# Consistency and Coherency Enhanced Story Generation

**Wei Wang**[1,*], **Piji Li**[2], **Hai-Tao Zheng**[1]

[1]Shenzhen International Graduate School, Tsinghua University
[2]Tencent AI Lab
w-w16@mails.tsinghua.edu.cn
pijili@tencent.com
zheng.haitao@sz.tsinghua.edu.cn

## Abstract

Story generation is a challenging task, which demands to maintain consistency of the plots and characters throughout the story. Previous works have shown that GPT2, a large-scale language model, has achieved good performance on story generation. However, we observe that several serious issues still exist in the stories generated by GPT2 which can be categorized into two folds: consistency and coherency. In terms of consistency, on one hand, GPT2 cannot guarantee the consistency of the plots explicitly. On the other hand, the generated stories usually contain coreference errors. In terms of coherency, GPT2 does not take account of the discourse relations between sentences of stories directly. To enhance the consistency and coherency of the generated stories, we propose a two-stage generation framework, where the first stage is to organize the story outline which depicts the story plots and events, and the second stage is to expand the outline into a complete story. Therefore the plots consistency can be controlled and guaranteed explicitly. In addition, coreference supervision signals are incorporated to reduce coreference errors and improve the coreference consistency. Moreover, we design an auxiliary task of discourse relation modeling to improve the coherency of the generated stories. Experimental results on a story dataset show that our model outperforms the baseline approaches in terms of both automatic metrics and human evaluation.

## 1 Introduction

A story is anything which is narrated in the form of a causally/logically linked set of event plots involving some shared characters (Mostafazadeh et al., 2016). Story generation aims at automatically generating stories with such attributes. Different from other text generation tasks, story generation is

---

*Work was done during internship at Tencent AI Lab.

**Prompt**: The zombie apocalypse is upon us but all humans are immune to the virus . All other mammals are susceptible

**Story**: As I sat there I tried to remember how many days it had been since the last one died . [My wife and daughter were dead , I knew they would n't be alive any time soon .]₁ₐ The food ran out quickly , they had to be taken away to stay in the house , but that was no reason to be angry if there were a few people left living there . After a month of food being gone a little girl sat next to me . She was the only other survivor who was n't bitten , as it was hard for us to eat them , we were all sickly . That would be another month to go and if the other animals were still alive in the food , they would n't be able to eat her . " Hey little one . " I looked up and saw her . [I had a big smile on my face at this point .]₂ₐ [My wife was the only one who noticed , but her smile was all that was visible of me , and not of her .]₁ᵦ [I looked at her , and smiled sadly .]₂ᵦ She looked sad , but did n't care . She never cared for me .

Table 1: A story generated by GPT2.

more challenging because it is restricted to several constraints: (1) it must maintain **consistent plots** to form a reasonable story; (2) it must guarantee the **consistency** of the characters throughout the story; (3) the **coherency** of the text units such as the clauses or sentences should be concerned. In order to improve the quality of plot consistency, some previous works focus on the perspective of plot planning and then merge the text units according to the order of plots (Lebowitz, 1987; PÉrez and Sharples, 2001; Porteous and Cavazza, 2009; Riedl and Young, 2010; Li et al., 2013). We observe that those approaches rely heavily on human annotations and are restricted to the abstract story representation level without surface realization in natural language, such as producing event verb sequence and sentence segments. Therefore, these methods need to work with sentence templates or rules to generate stories.

In the past few years, several end-to-end approaches based on Sequence-to-Sequence (Seq2Seq) models (Sutskever et al., 2014; Bahdanau et al., 2014) are proposed, which can

generate a story at a stroke in a left-to-right manner (Jain et al., 2017; Clark et al., 2018; Fan et al., 2018). These methods are data-driven and can directly generate stories in natural language form instead of other abstract representation. However, these methods struggle to capture the high-level interactions between the plot points and maintain consistent plots throughout the story. Thus, several two-stage models for story generation have recently been proposed (Martin et al., 2018; Xu et al., 2018; Yao et al., 2019; Fan et al., 2019; Chen et al., 2019). These models usually decompose story generation into two stages: generating middle form first and then generating the final story. Different middle forms are applied in these methods, such as keywords, sentences and event tuples.

Recently, the OpenAI GPT2/3 language model (Radford et al., 2019; Brown et al., 2020) achieves strong performance on several language generation tasks. (See et al., 2019) and (Guan et al., 2020) verify the performance of GPT2 on story generation and GTP2 outperforms both end-to-end methods and two-stage methods. However, after analyzing the generated stories carefully, we observe that there are still some serious issues in the generated stories by GPT2. Take a story generated by GPT2 as shown in Figure 1 for example. The story is about survivors in the end of the world. First, plots consistency cannot be guaranteed among multiple sentences of a story, such as blue sentences in Figure 1. The sentence $1a$ describes "My wife and daughter were dead ". But the sentence $1b$ talks about "My wife" again. It is contradictory. There is the same problem in the sentence $2a$ and $2b$. Second, there are still coreference errors in generated stories, such as red text in Figure 1. It is not clear who *they* and *them* refer to. Moreover, Top-k sampling (Radford et al., 2019; See et al., 2019; Brown et al., 2020) is usually utilized as the decoding strategy in long text generation. The random operation in sampling will disturbance the generation procedure by producing improper tokens which will decrease the quality. This phenomenon is more pronounced at the border of sentences, therefore we can sometimes observe the bad performance in discourse coherency.

To solve the aforementioned problems, we propose a two-stage generation model based on Transformer-based auto-regressive language models to improve consistency and coherency of stories.

Specifically, the first stage is to organize the story outline which depicts the story plots and events, and the second stage is to expand the outline into a complete story. Therefore the plots consistency can be controlled and guaranteed explicitly. In addition, coreference supervision signals are incorporated to reduce coreference errors and improve the coreference consistency. Moreover, we design an auxiliary task of discourse relation modeling to enhance the discourse coherency of the generated stories. Both the backbone models in the two states are designed based on Transformer-based language models. Thus, on one hand, the framework can still inherit the superior performance of GPT2, on the other hand, it can guarantee the plot consistency, coreference consistency, as well as discourse coherency.

The main contributions of this paper are summarized as follows:

- We propose to improve the plot and coreference consistency as well as the discourse coherency for the task of story generation.
- A two-stage framework based on Transformer-based language models is designed to control the plots and improve consistency of generated stories.
- A coreference constraint is applied to improve the coreference consistency of generated stories.
- We design a discourse relation modeling component as an auxiliary task during training to enhance the performance of discourse coherency.
- Experiments on a story dataset from Reddit demonstrate that our model outperforms the baseline methods in terms of both automatic metrics and human evaluation.

## 2 Methodology

### 2.1 Overview

To begin with, we state the problem of story generation as follows: given a prompt context $\mathbf{X} = \{x_1, ..., x_i..., x_k\}$ where $x_i$ denotes each word in the prompt, the model needs to generate a story $\mathbf{Y} = \{y_1, ..., y_i, ..., y_n\}$ following the prompt $\mathbf{X}$ by maximizing the conditional probability $p(\mathbf{Y}|\mathbf{X})$.

As show in Figure 1, to enhance the consistency and coherency of generated stories, we propose a two-stage framework for story generation. The first stage is story outline generation which can generate the plot outline based on the given prompt.
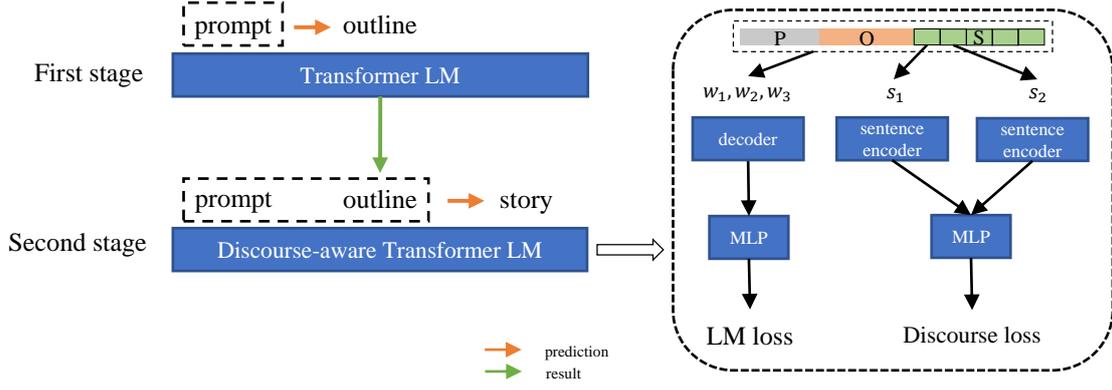
Figure 1: The framework of our model for story generation.

Then in the second stage, the whole story is completed by embellishing the outline generated in the first stage. Transformer-based language models are introduced as the backbone models for those two stages respectively. Moreover, a component of discourse relation classification is incorporated to the language model as an auxiliary task to further improve the coherency of the generated stories. To further improve the consistency, we design a coreference supervision component to encourage the language model to attend on correct entities when generating pronouns by maximizing the attention weights of the corresponding entities.

## 2.2 Transformer-based Language Model

Inspired by the popular pre-trained language models for text generation such as GPT2 (Radford et al., 2019), XLNET (Yang et al., 2019) and GPT3 (Brown et al., 2020), we also employ the Transformer-based auto-regressive language models as our backbone frameworks.

Transformer-based language models only contain a decoder. The decoder consists of $N$ identical self attention blocks and each block contains two sub-layers: a self multi-head attention layer and a feed-forward layer. A add & norm layer is employed around each of two sub-layers. Formally, given the input $\mathbf{H}^{n-1}$, the output $\mathbf{H}^n$ of each decoder block is computed as follows:

$$\mathbf{C}^n = \mathrm{LN}\left(\mathrm{SELF\text{-}ATT}\left(\mathbf{H}^{n-1}\right) + \mathbf{H}^{n-1}\right) \quad (1)$$
$$\mathbf{H}^n = \mathrm{LN}\left(\mathrm{FFN}\left(\mathbf{C}^n\right) + \mathbf{C}^n\right) \quad (2)$$

where SELF-ATT(·), LN(·), and FFN(·) are respectively self-attention mechanism, layer normalization, and feed-forward network with ReLU activation in between. SELF-ATT(·) computes atten-

tion over the input $\mathbf{H}^{n-1}$ as follows:

$$\mathrm{SELF\text{-}ATT}\left(\mathbf{H}^{n-1}\right) = \mathrm{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_k}}\right)\mathbf{V}$$
$$(3)$$

where $\{\mathbf{Q}, \mathbf{K}, \mathbf{V}\}$ are query, key and value vectors that are transformed from the input $\mathbf{H}^{n-1}$. $\sqrt{k}$ is the scaling factor where the $d_k$ is the dimension size of the query and key vectors. Given the word embeddings $\mathbf{E} = \{e_1, e_2, ..., e_m\}$ and corresponding positional embeddings $\mathbf{P} = \{p_1, p_2, ..., p_m\}$, the first block input $\mathbf{H}^0 = \mathbf{E} + \mathbf{P}$.

Finally, a linear function with softmax activation is used to compute the probability of next word $x_t$ via:

$$p\left(x_t | x_{\leq t-1}\right) = \mathrm{softmax}\left(g\left(h_t\right)\right) \quad (4)$$

We calculate negative log-likelihood loss for model training:

$$\mathcal{L}_{\mathrm{lm}} = -\frac{1}{T}\sum_t \log p\left(x_t | x_{\leq t-1}\right) \quad (5)$$

## 2.3 Two-stage Generation

### Outline Preparation

In order to regard the outline generation task as a supervised learning problem, we must construct a high-quality training dataset including sufficient prompt-outline pairs. As pre-mentioned, outline implies the story plots, therefore the quality of outline can affect the performance of story generation directly. If the outline contains too much information, the story generator will directly learn to copy from the outline and restrain the imagination and creativity. On the contrary, if the outline ignores

| S1 | marker | S2 |
|---|---|---|
| Her eyes flew up to his face. | and | Suddenly she realized why he looked so different. |
| The concept is simple. | but | The execution will be incredibly dangerous. |
| You used to feel pride. | because | You defended innocent people. |
| Belter was still hard at work. | when | Drade and barney strolled in. |
| I' ll tell you about it. | if | You give me your number. |
| We plugged bulky headsets into the dashboard. | so | We could hear each other when we spoke into the microphones. |
| It was mere minutes or hours. | before | He finally fell into unconsciousness. |
| And then the cloudy darkness lifted. | though | The lifeboat did not slow down. |

Table 2: Example pairs from Books 8 dataset.

the key-point information, the informativeness and consistency of stories will be decreased.

In this work, we investigate two forms of outline: keyword and abstract. These two forms retain the important information of the story and ignore some details and commonly used in two-stage based methods (Yao et al., 2019; Fan et al., 2019; Chen et al., 2019). Our motivation is to use two-stage generation to improve performance of GPT2 , so we do not design a new middle form. Specifically, we use the RAKE algorithm (Rose et al., 2010) [1] to extract keywords of story. According to (Yao et al., 2019) and the average lengths of stories in our corpus, we extract 10 keywords for each story. We use a variation of the TextRank algorithm (Barrios et al., 2016) [2] to extract abstract of story. In order to retain important information and ignore some detail information, we keep 30% sentences of each story as abstract. Thus, we can get (prompt, outline, story) pairs automatically to train the two-stage model.

**Prompt to Outline Generation**

A Transformer-based language model based decoder is used to generate outlines. Specifically, we concatenate prompt $\mathbf{X}$ and outline $\mathbf{Z}$ with `<SEP>` token to get a sequence $\mathbf{X}'$. For training, we compute cross entropy of all tokens in $\mathbf{X}'$ as normal language model. When testing, given the prompt tokens as context, the decoder generates outline tokens.

**Prompt and Outline to Story Generation**

Another decoder with the same architecture is used to generate stories. We concatenate prompt $\mathbf{X}$, outline $\mathbf{Z}$ and story $\mathbf{Y}$ with `<S>` and `<SEP>` token to get a sequence $\mathbf{X}''$. For training, we compute cross entropy of prompt and story tokens in $\mathbf{X}''$. Note that we don't calculate the loss of the outline tokens. Because, the tokens come from the story and we avoid computing loss of these tokens twice. When testing, given the prompt and the outline tokens as context, the decoder generates story tokens. Next, two components are incorporated in this stage to enhance discourse coherency and coreference consistency.

## 2.4 Discourse Coherency Enhancement

In order to improve discourse representation of Transformer-based language model, we design a discourse relation classification task as an auxiliary task. Discourse relations describe how two segments (e.g. clauses, sentences, and larger multi-clause groupings) of discourse are logically connected. These relations can be used to describe the high-level organization of text. Thus, discourse relation is an important aspect of story coherence. In this work, we only consider shallow discourse relations between adjacent sentences as many research on discourse relation classification do (Chen et al., 2016; Lan et al., 2017; Bai and Zhao, 2018).

**Discourse Information Preparation**

In order to get discourse label of adjacent sentences in stories, we need to train a golden discourse relation classification model. However, there is limited annotation corpus of implicit discourse relations and explicit discourse relations. For example, the commonly used dataset Penn Discourse Treebank 2.0 (Prasad et al., 2008) contains about 10k pairs. Following (Nie et al., 2019), we use discourse markers as replace of discourse relations. Because we are able to automatically curate a sizable training set of sentence pairs with discourse markers. We use discourse marker dataset Book 8 from (Nie et al., 2019), which contains 3.6M sentence pairs and each pair is labeled with one connective of 8 connectives as discourse label. Several sentence pairs and corresponding discourse markers are shown in Table 2.

---

[1]https://pypi.org/project/rake-nltk/
[2]https://radimrehurek.com/gensim/

We fine tuning BERT ([Devlin et al., 2018](#)) [3] in this dataset to get a golden discourse marker prediction model. Then we use this model to tag discourse relation label of sentence pairs in our story corpus. Considering that this automatic tagging may produce large errors, we only keep labels with high classification probability, and labels with lower probability are replaced with the ninth label, *unknown*. The sentence pairs with labels belonging to 8 connectives are used to train our discourse relation classification component.

**Discourse-aware Story Generation**

The discourse relation classification component contains a sentence encoder and a two-layers MLP. The encoder is used to extract sentence semantic feature and the MLP is used to convert feature into classification probability. The sentence encoder shares parameters with the story decoder exclude the output layer. For a story $\mathbf{Y}$ contains several sentence $\{\mathbf{S_1}, \mathbf{S_i}, \mathbf{S_p}\}$ and each sentence contains several words $\mathbf{S_i} = \{y_{i1}, y_{ij}, y_{iq}\}$, we get output $h_{ij}^w$ of encoder as word representation and use max pooling operation on words of this sentence to get sentence representation $h_i^s$:

$$\mathbf{H}_i^s = \text{encoder}(\mathbf{S_i}) \tag{6}$$
$$h_i^s = \max(\mathbf{H}_i^s) \tag{7}$$

Then the MLP is used to classify adjacent sentences as follows:

$$f = \tanh(\mathbf{W_f}[h_i^s, h_{i+1}^s] + b_f) \tag{8}$$
$$p(dis|\mathbf{S_i}, \mathbf{S_j}) = \text{softmax}(\mathbf{W_o}f + b_o) \tag{9}$$

The loss function $\mathcal{L}_{\text{dis}}$ of this component is the cross-entropy of discourse label. Then a joint loss function is applied to train the second stage model:

$$\mathcal{L} = \mathcal{L}_{\text{lm}} + \lambda_1 \mathcal{L}_{\text{dis}} \tag{10}$$

where $\lambda_1$ is a hyperparameter to balance two tasks.

## 2.5 Coreference Consistency Enhancement

Although Transformer-based language model has the ability of long-distance dependence, there are still some coreference errors in the generated stories. In order to encourage model to attend correct entities, we add a supervision on attention weight of entity mention tokens. We use Stanford's CoreNLP tool [4] to extract coreference annotation of stories.

Specifically, for a story $\mathbf{Y}$ we get $p$ coreference clusters and each cluster contains $q$ entity mentions. We assign each entity mention token $y_i^c$ in subsequence $\mathbf{Y}^c = \{y_1^c, y_i^c, y_{pq}^c\}$ a cluster label $\mathbf{C} = \{c_1, c_i, c_{pq}\}$. During training, for a entity mention token $y_i^c$, we get attention weights between current token and previous tokens $\{y^c \leq i-1\}$ in last self-attention layer of decoder, the sum of which is 1:

$$\sum_{k=1}^{i-1} \alpha_{ik} = 1 \tag{11}$$

We design a coreference loss to maximize attention weights of tokens in the same cluster as follows:

$$\mathcal{L}_{\text{coref}} = -\frac{1}{pq} \sum_{i=1}^{pq} \frac{1}{N_i} \sum_{k=1}^{i-1} \mathbb{1}(c_k = c_i) \log \alpha_{ik} \tag{12}$$

where $N_t$ is the number of entity mentions in the same cluster $c_i$. Considering these two components, the loss function for the second stage model is as follows:

$$\mathcal{L} = \mathcal{L}_{\text{lm}} + \lambda_1 \mathcal{L}_{\text{dis}} + \lambda_2 \mathcal{L}_{\text{coref}} \tag{13}$$

# 3 Experimental Setup

## 3.1 Settings and Data Set

For two Transformer decoders, we apply the same model size as GPT2-117M ([Radford et al., 2019](#)). Thus we can analysis the effect of pre-training weight of GPT2. Specifically, the dimension of word embedding and the dimension of hidden vectors are set to 768. The number of self attention block is set to 12 and 12 heads are used in self multi-head attention. We train the model using Adam ([Kingma and Ba, 2014](#)) with learning rate 0.0005. The dropout rate is set to 0.3 for regularization. $\lambda_1$ and $\lambda_2$ are set to 0.1 and 0.3 according to the performance in valid set. Following ([Fan et al., 2018](#)) we generate stories with random top $k$ sampling, where next words are sampling from the top $k = 20$ candidates rather than the entire vocabulary distribution.

We use writing prompts dataset from ([Fan et al., 2018](#)), which is collected from Reddit's WRIT-INGPROMPTS forum [5]. WRITINGPROMPTS is a community where online users inspire each other to write by submitting story prompts. Each prompt can have multiple story responses. The prompts have a large diversity of topic, length, and detail.

---

[3]https://github.com/huggingface/transformers
[4]https://stanfordnlp.github.io/CoreNLP/

[5]https://www.reddit.com/r/WritingPrompts/

There are 300k stories and the dataset is split into TRAIN, VAL and TEST (90%/5%/5%). For our experiments, we limit the length of the stories to 500 words maximum. We use the GPT2's BPE vocabulary with size of 50,527 in our model.

## 3.2 Evaluation Metrics

**Automatic Evaluation.** Many commonly used metrics based on n-gram overlap between the generated text and the human text, such as BLEU (Papineni et al., 2002), are not useful in story generation, which is also observed by previous work (Martin et al., 2018; Fan et al., 2018). Because we do not aim to generate a specific story; we want to generate viable and novel stories.

In order to evaluate different aspect of stories we use four type metrics. We use **Perplexity** to evaluate the fluency of stories. Perplexity is commonly used to evaluate the quality of language models, and it reflects how fluently the model can produce the correct next word given the preceding words. What's more, in order to evaluate the diversity of stories we compute **Distinct-1/2** (Li et al., 2016), which is the percentage of distinct n-grams in all generated stories and is widely used in conversation generation.

In order to evaluate the discourse coherency of the stories, we reuse the fine-tuned BERT for evaluation. Specifically, we use BERT to tag discourse label for sentence pairs in generated stories in the same way as the tagging process of training set in Section 2.4. We compute the percentage of sentence pairs with **Unknown** labels in generated stories. The less sentence pairs with unknown labels the model generates, the better the coherency of stories are. In order to evaluate the coreference coherence, we compute the averaged **Coreference Chains** in each story. Specifically, we use Stanford's CoreNLP tool [6] to extract coreference chains of generated stories.

**Human Evaluation.** To further evaluate the quality of generated stories, we conduct pair-wise comparisons with two strong baseline models (FConvS2S and GPT2P). We evaluate the models from the following three perspectives: **Relevance** to indicate whether a story is relevant to the given prompt, **Grammaticality** to indicate whether a story is natural and fluent, and **Logicality** to indicate whether a story is consistent and coherent in terms of causal dependencies in the context. Three

---

[6] https://stanfordnlp.github.io/CoreNLP/

aspects are independently evaluated. We randomly sample 100 stories from the test set and obtain 300 stories from three models. For each pair of stories (one by our model and the other by a baseline, along with the prompt), three annotators are asked to give a preference (win, lose, or tie) in terms of three metrics respectively. Majority voting is used to make final decisions among the three annotators.

## 3.3 Comparison Methods

**Conv Seq2Seq with self-attention (ConvS2S).** We replicate the model proposed by (Fan et al., 2018) using their source code, which applies a convolutional sequence-to-sequence model with gated self-attention to generate stories from prompts.

**Fusion of Conv Seq2Seq with self-attention (FConvS2S).** The model is also proposed by (Fan et al., 2018), which utilizes a fusion mechanism to integrate two **ConvS2S**.

**GPT2.** The model only contains a Transformer-based decoder and has the same model size as GPT2-117M (Radford et al., 2019). We train the model from scratch.

**GPT2 with Pre-training (GPT2P).** We first load pre-training weights of GPT2-117M and then fine tune the model on the used dataset.

**Ours.** Our overall model contains two-stage generation, discourse relation classification and coreference supervision. In order to evaluate the upper bound of two-stage generation, we use different percentages of tokens of ground truth outlines as contexts to generate stories. Ours(0%) means using own generated outlines as contexts in the second stage to generate stories. It is our final model. Ours(100%) means all tokens of ground truth outlines are used as contexts. It is the upper bound model.

## 4 Results and Discussions

### 4.1 Automatic Evaluation and Human Evaluation

As shown in Table 3, we compute four types metrics for these methods. We can see that GPT2 outperforms FConvS2S and ConvS2S in all metrics. This indicates that the self-attention based model is superior to the convolutional based model in story generation. Although FConvS2S and ConvS2S is enhanced with a self-attention mechanism, their ability to capture long-distance dependence is still weaker than GPT2. Compared to GPT2, GPT2P improves the perplexity and distinct significantly.

| Method | Perplexity↓ | Dis-1(%)↑ | Dis-2(%)↑ | Unknown(%)↓ | Coref Chains↑ |
|---|---|---|---|---|---|
| ConvS2S | 34.61 | 0.400 | 5.191 | 76.01 | 5.52 |
| FConvS2S | 33.97 | 0.482 | 6.271 | 75.60 | 5.43 |
| GPT2 | 29.50 | 0.474 | 6.796 | 74.95 | 5.67 |
| GPT2P | 25.64 | 0.493 | 7.333 | 73.61 | 5.61 |
| Ours(0% ground truth outline) | 30.84 | 0.531 | 7.379 | 75.19 | 5.98 |
| Ours(50% ground truth outline) | 19.21 | 1.311 | 13.253 | 75.15 | 5.97 |
| Ours(100% ground truth outline) | 10.32 | 1.509 | 15.266 | 74.97 | 5.80 |

Table 3: Automatic evaluation results on TEST set.

| Method | Relevance | | | Grammaticality | | | Logicality | | |
|---|---|---|---|---|---|---|---|---|---|
| | Win(%) | Tie(%) | Lose(%) | Win(%) | Tie(%) | Lose(%) | Win(%) | Tie(%) | Lose(%) |
| Ours vs. FConvS2S | **23** | 66 | 11 | **28** | 53 | 19 | **40** | 33 | 27 |
| Ours vs. GPT2P | **21** | 60 | 19 | **17** | 69 | 14 | **31** | 47 | 22 |

Table 4: Human evaluation results on TEST set.

| Method | Perplexity↓ | Dis-1(%)↑ | Dis-2(%)↑ | Unknown(%)↓ | Coref Chains↑ |
|---|---|---|---|---|---|
| **First stage** | | | | | |
| keyword | 74.46 | 0.964 | 7.132 | / | / |
| abstract | 35.53 | 0.776 | 10.060 | / | / |
| **Second stage** | | | | | |
| story with keyword | 17.82 | 0.461 | 6.188 | 74.26 | 5.67 |
| story with abstract | 10.65 | 0.512 | 7.358 | 74.54 | 5.81 |

Table 5: Comparison of different outlines.

GPT2P also generates least sentence pairs with *unknown* discourse relation. This shows that pre-training weights contributions to generating more fluent, diverse and coherent stories. Compared to these methods, our model (Ours(0%)) achieves best diversity and coreference performance. This demonstrates the effectiveness of our overall model. The upper bound model (Ours(100%)) achieves best perplexity score. This indicates that our model sacrifices part of fluency for the plot control. What's more, we can see that all two-stage models has a lower *unknown* score compared with GPT2 and GPT2P. We claim that two-stage generation and discourse relation component may repel each other. Next, we conduct ablation experiment to evaluate each component of our method.

Table 4 reports human evaluation results. Our method achieves best scores in three metrics. Specifically, our method mainly improve scores on Logicality. This shows that our method can generate more coherent stories by utilizing discourse and coreference supervision. Our method performs similar to GPT2P in term of Relevance and Gram-

maticality. Because both two methods use Transformer as the decoder and our model dose not design a component to improve the relevance to the prompt.
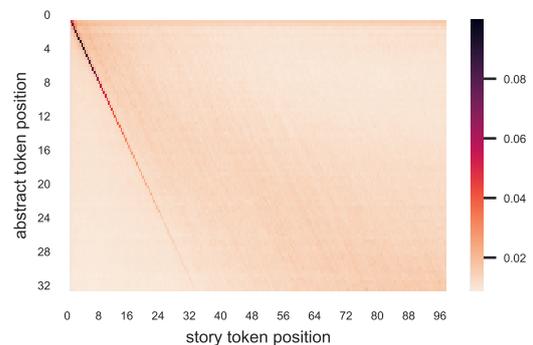
## 4.2 Outline Analysis



Figure 2: The attention weight distribution of story tokens in different positions.

We compare the performance of keyword and abstract as outlines. As shown in Table 5, in first stage keyword is more difficult to generate than

abstract, for that keyword gets a higher perplexity. From second stage, we can see that stories using abstract as outline get better scores in four metrics. This indicates that the abstract contributes to generating stories with better diversity and consistency. Therefore, we take abstract as outline in our model. In order to evaluate whether the stories are generated following the plot order of abstract, we plot story tokens' attention weight distributions on abstract tokens. The attention weight distributions are computed by averaging 2,000 generated stories. Because of the limited space, we only list tokens of the abstract and the story in the front positions. The result is shown in Figure 2. There are several lines with darker colors in the diagonal direction of the figure. This demonstrates that the story's focus follows the plot order of the abstract and our two-stage model can control the plots of the story well.

## 4.3 Discourse Relation Classification

| TLM+Discourse | And(%)↑ | When(%)↑ | Unknown(%)↓ |
|---|---|---|---|
| 0.1 | 11.43 | 2.90 | 72.94 |
| 0.3 | 11.38 | 2.80 | 73.60 |
| 0.5 | 10.91 | 2.72 | 73.78 |

Table 6: The percentages of discourse relations with different $\lambda_1$.

| Method | And(%)↑ | When(%)↑ | Unknown(%)↓ |
|---|---|---|---|
| ConvS2S | 8.52 | 2.45 | 76.01 |
| FConvS2S | 8.67 | 2.41 | 75.59 |
| GPT2 | 10.14 | 2.50 | 74.96 |
| GPT2P | 10.96 | 2.60 | 73.61 |
| TLM+Discourse | 11.43 | 2.90 | 72.94 |

Table 7: The percentages of discourse relations of different methods.

We evaluate the effect of discourse relation classification with **TLM+Discourse**, which enhances Transformer based language model with our proposed discourse relation classification component. We show the percentages of *and*, *when* and *unknown* for that the percentages of other relations are small. We first try several different values for $\lambda_1$ in loss function and the result is shown in Table 6. When $\lambda_1$ is high, there is no gain in three discourse metrics. When $\lambda_1$ is 0.1, we get most gain in three discourse metrics. This indicates that large weight affects the main task and damages the performance of language model. Therefore, we use 0.1 in our

final model. Then we compare **TLM+Discourse** with other baselines. The percentages of discourse relations are shown in Figure 7. We can see that GPT2 and GPT2P performs better than ConvS2S and FConvS2S. TLM+Discourse generates least sentence pairs with *unknown* relation and achieves best score among all methods. This shows that discourse relation classification component is helpful to generate more sentence pairs with certain discourse relations and further improve the coherency of generated stories.

## 4.4 Case Study

**Prompt**: The zombie apocalypse is upon us but all humans are immune to the virus . All other mammals are susceptible .

**GPT2P**: As I sat there I tried to remember how many days it had been since the last one died . [My wife and daughter were dead , I knew they would n't be alive any time soon .]$_{1a}$ The food ran out quickly , they had to be taken away to stay in the house , but that was no reason to be angry if there were a few people left living there . After a month of food being gone a little girl sat next to me . She was the only other survivor who was n't bitten , as it was hard for us to eat them , we were all sickly . That would be another month to go and if the other animals were still alive in the food , they would n't be able to eat her . " Hey little one . " I looked up and saw her . [I had a big smile on my face at this point .]$_{2a}$ [My wife was the only one who noticed , but her smile was all that was visible of me , and not of her .]$_{1b}$ [I looked at her , and smiled sadly .]$_{2b}$ She looked sad , but did n't care . She never cared for me .

**Our (generated abstract)**: He was the last to see the end of this world . The only reason John had survived was by his own hand , and it took much more than he could bear with a few months to heal his wounds . It took a few months for him to be able to see the end and get his arm back from the dead , and he did n't care for a few more years without having to go back .

**Our (generated story)**: When he had been first laid to rest the first few days he did n't care , and he had thought it was just a small thing , but it started getting worse . He was the last to see the end of this world . The only reason John had survived was by his own hand , and it took much more than he could bear with a few months to heal his wounds . John had lived the rest of his life on Earth without anything else but the clothes on his back , and the rifle he kept at home when he was alone . The first few days were easy . A little bit of food , a bit of protection . He spent a few days walking around the city and drinking a beer , thinking he was safe , but the next few months were hell . It took a few months for him to be able to see the end and get his arm back from the dead , and he did n't care for a few more years without having to go back . It was better to go back , to be safe , so he would be safe for a while , and so he would n't get infected .

Table 8: Comparison of different methods.

We analyze some generated stories to evaluate

our overall model. An example is shown in Table 8. Stories generated by GPT2P have poor plot consistency and have some coreference errors, such as blue sentences and red words in Table 8. Compared with GPT2P, our model can effectively control the plot consistency of the story through the abstract. Therefore, stories generated by our model have better plot consistency. In addition, our model has less coreference errors than GPT2P and generates stories with better coreference consistency. What's more, the coherency between sentences is also better than GPT2P.

## 5 Conclusion

In this paper, we propose a two-stage generation model to improve consistency and coherency of generated stories. The first stage is to build the story outline, and the second stage is to expand the outline into a complete story. What's more, we design a supplementary task of discourse relation classification to improve the discourse representation ability of the model. In addition, we enhance model with coreference supervision to improve coreference consistency in generated stories. Experimental results on a story dataset show that our method is superior to baseline methods.

## References

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Hongxiao Bai and Hai Zhao. 2018. Deep enhanced representation for implicit discourse relation recognition. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 571–583, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Federico Barrios, Federico López, Luis Argerich, and Rosa Wachenchauzer. 2016. Variations of the similarity function of textrank for automated summarization. *arXiv preprint arXiv:1602.03606*.

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 Mc-Candlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Gang Chen, Yang Liu, Huanbo Luan, Meng Zhang, Qun Liu, and Maosong Sun. 2019. Learning to predict explainable plots for neural story generation. *arXiv preprint arXiv:1912.02395*.

Jifan Chen, Qi Zhang, Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. 2016. Implicit discourse relation detection via a deep architecture with gated relevance network. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1726–1735, Berlin, Germany. Association for Computational Linguistics.

Elizabeth Clark, Yangfeng Ji, and Noah A Smith. 2018. Neural text generation in stories using entity representations as context. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2250–2260.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898.

Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2650–2660, Florence, Italy. Association for Computational Linguistics.

Jian Guan, Fei Huang, Zhihao Zhao, Xiaoyan Zhu, and Minlie Huang. 2020. A knowledge-enhanced pre-training model for commonsense story generation. *Transactions of the Association for Computational Linguistics*, 8:93–108.

Parag Jain, Priyanka Agrawal, Abhijit Mishra, Mohak Sukhwani, Anirban Laha, and Karthik Sankaranarayanan. 2017. Story generation from sequence of independent short descriptions. *arXiv preprint arXiv:1707.05501*.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Man Lan, Jianxiang Wang, Yuanbin Wu, Zheng-Yu Niu, and Haifeng Wang. 2017. Multi-task attention-based neural networks for implicit discourse relationship representation and identification. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1299–1308, Copenhagen, Denmark. Association for Computational Linguistics.

Michael Lebowitz. 1987. Planning stories. In *Proceedings of the 9th annual conference of the cognitive science society*, pages 234–242.

Boyang Li, Stephen Lee-Urban, George Johnston, and Mark Riedl. 2013. Story generation with crowd-sourced plot graphs. In *Twenty-Seventh AAAI Conference on Artificial Intelligence*.

Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and Bill Dolan. 2016. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 994–1003. Association for Computational Linguistics.

Lara J Martin, Prithviraj Ammanabrolu, Xinyu Wang, William Hancock, Shruti Singh, Brent Harrison, and Mark O Riedl. 2018. Event representations for automated story generation with deep neural nets. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation 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*, pages 839–849, San Diego, California. Association for Computational Linguistics.

Allen Nie, Erin Bennett, and Noah Goodman. 2019. DisSent: Learning sentence representations from explicit discourse relations. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4497–4510, Florence, Italy. 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.

Rafael PÉrez Ý PÉrez and Mike Sharples. 2001. Mexica: A computer model of a cognitive account of creative writing. *Journal of Experimental & Theoretical Artificial Intelligence*, 13(2):119–139.

Julie Porteous and Marc Cavazza. 2009. Controlling narrative generation with planning trajectories: the role of constraints. In *Joint International Conference on Interactive Digital Storytelling*, pages 234–245. Springer.

Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Miltsakaki, Livio Robaldo, Aravind Joshi, and Bonnie Webber. 2008. The Penn discourse TreeBank 2.0. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco. European Language Resources Association (ELRA).

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8).

Mark O Riedl and Robert Michael Young. 2010. Narrative planning: Balancing plot and character. *Journal of Artificial Intelligence Research*, 39:217–268.

Stuart Rose, Dave Engel, Nick Cramer, and Wendy Cowley. 2010. Automatic keyword extraction from individual documents. *Text mining: applications and theory*, 1:1–20.

Abigail See, Aneesh Pappu, Rohun Saxena, Akhila Yerukola, and Christopher D. Manning. 2019. Do massively pretrained language models make better storytellers? In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 843–861, Hong Kong, China. Association for Computational Linguistics.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.

Jingjing Xu, Xuancheng Ren, Yi Zhang, Qi Zeng, Xiaoyan Cai, and Xu Sun. 2018. A skeleton-based model for promoting coherence among sentences in narrative story generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4306–4315.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*.

Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7378–7385.