Data-to-Text Generation with Style Imitation

Shuai Lin^{1,2}, Wentao Wang², Zichao Yang², Xiaodan Liang¹, Frank F. Xu² Eric P. Xing^{2,3}, Zhiting Hu^{2,4*}

¹Sun Yat-sen University, ²Carnegie Mellon University, ³Petuum Inc., ⁴UC San Diego {shuailin97,xdliang328,zhitinghu}@gmail.com, {zichaoy,fangzhex,epxing}@cs.cmu.edu

Abstract

Recent neural approaches to data-to-text generation have mostly focused on improving content fidelity while lacking explicit control over writing styles (e.g., word choices, sentence structures). More traditional systems use templates to determine the realization of text. Yet manual or automatic construction of highquality templates is difficult, and a template acting as hard constraints could harm content fidelity when it does not match the record perfectly. We study a new way of stylistic control by using existing sentences as "soft" templates. That is, the model learns to imitate the writing style of any given exemplar sentence, with automatic adaptions to faithfully describe the content record. The problem is challenging due to the lack of parallel data. We develop a neural approach that includes a hybrid attention-copy mechanism, learns with weak supervisions, and is enhanced with a new content coverage constraint. We conduct experiments in restaurants and sports domains. Results show our approach achieves stronger performance than a range of comparison methods. Our approach balances well between content fidelity and style control given exemplars that match the records to varying degrees.¹

1 Introduction

Recent years have seen remarkable progress in *neural* natural language generation to produce well-formed coherent text (Sutskever et al., 2014; Vaswani et al., 2017). Yet, controllability over various text properties, as an essential demand to ensure the utility of generations in real-world applications, has not attained the same level of advancement. Data-to-text generation is one of such applications with ubiquitous practical use, in which

natural language text is generated to describe a given data record such as a box score of a sports player or an infobox table of a restaurant.

Though current data-to-text neural approaches with encoder-decoder models could produce fluent text with high fidelity to content ("what to say"), they largely lack control over the writing style, such as sentence structures and word choices ("how to say"). Many efforts have been made to promote the overall diversity in data-to-text generation through, e.g., latent variables (Ye et al., 2020) or customized model architectures (Jagfeld et al., 2018; Deriu and Cieliebak, 2018). Yet fine-grained style manipulation is not permitted. This contrasts with the traditional text generation systems which separate content planning and surface realization (Reiter and Dale, 1997), and usually determine the realization with explicit templates (Kukich, 1983; McRoy et al., 2000) or based on syntactic grammars (Robin and McKeown, 1996; Power et al., 2003).

Controlling writing style with "hard" templates could suffer from unscalable template creation and lack of generation flexibility. Though previous work (Wiseman et al., 2018; Dou et al., 2018; Angeli et al., 2010) has enabled automatic template extraction, the templates usually act as hard constraints and could harm the content fidelity of generations when the template does not exactly match the content in a record.

In this paper, we study a new way of stylistic control in data-to-text generation by using any existing sentences as "soft" templates. That is, we learn to *imitate* the writing style of a given exemplar sentence. The goal is two-fold: to generate text that not only faithfully describes all content in the record, but also inherits as many of the exemplar's stylistic characteristics as possible (Figure 1). The new paradigm sidesteps the restrictions with traditional dedicated templates and allows us to use arbitrary exemplar sentences that could be describ-

^{*}corresponding authors

¹Data and code are publicly available at https://github.com/ha-lins/DTG-SI

Data Record	Name	Food	Area	Price	Near				
	Loch Fyne	Italian	Riverside	£20-25	Strada				
Exemplar 1 Generation 1	Zizzi is a pub providing fine French dining but with an expensive price, located near Cocum in the city center. Loch Fyne provides fine Italian dining with a £20-25 price, located near Strada at the riverside.								
Exemplar 2 Generation 2	Located near the Blue Spice, there is a highly-rated place, the Mill, as a choice that frugally priced. Located near Strada by the river, there is a place with Italian foods, Loch Fyne, as a choice that priced £20-25.								
Exemplar 3 Generation 3	With a family-friendly atmosphere and a 5-star rating, Aromi is a pub in the city center. With Italian foods and a moderate price range, Loch Fyne is near Strada at the riverside.								

Figure 1: An example of generating sentences that describe the data record and imitate the style of given exemplar sentences (i.e., soft templates). The generations *adaptively* inherit the structural and phrasing characteristics (highlighted with cyan boxes) of the exemplars. For instance, exemplar 2 does not match the record content perfectly (e.g., it does not describe the food). The generation adapts the structure to add "with Italian foods". All such automatic adaptions are highlighted in orange. Note that the word "providing" in exemplar 1 is also adapted to "provides" for grammar correction.

ing distinct content. As shown in Figure 1, the model automatically adapts the soft templates to varying extents based on how well they match the record, and precisely expresses the desired content.

To this end, we develop a neural approach that balances well between content fidelity and style imitation. A key learning challenge is the lack of parallel data, i.e., triples of (record, exemplar sentence, target description). Instead, we usually only have access to abundant record-description pairs². The proposed approach learns with rich weak supervisions derived from the record-description pairs. Architecture-wise, we develop a hybrid attention-copy mechanism that offers differentiated treatments of the content and style sources. Further, based on the structural nature of data records, we devise a new content coverage constraint for the balanced embodiment of both content and style in the generation.

We conduct empirical studies on corpora from two domains, including restaurant recommendation (Dušek et al., 2019) and NBA reports (Wiseman et al., 2017). Experiments show our models strongly improves over a diverse set of comparison methods in terms of both automatic and human evaluations. In particular, given exemplar sentences that match data records to varying degrees, our approach retains a good content-style balance.

2 Related Work

Data-to-Text Generation Many efforts have been made to improve the fidelity of generated text to the record content, through sophisticated neural architectures (Wiseman et al., 2017; Gehrmann et al., 2018; Puduppully et al., 2019; Iso et al., 2019), hybrid retrieval and generation (Hashimoto et al., 2018; Weston et al., 2018; Cao et al., 2018; Pandey et al., 2018; Peng et al., 2019), and others. These approaches do not have the additional goal of style control as ours, and usually perform supervised learning based on record-description pairs. Traditional data-to-text generation systems implement a pipeline architecture consisting of separate components, including content planning, sentence planning, and surface realization (e.g., Reiter and Dale, 1997; Kukich, 1983; McRoy et al., 2000; Kondadadi et al., 2013). Recent work (Wiseman et al., 2018) integrates the template use in a more end-to-end neural model. Rather than treating templates as hard constraints as in the previous work, we study the new setting of using existing sentences as exemplars, allowing the model to adaptively imitate the style while ensuring content fidelity.

Text Style Transfer There has been growing interest in text style transfer (Hu et al., 2017; Shen et al., 2017; Yang et al., 2018; Subramanian et al., 2019, etc) which assumes an existing sentence of certain content, and modifies single or multiple textual attributes (e.g., sentiment) of the sentence without changing the content. Our problem differs

²This highlights the difference from the recent retrievaland-generation work (e.g., Hashimoto et al., 2018; Weston et al., 2018; Cao et al., 2018; Peng et al., 2019) which focuses only on content fidelity and thus is a supervised learning problem given the record-description pairs.

in important ways in that we assume the abstract writing style is encoded in an exemplar sentence and attempts to modify its concrete content to express the new information in a structured record (we thus can call our setting text content rewriting). The different settings can lead to different application scenarios in practice, and pose varying technical challenges. In particular, though the recent style transfer research (Subramanian et al., 2019; Logeswaran et al., 2018) has controlled multiple categorical attributes which are largely independent or loosely correlated to each other, a data record in our task, in comparison, can contain a varying number of fields, have many possible values, and are structurally coupled. Our empirical studies (sec 5) show the recent models designed for style transfer fail to perform well on the problem under study. We also note recent work of syntacticallycontrolled paraphrase generation based on either constituency parse (Iyyer et al., 2018) or reference sentences (Chen et al., 2019). The problem nature of data-to-text generation in this work leads to a solution with very different architectures and learning approaches.

Controlled Generation without Parallel Data

Controlling different aspects (e.g., content, style, discourse structures) in text generation requires grasping the intrinsic mapping between the aspects and the surface text. The lack of parallel data often poses challenges in learning the mapping, making it necessary to incorporate other forms of experiences (supervisions) (Hu and Xing, 2020). For example, the style transfer work (Hu et al., 2017; Shen et al., 2017; Yang et al., 2018) used auxiliary models such as attribute classifiers and language models for supervision signals. Tang et al. (2019) learned guided conversation flow using standard conversation data combined with logical control. Tan et al. (2020) created weak supervision labels from knowledge bases for aspect-based summarization. This work devises competing training objectives based on common record-description pairs. Joint optimization of the competing objectives drives the model to learn desired behaviors.

3 The Task: Data-to-Text Generation with Style Imitation

For clarity, we first formally describe the problem of data-to-text generation with style imitation. We also establish the key notations used in the paper.

Consider a data record x which consists of a set

of *fields* and their values (e.g., field "Food" and its value "Italian" in Figure 1). Note that different records can include different fields. For example, the field "Customer Rating" is included in some records but not the one in Figure 1. Data-to-text generation aims to produce a sentence to describe the content in the record. We are additionally given an exemplar sentence y_e which could be describing distinct content in the same domain. The goal of the task is thus to generate a new sentence y that achieves (1) content fidelity by describing the content in x accurately and completely, and (2) style embodiment by retaining as much of the writing style (e.g., sentence structure, word choice, etc) of y_e as possible.

A solution to the problem is required to balance well between the two objectives, by adaptively rewriting necessary portions of the reference y_e to express the desired content in a correct and fluent way, while at the same time editing y_e to a minimum extent to inherit its style. The demand for adaptive trade-off necessitates developing learning approaches for flexible imitation and generation.

To the best of our knowledge, there is no large data containing the desired (x, y_e, y) triples for supervised learning. Instead, we often only have access to *pairs* of record and its description which was originally written without following any designated style. In the next section, we develop a neural approach that learns style imitation given only the paired data.

4 The Approach

Denote the proposed neural model as $p_{\theta}(\boldsymbol{y}|\boldsymbol{x},\boldsymbol{y}_{e})$. The model has a hybrid attention-copy mechanism (sec 4.1) for differentiated treatment of source content and style exemplar. We learn the model by constructing weak supervisions from the available non-parallel data (sec 4.2), and further encourage accurate content description with a content coverage constraint (sec 4.3). Figure 2 presents an overview of the approach.

4.1 Hybrid Attention-Copy Architecture

The overall architecture of the neural model consists of two encoders and one decoder. The two encoders extract the representation of the data record \boldsymbol{x} and exemplar \boldsymbol{y}_e , respectively. Concretely, for each field in \boldsymbol{x} , we concatenate the embedding vectors of the field and its value, and feed the sequence of field-value embeddings to the encoder.

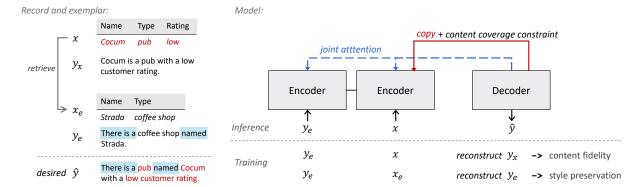


Figure 2: A (simplified) data example and retrieval (**left**) and the model overview (**right**). The proposed approach uses a hybrid attention-copy mechanism, and is learned with weak supervisions and a content coverage constraint.

The decoder generates the output sentence with a hybrid attention-copy mechanism. In particular, the decoder applies joint attention over both y_e and x, and uses a copy mechanism (Gu et al., 2016) only on the field values in the record x. More concretely, at each step t, the decoder first attends jointly to the hidden states of both encoders, and obtains a decoding hidden state h_t . The final output distribution is the weighted-sum of two distributions:

$$\boldsymbol{P}_{out}^{(t)} = g_t \cdot \boldsymbol{P}_V^{(t)} + (1 - g_t) \cdot \boldsymbol{P}_{\boldsymbol{x}}^{(t)}$$
(1)

where g_t is the probability of generating a token from the vocabulary; $P_V^{(t)}$ is the generation distribution over the whole vocabulary; $P_x^{(t)}$ is the copy distribution over the field values in the record.

4.2 Learning with Weak Supervisions

The two problem goals, namely content fidelity and style embodiment, are complementary and to some extent competitive. We derive weak forms of supervisions for each of them, respectively, based on the corpus of record-description pairs available.

Exemplar Retrieval First, for each record x, we automatically construct the exemplar y_e through retrieval. Specifically, we use x to retrieve another record x_e based on their *distance*, and use the description associated with x_e as the exemplar sentence y_e in training. We define the distance between y and y_e as follows:

$$\mathcal{D}(y, y_e) = \#[\mathcal{T}(x) \cup \mathcal{T}(x_e)] - \#[\mathcal{T}(x) \cap \mathcal{T}(x_e)]. (2)$$

where $\mathcal{T}(\cdot)$ is the set of all fields in the record; $\#[\cdot]$ represents the number of fields in the set. Figure 2 gives an illustration of retrieved exemplar (with distance =1). We study the effect of training with exemplars of varying distances in the experiments.

Content Objective Given the retrieved results, we next tackle content fidelity. Consider the description associated with x, which, though not following the desired style of y_e , has accurately presented the content in x. Denote the description as y_x . We thus devise the first learning objective that reconstructs y_x given (x, y_e) , in order to provide the model with the hints on how the x content can be presented in natural language:

$$\mathcal{L}_{content}(\boldsymbol{\theta}) = \log p_{\theta}(\boldsymbol{y}_{x}|\boldsymbol{x}, \boldsymbol{y}_{e}). \tag{3}$$

Style Objective For the second goal of style embodiment, we want to encourage the model to generate sentences in a similar form of y_e . To this end, we notice that, if we feed the model with the exemplar sentence y_e and its corresponding record x_e , then by definition the desired output would be y_e itself. We thus devise the second learning objective that reconstructs y_e given (x_e, y_e) :

$$\mathcal{L}_{style}(\boldsymbol{\theta}) = \log p_{\theta}(\boldsymbol{y}_e | \boldsymbol{x}_e, \boldsymbol{y}_e). \tag{4}$$

The objective essentially treats the exemplar sentence encoder and the decoder together as an autoencoding module, which effectively drives the decoder to reproduce the exemplar's characteristics.

Joint Training The above two learning objectives are competitive with each other such that, by combining them and optimizing jointly, the model is encouraged to learn to balance between content fidelity and style embodiment. A similar learning strategy of dividing a learning problem into multiple competitive objectives has also been used in previous work such as text style transfer (Hu et al., 2017; Shen et al., 2017). More formally, the above two objectives are coupled together to train the model as follows:

$$\mathcal{L}_{joint}(\boldsymbol{\theta}) = \lambda \mathcal{L}_{content}(\boldsymbol{\theta}) + (1 - \lambda) \mathcal{L}_{style}(\boldsymbol{\theta}), \quad (5)$$

	Restau	rant Reco	ommend.	NBA Reports		
	Train	Dev	Test	Train	Dev	Test
#Instances	29,486	6,299	6,273	31,444	6,765	6,930
#Tokens	0.54M	0.12M	0.12M	7.88M	1.69M	1.75M
Avg Text Length	18.36	18.34	18.35	25.07	25.10	25.32
#Unique Fields	8	8	8	34	34	34
Avg #Fields	5.38	5.38	5.35	4.32	4.31	4.35

Table 1: Statistics of the two datasets.

where $\lambda \in (0,1)$ is the balancing weight.

4.3 Content Coverage Constraint

As shown in the empirical study (section 5), the above learning performs well in general yet sometimes still fall short of expressing the record accurately. We thus devise an additional learning constraint to enhance content fidelity. The intuition is that, given the copy mechanism over the record x, each field value in x should be copied exactly once. We thus minimize the following L2 constraint that encourages the temporally aggregated copy probability of each field value in x to be 1:

$$C(\boldsymbol{\theta}) = \left\| \sum_{t} P_{\boldsymbol{x}}^{(t)} - 1 \right\|^{2} \tag{6}$$

where $P_x^{(t)}$, as defined in Eq.(1), denotes the copy distribution over all field values at decoding step t; and 1 is a vector with all ones.

The full model training objective with the constraint is thus written as:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathcal{L}_{joint}(\boldsymbol{\theta}) - \eta \cdot \mathcal{C}(\boldsymbol{\theta}) \tag{7}$$

where $\eta \geq 0$ is the weight of constraint.

5 Experiments

We study on two datasets in the restaurant recommendations and NBA reports domains, respectively. We conduct both automatic and human evaluations to assess model performance. Experiment results validate the proposed approach in learning an effective, balanced control of content and style.

5.1 Datasets

We derived and processed the two existing popular corpora as below. As defined in section 3, each resulting dataset contains record-description pairs. Table 1 shows the data statistics.

Restaurant Recommendations The dataset is extracted from the E2E NLG challenge (Dušek et al., 2019). A restaurant record can contain a subset of 8 fields, such as *Eat Type*, *Price Range*,

and others. See Figure 1 for an example record and the different possible ways of description.

NBA Reports We extract the dataset from the NBA game corpus developed in (Wiseman et al., 2017). The original corpus consists of box-score tables of NBA matches and the corresponding full-length match reports. We first split each report into individual sentences and extract the associated information from the box-score table as the data record. The data contains 34 unique fields, such as *Points*, *Rebounds*, *Field-Goal Percentage*, etc. Though the recorded fields look regular, the natural language descriptions are rich with variation. For example, for a field value *Points: 14*, one could say "contributed 18 points", "reached double figures", or, fusing with other fields, "scored an amazingly efficient 18 points on 7-of-8 shooting", etc.

5.2 Setup

Comparison Approaches

We compare with diverse approaches for a comprehensive analysis of the task and proposed approach:

- Reference for Content Fidelity: AttnCopy-S2S. We first consider a conventional data-totext model designed for only expressing the content. As style imitation is omitted, the method is expected to excel on content fidelity but fail on style control. Specifically, we use a sequenceto-sequence model (Sutskever et al., 2014) augmented with the proposed attention-copy mechanism (Section 4.1), which is trained supervisedly on the record-description pairs.
- Reference for Style Embodiment: Slot-filling. The second approach serving as a reference is a traditional slot-filling method that first removes the content words in the exemplar sentence y_e to make a template, and fills in the slots with respective values in the record x. As all content-independent tokens in y_e are preserved, the method is expected to perform well on style embodiment, but fail on content fidelity due to the possible mismatch between the exemplar sentences and desired content x. We manually crafted a large set of slot-filling rules for each of the two datasets respectively.
- Multi-Attribute Style Transfer (MAST) (Subramanian et al., 2019). We compare with a recent style transfer approach capable of manipulating multiple attributes. To apply to our task, we treat

		Restaura	nt Recommen	NBA Reports			
	Method	Con % Inclnew	tent %Exclold	Style m-BLEU	Cor Precision	ntent Recall	Style m-BLEU
Reference	AttnCopy-S2S Slot-filling	78.88±2.08 61.23	99.71±0.06 66.2	13.95±0.52 100	81.62±3.25 56.69	75.65±7.42 71.34	45.5±0.71 100
Baselines	MAST AdvST	36.28±0.25 51.64±4.45	37.06±0.16 57.06±4.44	91.76 ± 0.28 76.02±5.27	23.06±3.90 67.37±0.66	27.37 ± 3.88 66.79 ± 1.43	95.43 ± 2.71 64.67±4.81
Ours	Transformer w/o Coverage + Coverage	60.03±2.16 61.84±1.31	74.65±2.69 81.14±2.73	$77.81 {\scriptstyle \pm 3.83} \atop 80.29 {\scriptstyle \pm 0.35}$	62.58±2.88 67.74±0.79	$70.22{\pm}3.58$ 74.35 ${\pm}1.22$	81.75±2.32 81.97±2.87
	LSTM w/o Coverage + Coverage	60.83±1.29 65.02±4.16	81.45±1.10 82.53 ± 0.70	$78.91{\scriptstyle\pm1.05\atop82.92{\scriptstyle\pm3.18}}$	68.74±3.07 69.54±1.16	69.35±3.30 73.27±1.18	79.88±2.44 80.66±1.89

Table 2: Results of automatic evaluation, averaged over 3 runs \pm one standard deviation. The distance between the record and the exemplar is set to ≤ 5 for exemplar retrieval (see the text). Methods in the first block are two reference approaches (Section 5.2), i.e., AttnCopy-S2S for content fidelity and Slot-filling for style embodiment. For our method, we evaluate the variants with and without the coverage constraint (Section 4.3). The table highlights the best results in the blocks of Baselines and Ours under different metrics.

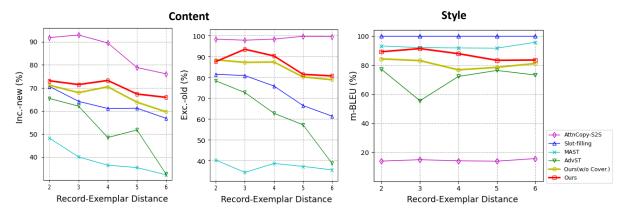


Figure 3: Effect of record-exemplar distance on model performance on the restaurant dataset. **Left**: Content fidelity performance, including "%Inc-new" and "%Exc-old". **Right**: Style embodiment performance by "m-BLEU".

the field values in record x as separate attributes. The method is based on back-translation (Sennrich et al., 2015) that first generates a target sentence \hat{y} conditioning on (x, y_e) , and then treat it as the reference to reconstruct y_e conditioning on (x_e, \hat{y}) . Auxiliary sentence y_x is used in an extra auto-encoding loss.

Adversarial Style Transfer (AdvST) (Logeswaran et al., 2018). As another style transfer approach for multiple attributes, the model incorporates back-translation with adversarial training to disentangle content and style representations.

Model Configurations

We studied both LSTM (Hochreiter and Schmidhuber, 1997) and Transformer (Vaswani et al., 2017) architectures. For LSTM, we use a single layer with the Luong attention (Luong et al., 2015) and copy mechanism (Gu et al., 2016). For Transformer, use the recent copy-augmented variant following (Su

et al., 2019) with 3 blocks. During training, we first set $(\lambda = 0, \eta = 0)$ to pre-train the model so that it captures the full characteristics of the exemplar sentence. We then switch to $(\lambda = 0.2, \eta = 1.0)$ for full training. Adam optimization (Kingma and Ba, 2014) is used with an initial learning rate of 0.001. At inference time, we use beam search with the width 5 and the maximum decoding length 50.

5.3 Automatic Evaluation

Metrics

Automatic evaluation of the task is an open and challenging problem. We use several quantitative metrics for the two goals of the task, namely content fidelity and style embodiment.

• Content fidelity. For the NBA data, we follow the original work (Wiseman et al., 2017) and use information extraction (IE) to measure content fidelity. Given a generated sentence \hat{y} and the input data record x, we extract field values from

	Restaura	ant Recomme	endations	Ī	NBA Reports	
Model	Content Fidelity	Style Embody	Fluency	Content Fidelity	Style Embody	Fluency
Slot-filling AdvST	3.36 3.56	5.00 4.24	4.70 4.02	2.79 2.88	5.00 4.00	4.86 4.09
Ours, LSTM w/o Coverage Ours, LSTM	3.91 4.28	4.38 4.73	4.58 4.54	3.43 3.88	4.13 4.53	4.59 4.52
	Ours Better	No Prefer.	Ours Worse	Ours Better	No Prefer.	Ours Worse
Slot-filling AdvST Ours, LSTM w/o Coverage	64.1 % 70.4 % 52.0 %	18.6% 14.3% 26.7%	17.3% 15.2% 21.3%	67.5% 68.8% 51.3%	17.5% 17.5% 32.5%	15.0% 13.8% 16.3%

Table 3: Results of human evaluation. Each metric achieves an average Pearson correlation coefficient \geq 0.73, showing a reasonable inter-annotator agreement. Our improvement in terms of mean annotator ratings is statistically significant (p<0.01, t-test). **Top:** Scoring three aspects on a 5-point Likert scale. **Bottom:** Ranking the generations from pairs of models. We use our LSTM-based full model to compare with other methods.

 \hat{y} with an IE tool and compute the precision and recall against x. We use the IE model provided in (Wiseman et al., 2017), which achieves 81% precision and 86% recall on the test set.

We found IE on the restaurant data is too difficult to serve as a reliable metric, because the descriptions are less structured. We thus instead train a BERT-based binary classifier to evaluate whether a field value is expressed in the generated sentence, which achieves 94% classification accuracy on the test set. We apply the classifier and compute both the percentage of desired \boldsymbol{x} field values expressed in the generation (%Incl.new) and the percentage of original content in \boldsymbol{y}_e (or equivalently, \boldsymbol{x}_e) removed from the generation (%Excl.-old). The higher both numbers, the more faithfully the generation describes \boldsymbol{x} .

• Style embodiment. Imitating the exemplar style involves inheriting the sentence structure, word choices, and other surface forms of y_e . Inspired by the text style transfer literature (Subramanian et al., 2019; Yang et al., 2018), we measure the BLEU score between the generated and the exemplar sentences. To reduce the influence of the change of content tokens, we mask in both sentences all obvious content tokens, e.g., player/team names and numbers, by replacing them with a special token <M>. We denote the metric as m-BLEU. This guarantees the reference approach, namely the slot-filling method, achieves an m-BLEU score of 100.

Study: Balance between Content and Style

Table 2 shows the automatic evaluation results on the two datasets. In this study, for exemplar retrieval (Section 4.2), we set the distance between a record and an exemplar to be no larger than 5 both during training and when constructing test cases. That is, the record and the exemplar sentence can have 5 mismatched fields, which thus requires strong flexibility of the generation model to be able to automatically adapt the exemplar in order to describe the record accurately.

As expected, the reference methods excel only in one of the two aspects, respectively. Specifically, AttnCopy-S2S expresses the desired content well, yet is incapable of embodying the designated style (e.g., m-BLEU=13.95). On the contrary, the Slot-filling method achieves perfect style m-BLEU by definition, but falls short of adaptively described the desired content in an accurate way, as shown by the low content scores. The two style transfer approaches (MAST and AdvST) also fail in terms of content fidelity performance. This is partly because these models are built on a different task assumption (i.e., modifying independent textual attributes) and are incompetent in manipulating the structured content well.

Our proposed approach is able to better balance between content fidelity and style embodiment. For example, in terms of content fidelity, our approach with an LSTM architecture improves over the Slot-filling results by 16.3 on NBA content precision and 12.9 on Restaurant content %Excl.-old. The approach meanwhile keeps a high style m-BLEU score of over 80. Regarding the ablation study, the results show the proposed content coverage constraint (Section 4.3) consistently improves both the content and style performance by a large margin. We note that the LSTM and transformer architectures perform comparably, with LSTM slightly better on the restaurant

Content Record	Name I Cocum o	EatType coffee shop	Food Italian	PriceRan £20-25	ge CustomRat	ng FamilyFriendly family friendly		
Exemplar 1	Looking fo	r French foo	d <mark>near</mark> Ziz	zzi? Come t	ry Strada, which ha	as a 3-star customer rating and priced lowly		
Slot filling	Looking fo priced £20-		food near	r Zizzi? Cor	ne try [] Cocum,	which has a high customer rating and		
AdvST	For Italian [] place near Zizzi? Come try [] Cocum, which has a high customer rating with priced £20-25.							
Ours	Looking for an Italian coffee shop? Come try family-friendly Cocum, which has a high customer rating and priced £20-25.							
Exemplar 2	Along the riverside near Cafe Rouge, there is a Japanese food place called The Golden Curry. It has an average customer rating since it is not a family-friendly environment.							
Slot-filling	Along the riverside near Cafe Rouge [], there is a Italian food [] place called Cocum. It has an high customer rating since it is not a family-friendly environment.							
AdvST	Along the riverside near the Ranch [], there is a Italian food [] place called Cocum. It has [] high customer rating since it is not a family-friendly environment.							
Ours	Priced £20-25, there is an Italian food coffee shop called Cocum. It has a high customer rating since it is a family-friendly environment.							
Content Record	PLAYER Patrick	PLAYER Dwight Ho		PLAYER Harden	PTS 10			
Exemplar	Both J.J. Hickson and Timofey Mozgov reached double - figures , scoring 10 and 15 points.							
Slot-filling	Both Patrick [] and Dwight Howard reached double - figures , scoring 10 and 15 points.							
AdvST	Both J.J. Hickson [] and Dwight Howard reached double - figures, scoring 10 and 10 points.							
Ours	Patrick , Dwight Howard and Harden reached double - figures , scoring 10 points.							

Table 4: Example outputs by different models given various exemplar sentences. Text of erroneous content and syntax are highlighted in red, where [...] indicates desired content that is missing. Text portions about the writing style in both exemplars and the generated sentences by our model are highlighted in blue.

dataset. We speculate that the copy mechanism of LSTM (Gu et al., 2016) is slightly more effective than that of transformer (Su et al., 2019).

Study: Effect of Record-Exemplar Distance

We then study how well the different methods would perform when given exemplars of varying distances (mismatchness) to the records. Figure 3 show the content and style results under different distances. We can see that, as the exemplars deviate more from the structure of the records, the model performance drops since it is getting harder to automatically adapt the exemplars to express the desired content. For example, the "%Excl.old" score (middle panel) of the methods Slot-filling and AdvST decreases quickly. Our approach maintains a more stable performance and keeps a better content-style balance. The results also show the proposed content coverage constraint consistently offers enhanced performance.

5.4 Human Evaluation

We also perform human evaluation for a more thorough and accurate comparison. Following the experimental settings in prior work (Subramanian et al., 2019; Logeswaran et al., 2018; Shen et al., 2017), we undertake two types of human evaluation: (1) We ask three human annotators to score generation results in three aspects, namely content fidelity, style embodiment, and sentence fluency, on a 5-point Likert scale. (2) We present to each annotator a pair of generated sentences, one from our model and the other from a comparison method, then ask the annotator to rank the two sentences by considering the above criteria jointly. Annotators can also choose "no preference" if the sentences are equally good or bad. For each study, we evaluated on 80 test instances. We use the LSTM architecture as it outperforms the transformer slightly in the automatic evaluation. We compare with the Slot-filling method, AdvST (which is better than MAST in automatic evaluation), and our variant without the coverage constraint.

Table 3 shows the results. From the top block, as discussed above, the Slot-filling method performs well in terms of style embodiment and fluency. However, its content fidelity is extremely weak. In contrast, our model achieves a better balance across the three criteria, by obtaining the best performance on content fidelity and reasonably

high scores on both style embodiment and fluency. The fluency of our full model is slightly inferior to the variant without the coverage constraint, which is not unexpected since the full model modifies more portions of the exemplar sentences, which would result in minor language mistakes.

The bottom block of Table 3 shows the human ranking results. We can see that our model consistently outperforms the comparison methods with over 50% wins on both datasets.

5.5 Qualitative Study

Table 4 shows samples on two test cases. We can see that the proposed full model performs superior to other approaches in effectively retaining the desired style and describing the content. For example, in the first two examples, other approaches often fail to remove the redundant content (e.g., "near Zizzi" or "riverside") from the generation while neglecting desired fields in the record. The proposed model performs better by adaptively adding and deleting text portions for accurate content description. Similarly, in the third case, both Slot-filling and AdvST fail to convey the new field value "Harden" given the exemplar, and leave in irrelevant information given the second one due to the different record structures between x and x_e . In contrast, our full model generates the desired sentence.

6 Conclusion

We have studied the new problem of data-to-text generation with style imitation. We developed a new approach with an attention-copy mechanism, weakly supervised learning, and a content coverage constraint. Experiments show the approach achieves a good balance between content fidelity and style control, and is flexible to adapt exemplars that do not match the record perfectly. We are interested in applying the style imitation approach to control longer paragraphs given full data tables.

References

- Gabor Angeli, Percy Liang, and Dan Klein. 2010. A simple domain-independent probabilistic approach to generation. In *EMNLP*, pages 502–512.
- Ziqiang Cao, Wenjie Li, Sujian Li, and Furu Wei. 2018. Retrieve, rerank and rewrite: Soft template based neural summarization. In *ACL*, pages 152–161.
- Mingda Chen, Qingming Tang, Sam Wiseman, and Kevin Gimpel. 2019. Controllable paraphrase generation with a syntactic exemplar. In *ACL*.

- Jan Milan Deriu and Mark Cieliebak. 2018. Syntactic manipulation for generating more diverse and interesting texts. In *INLG*, pages 22–34.
- Longxu Dou, Guanghui Qin, Jinpeng Wang, Jin-Ge Yao, and Chin-Yew Lin. 2018. Data2text studio: Automated text generation from structured data. In *EMNLP*, pages 13–18.
- Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2019. Evaluating the state-of-the-art of end-to-end natural language generation: The E2E NLG Challenge. *arXiv preprint arXiv:1901.11528*.
- Sebastian Gehrmann, Falcon Z Dai, Henry Elder, and Alexander M Rush. 2018. End-to-end content and plan selection for data-to-text generation. *arXiv* preprint arXiv:1810.04700.
- Jiatao Gu, Zhengdong Lu, Hang Li, and Victor OK Li. 2016. Incorporating copying mechanism in sequence-to-sequence learning. ACL.
- Tatsunori B Hashimoto, Kelvin Guu, Yonatan Oren, and Percy S Liang. 2018. A retrieve-and-edit framework for predicting structured outputs. In *NeurIPS*, pages 10073–10083.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*.
- Zhiting Hu and Eric P Xing. 2020. Learning from all types of experiences: A unifying machine learning perspective. In *KDD*.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P Xing. 2017. Toward controlled generation of text. In *ICML*.
- Hayate Iso, Yui Uehara, Tatsuya Ishigaki, Hiroshi Noji, Eiji Aramaki, Ichiro Kobayashi, Yusuke Miyao, Naoaki Okazaki, and Hiroya Takamura. 2019. Learning to select, track, and generate for data-to-text. In *ACL*, pages 2102–2113, Florence, Italy. Association for Computational Linguistics.
- Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. Adversarial example generation with syntactically controlled paraphrase networks. In *NAACL*.
- Glorianna Jagfeld, Sabrina Jenne, and Ngoc Thang Vu. 2018. Sequence-to-sequence models for data-to-text natural language generation: Word- vs. character-based processing and output diversity. In *INLG*, pages 221–232.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Ravi Kondadadi, Blake Howald, and Frank Schilder. 2013. A statistical NLG framework for aggregated planning and realization. In *ACL*, pages 1406–1415.

- Karen Kukich. 1983. Design of a knowledge-based report generator. In *ACL*, pages 145–150.
- Lajanugen Logeswaran, Honglak Lee, and Samy Bengio. 2018. Content preserving text generation with attribute controls. In *NeurIPS*, pages 5108–5118.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. *arXiv* preprint *arXiv*:1508.04025.
- Susan W. McRoy, Songsak Channarukul, and Syed S. Ali. 2000. YAG: A template-based generator for real-time systems. In *INLG*, pages 264–267.
- Gaurav Pandey, Danish Contractor, Vineet Kumar, and Sachindra Joshi. 2018. Exemplar encoder-decoder for neural conversation generation. In *ACL*, pages 1329–1338.
- Hao Peng, Ankur P Parikh, Manaal Faruqui, Bhuwan Dhingra, and Dipanjan Das. 2019. Text generation with exemplar-based adaptive decoding. In *NAACL*.
- Richard Power, Donia Scott, and Nadjet Bouayad-Agha. 2003. Generating texts with style. In *International Conference on Intelligent Text Processing and Computational Linguistics*, pages 444–452. Springer.
- Ratish Puduppully, Li Dong, and Mirella Lapata. 2019. Data-to-text generation with entity modeling. In *ACL*, pages 2023–2035.
- Ehud Reiter and Robert Dale. 1997. Building applied natural language generation systems. *Natural Language Engineering*, 3(1):57–87.
- Jacques Robin and Kathleen McKeown. 1996. Empirically designing and evaluating a new revision-based model for summary generation. *Artificial Intelligence*, 85(1-2):135–179.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Improving neural machine translation models with monolingual data. In *ACL*.
- Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2017. Style transfer from non-parallel text by cross-alignment. In *NeurIPS*, pages 6830–6841.

- Hui Su, Xiaoyu Shen, Rongzhi Zhang, Fei Sun, Pengwei Hu, Cheng Niu, and Jie Zhou. 2019. Improving multi-turn dialogue modelling with utterance ReWriter. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 22–31, Florence, Italy. Association for Computational Linguistics.
- Sandeep Subramanian, Guillaume Lample, Eric Michael Smith, Ludovic Denoyer, Marc'Aurelio Ranzato, and Y-Lan Boureau. 2019. Multiple-attribute text rewriting. In *ICLR*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *NeurIPS*, pages 3104–3112.
- Bowen Tan, Lianhui Qin, Eric P Xing, and Zhiting Hu. 2020. Summarizing text on any aspects: A knowledge-informed weakly-supervised approach. In *EMNLP*.
- Jianheng Tang, Tiancheng Zhao, Chengyan Xiong, Xiaodan Liang, Eric P Xing, and Zhiting Hu. 2019. Target-guided open-domain conversation. In *ACL*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*, pages 5998–6008.
- Jason Weston, Emily Dinan, and Alexander H Miller. 2018. Retrieve and refine: Improved sequence generation models for dialogue. arXiv preprint arXiv:1808.04776.
- Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2017. Challenges in data-to-document generation. In *EMNLP*.
- Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2018. Learning neural templates for text generation. In *EMNLP*.
- Zichao Yang, Zhiting Hu, Chris Dyer, Eric Xing, and Taylor Berg-Kirkpatrick. 2018. Unsupervised text style transfer using language models as discriminators. In *NeurIPS*.
- Rong Ye, Wenxian Shi, Hao Zhou, Zhongyu Wei, and Lei Li. 2020. Variational template machine for data-to-text generation. In *ICLR*.