as *Antonym*, *NotDesires*, *NotCapableOf*, and *NotHasProperty* with it and insert them as neighbor plots. In the navy blue box of Figure 2, *purse* has been added before *wallet* as its antonym. This guides the generation model to include consecutive contradictory elements in the generated text that harms the coherence of sections and makes it difficult to interpret.

Repetition Insertion. Repetition is one of the common issues that many generative models suffer from. Recently proposed top-k (Fan et al., 2018) and top-p (Holtzman et al., 2020) sampling techniques partially mitigated but not completely solved this issue. Guan and Huang (2020) proposed to replicate this problem in negative implausible text construction by repeating N-grams in consecutive positions. These heuristically constructed outputs only mirror local repetition issues, while the state-of-the-art generative models produce more complex and subtle repetitions throughout the whole text.

Manipulated Plots	Manipulation Type	Generated Implausible Stories
1. <i>waiting airport cousin jimmie</i> 2. <i>ended</i> practiced month	Random_Subst. Repetition	1. tobias was waiting at the airport to see his cousin jimmie. 2. although his parents ended up being unwanted, he practiced with his dad every day for a month.
3. <i>ended unwanted</i> wanted parents proud 4. ended losing	Contr_Insert. -	3. his parents ended up being unwanted, but he really wanted to make his parents proud. 4. he ended up losing to his cousin jimmie because he didn't practice enough.
5. parents ended proud	-	5. but his parents ended up being proud of him anyway.
1. finished <i>friend turned</i> assignment 2. friend	Repetition -	1. i finished my job , and my friend and i turned in our homework assignment. 2. my friend and i went to the mall.
3. <i>vendors games rides</i> 4. turned	Random_Subst. -	3. we went to vendors, played many games, and had rides. 4. when we got home, i turned in the assignment.
5. <i>went to class later</i> teacher called <i>home</i> office	Logic_Reorder. Contr_Insert.	5. when we went to class later teacher called us home from the office.
1. <i>made gina nervous perform</i> 2. <i>pained</i> ended <i>scared</i> missing	Random_Subst. Repetition Logic_Reorder.	1. gina's mom made gina get nervous about having to perform. 2. she was pained that she ended up scared of missing her bus
3. bus <i>scared allie</i> 4. one	Repetition -	3. she was on the bus and scared of allie. 4. no one was on the bus and she didn't know where they were.
5. <i>allie</i> scared	Repetition	5. therefore gina and allie were too scared to ride on the bus together.
1. billy noticed <i>billy</i> 2. <i>grandpa loved recall</i>	Repetition Random_Subst.	1. billy noticed that his buddy billy was out of gas. 2. billy's grandpa had just loved to recall a car recall he didn't want to recall.
3. <i>billy</i> finished filling <i>drove</i> 4. <i>unsuddenly</i> suddenly	Repetition Contr_Insert.	3. billy got in his car, finished filling it up, and drove away. 4. suddenly, suddenly, suddenly, billy's car was out of gas.
5. <i>billy noticed</i> billy driven	Repetition	5. billy noticed later that billy had driven off with that car recall.

Table 1: Examples of implausible stories generated based on manipulated plots. Bold italic keywords represent manipulated plots resulted from different proposed manipulations shown in the middle column.

We propose to repeat random plots of each text in various positions that would force the language model to duplicate them throughout the text and exhibit more realistic machine-generated repetitive examples. In Figure 2, the repetition of *floor* and *jake decided* compels the model to generate boring and repetitive sentences.

Random Substitution. Random sentence substitutions employed by many evaluation models amplify the implausibility sources in the text by inserting completely off-topic sentences that could potentially result in topical inconsistency throughout the text. Such scenarios are less likely for state-of-the-art high-quality generation models that use encoded context to generate tokens.

Once again, we propose to do the replacement at the plot level. Within our approach, even though the inserted random plots are completely irrelevant, the model would attempt to incorporate them into the text as much as possible by using encoded context sentences. This can be seen in the third sentence of Figure 2. Even if this sentence's plots are randomly inserted, the model is able to generate a sentence that does not have significant topical

inconsistency, thanks to the contextualized nature of the generative process.

Table 1 depicts four different machine-generated stories, each containing five sentences that are conditioned on the manipulated plots. Bold italic keywords represent manipulated plots resulted from the proposed approaches shown in the middle column.

3.2.2 Adversarial Filtering

Adversarial filtering (AF) technique was originally proposed to generate high-quality negative examples for a grounded commonsense inference task (Zellers et al., 2018). AF uses a committee of trained models to identify more appropriate negative endings from a pool of candidate samples generated for a given context. For each human-written text, there are N machine-generated endings. The goal is to select the most unbiased subset (A) of generated endings with similar stylistic features to the human-written ones.

AF starts by randomly specifying the best endings in the assignment set (A) from all N endings of each context (Zellers et al., 2018). In each iteration, the data is divided into two parts. The first part is

used for training a classifier to distinguish high/low quality endings, and the second part is used for replacing easy endings in A with adversarial endings from N . Easy endings are the ones that a trained classifier assigns a much lower score compared to human-written texts, e.g., due to their significantly different writing styles. Adversarial texts have a higher positive probability than easy texts indicating the challenge for a classifier to distinguish them from human-written texts. The replacement of easy texts with adversarial ones maximizes the empirical error of the trainable classifier. The steps outlined above are repeated till the assignment set is filled with high-quality endings for each context.

We use AF on top of the plot-based manipulations for generating implausible texts (briefly call AF_MANPLTS). Our approach for negative texts construction has two main stages: 1) generate a set of N implausible texts conditioned on manipulated plots 2) pick out the A most challenging high-quality implausible texts without stylistic biases based on applied adversarial filtering technique to increase the quality of negative samples.

4 Learnable Evaluation Models

We assess the plausibility of a text by training a classification model on the data that consists of human-written texts (positive examples) and constructed implausible stories (negative examples). Binary classifiers trained on this data can produce the probability of plausible/implausible labels for each text. The predicted probability of the positive class is interpreted as the text’s plausibility score.

4.1 Fine-tuning Language Models

The effectiveness of large pretrained language models has been proven in NLP downstream tasks (Devlin et al., 2019; Liu et al., 2019; Yang et al., 2019; Beltagy et al., 2020). RoBERTa introduced by Liu et al. (2019) is one of these models achieving impressive performances on text classification. We employ RoBERTa for our plausibility classification task. We start from pretrained RoBERTa parameters and fine-tune them on the constructed evaluation dataset to predict plausibility scores.

One of the main limitations of RoBERTa is its length requirement of at most 512 tokens. Recently, this limitation was addressed by considering a sparser set of attention mechanisms such as locality-sensitive hashing and sliding window attentions, which reduce the computation complexity

from $\mathcal{O}(n^2)$ to $\mathcal{O}(n \log n)$ and $\mathcal{O}(n)$ respectively (Kitaev et al., 2020; Beltagy et al., 2020). In this work, we broaden the scope of the text plausibility evaluation to cover not only short but also long texts with more than 512 tokens. To this end, we examine and evaluate the quality of long texts using Longformer (Beltagy et al., 2020) that has linear complexity in terms of the number of tokens in a text. We fine-tune the pretrained Longformer for long text plausibility evaluation.

4.2 Baselines

We benchmark both fine-tuned classifiers on the manipulated data with the two following baselines.

UNION. Recently, Guan and Huang (2020) proposed an automatic evaluation metric by training a BERT model (Devlin et al., 2019) with an auxiliary reconstruction objective which helps to recover the perturbation from a negative sample. The proposed model is trained on negative implausible texts constructed by adopting repetition, substitution, reordering, and negation sampling techniques. This model and its proposed approach for data construction were compared with previously proposed methods and shown to be more efficient.

SENTAVG. We complete our investigation by selecting SENTAVG (Lai and Tetreault, 2018) as another baseline model for the plausibility evaluation task. SENTAVG leverages LSTM to get sentence representation from their words GloVe embeddings. All the sentences vectors are averaged to form the representation for the whole text and this vector is passed to a hidden layer. A softmax layer at the end computes the probability distribution of texts over positive and negative labels.

5 Experiments

We investigate the effectiveness of our proposed approach versus heuristic negative sampling techniques by focusing on the evaluation of open-domain story generation models in two datasets with short and long stories. We show the generalizability of metrics trained on our proposed plot manipulation data. We also separately assess the impact of each manipulation technique on the metric accuracy.

5.1 Datasets

We conduct our experiments on two English stories datasets that are significantly different in terms of length and topic; ROCStories (shortly ROC) and

Dataset	Train/Valid/Test
HEUR_SENT_SUB	47.1k/5.9k/5.9k
HEUR_KEY_SUB	47.1k/5.9k/5.9k
UNION_DATA	47.1k/5.9k/5.9k
MANPLTS	47.1k/5.9k/5.9k
AF_MANPLTS	94.2k/11.8k/11.8k

Table 2: Plausibility evaluation datasets for ROC stories using different negative sampling techniques.

Writing Prompt (briefly WP) datasets including on average 49.4 and 734.5 tokens in each story.

ROCStories. ROCStories is a resource of five-sentence commonsense stories collected via crowdsourcing (Mostafazadeh et al., 2016) covering a logically linked set of daily events. We follow the approach proposed by Yao et al. (2019) to extract story plots (storylines) for the stories and manipulate them to guide conditional language models to generate negative samples.

Writing Prompt. Writing Prompt dataset contains abstract high-level prompts and their corresponding long human-written stories from an online forum (Fan et al., 2018). To apply the plot manipulation technique for implausible text construction, we follow the procedure proposed by Fan et al. (2019) to extract the plots with verb and argument type role labeling tags.

Data Preparation. We split the stories from both datasets into two subsets for training generation and evaluation models, respectively. We use 70 percent of stories in ROC (ROC_LM) and WP (WP_LM) for fine-tuning GPT2 (Radford et al., 2019) language model with batch size of 4.² After 3 epochs of fine-tuning, the perplexity on the validation set of ROC and WP datasets are 8.28 and 25.04, respectively.

The remaining 30 percent of stories from ROC (ROC_Eval) and WP (WP_Eval) are used for training and evaluating the evaluation models. All stories in the original dataset represent plausible texts. We apply approaches from Section 3 to augment negative samples. Table 2 and Table 3 summarize the resulting datasets for ROC and WP. In HEUR_SENT_SUB, we extract all stories with at least 2 sentences and replace 50% of their sentences with random ones. For HEUR_KEY_SUB, we do random substitution of 15% of keywords with their corresponding antonyms extracted from ConceptNet and ignore stories without substitutable keywords. The UNION_Data is resulted by following rules from Guan and Huang (2020) and is applied

²We fine-tune GPT2 language model using <https://github.com/huggingface/transformers>.

Dataset	Train/Valid/Test
HEUR_SENT_SUB	163.1k/9.3k/9.1k
HEUR_KEY_SUB	162.7k/9.3k/9.0k
UNION_DATA	161.8k/9.2k/9.0k
MANPLTS	84.5k/4.7k/4.7k
AF_MANPLTS	107.2k/35.7k/35.7k

Table 3: Plausibility evaluation datasets for WP stories using different negative sampling techniques.

Data	Texts	Annotators	Kappa
ROC	300	27	0.61
WP	300	75	0.56

Table 4: Statistics and inter-annotator agreement of AMT annotations for plausibility metrics evaluation.

to stories with at least four sentences.

To create MANPLTS dataset, we first fine-tune the BART model (Lewis et al., 2019) with a batch size of 8 for three epochs on pairs of ground-truth plots and stories from ROC_LM and WP_LM data with the resulting perplexity of 3.44 and 6.79 for the validation sets. Afterward, 15% of plots are employed and two up to four proposed manipulation techniques in Section 3.2 are randomly selected and applied. We leverage the fine-tuned BART model and use the top-50 sampling technique with a temperature of 0.8. We specify the maximum length of 200 for ROC dataset and 1024 for WP dataset to generate implausible texts on manipulated plots.

In the AF_MANPLTS dataset, we apply the adversarial filtering technique on top of six generated implausible stories using the fine-tuned BART model conditioned on the manipulated plots. The output contains each human-written story and its three most challenging implausible samples.

5.2 Human Annotations

The performance of automatic evaluation metrics is assessed based on their correlations with human judgments. To this end, we gather human evaluations and examine the Spearman (ρ) and Kendall (τ) correlations with metrics predicted scores (Newman et al., 2010; Lai and Tetreault, 2018; Guan and Huang, 2020). Spearman and Kendall are beneficial in estimating monotonic associations for not normally distributed and ranked scores.

We collect human judgments through Amazon Mechanical Turk (AMT) experiments. We randomly choose 150 human-written stories from ROC_Eval and WP_Eval test sets and 150 machine-generated texts by the fine-tuned GPT2 models. Five distinct participants are asked to rate each story on a scale of 0 to 5 (from *not at all plausible*

to *completely plausible*). We prepare an attention check test to guarantee the accuracy of human annotations and recollect evaluations for users who do not pass the test. The average score of the five annotators is treated as the final human score for each text. We normalize human scores to be in the same range of 0-1 as the model’s output scores are. Table 4 shows the statistics and agreements in the conducted experiments.

5.3 Experimental Setup

We conduct a comprehensive set of experiments to examine and show the importance of training data in the plausibility evaluation task. We train both evaluation and language models on a machine with a GeForce RTX 2080 Ti GPU.

In our experiments, we have SENTAVG as the baseline model. We compare SENTAVG across more powerful classifiers – RoBERTa for ROC stories and Longformer for WP stories (FT_LM). We fine-tune pretrained RoBERTa-base model with the learning rate of 2e-5 and batch size 8 for three epochs and process the ROC stories with a maximum of 128 tokens. To evaluate WP with lengthy stories, we fine-tune pretrained Longformer-base model with the learning rate of 2e-5 and batch size 3 by encoding texts with at most 1024 tokens for three epochs.³

We complete the models’ comparisons by incorporating the recently proposed UNION model (Guan and Huang, 2020) to our experiments. We retrain it on the ROC_Eval and WP_Eval sets with the same hyper-parameters stated in their paper.

5.4 Experimental results

Table 5 depicts the quantitative results of correlation analysis between human and automatic evaluation metrics. For almost all constructed datasets for evaluation, the RoBERTa and Longformer in the case of short and long stories surpass the baseline models that show the impact of large transformer-based models in this evaluation task. The models trained on heuristically generated implausible samples by random sentence/keyword substitutions show the lowest correlations. The main reason for such weakness is the huge dissimilarity of heuristically generated training data and machine-generated test data, which has a significant negative impact on the model’s performance. The positive

³We fine-tune RoBERTa and Longformer models using <https://github.com/huggingface/transformers>.

Dataset	Model	ROC		WP	
		ρ	τ	ρ	τ
HEUR_SENT_SUB	SENTAVG	0.04	0.03	-0.13	-0.10
	FT_LM	0.10	0.07	0.12	0.10
HEUR_KEY_SUB	SENTAVG	-0.04	-0.03	-0.26	-0.18
	FT_LM	0.31	0.22	0.08	0.06
UNION_DATA	SENTAVG	0.11	0.08	-0.22	-0.15
	FT_LM	0.46	0.34	0.49	0.32
	UNION	0.22	0.15	0.19	0.15
MANPLTS	SENTAVG	0.24	0.16	0.22	0.20
	FT_LM	0.50	0.37	0.71	0.48
AF_MANPLTS	SENTAVG	0.22	0.16	0.25	0.23
	FT_LM	0.56	0.41	0.74	0.52

Table 5: Higher correlations of plausibility evaluation models trained on manipulated plots and adversarially filtered negative samples with human judgments versus heuristically constructed negative samples. Ft_LM represents fine-tuned RoBERTa and Longformer models for ROC and WP datasets, respectively.

impact of UNION_Data is visible in Table 5. It demonstrates that the construction of implausible stories based on a more complete set of heuristic alterations yields better training data but still has its own shortcomings. This could be due to fact that text-level manipulations introduce artifacts that break the naturalness of the texts and have quite different styles compared to machine-generated implausible texts.

The superiority of RoBERTa and Longformer models trained on MANPLTS and AF_MANPLTS datasets show the effectiveness of our proposed plot manipulation technique in enhancing the similarity between the training and test data. Adversarial filtering technique further helps to increase the quality of negative samples and generate better implausible machine-generated texts, which consequently improves the accuracy of evaluation. By applying hypothesis testing to compare the metrics correlations with human scores (Diedenhofen and Musch, 2015), we verify that these improvements are statistically significant ($p < .05$). We also note that the correlations between plot manipulation-based metrics and human evaluation are much higher in WP dataset. This could result from the limited ability of the current generative models to generate plausible long stories, thus making them easily distinguishable both by humans and automated metrics.

One of the desirable features of automated evaluation metrics for story generation is their generalizability or robustness to different datasets (Sellam et al., 2020; Guan and Huang, 2020). The dataset

Dataset	ROC \rightarrow WP		WP \rightarrow ROC	
	ρ	τ	ρ	τ
UNION_DATA	0.17	0.15	0.12	0.07
MANPLTS	0.57	0.39	0.23	0.16
AF_MANPLTS	0.60	0.42	0.26	0.18

Table 6: Correlation of plausibility metrics with human judgements. Arrow shows the train and test data used for examining metrics robustness.

Dataset	ρ	τ
MANPLTS-REORDER	0.65	0.42
MANPLTS-CONTINSER	0.69	0.45
MANPLTS-REPEATINSER	0.67	0.43
MANPLTS-RANDSUB	0.68	0.45

Table 7: Correlations of Longformer model fine-tuned on plot-level manipulated datasets with one specific excluded technique.

shifting robustness shows the metric’s success in accurately evaluating texts in different datasets. We examine the robustness of metrics by leveraging ROC and WP as two distributionally different types of stories datasets. We train models on various training data constructed from negative sampling techniques in ROC dataset and test them on human scores collected through AMT experiments conducted on WP dataset (**ROC \rightarrow WP**) and vice versa (**WP \rightarrow ROC**). In Table 6, we show the robustness of fine-tuned language models trained on the last three datasets of Table 5 as the best performing models in comparison to models trained on sentence and keyword substitutions. According to Table 6, the correlation drops due to the quite different structure of two datasets. RoBERTa/Longformer models fine-tuned on AF_MANPLTS in ROC/WP datasets and subsequently tested on WP/ROC dataset have the highest correlations with human judgments and can be generalized well on two datasets. The data shifting from ROC to WP better preserves the performance of metrics rather than the counterpart shifting. The reason for correlation decline of models trained on WP and tested on ROC could be the format of implausible texts in WP that could not be found in ROC data since the stories are shorter in this data and the reason for implausibility is fewer.

5.4.1 Ablation Study

The positive impact of plot-level manipulations in precisely evaluating the plausibility can be assessed with regard to the four different manipu-

lation techniques. We conduct an ablation study on WP dataset to examine each manipulation technique’s impact separately. We construct different training data each time by excluding one of the manipulation techniques and generating a new set of negative samples. Then we fine-tune Longformer on all these training datasets with different negative samples and compute the correlation of the fine-tuned Longformer as the evaluation metric with human judgments.

The lower correlations shown in Table 7 in comparison to Table 5 illustrate the harms that the elimination of each of the proposed approaches from the construction of training data could cause. This attests to the effectiveness of all proposed manipulation techniques in the generation of higher quality training data and subsequently resulting in more accurate evaluation metrics.

As this table demonstrates, the correlation drops the most by ablating the reordering and repeating plots, which shows that they are the major problems in generating long texts by language models and have the most significant role in constructing high-quality implausible samples and consequently accurate evaluation metrics.

6 Conclusion

Automatic plausibility evaluation models that are trained on heuristically generated data show low correlation with human judgement. We address this issue by creating a better quality set of implausible texts. In contrast to existing methods that modify text at token level, our approach introduces incoherence sources at a more abstract plot level, which helps to guide the generative model conditioned on those manipulated plots to generate negative samples that are more similar to machine-generated incoherent texts. We further improve the data quality by applying adversarial filtering to select more challenging and refined negative samples. Our experiments demonstrate that negative examples generated according to the proposed method result in more realistic implausible texts and consequently lead to more accurate evaluation metrics that have higher correlation with human judgement.

7 Ethics

All co-authors of this work totally understand and agree with *ACM Code of Ethics* and its importance in expressing the conscience of the profession. We ensure this work is compatible with the provided

code, specifically in the terms of providing non-offensive dataset construction.

1) training data construction In our approach, we use BART model conditioned on manipulated story plots to construct implausible samples that better reflect the implausibility in generation models. The main concern that arises here is the probability of generating abusive language samples from manipulated plots. Indeed, these plots origin from human-written stories without abusive languages provided by (Mostafazadeh et al., 2016; Fan et al., 2018) where users are not allowed to write profanity and inappropriate content. Accordingly, our manipulated version of plots and the BART model conditioned on them generate samples unlikely to contain strong biases or abusive content. It is noteworthy to mention that even the source plots are relatively benign, the process of altering them would have the possibility of creating objectionable texts. Other potential attack could be the dual-usage of metrics by augmenting offensive language texts as plausible samples. This would harm the underlying tasks such as story generation models to be encouraged to generate inappropriate stories. Such attacks can be identified and dissolved by security trended studies which are out of this work’s scope.

1) testing data collection We collect human judgments by conducting Amazon Mechanical Turk (AMT) experiments that are leveraged to compare the accuracy of trained metrics in terms of their correlations with human scores. The conducted AMT does not disrupt user privacy as we do not contain their personal information. This fades the possibility of any gender bias problems and IRB approval needs. Annotators were asked to rate the coherence of stories in each HIT page of AMT in the range of 0 up to 5. We fairly compensated annotators. The average time of annotating each HIT in AMT was 25 minutes (including three stories for evaluation and their explanations), and according to the per hour wage of \$13, we fairly paid them \$6 per HIT.

This work targets the NLP open-domain generation community. Our metrics establish the main basis to achieve higher-quality generations by automatically assess the outputs and save time, cost, and human efforts. We don’t anticipate specific failure modes in our work since the provided approach’s success has been investigated through a comprehensive set of comparisons with other existing metrics.

Acknowledgment

This work is supported by the CwC program under the Contract W911NF-15-1-0543 with the US Defense Advanced Research Projects Agency (DARPA). We would like to thank the anonymous reviewers for their helpful comments and the members of PLUSlab from USC/UCLA, Shushan Arakelyan, and Ninareh Mehrabi for their constructive feedback.

References

- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. Comet: Commonsense transformers for automatic knowledge graph construction. In *Association for Computational Linguistics (ACL)*.
- Faeze Brahman, Alexandru Petrusca, and Snigdha Chaturvedi. 2020. Cue me in: Content-inducing approaches to interactive story generation. In *Asia-Pacific Chapter of the Association for Computational Linguistics (AACL)*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*.
- Birk Diedenhofen and Jochen Musch. 2015. cocor: A comprehensive solution for the statistical comparison of correlations. *PloS one*.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. *arXiv preprint arXiv:1805.04833*.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. In *Association for Computational Linguistics (ACL)*.
- Seraphina Goldfarb-Tarrant, Tuhin Chakrabarty, Ralph Weischedel, and Nanyun Peng. 2020. Content planning for neural story generation with aristotelian rescoring. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Seraphina Goldfarb-Tarrant, Haining Feng, and Nanyun Peng. 2019. Plan, write, and revise: an interactive system for open-domain story generation. In *2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT 2019), Demonstrations Track*, volume 4, pages 89–97.

- Jian Guan and Minlie Huang. 2020. Union: An unreference metric for evaluating open-ended story generation. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *ICLR*.
- Ari Holtzman, Jan Buys, Maxwell Forbes, Antoine Bosselut, David Golub, and Yejin Choi. 2018. Learning to write with cooperative discriminators. In *Association for Computational Linguistics (ACL)*.
- Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. 2020. Reformer: The efficient transformer. In *ICLR*.
- Alice Lai and Joel Tetreault. 2018. Discourse coherence in the wild: A dataset, evaluation and methods. *arXiv preprint arXiv:1805.04993*.
- Mirella Lapata and Regina Barzilay. 2005. Automatic evaluation of text coherence: Models and representations. In *IJCAI*, volume 5, pages 1085–1090.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Association for Computational Linguistics (ACL)*.
- Jiwei Li and Dan Jurafsky. 2016. Neural net models for open-domain discourse coherence. *arXiv preprint arXiv:1606.01545*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Eleni Miltsakaki and Karen Kukich. 2004. Evaluation of text coherence for electronic essay scoring systems. *Natural Language Engineering*, 10(1):25–55.
- N Mostafazadeh, N Chambers, X He, D Parikh, D Batra, L Vanderwende, P Kohli, and J Allen. 2016. A corpus and cloze evaluation framework for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL HLT 2016)*, San Diego, CA, USA. Association for Computational Linguistics.
- David Newman, Jey Han Lau, Karl Grieser, and Timothy Baldwin. 2010. Automatic evaluation of topic coherence. In *Human language technologies: The 2010 annual conference of the North American chapter of the association for computational linguistics*, pages 100–108. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Nanyun Peng, Marjan Ghazvininejad, Jonathan May, and Kevin Knight. 2018. Towards controllable story generation. In *Proceedings of the First Workshop on Storytelling*, pages 43–49.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- Hannah Rashkin, Asli Celikyilmaz, Yejin Choi, and Jianfeng Gao. 2020a. Plotmachines: Outline-conditioned generation with dynamic plot state tracking. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Hannah Rashkin, Asli Celikyilmaz, Yejin Choi, and Jianfeng Gao. 2020b. PlotMachines: Outline-conditioned generation with dynamic plot state tracking. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. In *Association for Computational Linguistics (ACL)*.
- Robyn Speer and Catherine Havasi. 2012. Representing general relational knowledge in conceptnet 5. In *LREC*, pages 3679–3686.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5754–5764.
- Lili Yao, Nanyun Peng, Weischedel Ralph, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In *The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)*.
- Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. Swag: A large-scale adversarial dataset for grounded commonsense inference. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Yizhe Zhang, Siqu Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. Dialogpt: Large-scale generative pre-training for conversational response generation. In *Association for Computational Linguistics (ACL)*.