

KGPT: Knowledge-Grounded Pre-Training for Data-to-Text Generation

Wenhu Chen, Yu Su, Xifeng Yan, William Yang Wang

University of California, Santa Barbara, CA, USA

{wenhuchen, xyan, william}@cs.ucsb.edu, su.809@osu.edu

Abstract

Data-to-text generation has recently attracted substantial interests due to its wide applications. Existing methods have shown impressive performance on an array of tasks. However, they rely on a significant amount of labeled data for each task, which is costly to acquire and thus limits their application to new tasks and domains. In this paper, we propose to leverage pre-training and transfer learning to address this issue. We propose a knowledge-grounded pre-training (KGPT), which consists of two parts, 1) a general knowledge-grounded generation model to generate knowledge-enriched text. 2) a pre-training paradigm on a massive knowledge-grounded text corpus crawled from the web. The pre-trained model can be fine-tuned on various data-to-text generation tasks to generate task-specific text. We adopt three settings, namely fully-supervised, zero-shot, few-shot to evaluate its effectiveness. Under the fully-supervised setting, our model can achieve remarkable gains over the known baselines. Under zero-shot setting, our model without seeing any examples achieves over 30 ROUGE-L on WebNLG while all other baselines fail. Under the few-shot setting, our model only needs about one-fifteenth as many labeled examples to achieve the same level of performance as baseline models. These experiments consistently prove the strong generalization ability of our proposed framework¹.

1 Introduction

Data-to-text generation, i.e., generating textual description from structured data, is an important task with many real-world applications such as generating weather reports (Liang et al., 2009), sports news (Wiseman et al., 2017), dialog response (Wen et al., 2016; Dušek et al., 2019), etc. Neural gener-

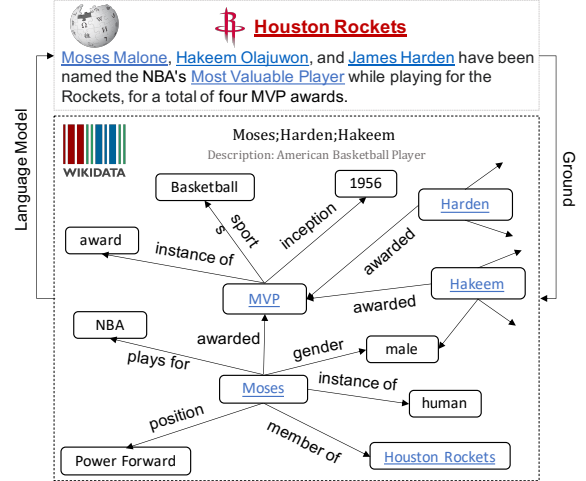


Figure 1: An example from the constructed KGTEXT, which pairs a hyperlinked sentence from Wikipedia with a knowledge subgraph from WikiData.

ation models based on different strategies like soft-template (Wiseman et al., 2018; Ye et al., 2020), copy-mechanism (See et al., 2017), content planning (Reed et al., 2018; Moryossef et al., 2019), and structure awareness (Liu et al., 2018; Colin and Gardent, 2019) have achieved impressive results. However, existing studies are primarily focused on fully supervised setting requiring substantial labeled annotated data for each subtask, which restricts their adoption in real-world applications.

In this paper, we are interested in developing a general-purpose model that can easily adapt to different domains/tasks and achieve strong performance with only a small amount or even zero annotated examples. Our model draws inspiration from the recent wave of pre-trained language model (Devlin et al., 2019; Radford et al., 2019; Dai et al., 2019) to exploit large-scale unlabeled data from the web for pre-training. The data pairs are constructed through the following procedure. We first crawl sentences with hyperlinks from Wikipedia, and then link the hyperlinked entities to Wiki-

¹<https://github.com/wenhuchen/KGPT>

Data (Vrandečić and Krötzsch, 2014) to find their 1-hop knowledge triples. Finally, we build a subgraph based on the linked triples. Such automatic alignment between knowledge graph and texts provides distant supervision (Mintz et al., 2009) for pre-training but it is bound to be noisy. Therefore, we design a selection strategy and only retain plausible alignments with high semantic overlap. The harvested knowledge-grounded corpus KGTEXT consists of over 1.8M (knowledge subgraph, text) pairs, as depicted in Figure 1.

We unify the input of KGTEXT and downstream data-to-text tasks into a generalized format and design a novel architecture KGPT to encode it. We use KGTEXT to first pre-train KGPT and then fine-tune it on downstream data-to-text tasks like WebNLG (Shimorina and Gardent, 2018), E2ENLG (Dušek et al., 2019) and WikiBio (Liu et al., 2018). Experimental results demonstrate KGPT’s several advantages: 1) with full downstream dataset, KGPT can achieve remarkably better performance than known competitive baselines, 2) with zero training, KGPT can still achieve a reasonable score on WebNLG. 3) with a few training instances, KGPT can maintain a high BLEU score while the non-pre-trained baselines only generate gibberish text. A quantitative study shows that our pre-training scheme can reduce annotation costs by roughly 15x to achieve a decent BLEU score of 30. Our contribution is summarized as follows:

- i). We design a distantly supervised learning algorithm to exploit large-scale unlabeled web text to pre-train data-to-text models.
- ii). The proposed pre-training algorithm can bring significant performance under different settings, especially zero-shot and few-shot scenarios.

2 Related Work

Data-to-Text Generation Data-to-text is a long-standing problem (Kukich, 1983; Reiter and Dale, 1997), which involves generating natural language surface form from structured data. The traditional system is primarily built on a template-based algorithm. Recently, with the development of deep learning, attention has been gradually shifted to end-to-end neural generation models, which achieve significant performances on existing large-scale datasets like WebNLG (Shimorina and Gardent, 2018), E2ENLG (Dušek et al., 2019), WikiBio (Lebret et al., 2016), ROTOWIRE (Wiseman et al., 2017), TOTTO (Parikh et al., 2020), Log-

icNLG (Chen et al., 2020a), etc. However, these neural generation models are mainly focused on fully supervised learning requiring a huge amount of human annotation for the specific task. Our paper focuses on building a more generalized model architecture, which can adapt to specific tasks well with only a handful of training instances.

Knowledge-Grounded Language Modeling It is of primary importance to ground language models on existing knowledge of various forms. The neural language models (Bengio et al., 2003) have been shown to well capture the co-occurrences of n-grams in the sentences, but falls short to maintain the faithfulness or consistency to world facts. To combat such an issue, different knowledge-grounded language models (Ahn et al., 2016; Hayashi et al., 2020; Logan et al., 2019) have been proposed to infuse structured knowledge into the neural language model. These models are mainly focused on enhancing the factualness of unconditional generative models. Inspired by these pioneering studies, we explore the possibility to connect the unconditional generative model with downstream conditional generation tasks. The most straightforward knowledge-intensive conditional generative task is the data-to-text generation, which aims to verbatim given knowledge into lexical format. We demonstrate great potential of the knowledge-grounded pretraining in enhancing the model’s factualness on these down-stream data-to-text tasks and believe such language models can be applied to broader range of NLP tasks requiring knowledge understanding.

Pre-trained Language Model Recently, the research community has witnessed the remarkable success of pre-training methods in a wide range of NLP tasks (Devlin et al., 2019; Radford et al., 2018, 2019; Dai et al., 2019; Yang et al., 2019; Liu et al., 2019b; Kesar et al., 2019; Lan et al., 2020; Lewis et al., 2019; Raffel et al., 2019). These models trained on millions or billions of data unlabeled data demonstrate unprecedented generalization ability to solve related down-stream tasks. However, the existing pre-trained text generation models (Radford et al., 2019; Kesar et al., 2019; Raffel et al., 2019) are initially designed to condition on text input, thus lacking the ability to encode structured inputs. The work closest to our concept is Switch-GPT-2 (Chen et al., 2020b), which fits the pre-trained GPT-2 model as the decoder part

to perform table-to-text generation. However, their knowledge encoder is still trained from scratch, which compromises the performance. In this paper, we follow the existing paradigm to construct an unlabeled web data for LM pre-training.

3 Dataset Construction

The construction process has two stages, namely the crawling stage and the selection stage:

3.1 Hyperlinked Sentence Crawling

We use English Wikidump² as our data source. For each Wikipedia page, we split the whole paragraphs into an array of sentences and then tokenize with the nltk toolkit (Loper and Bird, 2002). We loop through each sentence to keep the sentences with more than 2 Wikipedia anchor links and within the length of 10 and 50. For each candidate sentence, we use its Wikipedia hyperlink to query WikiData (Vrandečić and Krötzsch, 2014) and obtain its corresponding entity page³. We retrieve the neighboring knowledge triples from these entity pages to construct a local 1-hop graph for each entity. The knowledge triples are divided into two types: 1) the object of the triple is also an entity like ‘(Roma F.C., country, Italy)’, 2) the object of the triple is in plain text like ‘(Roma F.C., inception, 7 June 1927)’. In the first case, if the object entity also appears in the sentence, we use it as the bridge to build a multi-hop graph like Figure 2. After this step, we collected roughly 4 million pairs in the form of (subgraph, sentence) as the candidate for the following step.

3.2 Data Selection

We observe that the collected pairs are overly noisy with many sentences totally irrelevant to their paired subgraphs. Apparently, these pairs cannot serve our goal to build a knowledge-grounded language model. Therefore, we propose a data selection step to suppress the noise and filter out the data pairs of our interests. An example is depicted in Figure 2, the first sentence does not rely on any information provided by the knowledge graph, while the second sentence has a tight connection to the facts presented in the knowledge graph. Ideally, our proposed strategy should favor the second sentence over the first one.

To achieve this, we propose a simple lexical-based selection strategy to perform data selection. For example, the sentence ‘He was born ...’ in Figure 2 has two query words ‘Italy’ and ‘Germany’, we will conduct two rounds of lexical matching. In the first round, we use ‘Italy’ to query its surrounding neighbors in WikiData to the neighboring unigram, i.e. ‘(Rome, capital, Europe, Continent, Country, Roma F.C)’. We compute the unigram overlap with the original sentence ‘(He, was, ...)’, which is still 0%. In the second round, we use ‘Germany’ to do the same computation and calculate the lexical overlap, which is still 0%. So the final averaged grounding score of two rounds is 0%. We can follow the same procedure to compute the grounding score for the second sentence in Figure 2 with four rounds ‘(AS Roma, FB, Rome, Italy)’. The grounding score is above 30%, which indicates that the sentence is highly grounded on WikiData subgraph. In this paper, we use a threshold of 0.13, which selects the top 7M ‘good’ sentences from the original 12M Wikipedia corpus.

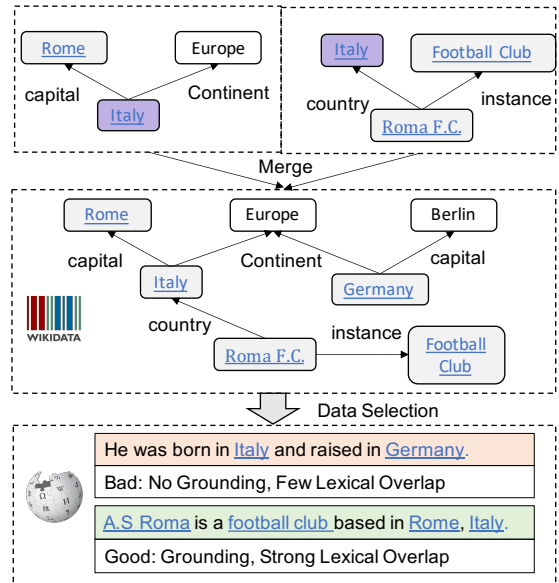


Figure 2: Data denoising procedure for the KGTEXT.

After the selection step, we obtain a denoised knowledge-grounded corpus KGTEXT for pre-training. However, there still exist noisy false positives in the corpus, for example, a subgraph contains triple ‘(Roma F.C., country, Italy)’, which is associated with the text ‘An Italian player plays for A.S. Roma’. Though the two entities co-occur, they are not meant to describe the fact triple. By applying more strict rules, we can suppress such false positives, but the data capacity could significantly drop consequently. We experimented with differ-

²<https://dumps.wikimedia.org/>

³<https://www.wikidata.org>

ent thresholds to balance noise and data capacity and finally decide on a threshold with an acceptable noise degree. The detailed statistics of the KGTEXT is listed in Table 1. We held-out 10,000 sentences for both validation and testing to evaluate the pre-trained model.

#Sent	Length	#Ent	#Pred	#Triple	#Ent/Sent
7M	20.2	1.8M	1210	16M	3.0

Table 1: Statistics of collected KGText dataset

4 Model

We formally define the problem setting and KGPT’s architectures in this section.

4.1 Problem Setting

In this paper, we consider inputs from structured data with diverse formats, like knowledge subgraph in KGTEXT, dialog act in E2E (Dušek et al., 2019), RDF triples in WebNLG (Shimorina and Gardent, 2018) and tables in WikiBio (Lebret et al., 2016). Here we unify them into a generalized dictionary format, which uses keys to represent subjects and values to denote the predicate-object pairs following the subject. We showcase the conversion criteria from structured inputs in different data-to-text datasets into our generalized format in Figure 3. The generalized input is denoted as X , and the output is denoted as y . Our model encodes X into a sequence of dense vectors, and then uses the decoder to attend and generate y .

4.2 Encoder

The encoder network is crucial to our model to capture the highly structured graph input. We mainly experiment with two types of encoders:

Graph Encoder This encoder is mainly based on graph attention network (Li et al., 2016; Kipf and Welling, 2017; Veličković et al., 2018) to explicitly encode the structure information. Specifically, we view each object, predicates, and subjects as the leaf nodes, and add [ENT], [TRIPLE] as pseudo nodes for message passing purposes. The built graph is depicted in Figure 4.

First of all, we initialize the node representation with the averaged embedding of its subword units. For example, the node ‘Moses Malone’ has a representation of $(E[\text{Mos}] + E[\text{es}] + E[\text{Ma}] + E[\text{lone}]) / 4$ with E denoting the embedding. After we obtain the initial node representation, we use

message propagation to update the node representations based on neighboring information.

In the first layer, we exchange the information between nodes inside a triple, e.g., ‘Moses Malone’ receives message from siblings ‘Gender’ and ‘Male’. In the second layer, we aggregate information from sub/pred/obj nodes to the [TRIPLE] node, e.g., [‘TRIPLE1’] receives message from children ‘Moses, Gender, Male’. In the third layer, we aggregate the information from different [TRIPLE] to the [ENT] node. In the fourth layer, we exchange information between different [ENT] nodes to enhance cross-entity interactions. Formally, we propose to update the representation of the i -th node $g_i \in \mathbb{R}^D$ with the multi-head attention network, which aggregates information from neighboring nodes $g_j \in \mathcal{N}_i$ as follows:

$$\begin{aligned} \alpha_j^m &= \frac{e^{(W_Q^m g_i)^T (W_K^m g_j)}}{\sum_{j \in \mathcal{N}_i} e^{(W_Q^m g_i)^T (W_K^m g_j)}} \\ v &= \text{concat}[\sum_{j \in \mathcal{N}_i} \alpha_j^m W_V^m(g_j)] \\ \hat{g}_i &= \text{LayerNorm}(\text{MLP}(v + g_i)) \end{aligned} \quad (1)$$

where m denotes the m -th head in the attention layer, $W_Q^m, W_K^m, W_V^m \in \mathbb{R}^{D \times D}$ are the matrices to output query, key, value vectors for m -th head. The attention output v and the residue connection from g_i are fed through the final MLP and LayerNorm to update i -th node representation as \hat{g}_i . The output of graph encoder is denoted as $G \in \mathbb{R}^{n \times D} = \{g_1, \dots, g_n\}$ with n nodes.

Sequence Encoder This encoder is mainly based on transformer (Vaswani et al., 2017) with special embedding as an auxiliary input to infuse the structure information to the sequence model. The concept of special embedding was initially proposed by BERT (Devlin et al., 2019), more recently, it has been adopted by Herzig et al. (2020) to infuse structural information. We visualize the embedding layer in Figure 5, where we leverage additional entity embedding, triple embedding, and property embedding to softly encode the structure of the subgraph as a linearized sequence. For example, the entity embedding can inform the model which entity the current token belongs to, while the triple embedding can indicate which triple the current token belongs to and the property embedding indicates whether the token is a subject, predicate, or a subject. Such an encoding mechanism is designed to softly encode the graph structure into the embedding space for further self-attention. Compared

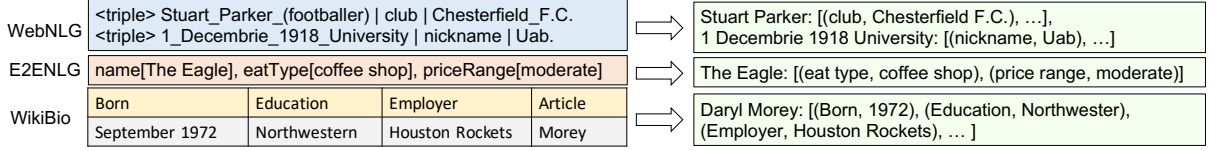


Figure 3: The conversion criterion to unify different structured data input into our generalized format.

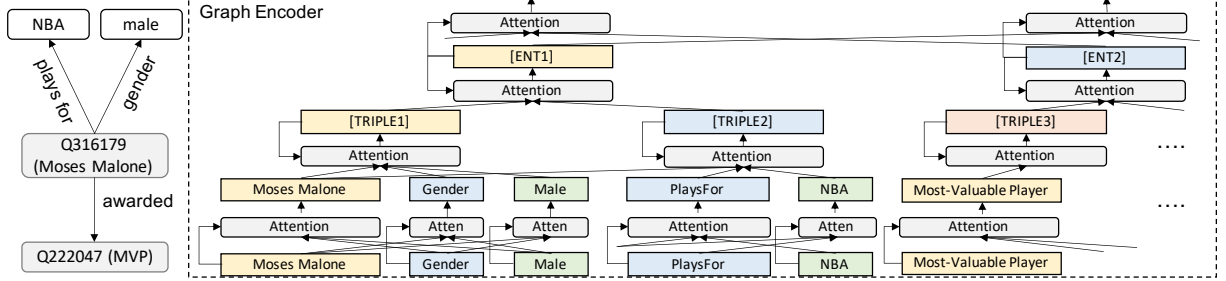


Figure 4: Graph Encoder with hierarchical propagation, where we propagate the information from bottom to top.

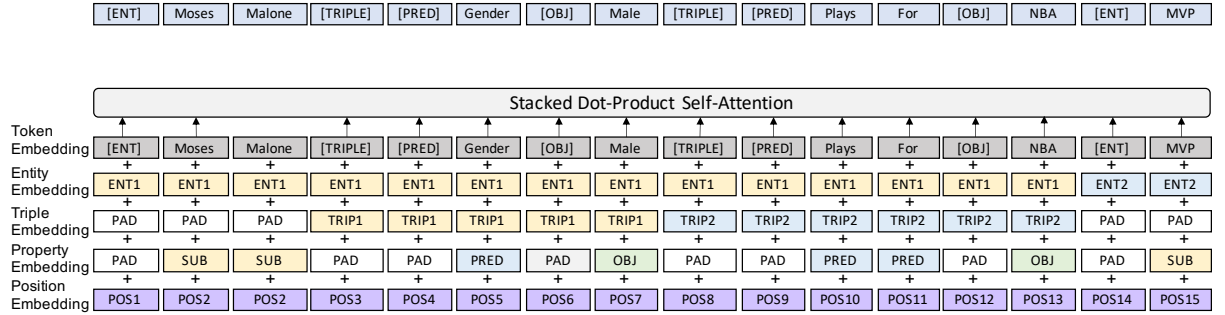


Figure 5: Encoding of the knowledge graph as a sequence using special embedding.

to the graph encoder, the sequence encoder does not enforce the structure as a hard constraint and allows more flexibility for the model to perform cross-triple and cross-entity interactions. Formally, the dot-product self-attention follows the definition of Transformer (Vaswani et al., 2017):

$$f_{att}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V$$

$$G_m = f_{att}(QW_Q^m, KW_K^m, VW_V^m)$$

$$G = \text{MLP}(\text{Concat}(G_1, \dots, G_m))$$
(2)

where Q, K, V are the computed from the input embedding, m represents m -th head and f_{att} is the core attention function, the final output is denoted as $G \in \mathbb{R}^{n \times D}$ with n denoting the sequence length.

4.3 Decoder

Our decoder architecture is mainly based on Transformer (Vaswani et al., 2017) and copy mechanism (See et al., 2017). At each decoding time step, the model has a copy gate p_{gen} to select y_i should be generated from the vocabulary $w \in \mathcal{V}$ or

copied from the input tokens x :

$$\alpha_j = \frac{e^{o_i^T G_j}}{\sum_{j'} e^{o_i^T G_{j'}}}, \quad p_{gen} = \sigma(\text{MLP}(o_i))$$
(3)

$$P(y_i = w) = p_{gen} P_{voc}(w) + (1 - p_{gen}) \sum_{j: x_j = w} \alpha_j$$

where o_i is the last layer hidden state of the decoder at i -th time step, α_j is the copy probability over the whole input token sequences x .

4.4 Optimization

As we have defined our encoder-decoder model, we will simply represent it as $p_{encdec}(x)$ to output a distribution over word $y_i \in \mathcal{V}$ at the i -th time step. During pre-training, we optimize the log-likelihood function on D_{KGText} . After pre-training, we convert the downstream task's input into the defined dictionary format and denote the dataset as D_{down} , and then further optimize the log-likelihood objective with θ initialized from the pre-training stage.

The pre-train and fine-tuning procedure is displayed in Figure 6, where we first use KGTEXT to pre-train KGPT, and then fine-tune with different

types of inputs using the standard auto-regressive log-likelihood objective.

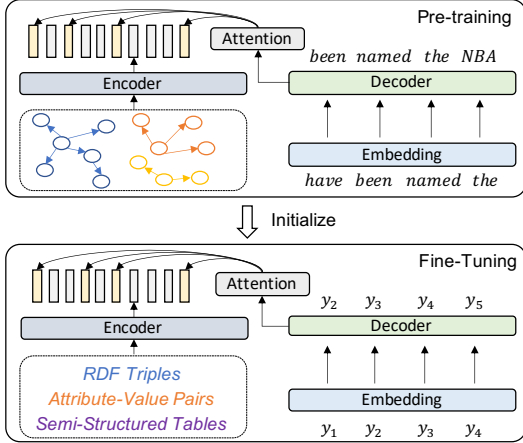


Figure 6: Overall pre-training and fine-tuning procedures for KGPT. The downstream knowledge data formats are converted into the generalized format.

5 Experiments

We experiment with three different down-stream tasks, which covers various table-to-text applications to verify the generalization capability of KGPT. Besides the fully supervised learning, we also evaluate zero-shot and few-shot learning.

5.1 Datasets

We use WebNLG (Shimorina and Gardent, 2018), E2ENLG (Dušek et al., 2019) and WikiBio (Lebret et al., 2016) to evaluate the performance of KGPT. Their basic statistics are listed in Table 2. WebNLG and E2ENLG are both crowd-sourced by human annotator while WikiBio is from the Web.

Dataset	Train	Val	Test	Input
WebNLG	34,338	4,313	4,222	RDF Triple
E2ENLG	42,061	4,672	4,693	Dialog Act
WikiBio	582,657	72,831	72,831	Table

Table 2: Statistics of different data-to-text datasets

WebNLG This dataset (Shimorina and Gardent, 2018) aims to convert RDF triples into a human annotated textual description. We use the recent release 2.0 from GitLab⁴. It contains sets with up to 7 triples each along with one or more references. The number of KB relations modeled in this scenario is potentially large and generation involves solving various subtasks (e.g. lexicalisation

⁴<https://gitlab.com/shimorina/webnlg-dataset>

and aggregation). As the input RDF triples were modified from the original triples in DBPedia, we first need to check whether there are seen triples in pre-training dataset KGTEXT. We verify that there is zero RDF triple seen during pre-training though 31% entities are seen. Therefore, we can confirm the comparison with other baselines is still fair given no information from test/dev is leaked.

E2ENLG This dataset (Dušek et al., 2019) aims to convert dialog act-based meaning representation into a spoken dialog response. It aims to provide higher-quality training data for end-to-end language generation systems to learn to produce more naturally sounding utterances. In this dataset, each meaning representation is associated with on average with 8.65 different reference utterances.

WikiBio This dataset (Lebret et al., 2016) aims to generate the first sentence of biography description based on a Wikipedia infoboxes table, with each table associated with only one reference. Unlike the previous two human-annotated datasets from different domains, WikiBio is also scraped from Wikipedia. Therefore, we filtered out the instances of KGTEXT from the first paragraph of the biography domain to ensure no overlap or leakage about Wikibio’s dev/test set.

5.2 Experimental Setup

We apply the standard GPT-2 (Radford et al., 2019) tokenizer from Huggingface Github⁵ to tokenize the text input, which has a vocabulary of over 50K subword units. We test with both graph encoder and sequence encoder. We set their hidden size to 768 and stack 6 layers for both encoder and decoder with 8 attention heads. During pre-training, we run the model on KGTEXT on 8 Titan RTX GPUs with a batch size of 512 for 15 epochs using Adam (Kingma and Ba, 2015) optimizer with a learning rate of 1e-4. The pre-training procedure takes roughly 8 days to finish. We use a held-out validation set to select the best checkpoint. During fine-tuning, we use a learning rate of 2e-5.

In our following experiments, we compare with the known best models from different datasets. As none of these models are pre-trained, we also add Template-GPT-2 (Chen et al., 2020a) and Switch-GPT-2 (Chen et al., 2020b) as our pre-trained baselines. Both models apply GPT-2 (Radford et al.,

⁵<https://github.com/huggingface/transformers>

2019) as the generator to decode description from a table. For the ablation purposes, we list the performance of all non-pre-trained KGPT to see the performance gain brought by pre-training alone. All the best models are selected based on the validation set score, and the numbers are reported in the following tables are for test split. For evaluation, we report the performance with BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005) and ROUGE-L (Lin, 2004) using e2e-metric⁶. It’s worth noting that we perform comprehensive data contamination studies in the following experiments to make sure the pre-training data contains very little overlap with the test split in downstream tasks. We filter out potentially information-leaking pages during the data crawling process.

5.3 Preliminary Study on KGTEXT

In the preliminary study, we evaluate our pre-trained model’s performance on the held-out set of KGTEXT to conduct ablation study over KGPT. Specifically, we investigate 1) which encoding mechanism is better, 2) whether we need copy mechanism or copy supervision. As demonstrated in Table 3, we observe that the trivial difference between two encoder designs. With the copy mechanism, KGPT can greatly decrease the perplexity. However, supervising the copy attention does not have much influence on the performance. Therefore, in the following experiments, we will run experiments for both encoding schemes with a copy mechanism without copy loss.

Model	BLEU-4	Perplexity
KGPT-Graph	24.71	4.86
KGPT-Graph + Copy Loss	24.77	4.91
KGPT-Graph w/o Copy	22.69	7.23
KGPT-Seq	24.49	4.95
KGPT-Seq + Copy Loss	24.31	4.93
KGPT-Seq w/o Copy	22.92	7.11

Table 3: Ablation Study on held-out set of KGTEXT.

5.4 Fully-Supervised Results

We experiment with KGPT under the standard fully-supervised setting to compare its performance with other state-of-the-art algorithms.

WebNLG Challenge We list WebNLG’s experimental results in Table 4, here we compare with

⁶<https://github.com/tuetschek/e2e-metrics>

the known models under the unconstrained setting. The baseline models (Shimorina and Gardent, 2018) uses sequence-to-sequence attention model (Luong et al., 2015) as the backbone and propose delexicalization and copy mechanism to enhance model’s capability to handle rare items from the input. The GCN model (Marcheggiani and Perez-Beltrachini, 2018) uses graph convolutional neural encoder to encode the structured data input. Its implementation is from Github⁷. As can be seen, KGPT without pre-training already achieves better performance than the GCN baseline. With pre-training, the performance is further boosted by 1-2 BLEU-4, which reflects the effectiveness of our method.

Model	BLEU	METEOR	ROUGE
Seq2Seq [†]	54.0	37.0	64.0
Seq2Seq+Delex [†]	56.0	39.0	67.0
Seq2Seq+Copy [†]	61.0	42.0	71.0
GCN	60.80	42.76	71.13
KGPT-Graph w/o Pre	62.30	44.33	73.00
KGPT-Seq w/o Pre	61.79	44.39	72.97
KGPT-Graph w/ Pre	63.84	46.10	74.04
KGPT-Seq w/ Pre	64.11	46.30	74.57

Table 4: Experimental results on WebNLG’s test set, w/ Pre refers to the model with pre-training, otherwise it refers to the model training from scratch. [†] results are copied from Shimorina and Gardent (2018).

E2E Challenge We list E2ENLG’s experimental results in Table 5, here we compare with the state-of-the-art systems on the leaderboard of E2E challenge⁸. These baselines methods are based on neural template model (Wiseman et al., 2018), syntax-enhanced algorithms (Dušek and Jurcicek, 2016), slot alignment (Juraska et al., 2018) and controlling mechanism (Elder et al., 2018). As is seen from the table, KGPT can beat the SOTA systems by a remarkable margin. Overall, the improvement brought by pre-training is roughly 0.5-1.0 in terms of BLEU-4, which is less significant than WebNLG. Such a phenomena is understandable given that this dataset contains limited patterns and vocabulary in the input meaning representation, a full training set over 40K instances is more than enough for the generation model to memorize. In the following few-shot experiments, we will show the strength of KGPT to generate high-quality faithful descriptions with only 0.1% of training data.

⁷<https://github.com/diegma/graph-2-text>

⁸<http://www.macs.hw.ac.uk/InteractionLab/E2E/>

Model	BLEU	METEOR	ROUGE
NTemp	55.17	38.75	65.01
TGen	65.93	44.83	68.50
SLUG2SLUG	66.19	44.54	67.72
Adapt	67.37	45.23	70.89
KGPT-Graph w/o Pre	66.47	44.20	67.78
KGPT-Seq w/o Pre	67.67	45.33	70.39
KGPT-Graph w/ Pre	67.87	44.50	70.00
KGPT-Seq w/ Pre	68.05	45.80	70.92

Table 5: Experimental results on E2E’s test set. NTemp is from [Wiseman et al. \(2018\)](#), TGen is from [Dušek and Jurcicek \(2016\)](#), SLUG2SLUG is from [Juraska et al. \(2018\)](#) and Adapt is from [Elder et al. \(2018\)](#).

WikiBio Dataset We list WikiBio’s experimental results in Table 6 and compare with models like Table2Seq([Bao et al., 2018](#)), Order Planning ([Sha et al., 2018](#)), Field Gating ([Liu et al., 2018](#)), Background-KB Attention ([Chen et al., 2019](#)), Hybrid Hierarchical Model ([Liu et al., 2019a](#)) trained with multiple auxiliary loss functions. We also train Template-GPT-2 on this dataset to observe pre-trained model’s performance. As can be seen from the table, KGPT can achieve better results than the mentioned baseline models. Pre-training can yield an improvement of roughly 0.5 BLEU-4. As this dataset trainin/testing have similar table schema and the large number of training instances already teach the model to memorize the generation patterns, exploiting an external corpus of on par size (1.8M) does not bring a significant boost. So is the template-GPT-2 ([Chen et al., 2020a](#)), which performs on par with Field Gating ([Liu et al., 2018](#)). However, in the few-shot setting, we will show the 25+ BLEU gain brought by pre-training.

Model	BLEU
Table NLM (Lebret et al., 2016)	34.70
Table2Seq (Bao et al., 2018)	40.26
Order Planning (Sha et al., 2018)	43.91
Field-Gating (Liu et al., 2018)	44.71
KBAtt (Chen et al., 2019)	44.59
Hierarchical+Auxiliary Loss (Liu et al., 2019a)	45.01
Template-GPT-2	44.67
KGPT-Graph w/o Pre	44.64
KGPT-Seq w/o Pre	44.58
KGPT-Graph w/ Pre	45.10
KGPT-Seq w/ Pre	45.06

Table 6: Experimental results on WikiBio’s test set.

5.5 Few-Shot Results

The few-shot learning setting aims to study the potential of the proposed pre-training to decrease annotation labor in data-to-text generation tasks. Under this setting, we not only compare with non-pre-trained baselines to observe how pre-training

can benefit the model’s few-shot learning capability but also compare with other pre-trained LM ([Chen et al., 2020b,a](#)) to see the benefit of KGPT over existing pre-trained LM.

Model	0.5%	1%	5%	10%
Seq2Seq	1.0	2.4	5.2	12.8
Seq2Seq+Delex	4.6	7.6	15.8	23.1
KGPT-Graph w/o Pre	0.6	2.1	5.9	14.4
KGPT-Seq w/o Pre	0.2	1.7	5.1	13.7
Template-GPT-2	8.5	12.1	35.3	41.6
KGPT-Graph w/ Pre	22.3	25.6	41.2	47.9
KGPT-Seq w/ Pre	21.1	24.7	40.2	46.5

Table 7: Few-shot results on WebNLG’s test set.

Model	0.1%	0.5%	1%	5%
TGen	3.6	27.9	35.2	57.3
KGPT-Graph w/o Pre	2.5	26.8	34.1	57.8
KGPT-Seq w/o Pre	3.5	27.3	33.3	57.6
Template-GPT-2	22.5	47.8	53.3	59.9
KGPT-Graph w/ Pre	39.8	53.3	55.1	61.5
KGPT-Seq w/ Pre	40.2	53.0	54.1	61.1

Table 8: Few-shot results on E2ENLG’s’s test set.

WebNLG & E2ENLG Dataset In these two datasets, we use 0.1%, 0.5%, 1%, 5%, 10% of training instances to train the model and observe its performance curve in terms of BLEU-4.

For WebNLG challenge, the few-shot situation will pose a lot of unseen entities during test time. From Table 7, we can observe that the delexicalization mechanism can remarkably help with the few-shot situation. However, the improvement brought by delexicalization is much weaker than our proposed pre-training. Under the 5% setting, while the non-pre-trained baselines are only able to generate gibberish text, pre-trained KGPT can maintain a high BLEU score over 40.0 due to its strong generalization ability.

For E2E challenge, the task is comparatively simpler with rather limited items. From Table 8, we can observe that TGen ([Dušek and Jurcicek, 2016](#)) is achieving similar performance as our non-pre-trained KGPT, they both perform quite well even under 1% training instances. However, after we further reduce the training samples to roughly 0.1%, the baseline models fail while pre-trained KGPT still maintains a decent BLEU over 40.0.

WikiBio Dataset In this dataset, we adopt the same setting as Switch-GPT-2 ([Chen et al., 2020b](#)) and Pivot ([Ma et al., 2019](#)) to use 50, 100, 200 and 500 samples from the training set to train the generation model. From the results in Table 9, we observe that KGPT can achieve best scores and out-

perform both Template-GPT-2 and Switch-GPT-2 under most cases. Though Template-GPT-2 is getting slightly better score with 500 training samples, the overall performance on three datasets are remarkably lower than KGPT, especially under more extreme cases. It demonstrates the advantage of our knowledge-grounded pre-training objective over the naive LM pre-training objective.

Model	50	100	200	500
Field-Infusing	1.3	2.6	3.1	8.2
KGPT-Graph w/o Pre	0.2	1.1	3.8	9.7
KGPT-Seq w/o Pre	0.6	1.7	3.0	8.9
Pivot [†]	7.0	10.2	16.8	20.3
Switch-GPT-2 [†]	17.2	23.8	25.4	28.6
Template-GPT-2	19.6	25.2	28.8	30.8
KGPT-Graph w/ Pre	24.5	27.5	28.9	30.1
KGPT-Seq w/ Pre	24.2	27.6	29.1	30.0

Table 9: Few-shot results on Wikibio’s test set. [†] results are copied from [Chen et al. \(2020b\)](#).

Quantitative Study We further investigate how much sample complexity KGPT can reduce. Specifically, we specify a BLEU-4 score and vary the training data size to observe how much training samples are required to attain the performance. We specify BLEU=30 as our standard and display our results in [Table 10](#). We compute the ratio of

Model	WebNLG	E2ENLG	WikiBio
KGPT w/o Pre	~10000	~300	~8000
KGPT w/ Pre	~700	~20	~500
Ratio	14x	15x	16x

Table 10: Required number of training samples to reach designated BLEU on different dataset.

sample quantity to characterize the benefits from pre-training. Roughly speaking, pre-training can decrease the sample complexity for training by 15x, which suggests the great reduction rate the annotation cost with pre-trained KGPT to achieve the desired ‘promising’ performance.

5.6 Zero-Shot Results

We further evaluate KGPT’s generalization capability under the extreme zero-shot setting and display our results for WebNLG in [Table 11](#). As can be seen, all the non-pre-trained baselines and Template-GPT-2 fail under this setting, while KGPT can still manage to generate reasonable outputs and achieve a ROUGE-L score over 30. Given that no input knowledge triples in WebNLG were seen during pre-training, these results reflect KGPT’s strong generalization ability to cope with out-of-domain unseen knowledge inputs.

Model	BLEU	METEOR	ROUGE
All Baselines	0	0	1.2
Template-GPT-2	0.3	0.5	3.4
KGPT-Graph w/ Pre	13.66	19.17	30.22
KGPT-Seq w/ Pre	13.86	20.15	30.23

Table 11: Zero-shot results on WebNLG’s test set.

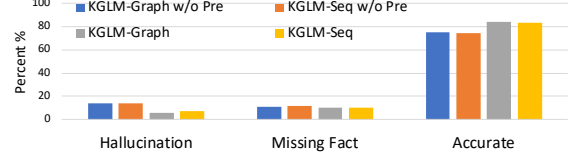


Figure 7: Human evaluation of the factual consistency of different models on WebNLG samples.

5.7 Human Evaluation

We conduct human evaluation to assess the factual accuracy of the generated sentences. Specifically, we sample 100 test samples from WebNLG and observe the model’s factual consistency with given fact triples. We use AMT to distribute each generated sentence to four high-quality workers (95% approval rate, 500+ approved jobs) to choose from the three ratings. The majority voted rating is the final rating. We compare four different systems, i.e., non-pre-trained and pre-trained KGPT. Conditioned on the fact triples, we categorize the generated samples into the following categories: 1) hallucinating non-existing facts, 2) missing given facts without hallucination, 3) accurate description of given facts. We visualize the results in [Figure 7](#), from which we observe that pre-trained KGPT are less prone to the known hallucination issue and generate more accurate text. The human evaluation suggests that pre-training can enhance the model’s understanding over rare entities, thus reducing the over-generation of non-existent facts.

5.8 Conclusion

In this paper, we propose a pre-training recipe to exploit external unlabeled data for data-to-text generation tasks. Our proposed model has achieved significant performance under zero-shot and few-shot settings. Such a framework provides a plausible solution to greatly reduce human annotation costs in future NLG applications.

Acknowledgement

The authors would like to thank the anonymous reviewers for their thoughtful comments. This research is sponsored in part by NSF IIS 1528175, we also want to thank their financial support.

References

- Sungjin Ahn, Heeyoul Choi, Tanel Pärnamaa, and Yoshua Bengio. 2016. A neural knowledge language model. *arXiv preprint arXiv:1608.00318*.
- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Junwei Bao, Duyu Tang, Nan Duan, Zhao Yan, Yuanhua Lv, Ming Zhou, and Tiejun Zhao. 2018. Table-to-text: Describing table region with natural language. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2003. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155.
- Shuang Chen, Jinpeng Wang, Xiaocheng Feng, Feng Jiang, Bing Qin, and Chin-Yew Lin. 2019. Enhancing neural data-to-text generation models with external background knowledge. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3013–3023.
- Wenhu Chen, Jianshu Chen, Yu Su, Zhiyu Chen, and William Yang Wang. 2020a. Logical natural language generation from open-domain tables. *ACL 2020*.
- Zhiyu Chen, Harini Eavani, Wenhu Chen, Yinyin Liu, and William Yang Wang. 2020b. Few-shot nlg with pre-trained language model. *ACL*.
- Émilie Colin and Claire Gardent. 2019. Generating text from anonymised structures. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 112–117.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Ondřej Dušek and Filip Jurcicek. 2016. Sequence-to-sequence generation for spoken dialogue via deep syntax trees and strings. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 45–51.
- 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*.
- Henry Elder, Sebastian Gehrmann, Alexander O’Connor, and Qun Liu. 2018. E2e nlg challenge submission: Towards controllable generation of diverse natural language. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 457–462.
- Hiroaki Hayashi, Zecong Hu, Chenyan Xiong, and Graham Neubig. 2020. Latent relation language models. *Thirty-Fourth AAAI Conference on Artificial Intelligence*.
- Jonathan Herzig, Paweł Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Martin Eisenschlos. 2020. Tapas: Weakly supervised table parsing via pre-training. *ACL*.
- Juraj Juraska, Panagiotis Karagiannis, Kevin Bowden, and Marilyn Walker. 2018. A deep ensemble model with slot alignment for sequence-to-sequence natural language generation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 152–162.
- Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. *ICLR 2015*.
- Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. *ICLR*.
- Karen Kukich. 1983. Design of a knowledge-based report generator. In *Proceedings of the 21st annual meeting on Association for Computational Linguistics*, pages 145–150. Association for Computational Linguistics.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations. *ICLR*.
- Rémi Lebrete, David Grangier, and Michael Auli. 2016. Neural text generation from structured data with application to the biography domain. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1203–1213.

- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.
- Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. 2016. Gated graph sequence neural networks. *ICLR*.
- Percy Liang, Michael I Jordan, and Dan Klein. 2009. Learning semantic correspondences with less supervision. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1*, pages 91–99. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. **ROUGE: A package for automatic evaluation of summaries**. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Tianyu Liu, Fuli Luo, Qiaolin Xia, Shuming Ma, Baobao Chang, and Zhifang Sui. 2019a. Hierarchical encoder with auxiliary supervision for neural table-to-text generation: Learning better representation for tables. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6786–6793.
- Tianyu Liu, Kexiang Wang, Lei Sha, Baobao Chang, and Zhifang Sui. 2018. Table-to-text generation by structure-aware seq2seq learning. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Robert Logan, Nelson F Liu, Matthew E Peters, Matt Gardner, and Sameer Singh. 2019. Barack’s wife hillary: Using knowledge graphs for fact-aware language modeling. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5962–5971.
- Edward Loper and Steven Bird. 2002. Nltk: The natural language toolkit. In *Proceedings of the ACL-02 Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics*, pages 63–70.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1412–1421.
- Shuming Ma, Pengcheng Yang, Tianyu Liu, Peng Li, Jie Zhou, and Xu Sun. 2019. Key fact as pivot: A two-stage model for low resource table-to-text generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2047–2057.
- Diego Marcheggiani and Laura Perez-Beltrachini. 2018. Deep graph convolutional encoders for structured data to text generation. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 1–9.
- Mike Mintz, Steven Bills, Rion Snow, and Dan Jurafsky. 2009. Distant supervision for relation extraction without labeled data. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2-Volume 2*, pages 1003–1011. Association for Computational Linguistics.
- Amit Moryossef, Yoav Goldberg, and Ido Dagan. 2019. Step-by-step: Separating planning from realization in neural data-to-text generation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2267–2277.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Ankur P Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. Totto: A controlled table-to-text generation dataset. *arXiv preprint arXiv:2004.14373*.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. [URL https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language_understanding_paper.pdf](https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language_understanding_paper.pdf).
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Lena Reed, Shereen Oraby, and Marilyn Walker. 2018. Can neural generators for dialogue learn sentence planning and discourse structuring? In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 284–295.

- Ehud Reiter and Robert Dale. 1997. Building applied natural language generation systems. *Natural Language Engineering*, 3(1):57–87.
- Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083.
- Lei Sha, Lili Mou, Tianyu Liu, Pascal Poupart, Sujian Li, Baobao Chang, and Zhifang Sui. 2018. Order-planning neural text generation from structured data. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Anastasia Shimorina and Claire Gardent. 2018. [Handling rare items in data-to-text generation](#). In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 360–370. Association for Computational Linguistics.
- 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 *Advances in neural information processing systems*, pages 5998–6008.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. *ICLR*.
- Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85.
- Tsung-Hsien Wen, Milica Gasic, Nikola Mrkšić, Lina M Rojas Barahona, Pei-Hao Su, David Vandyke, and Steve Young. 2016. Multi-domain neural network language generation for spoken dialogue systems. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 120–129.
- Sam Wiseman, Stuart Shieber, and Alexander Rush. 2017. [Challenges in data-to-document generation](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2253–2263, Copenhagen, Denmark. Association for Computational Linguistics.
- Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2018. Learning neural templates for text generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3174–3187.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5754–5764.
- Rong Ye, Wenxian Shi, Hao Zhou, Zhongyu Wei, and Lei Li. 2020. Variational template machine for data-to-text generation. *ICLR 2020*.

A Learning Curve

Here we observe the learning trend of both non-pre-trained and pre-trained models by evaluating the validation BLEU at each epoch end, here we show our findings in Figure 8. As can be seen from the figure, the pre-trained model converges much faster to the best score. More specifically, it only takes 20 epochs for the model to reach BLEU-4 over 60 while it takes 80-90 epochs for a non-pre-trained model to reach equivalent performance.

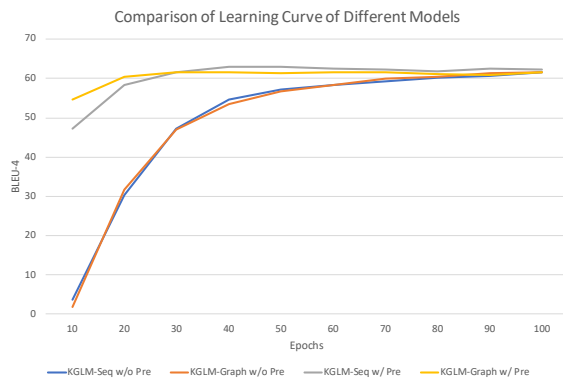


Figure 8: The learning curve of different models during training for the WebNLG dataset.

B Predicate Distribution

Here we demonstrate the most popular predicates in Figure 9. As can be seen, the most popular predicates are ‘instance of’, ‘occupation’, ‘country’, ‘located in’, etc. There are over 1000 predicates in our dataset, which covers the commonly seen categories in different domains like politics, athletics, music, news, etc.

C Case Study

Here we demonstrate some empirical study over the generated samples from our models in Figure 10. As can be seen, KGPT has developed a really strong generation capability to output fluent and coherent sentences. In the first line, the decoded sentence is mostly correct, just the name of ‘municipality’ should be ‘Belgrade’ rather than ‘Zemun’ itself according to <https://www.wikidata.org/wiki/Q189419>. In the second line, the sentence is mostly correct, the error comes from the end date of Annibale. The third sentence is completely correct. The fourth sentence also suffers from a factual error, the relationship should be ‘married’ rather than ‘daughter’.

From these sentences, it’s understandable that the model can achieve reasonable zero-shot performance on the WebNLG dataset given that WebNLG also comes from a similar domain. The case study reveals that our generation model though generates fluent and relevant sentences from the given knowledge triples, the groundedness is still questionable with quite an amount of hallucination issues.

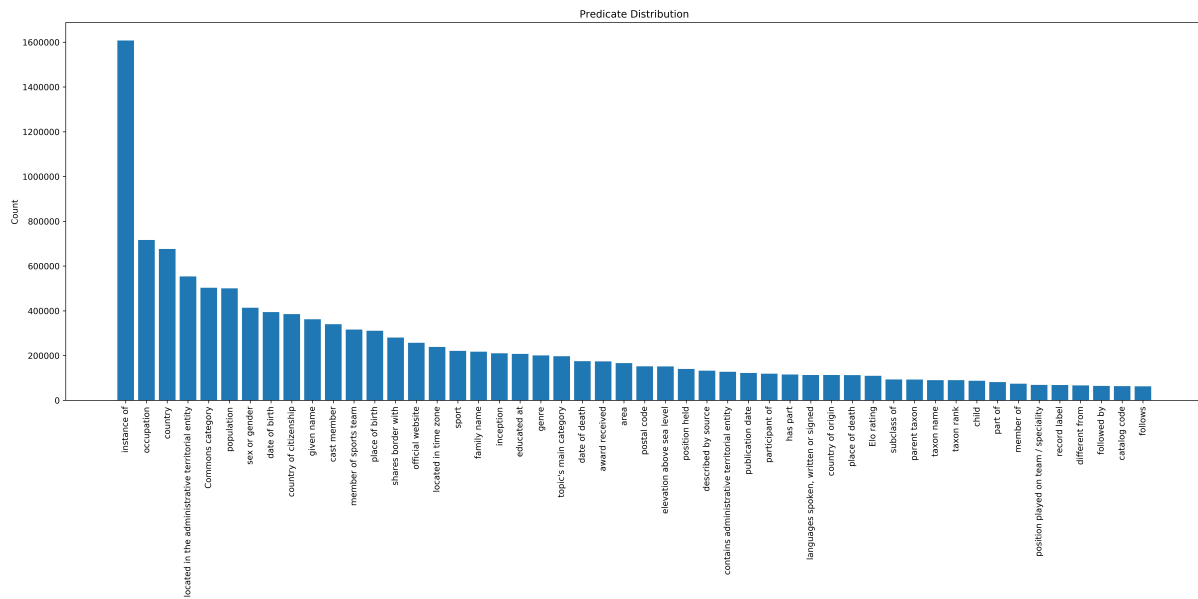


Figure 9: Predicate distribution over the knowledge triples in KGTEXT.

Input	'Q189419', 'Q403'
Decoded	Zemun is a village in the municipality of Zemun, Serbia.
Reference	Zemun Polje railway station is the rail station in Zemun Polje, Zemun, Serbia.
Input	'TITLE:::Q566851', 'Q548320', 'Q190353', 'Q50001', 'Q49757'
Decoded	Annibale Caro (6 June 177 – 11 November 1766) was an Italian poet.
Reference	Fra' Annibale Caro , K.M. , (6 June 150717 November 1566) was an Italian writer and poet .
Input	'TITLE:::Q7151519', 'Q50054', 'Q5925', 'Q99', 'Q286803'
Decoded	Watford was born in Garden Grove, Orange County, California, and graduated from Laguna Beach High School.
Reference	Watford was born in Garden Grove in Orange County , California and graduated from Laguna Beach High School in 1985.
Input	'TITLE:::Q4738713', 'Q6257160', 'Q578478', 'Q23129'
Decoded	She was the daughter of John Scudamore, sheriff of Herefordshire.
Reference	Alys is known to have married Sir John Scudamore a sheriff of Herefordshire .
Input	'Q5372', 'Q3741166'
Decoded	The Iran national basketball team represents Iran in international basketball and is controlled by the Islamic Republic of Iran Federation of Basketball Iran.
Reference	The Iranian national basketball team represents Iran in international basketball competitions , and is controlled by the IR Iran Basketball Federation .
Input	'Q5472010', 'Q16', 'Q7940062'
Decoded	Fort Selkirk volcanic field is a Canadian stratovolcano located on Volcano Mountain in the Yukon Territory.
Reference	The Fort Selkirk Volcanic Field in central Yukon is the northernmost Holocene volcanic field in Canada , including the young active cinder cone , Volcano Mountain .

Figure 10: Randomly generated samples from KGTEXT, where the inputs are the WikiData entities, you can search it online to see it information. For example, the entity 'Q403' and its fact triples can be seen from <https://www.wikidata.org/wiki/Q403>.