

# STYLE EXAMPLE-GUIDED TEXT GENERATION USING GENERATIVE ADVERSARIAL TRANSFORMERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We introduce a language generative model framework for generating a styled paragraph based on a context sentence and a style reference example. The framework consists of a style encoder and a texts decoder. The style encoder extracts a style code from the reference example, and the text decoder generates texts based on the style code and the context. We propose a novel objective function to train our framework. We also investigate different network design choices. We conduct extensive experimental validation with comparison to strong baselines to validate the effectiveness of the proposed framework using a newly collected dataset with diverse text styles. Both code and dataset will be released upon publication.

## 1 INTRODUCTION

Text generation refers to the task of generating a new text based on some user input. The task exists in many forms, but arguably the most common form concerns generating a coherent and consistent text based on an input context such as the first few sentences of the target output. This is often achieved by giving the context to a generative language model. Generative language models play a central role in machine learning and natural language processing (NLP). Not only they serve as the main mean for unsupervised feature representation learning but also find use in various applications, including question answering, dialogue agents, summarization, and content creation systems.

Thanks to the introduction of novel deep learning architectures and the availability of large-scale training corpora, the state-of-the-art text generation has advanced significantly in recent years. We can now train language models capable of generating fluent and coherent texts that people cannot tell them apart from those written by humans. However, despite the great achievement, existing generative models are limited and inflexible in the sense that a trained model is only capable in generating texts of one style. It can not be used to generate texts of different styles. For instance, a news generative model can only be used to generate news, and a lyric generative model can only be used to generate lyrics. In contrast, humans can compose texts in various styles.

To bridge the gap, we propose a style example-guided text generation framework that can generate styled texts based on the style of the example reference text. In our framework, the generator takes two inputs where one is the context input while the other is the style reference example. We use the style reference example to change the generation behavior of our generative model dynamically. For a fixed context, when the provided style reference is a paragraph sampled from a news article, it becomes a news generator. When the provided style reference is a review, it becomes a review generator. In other words, the same generator can generate texts of different styles based on the examples. In Tab. 1, we show example outputs of the proposed framework where we generate texts of dramatically different styles for the same input sentence.

The proposed style example-guided text generation framework is based on the generative adversarial networks (GANs), and we utilize the transformer in both the generator and discriminator design. We collect a large dataset containing documents with many different styles for training. Using a novel learning objective function, our network learns to generate styled texts based on the input style example in an unsupervised manner. We conduct extensive experimental validations with comparisons to strong baselines. We also investigate different ways of designing the generator and compare their performance. Through detailed quantitative and user study results, we prove the effectiveness of the proposed framework for the style example-guided text generation task.

<p><b>Context – Wiki Style (Machine Learning Intro.)</b> Machine learning (ML) is the scientific study of algorithms and statistical models that computer systems use to perform a specific task without using explicit instructions, relying on patterns and inference instead. It is seen as a subset of artificial intelligence. Machine learning algorithms build a mathematical model based on sample data, known as "training data", in order to make predictions or decisions without being explicitly programmed to perform the task.[1][2].</p>	
<p><b>Reference – News Style</b> (CNN)Remember that robot dog that went viral a few years ago? The one that can run uphill, climb stairs, and looks like the killer robot dogs in the TV show "Black Mirror"? It's now available for select customers.</p> <p>Boston Dynamics, the company behind the dog, which is named Spot, released a video on Tuesday advertising its launch. It shows Spot navigating challenging terrain, picking up construction objects, moving through bad weather, and picking itself up after a fall. The video and the product launch page also show off Spot's capabilities – for instance, it can run at 1.6 meters per second (about 5.2 feet per second), has 360-degree cameras that give it a full range of vision, is dust- and water-proof, and can operate in temperatures ranging from -20 to 45 degrees Celsius (-4 to 113 Fahrenheit).</p>	<p><b>Reference – Yelp Reviews Style</b> Yumm</p> <p>The line was AWFUL! Came by Saturday around 5:20pm and the line was CRAZY. I do feel like they're understaffed. Despite that, it was my first time and i got the chance to taste test two flavors, matcha and milk tea! Both were great!!</p> <p>I got the : - Matcha swirl with a custard fish, sprinkles and coconut flakes : The custard was so good! Paired great with the ice cream.</p> <p>Staff were so friendly! Especially the cashier (guy) and person who swirled the ice cream that shift! She was super patient as I, myself is very indecisive, haha! Most definitely coming by again!!</p>
<p><b>Generated</b> Machine Learning is an ancient concept that has entered the mainstream today, which means that it can survive in the face of such heightened competition. In its application, Flurry wants to improve the efficiency of AI driven computing machines, and with a large footprint. As is seen on me, Flurry has designed a large extension that (unfortunately) can be used in its development, but with lower speed lag. Now that the hybrid technology is here, it's there for the taking.</p> <p>There may also be a Something About Mary in the future. "So it's nice to be able to mount a cat across your desk without requiring any backup access."</p>	<p><b>Generated</b> Since the day machines started boosting up, dozens of companies have recently done so. In fact, I highly recommend RoboSports, an example of small-batch ice cream places that have taken up artificial advances in AI: their ice cream is the perfect size, slim and hot, it was great. If you are a connoisseur of ice cream, check out this wonderful ice cream shop.</p>

Table 1: Given the same context sentences but different style reference examples, our generation model can generate paragraphs of different styles. The output paragraph shares the same style as the reference example.

## 2 RELATED WORK

**Language modeling** has seen many advancements in recent years, which has resulted in significant improvements on various NLP tasks. Early language models focused on using n-grams to represent a text distribution. Bengio et al. (2003) introduced a neural language model in a shift from more traditional n-gram models. Many works later (Mikolov et al. (2013); Pennington et al. (2014)) focused on word embeddings as a way to represent tokens within the text. More recently, Peters et al. (2018) used bi-directional LSTMs to obtain deep contextualized word representation. However, RNNs can only represent a limited context. Vaswani et al. (2017) introduced the transformer networks which use the connections between long-distance word pairs embedded in attention mechanisms and can easily enable the learning of long-term dependency. Many later models (Devlin et al. (2019); Liu et al. (2019d); Dai et al. (2019); Yang et al. (2019)) used transformer model and obtained significant improvements on downstream tasks (Wang et al. (2019); Rajpurkar et al. (2016); Zellers et al. (2018)). Lately, (Radford et al. (2019)) introduced GPT-2, a generative left-to-right language model based on the transformer and showed that these models are able to generate coherent text when pre-trained on a large corpus. Shoeybi et al. (2019) further scaled up the GPT-2 model and demonstrated improved performance. Our work differs from the prior works because we aim for allowing user flexible control over the style of the generated text.

**Texts generation** includes review generation (Radford et al. (2018); Zang & Wan (2017)), sentiment texts generation (Wang & Wan (2018); Hu et al. (2017); Merity et al. (2017)), Wikipedia generation (Liu et al. (2018); Lebret et al. (2016)), fake news generation (Bakhtin et al. (2019); Zellers et al. (2019)), abstractive summarization (Li et al. (2018); Zhang et al. (2019); Pasunuru et al. (2017)), and conversation/dialogue system (Vinyals & Le (2015); Budzianowski & Vulić (2019)). Although many of them trained a transformer on large-scale corpora, their results were limited in their specific domain (e.g., reviews, news, etc.) because they either utilized domain-specific priors in their model design or were not designed to generate texts in many different domains or styles.

**Control on texts generation.** In addition, there are literature utilizing insertion-base (Stern et al. (2019); Chan et al. (2019)), GAN-based (Yu et al. (2017); d’Autume et al. (2019)), variational autoencoder-based (Xu et al. (2019)), normalizing flow-based (Tran et al. (2019)) approaches for general texts generation task. However, we instead focus on generating styled paragraphs conditioning on a context and a reference paragraph. A recent work by Keskar et al. (2019) is most related to ours. They propose a conditional transformer using a control code to perform language generation in a sequence-to-sequence manner. We demonstrate our method outperforms theirs by a large margin in the experiment section.

**Text style transfer** concerns transferring an input text of one style to a different style (Kerpedjiev (1992); Rao & Tetreault (2018); Xu (2017); Xu et al. (2012); Fu et al. (2018); Hu et al. (2017);

Prabhumoye et al. (2018); Shen et al. (2017); Li et al. (2019)). Our work is different since we do not aim for changing the style of a given text. Instead, we aim for a style-controllable way for generating texts from scratch. Also, rather than handling transferring between two styles (e.g., positive  $\leftrightarrow$  negative sentiments), our model can generate texts of many different styles. Finally, our model outputs paragraphs while existing text style transfer works mostly output sentences.

**Image Style transfer** is a popular topic in computer vision. There are many successful techniques, including iterative optimization on the gram matrix (Gatys et al. (2016)), perceptual loss (Johnson et al. (2016); Gupta et al. (2017)), feature transformation (Li et al. (2017)), adaptive instance-normalization (Dumoulin et al. (2017); Huang & Belongie (2017)), and GAN-based methods (Zhu et al. (2017); Kim et al. (2017)). Our proposed framework also gets inspiration from them.

### 3 PRELIMINARIES

Our framework is based on the transformer network (Vaswani et al. (2017)) and the GAN framework (Goodfellow et al. (2014)). In this section, we briefly review these two components.

**Transformer** is the state-of-the-art network for various natural language processing tasks. Different from RNNs (Hochreiter & Schmidhuber (1997); Bengio et al. (2003); Chung et al. (2014)), which consume a sequence token by token, in a transformer network, the entire sequence is fed into layers of transformer modules. The representation of a token at a layer is then computed by attending to the latent representations of all the other tokens in the preceding layer.

Variants of transformer networks are available. We build our model based on GPT-2 transformer network (Radford et al. (2019); Shueybi et al. (2019)), which train a deep transformer using a left-to-right language model:

$$p(\mathbf{w}) = \prod_t^T p(w_t | w_{t-1} \dots w_1), \quad (1)$$

where  $w_t$ 's denote the word tokens. Different from BERT-like transformer networks (Devlin et al. (2019); Liu et al. (2019d)), GPT-2 makes a casual assumption, i.e., the latent representation of a token is calculated using only the latent representations of the preceding tokens. Thus, during generation, GPT-2 can be directly applied to complete the text given the context sentence.

**GAN** defines a zero-sum game played by a generator  $F$  and a discriminator  $D$ . Under some nice conditions, the generator learns to convert a random noise vector to a realistic signal in a way that the discriminator cannot tell it apart from real signals. In this case, the distribution of the output signals produced by the generator converges to the distribution of signals observed in the real world.

We use a conditional GAN where  $F$  takes a context sentence and a style reference example as inputs. To avoid non-differentiability in text decoding (e.g., beam search), we use a latent GAN formulation (Achlioptas et al. (2017)). We first divide  $F$  into a feature extractor  $F_f$  and an output embedding layer  $F_o$ ; that is  $F \equiv F_o \circ F_f$ . Now, instead of using the output text from  $F_o$  as the discriminator input, we feed the latent representation computed by  $F_f$  to the discriminator. For real text, we use a pretrained trained GPT-2 model  $H$ . Again, we decompose  $H$  into a feature extractor  $H_f$  and an output embedding layer  $H_o$  ( $H \equiv H_o \circ H_f$ ). The GAN discriminator then takes features extracted by  $H_f$  as input for real texts. Using this latent GAN formulation, we aim for aligning the feature distribution of our generator to the feature distribution of the pretrained GPT-2 model.

### 4 STYLE EXAMPLE-GUIDED TEXT GENERATION

We propose a language generative model framework that allows us to control style of the output text using a style reference example. Given few context sentences  $\mathbf{w} = \{w_t\}_{t=1}^T$  and a reference text  $\mathbf{s}$ , our generator  $F$  generates output text  $\mathbf{y}$  that has the same style as the reference example  $\mathbf{s}$  given by

$$\mathbf{y} = F(\mathbf{w}, \mathbf{s}) \equiv F_o(F_f(\mathbf{w}, \mathbf{s})). \quad (2)$$

We divide the feature extractor  $F_f$  into a style encoder  $F_s$  and a text decoder  $F_g$  where the style encoder extracts a style representation from the style example,  $\mathbf{z} = F_s(\mathbf{s})$ , and the text decoder  $F_g$  consumes the style representation and the context sentences to compute a feature for  $F_o$  to generate

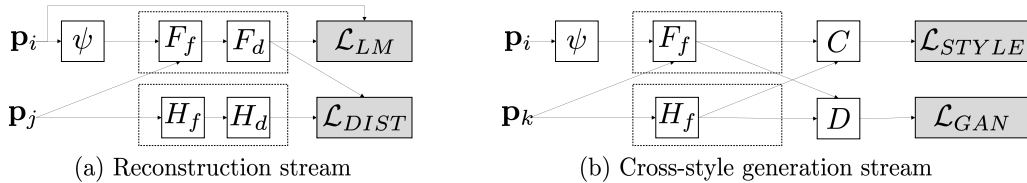


Figure 1: We employ two data streams to train our framework. While  $\mathbf{p}_i$  and  $\mathbf{p}_j$  have the same style,  $\mathbf{p}_i$  and  $\mathbf{p}_k$  do not. (a) The reconstruction stream is trained using the language modeling loss  $\mathcal{L}_{LM}$  and the distillation loss  $\mathcal{L}_{DIST}$ . (b) The cross-style generation stream is trained using the style loss  $\mathcal{L}_{STYLE}$  and the GAN loss  $\mathcal{L}_{GAN}$ . Note that we decompose each network into a feature extractor and an embedding layer.

the styled text  $\mathbf{y}$ . In this section, we will first introduce the data streams employed during training and our novel learning objective function. We will then discuss various generator design choices.

#### 4.1 LEARNING DATA STREAMS

Let  $\mathbb{D} = \{(\mathbf{d}_n, l_n)\}$  be a dataset of documents where  $\mathbf{d}_n$  is a document and  $l_n$  is its style label. We assume a finite set of style labels  $\mathcal{L} = \{1, 2, \dots, L\}$  where each integer represents a style class such as news, review, lyric, poem, novel, and children book. During training, our framework employs two data streams where the first one is called the reconstruction stream while the other is referred to as the cross-style generation stream. We note that such a two-stream processing pipeline is common in GAN-based image translation frameworks (Liu et al. (2017); Huang et al. (2018); Liu et al. (2019a)) but is less explored for language modeling.

**Reconstruction stream (RS).** For this stream, we first sample two documents with the same style from  $\mathbb{D}$ :  $(\mathbf{d}_i, l_i)$  and  $(\mathbf{d}_j, l_j)$  where  $l_i = l_j$ . We then sample two paragraphs<sup>1</sup>:  $\mathbf{p}_i \sim \mathbf{d}_i$  and  $\mathbf{p}_j \sim \mathbf{d}_j$ . We extract the first few sentences from  $\mathbf{p}_i$  as the input context  $\mathbf{w} = \psi(\mathbf{p}_i)$ , where  $\psi$  is the extraction function, and use  $\mathbf{p}_j$  for the style reference  $\mathbf{s}$ . Feeding  $\mathbf{w}$  and  $\mathbf{p}_j$  to the generator  $F$ , we expect  $F$  should be able to reconstruct  $\mathbf{p}_i$ :  $F(\psi(\mathbf{p}_i), \mathbf{p}_j) \approx \mathbf{p}_i$ .

**Cross-style generation stream (CS).** We first sample two documents  $(\mathbf{d}_i, l_i) \sim \mathcal{D}$  and  $(\mathbf{d}_k, l_k) \sim \mathcal{D}$  where  $l_i \neq l_k$ . We then sample paragraphs  $\mathbf{p}_i \sim \mathbf{d}_i$  and  $\mathbf{p}_k \sim \mathbf{d}_k$ . We again extract the first few sentences from  $\mathbf{p}_i$  as the input context  $\mathbf{w} = \psi(\mathbf{p}_i)$  and use  $\mathbf{p}_k$  for the style reference  $\mathbf{s}$ . As feeding  $\mathbf{w}$  and  $\mathbf{p}_k$  to the generator  $F$ , we expect  $F$  should output  $\mathbf{p}_{i \rightarrow k} = F(\psi(\mathbf{p}_i), \mathbf{p}_k)$  where  $\mathbf{p}_{i \rightarrow k}$  should have the same style as  $\mathbf{d}_k$ . Let  $C^*$  be an oracle style comparator function that outputs 1 if the two input texts have the same style and 0 otherwise. We aim for  $C^*(\mathbf{p}_{i \rightarrow k}, \mathbf{p}_k) = 1$ .

#### 4.2 LEARNING OBJECTIVE

We propose an objective function consisting of four carefully designed loss terms for training the proposed framework using the above two data streams. The objective function is given by

$$\mathcal{L} = \mathcal{L}_{LM} + \lambda_{DIST} \mathcal{L}_{DIST} + \lambda_{STYLE} \mathcal{L}_{STYLE} + \lambda_{GAN} \mathcal{L}_{GAN}, \quad (3)$$

where  $\mathcal{L}_{LM}$  is the language modeling loss,  $\mathcal{L}_{DIST}$  the distillation loss,  $\mathcal{L}_{STYLE}$  is a style comparison loss, and  $\mathcal{L}_{GAN}$  is the latent GAN loss. The scalars  $\lambda_{DIST}$ ,  $\lambda_{STYLE}$ , and  $\lambda_{GAN}$  are the hyper-parameters controlling relative importance of the terms. The values for these hyperparameters and the method for determining their values are discussed in Appendix A. We visualize training with the proposed objective function using the two data streams in Fig. 1.

**Language modeling loss  $\mathcal{L}_{LM}$**  formulates the probability distribution of a paragraph  $\mathbf{p}$  as the product of the conditional probability of each token  $w_t$  given the previous tokens  $\{w_t\}_t^{T-1}$  as shown in (1). We use  $\mathcal{L}_{LM}$  to supervise the training of the data reconstruction stream. It is given by

$$\mathcal{L}_{LM} = E_{(\mathbf{p}_i, \mathbf{p}_j) \sim \text{RS}} \left[ \frac{-1}{T} \sum_t \log \left( \frac{e^{F(w_t | w_{1:t-1}, \mathbf{p}_j)}}{\sum_v e^{F(w_v | w_{1:t-1}, \mathbf{p}_j)}} \right) \right], \quad (4)$$

where  $(\mathbf{p}_i, \mathbf{p}_j) \sim \text{RS}$  denotes that  $\mathbf{p}_i$  and  $\mathbf{p}_j$  are from the reconstruction stream. The variable  $T$  is the total number of tokens in  $\mathbf{p}_i$  and  $V$  is the size of the vocabulary.

<sup>1</sup>For the purpose of data augmentation, in our implementation, a paragraph we sample may not be the full paragraph in the nominal sense. It could start from the middle of a nominal paragraph.

**Distillation loss.** We use  $\mathcal{L}_{DIST}$  to regularize the learning as processing the data reconstruction stream. We pretrain a GPT-2 model using our dataset  $\mathbb{D}$  and use it as our distillation target. We denote the pretrained GPT-2 model as  $H$ . (Note that  $H$  does not have the desired style control capability.) By jointly optimizing  $\mathcal{L}_{LM}$  and  $\mathcal{L}_{DIST}$ , we train  $F$  to generate fluent texts (by minimizing  $\mathcal{L}_{LM}$ ) as well as behave similarly to  $H$  (by minimizing  $\mathcal{L}_{DIST}$ ). The distillation loss is calculated by minimizing the mutual information between output distributions of  $F$  and  $H$ , which is given by

$$\mathcal{L}_{DIST} = E_{(\mathbf{p}_i, \mathbf{p}_j) \sim \text{RS}} \left[ \frac{-1}{T} \sum_t \sum_v^{|V|} \frac{e^{H(w_{t=v}|w_{1:t-1})}}{\sum_v^{|V|} e^{H(w_v|w_{1:t-1})}} \log \left( \frac{e^{F(w_t|w_{1:t-1}, \mathbf{p}_j)}}{\sum_v^{|V|} e^{F(w_v|w_{1:t-1}, \mathbf{p}_j)}} \right) \right]. \quad (5)$$

We note that the distillation loss has been used in various tasks including model compression, transfer learning, life-long learning, etc (Hinton et al. (2015); Kim & Rush (2016); Liu et al. (2019c); Mirzadeh et al. (2019); Liu et al. (2019b); Hou et al. (2018)). In this paper, we extend its use to the style example-guided language generative model training task.

**Style loss**  $\mathcal{L}_{STYLE}$  helps ensure the output from the cross-style generation stream has the same style as the input reference. A pretrained style comparator  $C$  is used for computing the loss. The comparator takes two paragraphs as input and is trained to output 1 when the two paragraphs have the same style and 0 otherwise. We use  $\mathbb{D}$  for pretraining  $C$  since it contains style labels for each document. We pretrain  $C$  using the binary cross entropy loss. The comparator  $C$  is highly accurate. It achieves a classification accuracy of 87.8% to 98.8% in our held-out validation sets. After pretraining, we fix  $C$  and use it to train  $F$ . The style loss  $\mathcal{L}_{STYLE}$  is then given by

$$\mathcal{L}_C = E_{(\mathbf{p}_i, \mathbf{p}_k) \sim \text{CS}} \left[ -\log \left( C(H_f(\mathbf{p}_k), F_f(\psi(\mathbf{p}_i), \mathbf{p}_k)) \right) \right] \quad (6)$$

where  $(\mathbf{p}_i, \mathbf{p}_k) \sim \text{CS}$  denotes the pair is sampled from the cross-style generation stream.

Here, we would like to make two remarks. First, since  $C$  takes the latent feature from  $F_f$  as input, we avoid the non-differentiability of the text decoding mechanism and can directly train  $F_f$ . Second, despite that  $C$  is pretrained using feature extracted from  $H_f$ , we use the feature extracted from  $F_f$  as input. We can perform this operation not only because these two features have the same dimension but also because we enforce them to have a similar distribution via optimizing the GAN loss, discussed below.

**GAN loss**  $\mathcal{L}_{GAN}$  is used to match the distributions of the features generated by  $F_f$  and those generated by  $H_f$ , respectively, as processing the cross-style generation stream. We use a latent GAN formulation where we train a GAN discriminator  $D$  to differentiate features extracted from  $F_f$  to  $H_f$ . The GAN loss is given by

$$E_{\mathbf{p}_i, \mathbf{p}_k} \left[ -\log(D(H_f(\mathbf{p}_k))) - \log(1 - D(F_f(f(\mathbf{p}_i), \mathbf{p}_k))) \right]. \quad (7)$$

We realize the discriminator  $D$  using a -based transformer network.

### 4.3 GENERATOR DESIGN

We realize the style encoder  $F_s$  using a GPT-2-based transformer identical to  $H_f$ . After extracting a representation  $\mathbf{z}_t$  for each token  $t$  in  $\mathbf{s}$ , we utilize a 3-layer position-wise fully-connected network to obtain the final style code  $\mathbf{z}$  as illustrated in Fig. 2. The text decoder  $F_g$  is also a GPT-2-based transformer identical to  $H$ . We initialize the weights in  $F_s$  and  $F_g$  using the weights in the pretrained  $H$ . Next, we compare four different ways of injecting outputs from  $F_s$  into  $F_g$ , which represent different inductive biases and result in difference performances.

**Model A: style code as a bias to the input.** In this model, the style code  $\mathbf{z} = F_s(\mathbf{s})$  is directly summed up with the token-embedding and position embedding before inputting to the first transformer module in  $F_g$ . In other words, the input to the first transformer module in  $F_g$  is  $\mathbf{e}_t^w + \mathbf{e}_t^p + \mathbf{z}$  where  $\mathbf{e}_t^w$  denotes as the  $t$ th word embedding, and  $\mathbf{e}_t^p$  denotes as the  $t$ th position embedding.

**Model B: style code as a summarization token.** In this model, the computed style code  $\mathbf{z} = F_s(\mathbf{s})$  is treated as a special token that is inserted to the beginning of the input sequence and is directed fed in the first transformer module in  $F_g$ . That is the input sequence length becomes  $T + 1$ . This design is motivated by the traditional sequence-to-sequence modeling techniques (Chung et al. (2014); Cho et al. (2014); Sutskever et al. (2014); Bahdanau et al. (2016); Vinyals & Le (2015)).

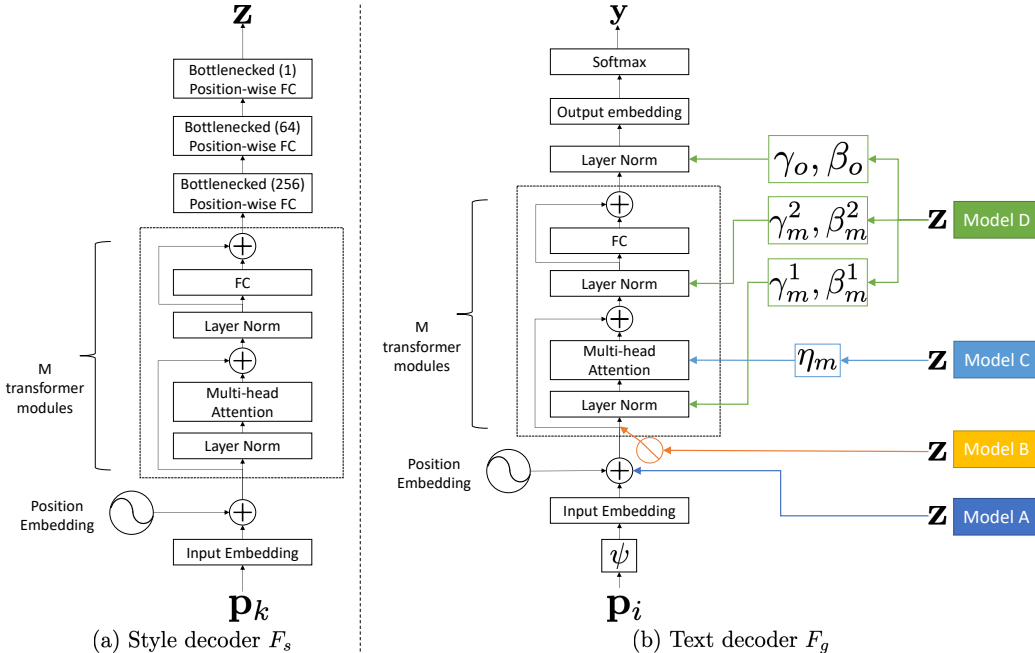


Figure 2: (a) Style encoder  $F_s$ . (b) Text decoder  $F_g$ . We consider 4 different ways of injecting the style code  $\mathbf{z}$  into  $F_g$  termed Model A, B, C, and D. Checkout main texts for more details.

**Model C: style-aware self-attention.** In this model, we input  $\mathbf{z}$  into each self-attention layer in  $F_g$  to influence its computation given by  $\text{Softmax}(\frac{\mathbf{q}_m \mathbf{k}_{m-1}^T}{\sqrt{B}}) \mathbf{v}_{m-1}$  where  $\mathbf{q}_m = \eta_m(\mathbf{z})$  which  $\eta_m$  denotes an affine transformation,  $\mathbf{k}_{m-1}$  and  $\mathbf{v}_{m-1}$  denotes the key and value embeddings from the  $(m - 1)$ th hidden layer, and  $B$  denotes the hidden dimension.

**Model D: adaptive layer normalization.** Inspired by the recent success in image generation tasks (Park et al. (2019); Karras et al. (2019)), we utilize the style code to modulate the hidden representations within the text decoder via normalization layers. Specifically, we replace the scale and bias parameters in the affine transformation step of the layer normalization (Ba et al. (2016)) with a style code determined scale and bias. That is

$$\gamma_{m,c}^a(\mathbf{z}) \frac{h_{m,c,t}^a - \mu_{m,t}^a}{\sigma_{m,t}^a} + \beta_{m,c}^a(\mathbf{z}), \tag{8}$$

where  $h_{m,c,t}^a$  denotes the  $c$ th hidden representation of the  $t$ th token at the  $m$ th transformer layer. We note  $a = \{1, 2\}$  since there are two layer normalization layers in each transformer in our implementation. The mean and deviation  $\mu_{m,t}^a$  and  $\sigma_{m,t}^a$  are computed across the channel dimension.

We illustrate how these models inject  $\mathbf{z}$  to  $F_g$  in Fig. 2. In Section 5, we compare the performance of these variants and show that Model D achieves the best style generation performance.

## 5 EXPERIMENTS

**Implementation.** We set the latent dimension  $B$  to 768, number of attention-heads to 16, number of transformer layers  $M$  to 16, number of tokens in a paragraph  $T$  to 512, and the vocabulary size  $V$  to 50257 using BPE-encoding (Sennrich et al. (2015)) vocabulary from Radford et al. (2019) throughout out all the models and experiments. We use a pretrained GPT-2 model  $H$  and a style comparator  $C$  in our framework. The training details of these two models are given in Appendix B. All of the experiments are conducted using an NVIDIA DGX1 machine.

**Datasets.** We compare competing methods using two newly composed datasets based on (Zhu et al. (2015); Zellers et al. (2019); Santiago (2015); See et al. (2017)).

**3-Style.** The dataset consists of documents from the RealNews dataset (Zellers et al. (2019)), the BookCorpus dataset (Zhu et al. (2015)), and the Reviews dataset (Yelp (2019); McAuley &

Leskovec (2013); Maas et al. (2011); Dataworld (2017); Liu (2017)). The 3 styles are news, book, and review. In detail, the news set has  $33M$  documents and  $113B$  words, the books set has  $50K$  documents and  $7.2B$  words, and the review set has  $4.8M$  documents and  $5.4B$  words after cleaning. The total dataset has  $37.85M$  documents and  $125.6B$  words. We hold out  $3.78M$  documents as the validation set and  $6K$  documents as the testing set.

**21-Style.** We build a dataset that contains 21 text styles. We first classify the documents in Real-News into 9 styles, including Sciences, Sport, Politics, Business, Technology, Entertainment, Opinion, Life, and News. Then, we divide the documents in BookCorpus into 8 different styles, which are Romance, Fantasy, Sciencefiction, Childrensbooks, Thriller, Adventure, Poetry, and Plays. We split the documents into multiple small documents by extracting the dialogues except for the Poetry and Plays. We divide the Review dataset into 3 styles, namely Yelp, Hotel, and Movie. Finally, we crawl  $0.77M$  lyrics from <http://www.azlyrics.com/>. The total dataset has  $35.5M$  documents. We hold out  $3.55M$  documents as the validation set and  $21K$  documents as the testing set.

**Auto-evaluation metrics.** We evaluate different models using fluency score, style score, style diversity score, and content novelty score. The fluency score measures whether the output paragraph reads like a human-written one. The style score checks whether the output text carries the target style. Our framework supports multimodal outputs (Huang et al. (2018)). For the same input context but different reference examples of the same style, our framework should produce different output texts but all with the same style. To measure how different these outputs are, we use the style diversity score. Finally, the content novelty score is used to measure the difference between the output and the reference example. A model that directly duplicates the reference to the output is undesirable. The details of these automatic evaluation metrics are available in Appendix C.

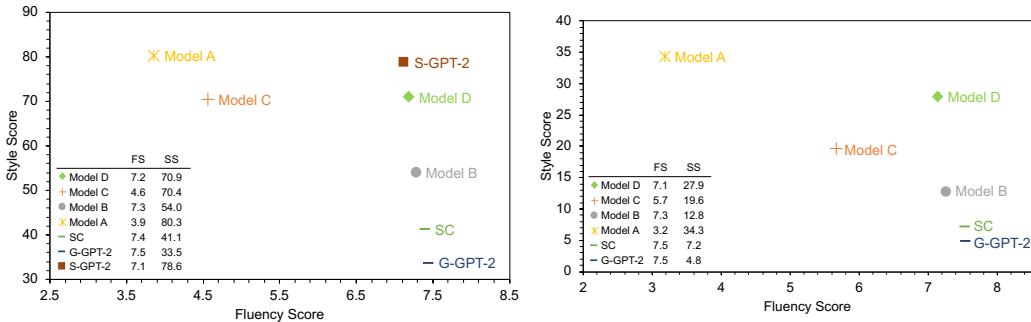
**Human study settings.** We use the Amazon Mechanical Turk (AMT) platform for user studies. We conduct two studies where one evaluates fluency of the generated paragraphs while the other verifies the style correctness. For the fluency study, we present a human-written text and a machine-generated text in random order and ask the worker to choose which one is written by a human. For this metric, the closer the preference score to 50%, the better the performance.

For the style study, we perform two tests. In one test, we present a worker a generated paragraph that supposes to be in the target style. We also give the worker two human-written reference paragraphs where one is with the target style while the other is not. We then ask the worker to choose which reference paragraph has a style more similar to the generated one. In the other test, we again present a worker a generated paragraph but this time with the style categorical labels to choose from instead of the reference paragraphs. We compute the frequency that the worker selects the right style. The higher the score, the better the performance. More details are in Appendix D.

**Strong baselines.** We compare our framework to three strong baselines, namely the general GPT-2 model (G-GPT-2), a baseline consists of multiple style-specialized GPT-2 models (S-GPT-2), and the style-code encoding (SC) method based on the description in Keskar et al. (2019). G-GPT-2 is trained on the entire dataset using  $\mathcal{L}_{LM}$ . It does not allow style control but can generate fluent texts. In S-GPT-2, we train a GPT-2 model per style. As training a GPT-2 model is costly, we only use this baseline for the 3-Style dataset evaluation. In SC, an one-hot encoding of the style class label is used as a special token for style-controllable paragraph generation. Unlike the proposed method that extracts the style code from the input paragraph, SC input the style label. The rest of the model is similar to our Model B without the style encoder.

## 5.1 RESULTS

In Fig. 3, we plot the fluency and style scores achieved by our models as well as those by the baselines on the 3-Style and 21-Style datasets. The closer the model to the top-right corner, the more superior the model is. From the figure, we found that among our models, Model D performs the best. As expected, G-GPT-2 achieves the best fluency score. However, since it does not support style control, it has a poor style score. On the other hand, S-GPT-2 achieves good fluency and style scores for the 3-Style dataset. This is understandable as it utilizes a GPT-2 model for each style. However, such an approach does not scale well as GPT-2 training is expensive. We also found that SC does not achieve good style score and is inferior to our models. We suspect this is because the one-hot style class code is largely ignored during inference. Since Model D performs the best



(a) Fluency vs. style on the 3-Style dataset. (b) Fluency vs. style on the 21-Style dataset.

Figure 3: Fluency and style scores achieved by the competing models on the 3-Style and 21-Style datasets.

	Model A	Model B	Model C	Model D
3-Style	<b>11.53</b>	11.61	10.58	11.13
21-Style	10.17	10.09	<b>11.32</b>	10.52

Table 2: Style diversity scores achieved by the computing methods. We note the lower bound and upper bound for the style diversity scores are 4.52 and 15.67, respectively.

(%)	Model D	SC	Random
3-Style	56	<b>54</b>	50
21-Style	<b>57</b>	63	50

(%)	Model D	SC	Random
3-Style by reference	<b>56</b>	52	50
3-Style by category	<b>65</b>	54	50
21-Style by reference	<b>66</b>	49	50
21-Style by category	<b>69</b>	50	50

Table 3: (Left): Human study results on fluency. (Right): Human study results on style control. Random denotes the accuracy for random guess. Model D performs favorably over the baseline SC.

Model D	Fluency Score	Style Score	Style Diversity Score	Content Novelty Score
$\mathcal{L}_{DIST}$	7.32	51.16	11.40	24.01
$\mathcal{L}_{STYLE}$	7.35	5.40	9.22	29.27
$\mathcal{L}_{GAN}$	6.85	28.67	10.35	26.77
All	7.14	27.90	10.52	25.85

Table 4: Ablation study on the various loss terms in the proposed objective function.

in our framework, for the rest of the paper, we use it as our representative model for performance comparison as well as ablation study.

In Tab. 1, we show several generation results from our Model D. We find that the output texts are fluent and respect the styles of the references. More output examples are available in Appendix E.

In Tab. 2, we show the style diversity scores achieved by our models. We found that all of our 4 models can generate diverse styled paragraphs conditioning on the same context and different reference examples with the same style.

**Human evaluation.** In Tab. 3, we report user study results on fluency and style control. We found that our model achieves great fluency on both of the datasets. Compared to SC, our model performs better in controlling the style in the output texts.

**Ablation study.** We conduct an ablation study on the loss terms in the proposed objective function and report the results in Tab. 4 using the 21-Style dataset. The results show that each term is important. Removing  $\mathcal{L}_{DIST}$  leads to a degraded content novelty score. Removing  $\mathcal{L}_{STYLE}$  leads to a degraded style score, though an improved fluency score and a content novelty score. Removing  $\mathcal{L}_{GAN}$  leads to both degraded fluency and style diversity scores.

## 6 CONCLUSION

We presented a language generative framework for style example-guided paragraph generation. To the best of our knowledge, we were the first to achieve such style-controllability on paragraph generation. We attributed the success to our carefully designed learning objective function, the generator network, and the newly composed large-scale dataset consisting of documents of various text styles.



## REFERENCES

- Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. *arxiv preprint arXiv:1707.02392*, 2017.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *International Conference on Learning Representations (ICLR)*, 2016.
- Anton Bakhtin, Sam Gross, Myle Ott, Yuntian Deng, Marc’Aurelio Ranzato, and Arthur Szlam. Real or fake? learning to discriminate machine from human generated text. *arxiv preprint arXiv:1906.03351*, 2019.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic language model. *Journal of Machine Learning Research (JMLR)*, 2003.
- Paweł Budzianowski and Ivan Vulić. Hello, it’s gpt-2—how can i help you? towards the use of pretrained language models for task-oriented dialogue systems. *arxiv preprint arXiv:1907.05774*, 2019.
- William Chan, Nikita Kitaev, Kelvin Guu, Mitchell Stern, and Jakob Uszkoreit. Kermit: Generative insertion-based modeling for sequences. *arxiv preprint arXiv:1906.01604*, 2019.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2014.
- Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arxiv preprint arXiv:1412.3555*, 2014.
- Zihang Dai, Zhilin Yang, Yiming Yang, William W Cohen, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- Dataworld. Hotel reviews dataset. In <https://data.world/datafiniti/hotel-reviews>, 2017.
- Cyprien de Masson d’Autume, Mihaela Rosca, Jack Rae, and Shakir Mohamed. Training language gans from scratch. In *Advances in Neural Information Processing Systems (NIPS)*, 2019.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019.
- Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. In *International Conference on Learning Representations (ICLR)*, 2017.
- Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, and Rui Yan. Style transfer in text: Exploration and evaluation. In *Association for the Advancement of Artificial Intelligence (AAAI)*, 2018.
- Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- Agrim Gupta, Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Characterizing and improving stability in neural style transfer. In *IEEE International Conference on Computer Vision (ICCV)*, 2017.

- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arxiv preprint arXiv:1503.02531*, 2015.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 1997.
- Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Lifelong learning via progressive distillation and retrospection. In *European Conference on Computer Vision (ECCV)*, 2018.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P Xing. Toward controlled generation of text. In *International Conference on Machine Learning (ICML)*, 2017.
- Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. *European Conference on Computer Vision (ECCV)*, 2018.
- Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision (ECCV)*, 2016.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- Stephan M. Kerpeldjiev. Generation of informative texts with style. In *International Conference on Computational Linguistics (COLING)*, 1992.
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. Ctrl: A conditional transformer language model for controllable generation. In *Einstein.ai*, 2019.
- Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jungkwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. In *International Conference on Machine Learning (ICML)*, 2017.
- Yoon Kim and Alexander M Rush. Sequence-level knowledge distillation. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2016.
- Rémi Lebret, David Grangier, and Michael Auli. Neural text generation from structured data with application to the biography domain. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2016.
- Chenliang Li, Weiran Xu, Si Li, and Sheng Gao. Guiding generation for abstractive text summarization based on key information guide network. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2018.
- Dianqi Li Li, Yizhe Zhang, Zhe Gan Gan, Yu Cheng, Chris Brockett, Ming-Ting Sun, and Bill Dolan. Domain adaptive text style transfer. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2019.
- Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Jason Liu. 515k hotel reviews data in europe. In <https://www.kaggle.com/jiashenliu/515k-hotel-reviews-data-in-europe>, 2017.
- Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In *IEEE International Conference on Computer Vision (ICCV)*, 2019a.

- Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. In *International Conference on Learning Representations (ICLR)*, 2018.
- Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Improving multi-task deep neural networks via knowledge distillation for natural language understanding. *arxiv preprint arXiv:1904.09482*, 2019b.
- Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Multi-task deep neural networks for natural language understanding. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019c.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arxiv preprint arXiv:1907.11692*, 2019d.
- Andrew L Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2011.
- Julian John McAuley and Jure Leskovec. From amateurs to connoisseurs: modeling the evolution of user expertise through online reviews. In *International World Wide Web Conference (WWW)*, 2013.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *International Conference on Learning Representations (ICLR)*, 2017.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- Seyed-Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, and Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant: Bridging the gap between student and teacher. *arxiv preprint arXiv:1902.03393*, 2019.
- Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially adaptive normalization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- Ramakanth Pasunuru, Han Guo, and Mohit Bansal. Towards improving abstractive summarization via entailment generation. In *Annual Meeting of the Association for Computational Linguistics (ACL) Workshop*, 2017.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. volume 14, pp. 1532–1543, 01 2014. doi: 10.3115/v1/D14-1162.
- Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2018.
- Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, and Alan W Black. Style transfer through back-translation. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2018.
- Alec Radford, Rafal Jozefowicz, and Ilya Sutskever. Learning to generate reviews and discovering sentiment. In *International Conference on Learning Representations (ICLR)*, 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. In *OpenAI Blog*, 2019.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2016.

- Sudha Rao and Joel Tetreault. Dear sir or madam, may i introduce the yafc corpus: Corpus, benchmarks and metrics for formality style transfer. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2018.
- Santiago. Wikipedia xml data. In *Amazon*, 2015.
- Abigail See, Peter J Liu, and Christopher D Manning. Get to the point: Summarization with pointer-generator networks. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2017.
- Stanislau Semeniuta, Aliaksei Severyn, and Sylvain Gelly. On accurate evaluation of gans for language generation. *arxiv preprint arXiv:1806.04936*, 2018.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2015.
- Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. Style transfer from non-parallel text by cross-alignment. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 6833–6844, 2017.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using gpu model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- Mitchell Stern, William Chan, Jamie Kiros, and Jakob Uszkoreit. Insertion transformer: Flexible sequence generation via insertion operations. In *International Conference on Machine Learning (ICML)*, 2019.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- Dustin Tran, Keyon Vafa, Kumar Krishna Agrawal, Laurent Dinh, and Ben Poole. Discrete flows: Invertible generative models of discrete data. In *International Conference on Learning Representations (ICLR) Workshop*, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Oriol Vinyals and Quoc Le. A neural conversational model. *arxiv preprint arXiv:1506.05869*, 2015.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations (ICLR)*, 2019.
- Ke Wang and Xiaojun Wan. Sentigan: Generating sentimental texts via mixture adversarial networks. In *International Joint Conferences on Artificial Intelligence (IJCAI)*, 2018.
- Peng Xu, Yanshuai Cao, and Jackie Chi Kit Cheung. Unsupervised controllable text generation with global variation discovery and disentanglement. *arxiv preprint arXiv:1905.11975*, 2019.
- Wei Xu. From shakespeare to twitter: What are language styles all about? In *Empirical Methods in Natural Language Processing (EMNLP) Workshop*, 2017.
- Wei Xu, Alan Ritter, Bill Dolan, Ralph Grishman, and Colin Cherry. Paraphrasing for style. *International Conference on Computational Linguistics (COLING)*, 2012.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *arxiv preprint arXiv:1906.08237*, 2019.
- Yelp. Yelp dataset challenge. In <https://www.yelp.com/dataset/challenge>, 2019.
- Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. Seqgan: Sequence generative adversarial nets with policy gradient. In *Association for the Advancement of Artificial Intelligence (AAAI)*, 2017.

- Hongyu Zang and Xiaojun Wan. Towards automatic generation of product reviews from aspect-sentiment scores. In *International Conference on Natural Language Generation (INLG)*, 2017.
- Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. Swag: A large-scale adversarial dataset for grounded commonsense inference. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. Defending against neural fake news. In *Advances in Neural Information Processing Systems (NIPS)*, 2019.
- Haoyu Zhang, Yeyun Gong, Yu Yan, Nan Duan, Jianjun Xu, Ji Wang, Ming Gong, and Ming Zhou. Pretraining-based natural language generation for text summarization. *arXiv preprint arXiv:1902.09243*, 2019.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.

## A HYPER-PARAMETERS TUNNING

We tune the hyper-parameters on a pilot-test dataset. This pilot-test dataset has  $3K$  training examples and 400 hold-out examples. We perform a grid search on log-scale. We utilize the Frechet Embedding Distance (Semeniuta et al. (2018)) to select best hyper-parameters. For  $\lambda_{DIST}$ ,  $\lambda_{STYLE}$ , and  $\lambda_{GAN}$ , we try  $\{1.0, 0.1, 0.1\}$ ,  $\{0.1, 0.1, 0.1\}$ , and  $\{1.0, 0.01, 0.01\}$ . For *betas* used in Adam Optimizer, we try  $\{0.9, 0.999\}$  and  $\{0.0, 0.9\}$ . For weight decay, we try 0.01 and 0.0. For the initial learning rate, we try 0.00025 and 0.000025. Eventually, we use 0.00025 as initial learning rate and  $\{1.0, 0.01, 0.01\}$  for  $\lambda_{DIST}$ ,  $\lambda_{STYLE}$ , as well as  $\lambda_{GAN}$  for all models except for the Model B. We use  $\{1.0, 0.01, 0.1\}$  for the Model B instead.

## B PRETRAINING OF $H$ AND $C$

**Pretraining of the GPT-2 model  $H$ .** We pretrain  $H$  on our collected dataset  $\mathbb{D}$  from scratch. We use Adam optimizer with 0.00015 initial learning rate,  $\beta_1$  and  $\beta_2$  are set by  $\{0.9, 0.999\}$ , cosine learning rate decay style, and 0.01 weight decay. The batch size is set to 512. The total training iterations is  $320K$  which takes 2 weeks.

**Pretraining of the style comparator  $C$ .** We pretrain the Style Comparator  $C$  using 256 and 512 batch size, respectively. The initial learning rate is 0.00015 with 0.01 weight decay and cosine learning rate decay style. The optimizer is also Adam with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Since the training converges very quickly, we do early stopping if the accuracy on validation set plateaus. Eventually, we use the checkpoint at  $28K$  and  $50K$  iterations to train on the 3-style and 21-style datasets, respectively. The accuracy on hold-out validation set are 98.8% and 87.8% on 3-styles and 21-styles datasets, respectively.

## C AUTO-EVALUATION METRICS

**Fluency score.** To ensure the generated paragraph is fluent and coherent, we utilize a pretrained GPT-2 model to measure the perplexity of the generated paragraph. We compute the generation likelihood over each token using the model and treat the generated paragraph  $y_{0:T-1}$  as inputs and  $y_{1:T}$  as labels. Given the input text, the pretrained GPT-2 returns the probability distribution of next token over the vocabulary. Then, we measure the perplexity by this probability distribution and label. Since our dictionary size is 50257, the random guess of the next token would result in 50257 perplexity ( $\mathcal{L}_{LM} = \log(50257) \approx 10.83$ ). Thus, we set 10.83 as an upper bound and define the fluency score of the generated paragraph as  $10.83 - \log(perplexity)$ . In this sense, a higher fluency score means lower perplexity.

**Style score.** We train 3/21 binary style classifiers (since we have 3/21 different styles in the 3-Style/21-Style dataset) by finetuning a GPT-2 network to automatically evaluate whether the generated text carries the style of a target class. These 3/21 classifiers achieve average classification accuracies of 99.1%/96.3%. During the testing phase, for a target style, if the corresponding style classifier correctly predicts 1 for the generated paragraph computed by a model, we count it as a successful trial. We compute the success rate over the test set and use the result as the style score for the model.

**Style diversity score.** We adopt the LPIPS distance (Zhang et al. (2018)) to measure the diversity of the generation outputs conditioning on the same context. To implement this metric, we first extract the feature representation from each token in a generated paragraph by a pretrained GPT-2 model. We compute the mean representation of the tokens in a paragraph as the paragraph-level representation. Then, we measure the  $L_2$  distance between two paragraph-level representations of two different paragraphs generated using the same context but two different references written in the same style. In this sense, a larger distance value implies the styles of the two generated paragraphs are more different.

To get an idea of the range of this metric, we compute an upper bound and a lower bound. We consider two paragraphs from two documents of different styles should have a high style diversity score. We hence sample 1000 paragraphs from each style and use the pretrained GPT-2 model to extract deep features. After taking average over the token-dimension, we obtain 1000 by 768

representation for each style. Then, we compute the  $\mathcal{L}_2$  distance between of these matrices divided by 1000. This gives us a matrix of size  $21 \times 21$  measuring the pairwise distance between two styles. We use the largest value in this matrix as our upper bound, which is 15.67.

For the lower bound, since two different paragraphs from the same document should have a low style diversity score, we use their scores for the lower bound computation. Specifically, we compute the average distance between two different paragraphs from the same document. We do this for each style and obtain 21 different values. We obtain the lower bound by taking average over these values, which is 4.52.

**Content novelty score.** To verify that our model is not simply duplicating the content from reference paragraph, we utilize LPIPS distance (Zhang et al. (2018)) to measure the difference between the generated paragraph and the input reference paragraph. We again use a pretrained GPT-2 model for extracting a feature representation for each token. To compute the distance between two paragraphs, we compute the bipartite matching cost between the tokens from the two paragraphs. Specifically, we first compute the  $L2$  distances between any token representation in one paragraph to all the token representations in the other paragraph. We then compute the minimum cost assignment by solving a bipartite matching problem. In order to get an idea about the range of the content novelty score, we compute an upper bound and a lower bound using a similar approach as the one used for the style diversity score. We find the upper bound value is 31.46 and the lower bound value is 15.13.

## D HUMAN EVALUATION

To participate in our user study, a worker has to be awarded the Master Qualification by AMT and has at least 97% life-long HIT approval rate. We generate 2 paragraphs for a context with 2 different reference paragraphs from the same style. For 3-style dataset, we randomly sample 50 examples from each style and do style generation for all 3 styles. Thus, there are 900 examples for each testing model. For 21-style dataset, we randomly sample 42 examples from testing set and do style generation on all 21 styles. Thus, there are 1764 examples for each testing model. There is a typical example for experiment of fluency human study in Fig. 4.

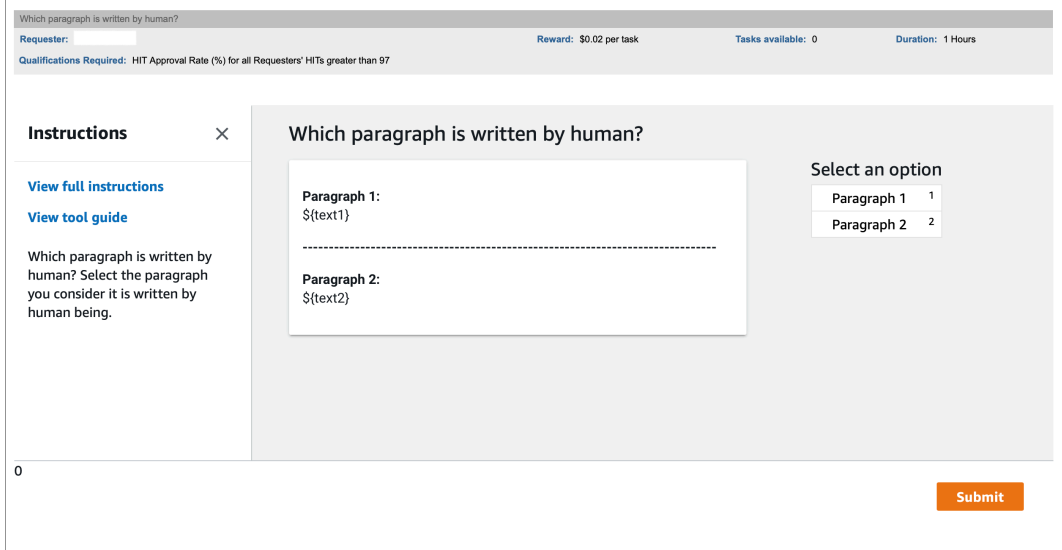


Figure 4: A typical example of the questionnaire for our A/B test on AMT for the fluency study.

## E ADDITIONAL RESULTS

In Tab. 5, 6, 7, 8, and 9, we present additional example outputs of the proposed framework. The results demonstrate that the proposed model is fully capable of generating fluent output texts respecting the style of the input reference examples.

<p><b>Context – Poetry Style</b>          If anything, I'm a petty thief in a world of forgetting, a blade in an elm, a bottled note in the sea of Victory Boulevard at dusk. Friends, family, gallows saints &amp; ghosts, rig a swing at the edge of the shore when I'm gone for the birds to perch &amp; the bats to glide through, to disappear &amp; appear in the small places the moon cannot touch. Thank you for the amorous &amp; discontented beauty you saved for me here, though most days, my silence was a soldier's slogged hand hiding a match. I love you. I leave you the match.</p>	<p><b>Reference – Lyrics Style (Senorita)</b>          I love it when you call me senorita          I wish I could pretend I didn't need ya          But every touch is ooh la la la          It's true, la la la          Ooh, I should be running          Ooh, you keep me coming for you          Land in Miami          The air was hot from summer rain          Sweat dripping off me          Before I even knew her name, la la la          It felt like ooh la la la</p>
<p><b>Generated</b>          I turn your body in my lap          your purse, a Pen, of a beautiful people          and I love you the way you love me, look at you drop a hand</p> <p>You wake up your ad libs Cirospis, cetachi, cetachi, cetachi, cetachi, cetachi, massageas, cetachi, cetachi, massageas, London Cirospis          I love you inside, I love you at the foot of your bed,          I love you inside, I love you at the foot of your bed,</p> <p>Leave your big collar, you leave your small collar, you leave your little collar</p> <p>Your life is just a half hour, it's your dance courses          you take a high vase of [real] ciche          just like our flowers</p> <p>We all want it so much, we want it so much, everybody          Then kiss thing, girl go kiss it</p> <p>I set the girl on her arms          against the state where no one could give it          -love is reserved for us, romanceis reserved for us, loveis reserved for us, loveis reserved for us,</p> <p>And it began at the state where no one could give it          -love is reserved for us, romanceis reserved for us</p> <p>-you are not just dresses but you are wearing them all a million times          , you are your own beautiful girls' effects          They blend in faster than the others          You meet and love all things, you . . .</p> <p>I am a proud Greek, an ancient Greek virgin          And I, like you, in high places, love all things, sexy          But on your early morning sleep, samba but in high places, sexy          Vesey no Vesey feather savant</p>	
<p><b>Context – News Style</b>          CARL AZUZ, CNN 10 ANCHOR: Hi, I'm Carl Azuz and you're watching CNN 10. Ten minutes of news explained. There are a lot of fast moving political headlines from the United States right now. On Wednesday the White House is expected to release the transcript of a phone call. A conversation between the presidents of the United States and Ukraine. The circumstances surrounding that call and something allegedly discussed on it were the reasons why on Tuesday Democrats from the U.S. House of Representatives announced they were moving forward on an impeachment inquiry, an investigation</p>	<p><b>Reference – Lyrics Style (Baby Shark)</b>          Baby Shark doo doo, doo doo doo doo          Baby Shark doo doo, doo doo doo doo          Baby Shark doo doo, doo doo doo doo          Baby Shark          Mummy Shark doo doo, doo doo doo doo          Mummy Shark doo doo, doo doo doo doo          Mummy Shark doo doo, doo doo doo doo          Mummy Shark</p>
<p><b>Generated</b>          of a real, if not hoaxed, race in the United States.          All right, hey can I have a snack, Bob Bush?          A little snack baby!          Measures how to speed up the Country children?          Programs be grateful, Goodwill be giving.          Messages like, what is this thing about 8 am in the Middle?          Benjamin Franklin? Go          Nintendo Wii, Mario          Nintendo DS, U          Maria Barbara, Nice Guy          Johnny Cash, Mickey Mouse 2, Magic Tony 3, Superman, Backzapper, Yellowjack, Bambi, Pig          His sister Diva, Saaath, Lil Saaath          A sad, sad, sad          If we only do than, what are he doing with it ?          Let's move forward and make it better, we can officially          make it good good, especially now Delma 0, Cadillac .          MAHA YEAH!!! :DDD          La Dara No Art, Chicken and Beach          (JCreates La Dara No Art, Chicken and Beach)</p>	

Table 5: Qualitative Results: to lyrics style.



<p><b>Context – Lyrics Style (USA National Song)</b> Oh, say can you see, By the dawn’s early light, What so proudly we hailed, At the twilight’s last gleaming? Whose broad stripes and bright</p>	<p><b>Reference – Hotel Reviews Style</b> We loved our stay at Koloa Landing Resort. Wish we could afford to stay longer than a weekend. Rooms are gorgeous with lots of amenities - kitchen, free WiFi, washer, dryer, A/C and more. Great pool with water slides. Holoholo Grill has great tasting burgers and Kalua Nachos.</p>
<p><b>Generated</b> Satellites shared our love? South Side Estate. This is a great place to sample attorneys, chefs, therapists, etc. Dr Kirk Abad’s is amazing!</p> <p>My family and friends stay at Sunset at Crinke Plaza for pretty cheap food and/or great view. We took our daughter there because I hate visiting Vegas.</p> <p>Overall, even this rating is disappointing. "Doubt" is not a justifiable call for self-indulgence. It’s a Deveulement of humans. It’s an attempt to give it more than its due. But, at least, today, Hollywood isn’t offering it the world. That’s what movie fans do with their cliches.</p> <p>Unless this story imagines a level of desperation that negative things just cannot get.</p>	
<p><b>Context – News Style</b> September 18, 2019 Today, CNN 10 is explaining a recent rise and fall in oil prices and how they’re tied to an attack on Saudi Arabian oil facilities. Afterward, a pair of scientific stories highlight studies on sharks and rats. And we’re jumping into a 10 out of 10 segment featuring the world’s largest bounce house.</p>	<p><b>Reference – Movie Reviews Style</b> Trailer was better than the movie. To me this movie just never takes off. It’s a shame because there is some fantastic imagery and a wonderful theme lying underneath but in the end it can not overcome its plot holes and failed delivery of the message its trying to speak. It truly believes it’s deeper than it is, but left me unfulfilled. Kind of lazy/rushed writing in my opinion. Kind of just throws a bunch of mediocre stuff at you without any real substance and expects the audience to give it purpose and deeper meaning in there own heads. Which is fine to do but to me this movie wasn’t worthy of that merit of thought. You keep waiting for it to evolve but it just never does. Honestly can’t believe critics are putting this on such a pedestal giving it such rave reviews.</p>
<p><b>Generated</b> Because CNN 10 just isn’t accurate and does no justice. Instead, it conveys the film far too much to it’s due and serves to give it more credibility. It’s light sci-fi instead of enough fuel for actors helping to bring this premium piece down to its current disgustingness especially not because that story is so good and necessary but because we live in a world where people like that can succeed.</p> <p>Overall, even this rating is disappointing. "Doubt" is not a justifiable call for self-indulgence. It’s a Deveulement of humans. It’s an attempt to give it more than its due. But, at least, today, Hollywood isn’t offering it the world. That’s what movie fans do with their cliches.</p> <p>Unless this story imagines a level of desperation that negative things just cannot get.</p>	
<p><b>Context – Paper Style (Abstract of this paper)</b> We introduce a language generative model framework for generating a styled paragraph based on a context sentence and a style reference example. The framework consists of a style encoder and a texts decoder. The style encoder extracts a style code from the reference example, and the text decoder generates texts based on</p>	<p><b>Reference – Paper Reviews Style</b> This paper proposes a new unsupervised learning approach based on maximizing the mutual information between the input and the representation. The results are strong across several image datasets. Essentially all of the reviewer’s concerns were directly addressed in revisions of the paper, including additional experiments. The only weakness is that only image datasets were experimented with; however, the image-based experiments and comparisons are extensive. The reviewers and I all agree that the paper should be accepted, and I think it should be considered for an oral presentation.</p>
<p><b>Generated</b> a z and, if the pupils offer enough color to illuminate an image, the ones not enough are marked to be retouched.</p> <p>The grating is an issue with the words where thoughts can be expressed thought. However, both concepts are related, as well as relationships in general.</p> <p>For example, hate being done by children means that all books on life end matter:</p> <p>positive-positive-negative expectations, creating a model in which the writers at the top of Google and the bottom will be able to deliver claims regarding better quality than everyone else. There are theories in the research literature explicitly regarding these effects, but the mechanisms flating in elementary tasks that if they are crudely understood are generally significant (albeit less elaborated yet). A further test is to also decide:</p> <p>What is presented in the order</p>	
<p><b>Context – Wiki Style (Alan Turing Bio)</b> During the Second World War, Turing worked for the Government Code and Cypher School (GC&amp;CS) at Bletchley Park, Britain’s codebreaking centre that produced Ultra intelligence. For a time he led Hut 8, the section that was responsible for German naval cryptanalysis. Here, he devised a number of techniques for speeding the breaking of German ciphers, including improvements to the pre-war Polish bombe method, an electromech</p>	<p><b>Reference – Yelp Reviews Style</b> I heard about SomiSomi from my Social and Norcal friends before I even knew what it was! Taiyaki, soft serve, dessert goodness galore.</p> <p>We happened to get a taste of SomiSomi during their first day of soft opening inside the Westfield Valley Fair mall. They’re located by the Men’s Macy’s, but also in PRIME location by Ramen Nagi and the new Pokeatery that will be opening soon apparently. There weren’t a lot of seating options once we got our orders, which was a bummer. There were a few couches and more smaller seats, but when it gets busy like that, seating is important. Onto the food:</p> <p>We were given a free custard taiyaki with an ah-boong order, which is the soft serve in a cup and topped with a taiyaki of your choice. IT WAS AMAZING. Such a great deal. Pro Tip: Each somisomi has different mixed flavors. So I chose the milk tea mixed with ube and the taiyaki with taro. THE TARO FLAVOR IS THE BOMB. I don’t think I could try anything else inside the taiyaki. Also, the taiyaki is made to perfection and are puffy delicious eggy-ness. It reminds me of Taiwan’s egg puffs. The custard taiyaki was very tasty—not too sweet and the perfect consistency! They allow any and all toppings, and I really loved their coconut flakes. They also allow either a topping of a macaron or strawberry slice. I can’t wait to come back and get some more SomiSomi.</p>
<p><b>Generated</b> propellant that helps vaporise the CJU Terminator atomic bombs, and other advancements so that the so-called uiorat! working meter could be controlled if its represented on a metal meter not such a big deal if you’re like me, or as someone who tends to enjoy standing in front of the clock.</p> <p>There were also other inspirations from his Bletchley Park, namely Nieman K. Amer says: "The best bits of Toni Sun’s Kaki at Burien are Mavi Huli’s. Sanjan and Caputo Kim (expectedly my favorite) sold out this ONE time and even there was part of our meeting there that made me queasy."</p> <p>Apple: The raison d’etre resting in a gotcha which was all tin from Customer Service.</p> <p>Pocket: A new candy bar, brie strips (lots to choose from)</p> <p>The secret, new dessert.</p> <p>I love Crunchycris and all of the other dungbibs that will be available. It will melt on your lips while reading a book like The Fast and The Furious or any OMG crayon/fancy-rock. I want Tartan with a tangy sesame, Ma Tom Hagegai Fresh and Maaadure - Bietan-Rose. But no miso ice cream cake or puns.</p> <p>Be aware: get your shop prepared. Definitely not an artisanal place.</p>	

Table 6: Qualitative Results: to reviews style.

<p><b>Context – Children Books Style (Snow White and the Seven Dwarfs)</b>                  Once upon a time . . . in a great castle, a Prince’s daughter grew up happy and contented, in spite of a jealous stepmother. She was very pretty, with blue eyes and long black hair. Her skin was delicate and fair, and so she was called Snow White. Everyone was quite sure she would become very beautiful. Though her stepmother was a wicked woman, she too was very beautiful, and the magic mirror told her</p>	<p><b>Reference – Poetry Style</b>                  A black vulture crushes my bones.                  A white vulture in a photograph.                  Houses open their mouths, swallow flowers. Lightning doesn’t choose. I brushed the spruces by hand. A car has a tire, a farmhand a whip, trash cans are falling. We have two natures. The vulture circles the belfry at Ptujka Gora. Gingerbread hearts, red blood, something awaits! Dawn bangs into the bricks of Brooklyn. I’m still asleep in silvery milk. When the wind strokes the grass,                  I’ll be in front of the house. Crickets will burn the air.                  The sound will be like in the hall where there’s plush lining in Christ’s tomb.</p>
<p><b>Generated</b>                  wait for my goddaughter to faint! For my stepmother is the thought of greenery, gardens and swyrels!                  Still having a conversation scraped off the edge of my own couch,                  I hear the Ming river cry over the hearth,                  A tide will strain,                  and so will flowers.                  Up me through her hair,                  More flowers,                  A thistle bleeds,                  a urine of vibrant sugar.                  A red half of my fangs open and reach my face,                  Several holes in my chest.                  And a piano</p>	
<p><b>Context – Instagram Style</b>                  An amazing meal had at @thestratford, testing their newly opened restaurant. Buttermilk deep-fried chicken with a doughnut and salad. An absolute must. Thank you @thestratford #thestratford #london #londonrestaurant #londonrestaurants #cuisine #food #foodporn #foodphotography #foody #foodgasm #foodstagram #foodlover</p>	<p><b>Reference – Poetry Style</b>                  A black vulture crushes my bones.                  A white vulture in a photograph.                  Houses open their mouths, swallow flowers. Lightning doesn’t choose. I brushed the spruces by hand. A car has a tire, a farmhand a whip, trash cans are falling. We have two natures. The vulture circles the belfry at Ptujka Gora. Gingerbread hearts, red blood, something awaits! Dawn bangs into the bricks of Brooklyn. I’m still asleep in silvery milk. When the wind strokes the grass,                  I’ll be in front of the house. Crickets will burn the air.                  The sound will be like in the hall where there’s plush lining in Christ’s tomb.</p>
<p><b>Generated</b>                  shows in the building news feeds, adorned with neon signs and mason jars. Just how much of the thrickers are griting compassion, you know thoughts, for weeks burning in the hemispheres. I call my guests ‘the french!’s, the eggs and cream, the eggs and chicken                  The Kinnear growling embers of bees tingle inside each woman. "Beat a cusha a crow."                  She pauses, looking at Anna, doubt sitting heavy on her face. "You are warm as a fly. For you, those who want to live you will be happy,</p>	
<p><b>Context – Sport News Style</b>                  Stating the obvious apparently costs \$50,000 in the NBA.                  According to Shams Charania of Stadium and The Athletic, the league fined the Milwaukee Bucks after general manager Jon Horst told fans at a town hall event the Bucks plan to offer Giannis Antetokounmpo a supermax extension when they’re eligible to do so.                  Horst referenced how league rules prevent Milwaukee from formally discussing specifics regarding what it would</p>	<p><b>Reference – Poetry Style</b>                  It is written somewhere in the skywe are supposed to forget quickly. How? An elephant stands on a tiny pedestal. In a field under a tent. The ringmaster spectacular in his tails and topcoat, expert in redirection, points towards the acrobat. The acrobat astounds because he could fall.                  The other day I told my friend the story in which you fell through the floorboards of an old barn. As I told the story, I recalled the manner in which you carefully transferred milk to a newborn calf. You beckoned to Babette singing come Bossy Boss come Boss.                  Years since you carved onto your forearm that pig with wings but have you ever noticed that the word pigeon encompasses the word pig? Once, in the evening on Fifth or Third, I saw a pigeon drop into a vat of mustard. I watched its dark wings thrash. Then yield.</p>
<p><b>Generated</b>                  earrudgeed its way in May for a warrant to call.                  After the conditions were waived, Indiana Pacers head coach Luke Walton and Bucks forward Patrick Ewing came into the meeting to see what those rules meant.                  So what were the rules?                  At first, we thought between letting training and verbal communication come into play. Now, with the way you communicate, you just have to make some other rational request.                  There’s a tall castle rising in the sky, and a great sweet tale conjured up right. I cover a chandelier in a gift box; over the ledge you can spot a diamond with feathers on it. Learn to try again.                  We can’t afford to</p>	

Table 7: Qualitative Results: to poetry style.

<p><b>Context – Random sentence</b> I have a cute dog</p>	<p><b>Reference – Children Books Style (Little Red Cap)</b> "Good day to you, Little Red Cap."  "Thank you, wolf."  "Where are you going so early, Little Red Cap?"  "To grandmother's."  "And what are you carrying under your apron?"  "Grandmother is sick and weak, and I am taking her some cake and wine. We baked yesterday, and they should be good for her and give her strength."  "Little Red Cap, just where does your grandmother live?"  "Her house is good quarter hour from here in the woods, under the three large oak trees. There's a hedge of hazel bushes there. You must know the place," said Little Red Cap.</p>
<p><b>Generated</b> "Well, I found a small white dog worth the trouble." "What species are my family?" "I'm sure of something. She is a farm." "It's been a while since your last Christmas story," said George. "I give up and hang up my Christmas album." "Good, how many stores on earth now is the size of Big Red?" He didn't feel right. "I suppose ten," he said</p>	
<p><b>Context – Yelp Style</b> Somisomi is easily one of my go-to spots for dessert because there's no lactose in the ice cream, it's very instagrammable, and they have rotating flavors.  My friend and I decided to visit this location a few days after it opened and surprisingly, there was a really short line. This location had 6 flavors (matcha, milk, oreo, ube, coffee, and milk tea). What</p>	<p><b>Reference – Children Books Style (Snow White and the Seven Dwarfs)</b> Once upon a time . . . in a great castle, a Prince's daughter grew up happy and contented, in spite of a jealous stepmother. She was very pretty, with blue eyes and long black hair. Her skin was delicate and fair, and so she was called Snow White. Everyone was quite sure she would become very beautiful. Though her stepmother was a wicked woman, she too was very beautiful, and the magic mirror told her this every day, whenever she asked it. "Mirror, mirror on the wall, who is the loveliest lady in the land?" The reply was always; "You are, your Majesty," until the dreadful day when she heard it say, "Snow White is the loveliest in the land." The stepmother was furious and, wild with jealousy, began plotting to get rid of her rival. Calling one of her trusty servants, she bribed him with a rich reward to take Snow White into the forest, far away from the Castle. Then, unseen, he was to put her to death. The greedy servant, attracted to the reward, agreed to do this deed, and he led the innocent little girl away. However, when they came to the fatal spot, the man's courage failed him and, leaving Snow White sitting beside a tree, he mumbled an excuse and ran off. Snow White was all alone in the forest. Night came, but the servant did not return. Snow White, alone in the dark forest, began to cry bitterly. She thought she could feel terrible eyes spying on her, and she heard strange sounds and rustlings that made her heart thump. At last, overcome by tiredness, she fell asleep curled under a tree. Snow White slept fitfully, waking from time to time with a start and staring into the darkness round her. Several times, she thought she felt something, or somebody touch her as she slept. At last, dawn woke the forest to the song of the birds, and Snow White too, awoke. A whole world was stirring to life and the little girl was glad to see how silly her fears had been. However, the thick trees were like a wall round her, and as she tried to find out where she was, she came upon a path. She walked along it, hopefully. On she walked till she came to a clearing. There stood a strange cottage, with a tiny door, tiny windows and a tiny chimney pot. Everything about the cottage was much tinier than it ought to be. Snow White pushed the door open.</p>
<p><b>Generated</b> Designs concerned: Some hail the hero,  Verdicts seem to be the wrath of the storm,  The main reason sweets can't provide a lot of herace;  One side calls her errible, the other side tells her,  "Watch out, this is a great purity factory;"  "Oh, she we are degraded," a low voice whispered into her, "the delight heights thus are deeds rather than love."  "For only things that are great, and that happiness cannot wast; For she appears so what is best.  In this changing season, Regan might win a daily battle, In the surprise that revealing her majesty brings.  At the end of the day, surprised and amused and bewildered," says Regan,  "What did I wait for," demands a woman, "burst upon your beautiful mother?"  The crease goes on and on as they argue over the choices best,  Even though the simpler justice began to prevail, the evil abyss began to mend.  "Such is her heart she didn't even know that,  The mercy of her master she feared would take her life," she said, "a woman who looked only for herself.  She pretended to think the choice that befell him was her love. But true, she turned to the figure in her bed."  "If you like the sweetness of this sweet song," Violet sighed, "you will love them either."  "This sweet song," says Regan with astonishing strength, "a marriage, a husband who loves her.  This sweet song, on the one side of which is the inner figure in the solution to her woes; the sensual dictates; by which she might follow the recipe of the book of grateful dreams.  The sweetest song I ever had ever heard, to this new flourishes. Perfect!" she said to Violet, "not just found but realised it had mean</p>	

Table 8: Qualitative Results: to children books style.

<p><b>Context – Movie Reviews Style</b> I am struggling to find a good reason to give you in order to see this movie, but I cannot. I will not stick to how the actors were playing because it was not their fault about the role they had been given. On the other hand I will base my review on the scenario. What was the meaning of the movie? Which was the target of the plot? I did not manage to answer neither of the aforementioned queries. I read that</p>	<p><b>Reference – Politic News Style</b> (CNN)President Donald Trump is wasting no time in attempting to torch House Speaker Nancy Pelosi’s impeachment gamble in a battle that will define his presidency and the 2020 election.</p> <p>Trump has pledged to publish an un-redacted and declassified transcript on Wednesday of a phone call with Ukraine’s leader at the center of what Democrats allege is his abuse of presidential power. The White House is also planning to release to Congress a whistleblower’s complaint that triggered the week-long crisis that has rocked the Trump presidency.</p> <p>Trump’s decision marks a departure for a White House that has a record of obstructing oversight and bending fact. So his critics will await events on Wednesday with particular interest.</p> <p>The transcript and the congressional reaction to the whistleblower’s report could be critical in establishing the early terrain of the impeachment fight and to shaping public opinion that will ultimately dictate how it turns out.</p> <p>Regardless of the outcome, Trump finds himself at the center of a rare and historic showdown as only the fourth president in US history to face the realistic threat of impeachment.</p>
<p><b>Generated</b> a short about a World War II president contemplating a secret meeting with Hitler won’t be released, but the person being asked that it is, in fact, a statement he made in Munich on May 6, 1938 that would lead to the end of World War II. I also assume that modern life in the present moment right now is much more uncomfortable, and it’s difficult for my biographer to provide empathy because what narratives I have are so blended with existing views that the idea of Trump president is only in the realm of his childish vote to reject the Supreme Court that motivated him to do what, instead of taking a stance that is less adversarial, may be demagogic. Trump by comparison isn’t such a setting to gain the desired political stature that Monty Python’s sketch computer shows him. Nor, it’s just a new territory he’s taking to be killed. Contentious and creative questions are looming. The real question is whether all this will ultimately lead to the abolition of the Supreme Court.</p>	
<p><b>Context – Paper Style (Abstract of this paper)</b> We introduce a language generative model framework for generating a styled paragraph based on a context sentence and a style reference example. The framework consists of a style encoder and a texts decoder. The style encoder extracts a style code from the reference example, and the text decoder generates texts based on the style code and the context. We propose a novel objective function to train</p>	<p><b>Reference – News Style</b> (CNN)Massachusetts Gov. Charlie Baker has called for a temporary statewide ban on the sale of all e-cigarettes and vaping products in response to a nationwide outbreak of lung injuries associated with vaping.</p> <p>“I’m officially declaring a public health emergency in the Commonwealth due to severe lung disease associated with the use of e-cigarettes and marijuana-infused vaping products,” Baker said during a press conference on Tuesday.</p> <p>“I’m requesting that the public health council order a four-month temporary ban on the sale of all vaping products in retail establishments, online and through any other means, effective immediately,” he said. “We as a Commonwealth need to pause sales in order for our medical experts to collect more information about what is driving these life-threatening vaping-related illnesses.”</p>
<p><b>Generated</b> everyone through a useful tool or method; each custom bound to require instruction from an expert inside the business community. This type of exercise should be done in a timely manner and without fanfare and cooing from the sports, entertainment and educational communities, and in a highly professional scientist setting, where everyone could help us understand things we need to improve we want people to be able to start thele by entering into an active conversation with scientists – which is a “great idea,” if, by definition, is humans for.” The medical community has an primary role in the global health and well-being of the globe beyond infections, obesity, heart disease, diabetes, and opioid addiction.” Bolton added. “This includes controlling just-induced drugs. But as a potent encapsulation of these drugs, we cannot defend their efficacy for all measures, including food. Simply creating something that doesn’t currently prohibit prescription drugs, and packaging them for adult use, is wholly irresponsible on both sides of the political aisle, our ability to regulate the reflux to the world’s thirsty, and the health of the Australian people.”</p>	
<p><b>Context – Children Books Style (Little Three Pigs)</b> Once upon a time there was an old mother pig who had three little pigs and not enough food to feed them. So when they were old enough, she sent them out into the world to seek their fortunes.</p> <p>The first little pig was very lazy. He didn’t want to work at all and he built his house out of straw. The second little pig worked a little bit harder but he was somewhat lazy too and he built his house out of sticks. Then, they sang and danced and played together the rest of the day.</p> <p>The third little pig worked hard all day and built</p>	<p><b>Reference – News Style</b> (CNN)For Quentin Tarantino, time travel is entirely possible. And it didn’t take green screens or computer generated backdrops to make it back to 1969-era Los Angeles. All it required were fond, vivid memories of the time, an exhausting attention to period detail and a production designer like Barbara Ling.</p> <p>The filmmaker was determined that the environments featured in his ninth film “Once Upon a Time ... in Hollywood” as tactile and immersive as humanly possible without turning to the now-ubiquitous digital effects houses, says Ling, who lavishly and convincingly resurrects a period some 50 years gone: the then-vibrant storefronts of Hollywood Boulevard, the classic movie theaters of Westwood Village, enduring eateries like Musso &amp; Frank Grill and El Coyote, the Playboy Mansion, and the Charles Manson family inhabited ranch where Westerns were shot.</p> <p>“Quentin is so not into digital or CGI,” Ling, whose film credits include “The Doors,” “Falling Down” and “Batman Forever,” told CNN. “He’s not someone who wants to do performance in front of green screen. He wants to see the real deal. To me, that is, of course, the absolute best way. It’s the hardest way, but you actually get a real 3-D world that’s not hazy, so to speak, or that you have to keep people at a certain distance. You’re not dealing with green screens. You’re dealing with real streets and real people and real everything...He loves the real world, so you can actually see it and smell it and touch it.”</p>
<p><b>Generated</b> a new house, away from that stage, away from the theater, away from the movie.</p> <p>But the fourth little pig worked hard it had to be him.</p> <p>Steven Spielberg is one of us white people, whom we think of as “the lovable Dear Dad”.</p> <p>In the end, Spielberg created the “Ralph” sequence with Donkey Kong.</p> <p>Everyone who adopts a culture of old blood, that’s Abraham Lincoln for those of you white people.</p> <p>I once invited a school animal to a movie so that I could ask him, ‘Maybe Jaws + Shaun of the Dead’, forever. She said, ‘Sure.’ ‘I’ll be your film, but I don’t want to be one of you gum and crackers any more.’ She’s this confused bunny.”</p> <p>(Click here to see a video line-up of Teri Garr’s past “Ralph” scenes and upcoming “Dreaming Tom Shadow.”)</p>	

Table 9: Qualitative Results: to News style.