

Neural Pipeline for Zero-Shot Data-to-Text Generation

Anonymous ACL submission

Abstract

In data-to-text (D2T) generation, training on in-domain data leads to overfitting to the data representation and repeating training data noise. We examine how to avoid finetuning the pretrained language models (PLMs) on D2T generation datasets while still taking advantage of surface realization capabilities of PLMs. Inspired by pipeline approaches, we propose to generate text by gradually transforming single-item descriptions with a sequence of modules trained on general-domain text-based operations: ordering, aggregation, and paragraph compression. We train PLMs for performing these operations on a synthetic corpus WIKIFLUENT which we build from English Wikipedia. Our experiments on two major triple-to-text datasets—WebNLG and E2E—show that our approach enables D2T generation from RDF triples in zero-shot settings.¹

1 Introduction

The aim of data-to-text (D2T) generation is to produce natural language descriptions of structured data (Gatt and Krahmer, 2018; Reiter and Dale, 1997). Although pipelines of rule-based D2T generation modules are still used in practice (Dale, 2020), end-to-end approaches based on PLMs recently showed superior benchmark performance (Ke et al., 2021; Chen et al., 2020a; Ferreira et al., 2020; Kale and Rastogi, 2020b; Ribeiro et al., 2020), surpassing pipeline systems (Ferreira et al., 2019) in both automatic and human evaluation metrics.

Finetuning PLMs on human-written references is widely accepted as a standard approach for adapting PLMs to the D2T generation objective and achieving good performance on a given benchmark (Agarwal et al., 2021; Ke et al., 2021). However, finetuning the model on the domain-specific data

¹The anonymized version of our code and data is available at <https://anonymous.4open.science/r/zeroshot-d2t-pipeline/>.

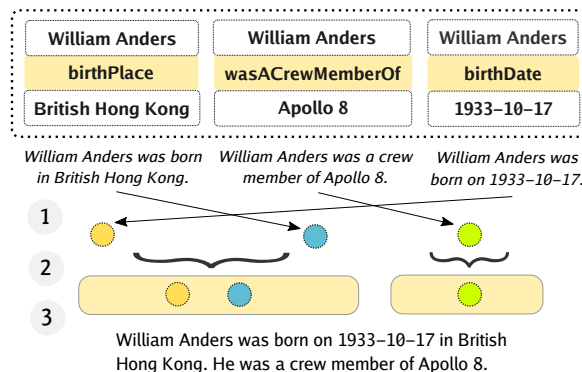


Figure 1: A scheme of our pipeline for zero-shot data-to-text generation from RDF triples: (1) ordering, (2) aggregation, (3) paragraph compression. Individual pipeline modules are trained on a large general-domain text corpus and operate over text in natural language. In-domain knowledge is included only in the simple hand-crafted templates for each predicate.

leads to overfitting on the particular benchmark, decreasing performance on out-of-domain data (Laha et al., 2020). Moreover, collecting a large set of references for a particular domain is costly and time-consuming, as the data are usually collected using crowdsourcing (Dušek et al., 2020). Few-shot approaches (Chen et al., 2020b; Ke et al., 2021; Su et al., 2021a) can operate with only several tens or hundreds of annotated examples, but their robustness is questionable—selecting a representative set of examples which would improve performance is difficult (Chang et al., 2021a), and the limited sample is often noisy, increasing the chance of hallucinations and omissions (Dušek et al., 2019; Harkous et al., 2020; Rebuffel et al., 2021).

In this paper, we provide an alternative to the traditional finetuning paradigm by formulating the D2T generation from RDF triples as a sequence of general-domain operations over text in natural language. We start by transforming individual triples to text using trivial templates, which we subsequently order, aggregate, and compress on the para-

graph level to produce the resulting description of the data. All the pipeline modules operate over natural language and are built upon PLMs trained on our new WIKIFLUENT corpus. WIKIFLUENT contains 934k examples of first paragraphs from the English Wikipedia, each supplied with a synthesized set of simple template-like sentences conveying the same meaning. Our approach allows generating natural language descriptions from triples with a minimum amount of domain-specific rules or knowledge and without using training data from the D2T datasets. We show that our approach can yield large improvements upon simple baselines and match older supervised systems on automatic metrics, while our error analysis suggests that our approach reduces the occurrence of omissions and hallucinations.

Our contributions are the following:

- (1) We propose an alternative D2T generation approach based on general-domain text-to-text operations (ordering, aggregation, and paragraph compression).
- (2) We introduce a synthetic WIKIFLUENT corpus containing 934k sentences based on English Wikipedia, providing training data for the operations in (1).
- (3) We apply our system on two D2T datasets and evaluate its performance both automatically and manually, including the contribution of individual pipeline modules.
- (4) We release our code, data, pretrained models, and system outputs to ease future research.

2 Related Work

D2T Generation with PLMs Large neural language models pretrained on self-supervised tasks (Lewis et al., 2020; Liu et al., 2019; Devlin et al., 2019) have recently gained a lot of traction in D2T generation research (Ferreira et al., 2020). Following Chen et al. (2020b), other works adopted PLMs for few-shot D2T generation (Chang et al., 2021b; Su et al., 2021a). Kale and Rastogi (2020b) and Ribeiro et al. (2020) showed that PLMs using linearized representations of data can outperform graph neural networks on graph-to-text datasets, recently surpassed again by graph-based models (Ke et al., 2021; Chen et al., 2020a). Although the models make use of general-domain pretraining tasks, all of them are eventually finetuned on domain-specific data.

Pipeline-based D2T Generation Until the recent surge of end-to-end approaches (Dušek et al., 2020), using several modules connected in a pipeline was a major approach for D2T generation (Gatt and Krahmer, 2018; Reiter, 2007; Reiter and Dale, 1997). Our approach is inspired by the pipeline approaches, in particular the pipelines utilizing neural-based modules (Ferreira et al., 2019). In contrast with these approaches, our pipeline works with unstructured data in natural language and it operates in zero-shot setting, i.e. without using the training data from the D2T datasets.

Laha et al. (2020) introduce a three-step pipeline for zero-shot D2T generation similar to ours. However, their approach is tailored for specific domains (using rule-based lexicalization conditioned on entity types) and they do not address content planning.

Content Planning in D2T Generation Content planning, i.e. the task of ordering input facts and aggregating them into individual sentences, is one of the steps of the traditional D2T pipeline (Gatt and Krahmer, 2018). As shown by Moryossef et al. (2019a,b) and confirmed by other works (Puduppully et al., 2019; Zhao et al., 2020; Trisedya et al., 2020; Su et al., 2021b), including a content plan improves the quality of outputs also in the neural D2T pipelines. Unlike the aforementioned planners, which use predicates or keys in the D2T datasets for representing the data items, our planner is trained on ordering sentences in natural language.

Sentence Ordering Sentence ordering is the task of organizing a set of natural language sentences to increase the coherence of a text (Barzilay et al., 2001; Lapata, 2003). Several neural methods for this task were proposed, using either interactions between pairs of sentences (Chen et al., 2016; Li and Jurafsky, 2017), global interactions (Gong et al., 2016; Wang and Wan, 2019), or combination of both (Cui et al., 2020). We base our ordering module (§5.2) on the recent work of Calizzano et al. (2021), who use a pointer network (Wang and Wan, 2019; Vinyals et al., 2015) on top of a PLM.

Fact Aggregation The compact nature of the target text description results in aggregating multiple facts in a single sentence. Previous works (Wiseman et al., 2018; Shao et al., 2019; Shen et al., 2020; Xu et al., 2021) capture the segments which correspond to individual parts of the input as latent variables. Unlike these works, we adopt a simpler

scenario using an already ordered sequence of facts, in which we selectively insert delimiters marking sentence boundaries.

Paragraph Compression We introduce *paragraph compression* as a final task in our D2T generation pipeline. Since we already work with natural language in this step, the focus of our task is on sentence fusion, rephrasing, and coreference resolution. Unlike text summarization or simplification (Zhang et al., 2020; Jiang et al., 2020), we aim to convey the complete semantics of the text without omitting any facts. In contrast to sentence fusion (Geva et al., 2019; Barzilay and McKeown, 2005) or sentence compression (Filippova and Altun, 2013), we operate in the context of multiple sentences in a paragraph. The task is the central focus of our WIKIFLUENT corpus (§4).

3 Method

We focus on the task of producing a natural language description Y for a set of n RDF triples $X = \{x_1, \dots, x_n\}$. Each triple $x_i = \{s_i, p_i, o_i\}$ consists of subject s_i , predicate p_i , and object o_i .

Our pipeline proceed as follows. Given a set of triples X on the input, we:

- (1) transform the set of triples to the set of *facts*, i.e. sentences in natural language,
- (2) sort the facts using an *ordering module*,
- (3) insert the sentence delimiters between the sorted facts using an *aggregation module*,
- (4) input the ordered sequence of facts with delimiters into the *paragraph compression* module, which generates the final description Y .

In the following sections, we provide a formal description of the individual steps: transforming individual triples to text (§3.1), ordering (§3.2), aggregation (§3.3), and paragraph compression (§3.4).

3.1 Transforming Triples to Facts

The first step in our pipeline involves transforming each of the input triples $x_i \in X$ into a fact $f_i \in F$ using a transformation $T : X \rightarrow F$. A fact f_i is a single sentence in natural language describing x_i . The transformation serves two purposes: (a) preparing the data for the subsequent text-to-text operations, (b) introducing in-domain knowledge about the semantics of individual predicates.

3.2 Ordering the Facts

We assume that the default order of triples X (and the respective facts F) is random. To maximize the

coherency of the resulting description, we apply an ordering model O to get an ordered sequence of facts: $F_o = \{f_{o_1}, \dots, f_{o_n}\}$, where $o_{1:n}$ is a permutation of indices. Grouping the related facts together (e.g. facts mentioning *birth date* and *birth place*) helps the paragraph compression model to focus only on fusing and rephrasing the neighboring sentences. We describe our ordering model in §5.2.

3.3 Aggregating the Facts

The aggregation model takes a sequence of ordered facts F_o as input and produces a sequence of sentence delimiters $A(F_o) = \{\delta_{o_1}, \delta_{o_2}, \dots, \delta_{o_{n-1}}\}$; $\delta_i \in \{0, 1\}$. The output $\delta_i = 1$ means that the neighboring facts are should be mentioned separately, serving as a hint for the paragraph compression model *not* to fuse the neighboring sentences. Conversely, $\delta_i = 0$ means that the facts should be aggregated and their corresponding sentences should be fused (see §5.3 and §5.4).

3.4 Paragraph Compression

The paragraph compression model (see §5.4) takes as input the ordered sequence of facts with delimiters $F_a = \{f_{o_1}, \delta_{o_1}, f_{o_2}, \dots, \delta_{o_{n-1}}, f_{o_n}\}$ and produces a resulting text Y . The objectives of the model are two-fold: (a) *fusing* related sentences, i.e., sentences i and j in between which $\delta_i = 0$, and (b) *rephrasing* the text to improve its fluency, e.g. fixing minor disfluencies in the templates, replacing noun phrases with referring expressions, etc. The focus is on minor rephrasing since the goal is to preserve the semantics of the original text.

4 WIKIFLUENT Corpus

In this section, we describe the process of building a large-scale synthetic corpus WIKIFLUENT. The corpus provides training data for the neural models which we use in our implementation of the ordering, aggregation, and paragraph compression modules (cf. §5).

In the corpus, we aim to cover a broad range of domains while capturing the sentence style in D2T generation with respect to both the input facts and the target descriptions. In other words, we aim to build a corpus in which (1) the input is a set of simple, template-like sentences, (2) the output is a fluent text in natural language preserving the semantics of the input. As we describe below in detail, we achieve that by applying *split-*

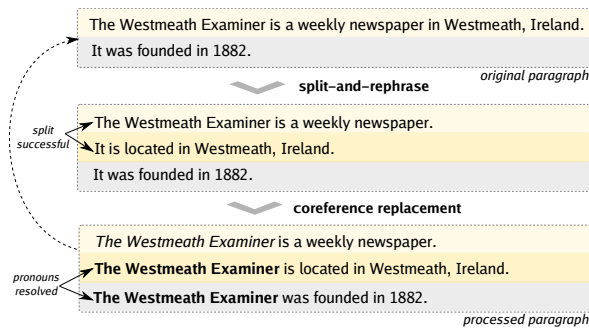


Figure 2: The building process of the WIKIFLUENT corpus. We apply a split-and-rephrase model on each sentence in the paragraph and resolve coreferences in the split sentences.

and-rephrase and coreference resolution models on a set of human-written paragraphs in English Wikipedia. For subsequent training, we consider the processed text as a source and the original text as the target. The process is illustrated in Figure 2; corpus statistics are included in Appendix A.

4.1 Data Source

For building the WIKIFLUENT corpus, we extracted 934k first paragraphs of articles from a Wikipedia dump² using WikiExtractor (Attardi, 2015). The paragraphs contain mostly concise, fact-based descriptions from a wide range of domains. We selected paragraphs with length between 30-430 characters; filtering out lists, disambiguations, repeated and malformed paragraphs. To balance the length of inputs, we selected 250k examples each from 4 equidistant length ranges (30-130 characters, etc.).

Wikipedia is commonly used for large-scale pre-training for D2T generation models (Jin et al., 2020; Chen et al., 2020a). Although the Wikipedia texts are not completely bias-free, they provide a more balanced sample of natural language use than typical D2T generation datasets.

4.2 Split-and-Rephrase

For generating the target set of sentences, we divide each paragraph into sentences using NLTK (Bird, 2006) and apply a *split-and-rephrase* model on each sentence. Split-and-rephrase is a task of splitting a complex sentence into a meaning preserving sequence of shorter sentences (Narayan et al., 2017). We train our model on the large-scale WikiSplit corpus by Botha et al. (2018), containing human-made sentence splits from Wikipedia

²enwiki-20210401-pages-articles-multistream

edit history. Following the setup in the rest of our experiments, we train the encoder-decoder PLM BART-base (Lewis et al., 2020) on the WikiSplit dataset in a sequence-to-sequence setting. We apply the trained split-and-rephrase model on each sentence, uniformly randomly choosing between 0-2 recursive calls to ensure that the splits are not deterministic. If the sentence cannot be meaningfully split, the model tends to duplicate the sentence on the output; in that case, we use only the original sentence and do not proceed with the splitting.

4.3 Coreference Replacement

Next, we concatenate the split sentences and apply a coreference resolution model (Gardner et al., 2018) in order to replace referring expressions with their antecedents (e.g., pronouns with noun phrases). This step is motivated by following the style of the facts in which the entities are fully verbalized (as each fact describes a single triple). Since this procedure replaces the referring expressions only in the source texts and keeps them in the target texts, it implicitly trains the paragraph compression module to generate the referring expressions for the entities in the final description.

4.4 Filtering

To assert that the generated sentences convey the same semantics as the original paragraph, we use a pretrained RoBERTa model³ (Liu et al., 2019) trained on the MultiNLI dataset (Williams et al., 2018) for checking the semantic accuracy of the generated text. Following Dušek and Kasner (2020), we test if the original paragraph entails each of the synthesized sentences (checking for omissions), and if the set of concatenated synthesized sentences entails the original paragraph (checking for hallucinations). In a filtered version of the WIKIFLUENT corpus, we include only the examples without omissions or hallucinations (as computed by the model), reducing it to 714k examples (approximately 3/4 of the original size).

5 Implementation

In this section, we describe how we implement our pipeline modules (§3) using simple template transformations (§5.1) and neural models trained on the WIKIFLUENT dataset (§5.2-5.4).⁴

³<https://huggingface.co/roberta-large-mnli>

⁴Our training setup details are included in Appendix C.

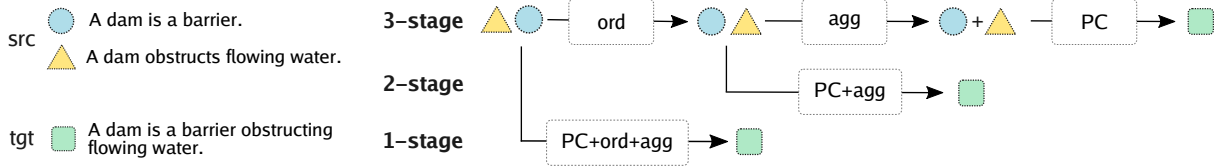


Figure 3: An example illustrating how the individual modules are trained and subsequently applied as the parts of the pipeline. See §5.2 for description of the ordering model (ORD), §5.3 for the aggregation model (AGG), and §5.4 for the versions of the paragraph compression model (PC, PC+AGG, PC+ORD+AGG).

dataset	predicate	template
WebNLG	instrument	<s> plays <o>.
	countryOrigin	<s> comes from <o>.
	width	<s> is <o> wide.
E2E	eatType	<s> is a <o>.
	food	<s> serves <o> food.
	area	<s> is in the <o>.

Table 1: Examples of templates for predicates in the WebNLG and E2E datasets with placeholders for the subject (<s>) and the object (<o>).

5.1 Templates

We transform triples to facts (§3.1) using a single-triple template t_i for each predicate. For example, if $p_i = \text{instrument}$; then $T(p_i) = "s_i \text{ plays } o_i"$ (cf. Table 1). We opt for this approach instead of automatic template generation used in Laha et al. (2020) for its simplicity and better applicability of our approach to arbitrary datasets. Simple hand-crafted templates have been also used in other works as an efficient way of introducing domain knowledge (Kale and Rastogi, 2020a; Kasner and Dušek, 2020).

Although this approach is sufficient in our case, we encourage future research on automatic template generation for making this step scalable and domain-independent. From another point of view, a more complex rule-based template generation engine (Heidari et al., 2021; Mehta et al., 2021) could ensure higher fluency and better controllability in case of practical deployment.

5.2 Ordering Model

For our ordering model (§3.2), we use the *Simple Pointer* model from Calizzano et al. (2021). The model is based on a pretrained BART-base extended with a pointer network from Wang and Wan (2019). We provide a short description of the model here; for details see Calizzano et al. (2021).

In the encoding phase, facts F are concatenated and tokenized. Each fact is surrounded by special tokens denoting the beginning (<s>) and the

end (</s>) of the fact. The sequence is processed by the BART encoder, generating a sequence of encoder states E for each end token </s> representing the preceding fact.

The decoding proceeds autoregressively. To bootstrap the decoding process, the pair of tokens <s></s> is fed into the decoder, producing the decoder state d_1 . The pointer network (attending to d_1 and E), selects the first ordered fact f_{o_1} , which is fed into the decoder in the next step. The process is repeated until the all the facts are decoded in a particular order.

The pointer network computes the probability of a fact to be on the j -th position, using the encoder output E and the decoder output d_j . The network is based on the scaled dot product attention, where d_j is the query and encoder outputs E_i are the keys:

$$Q = d_j W_Q$$

$$K = E W_K$$

$$P_j = \text{softmax} \left(\frac{QK^T}{\sqrt{b}} \right).$$

Here W_Q and $W_K \in \mathbb{R}^{b \times b}$, b is the dimension of BART hidden states, and $P_j \in \mathbb{R}^{n+1}$ is the probability distribution for the j -th position (i.e., P_{ji} is the probability that fact f_i is on the j -th position).

We train the model using the split sentences in the WIKIFLUENT corpus, randomly shuffling the order of the sentences and training the model to restore their original order.

5.3 Aggregation Model

We base our aggregation model (§3.3) on RoBERTa-large (Liu et al., 2019) with a token classification head.⁵ Similarly to the ordering model (§5.2), we input the sequence of facts F_o into the model, separating each pair of facts f_{o_i} with a special token </s> (used by the model as a separator).

⁵https://huggingface.co/transformers/model_doc/roberta.html#robertafortokenclassification

Subsequently, the token classification layer classifies each separator $\langle /s \rangle_i$ position into two classes $\{0, 1\}$ corresponding to the delimiter δ_i . We ignore the outputs for the non-separator tokens while computing the cross-entropy loss.

We create the training examples using the split sentences in the WIKIFLUENT corpus, in which we set $\delta_i = 0$ for the sentences $i, i + 1$ which were originally aggregated (i.e., are the result of splitting a single sentence) and $\delta_i = 1$ otherwise.

5.4 Paragraph Compression Model

We adopt BART-base for our paragraph compression model. We train the model in a sequence-to-sequence setting on the WIKIFLUENT corpus, concatenating the split sentences on the input. We add delimiters between sentences i and $i + 1$ where $\delta_i = 1$ using a special token $\langle \text{sep} \rangle$, which we add to the model vocabulary. As shown in Keskar et al. (2019), including control codes for training the model can steer the model towards producing certain outputs. We evaluate our model’s behavior with respect to ordering and aggregation in §7.3.

6 Experiments

We train our pipeline modules on the WIKIFLUENT corpus as described in §5. Next, we use the modules *without fine-tuning* for generating descriptions for RDF triples on two English D2T datasets, WebNLG and E2E. The datasets differ in domain, size, textual style, and number of predicates (see Appendix A for details):

- The **WebNLG** dataset (Gardent et al., 2017; Ferreira et al., 2020) contains RDF triples from DBpedia (Auer et al., 2007) and their crowdsourced descriptions. We use version 1.4 of the dataset for comparability to prior work. We hand-crafted templates for all 354 predicates, including unseen predicates in the test set.⁶
- The **E2E** dataset (Novikova et al., 2017; Dušek et al., 2020) contains restaurant recommendations in the form of attribute-value pairs. We use the cleaned version of the dataset (Dušek et al., 2019). Following previous work, we transformed the attribute-value pairs into RDF triples (using the restaurant name as a subject) and then applied the same setup as

for WebNLG. We created a template for each of the 8 attributes manually.

In order to evaluate individual components of our pipeline, we train three versions of the PC model (see §5.4). The models share the same architecture and targets, but differ in their inputs:

- **PC** – the model takes as an input ordered facts with delimiters (as described in §3.4),
- **PC+AGG** – the model takes as an input ordered facts *without* delimiters (i.e., the aggregation is left implicitly to the model),
- **PC+ORD+AGG** – the model takes as an input facts in *random* order and *without* delimiters (i.e., both ordering and aggregation are left implicitly to the model).

Subsequently, we test three versions of the pipeline in our **ablation study** (see Figure 3):

- **3-STAGE** – a full version of the pipeline consisting of the ordering model, the aggregation model and the PC model (following the full pipeline from §3),
- **2-STAGE** – a pipeline consisting of the ordering model and the PC+AGG model,
- **1-STAGE** – a single stage consisting of the PC+ORD+AGG model.

We evaluate all versions of the pipeline with PC models trained on the *full* and *filtered* versions of the WIKIFLUENT dataset (see §4).

7 Evaluation and Discussion

Our main aim is the evaluation of our pipeline on the downstream task of D2T generation.⁷ We evaluate outputs from the $\{1, 2, 3\}$ -STAGE variants of our pipeline using automatic metrics (§7.1) and we perform a detailed manual error analysis of the model outputs (§7.2). Furthermore, we specifically evaluate the performance of the content planning modules and the ability of the PC module to follow the content plan (§7.3).

7.1 Automatic Metrics

Following prior work, we use BLEU (Papineni et al., 2002) and METEOR (Banerjee and Lavie, 2005) to evaluate the outputs against the human references.⁸ We also evaluate the number of omission and hallucination errors (i.e., facts missing or added, respectively) using a metric from Dušek

⁷We also provide the results of our models on the test set of WIKIFLUENT in Appendix D.

⁸We use the implementation from <https://github.com/tuetschek/e2e-metrics>.

⁶See Appendix B for details on template creation.

Input	(Allen Forrest; background; solo singer), (Allen Forrest; genre; Pop music), (Allen Forrest; birthPlace; Dothan, Alabama)
Templ.	Allen Forrest is a solo singer. Allen Forrest performs Pop music. Allen Forrest was born in Dothan, Alabama.
Model	Allen Forrest is a solo singer who performs Pop music. He was born in Dothan, Alabama.
Human	Born in Dothan, Alabama, Allen Forrest has a background as a solo singer and was a pop artist.
Input	name[Wildwood], eatType[restaurant], food[French], area[riverside], near[Raja Indian Cuisine]
Templ.	Wildwood is a restaurant. Wildwood serves French food. Wildwood is in the riverside. Wildwood is near Raja Indian Cuisine.
Model	Wildwood is a restaurant serving French food. It is in the riverside near Raja Indian Cuisine.
Human	A amazing French restaurant is called the Wildwood. The restaurant is near the Raja Indian Cuisine in riverside. They love kids.

Table 2: Example outputs of our model (3-STAGE, filtered). See Appendix F for more examples.

		B	M	O	H
	UPF-FORGe*	38.65	39.00	0.075	0.101
	MELBOURNE*	45.13	37.00	0.237	0.202
	Ke et al. (2021) [†] *	66.14	47.25	-	-
	Laha et al. (2020) [†]	24.80	34.90	-	-
	COPY	37.18	38.77	0.000	0.000
<i>full</i>	3-STAGE	42.92	39.07	0.051	0.148
	2-STAGE	42.90	39.28	0.043	0.125
	1-STAGE	39.08	38.94	0.071	0.204
<i>filtered</i>	3-STAGE	43.19	39.13	0.152	0.073
	2-STAGE	43.49	39.32	0.146	0.096
	1-STAGE	42.99	38.81	0.202	0.093

Table 3: Automatic metrics on WebNLG. B = BLEU, M = METEOR, O = omissions / # facts, H = hallucinations / # examples. The systems marked with asterisk (*) are trained on the WebNLG dataset. Results for the systems marked with † are copied from the respective papers.

and Kasner (2020) based on a RoBERTa model (Liu et al., 2019) pretrained on natural language inference (NLI).⁹

We include a diverse set of baselines for comparison, including state-of-the-art (SotA) systems for both datasets. For WebNLG (see Table 3), we show the results of UPF-FORGe and MELBOURNE systems from the first run of WebNLG Challenge (Gardent et al., 2017) which are comparable in terms of automatic metrics and semantic errors, and the results of Ke et al. (2021), which is a SotA system with a structure-aware encoder and task-specific pretraining. Laha et al. (2020) is (to our knowledge) the only other zero-shot D2T generation system applied on WebNLG. For E2E, TGEN (Dušek and Jurčiček, 2015) is the baseline system for the E2E Challenge (Dušek et al., 2020) and Harkous et al. (2020) is a state-of-the art supervised system applied on the cleaned E2E (see Table 4).

⁹We additionally evaluated the outputs on the E2E dataset using the provided pattern-based slot error script. See Appendix E for the details.

		B	M	O	H
	TGEN*	40.73	37.76	0.016	0.083
	Harkous et al. (2020)*	43.60	39.00	-	-
	COPY	24.19	34.89	0.000	0.000
<i>full</i>	3-STAGE	36.04	36.95	0.001	0.001
	2-STAGE	35.84	36.91	0.001	0.001
	1-STAGE	30.81	36.01	0.009	0.122
<i>filtered</i>	3-STAGE	35.88	36.95	0.001	0.001
	2-STAGE	36.01	36.99	0.001	0.001
	1-STAGE	34.08	36.32	0.012	0.050

Table 4: Automatic metrics on E2E. B = BLEU, M = METEOR, O = omissions / # facts, H = hallucinations / # examples. The systems marked with asterisk (*) are trained on the E2E dataset. The results for Harkous et al. (2020) are copied from the paper.

For both datasets, COPY is the baseline of copying the templates verbatim.

The automatic evaluation suggests that although our system shows considerable improvements over the COPY baseline¹⁰ (e.g., ~12 BLEU points for E2E), it cannot be yet compared to the SotA supervised systems – for this reason, we focus on *error analysis* in our manual evaluation (§7.2). The 2-STAGE system is generally on par with the 3-STAGE system (or better), which indicates that implicit aggregation using the PC-AGG model may be sufficient. However, an advantage of having a separate aggregation module is the possibility to control the aggregation step explicitly. The filtered version of the dataset generally brings better results, although it brings also an increase in the number of omissions.

7.2 Manual Error Analysis

Since performance metrics do not provide insights into specific weaknesses of the system (van Mil-

¹⁰Also note that BLEU score of our COPY baseline is substantially better than the zero-shot system of Laha et al. (2020).

		WebNLG					E2E				
		H	I	O	R	G	H	I	O	R	G
<i>full</i>	3-STAGE	3	39	2	2	16	0	1	0	0	17
	2-STAGE	8	36	1	5	16	1	1	0	1	23
	1-STAGE	28	27	6	10	20	17	0	1	79	45
<i>filtered</i>	3-STAGE	2	37	2	1	15	0	0	0	0	17
	2-STAGE	5	32	1	2	14	0	0	0	0	11
	1-STAGE	8	40	6	6	16	11	2	1	41	22

Table 5: Number of manually annotated errors on 100 examples: H = hallucinations, I = incorrect fact merging, O = omissions, R = redundancies, G = grammar errors or disfluencies.

tenburg et al., 2021), we manually examined 100 outputs of the models regarding factual errors (hallucinations, omissions, incorrect fact merging, redundancies) and grammatical errors. The results are listed in Table 5.

The 1-STAGE model (which has to order the facts implicitly) tends to repeat the facts in the text (especially in E2E) and produces frequent hallucinations. These problems are largely eliminated with 2-STAGE and 3-STAGE models, which produce almost no hallucinations or omissions. However, the outputs on WebNLG for all systems suffer from semantic errors resulting from merging of unrelated facts. This mostly happens with unrelated predicates connected to the same subject/object (e.g. “X was born in Y”, “X worked as Z” expressed as “X worked as Z in Y”; see Appendix F for examples). As we discuss in §8, more research is needed for ensuring the output text consistency.

On the E2E data, which has a simpler triple structure (all predicates share the same subject), the outputs are generally consistent and the 2-STAGE and 3-STAGE models exhibit almost no semantic errors. The grammar errors and disfluencies stem mainly from over-eager paragraph compression or from artifacts in our templates and are relatively minor (e.g., missing “is” in “serves French food and family-friendly”).

7.3 Content Planning

Following Su et al. (2021b) and Zhao et al. (2020), we report the accuracy (Acc) and BLEU-2 score (B-2) of our **ordering model** on WebNLG against the human-generated plans from Ferreira et al. (2018). The results are listed in Table 6. RANDOM is the baseline of generating a random order. The results show that although our approach lacks behind state-of-the-art supervised approaches, it can outperform both the random baseline and the Transformer-

	B-2	Acc
Transformer (Ferreira et al., 2019) [†]	52.20	0.35
Step-by-step (Moryossef et al., 2019b) [†]	70.80	0.47
PLANENC (Zhao et al., 2020) [†]	80.10	0.62
Plan-then-generate (Su et al., 2021b) [†]	84.97	0.72
RANDOM	47.00	0.29
BART+ptr (Calizzano et al., 2021)	59.10	0.48

Table 6: Evaluation of our zero-shot ordering model based on Calizzano et al. (2021). The results marked with † are copied from the respective papers.

based approach from Ferreira et al. (2019) while not using any training examples from WebNLG.

We also evaluate the accuracy of our **aggregation model**, using triples ordered according to the plans from Ferreira et al. (2018) as input. The accuracy is 0.33 per example and 0.62 per sentence boundary (random baseline is 0.23 and 0.50, respectively). The results show that although our approach is better than the random baseline, further investigation regarding plausible fact aggregation schemes is needed.

Finally, we manually evaluate how the **PC model** follows the content plan using 100 randomly chosen examples with more than 1 triple on WebNLG and E2E. We find that the model follows the content plan in 95% and 100% of cases, respectively. The incorrect cases include a fact not properly mentioned and an extra boundary between the sentences without a separator. We can thus conclude that the pretraining task successfully teaches the PC model to follow a given content plan.

8 Conclusion and Future Work

We presented an approach using PLMs and a general domain corpus for D2T generation without using rules or domain-specific training data. Possible directions for extending our research include automatic generation of templates (cf. §5.1), using our approach as a task-specific pretraining for more efficient fine-tuning of D2T models (e.g., with a small amount of clean data), or extending the approach for more complex data inputs.

More research is also needed regarding the main shortcoming of our approach, i.e., the semantic errors stemming from merging of facts in improper ways. We suggest that explicitly controlling the semantics of sentence fusion (Ben-David et al., 2020) could help to mitigate this issue, while still keeping the advantages of a zero-shot approach.

9 Limitations and Broader Impact

We study zero-shot D2T generation with the focus on generating descriptions for RDF triples. Although the task of D2T generation has numerous applications, using neural models for D2T generation (especially in the zero-shot context) is still limited to experimental settings (Dale, 2020). Similarly to other recent approaches for D2T generation, our approach relies on PLMs, which are known to reflect the biases in their pretraining corpus (Bender et al., 2021; Rogers, 2021). Our system may therefore rely on spurious correlations for verbalizing e.g. gender or occupation of the entities. Since we cannot guarantee the factual correctness of the outputs of our system, the outputs should be used with caution.

On the flip side, our approach helps to reduce the number of omissions and hallucinations stemming from noise in human-written references. Our work thus contributes to the general aim of D2T generation in conveying the data semantics accurately and without relying on implicit world knowledge.

References

Oshin Agarwal, Heming Ge, Siamak Shakeri, and Rami Al-Rfou. 2021. Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3554–3565.

Giuseppe Attardi. 2015. Wikiextractor. <https://github.com/attardi/wikiextractor>.

Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. Dbpedia: A nucleus for a web of open data. In *The semantic web*, pages 722–735. Springer.

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, Ann Arbor, Michigan. Association for Computational Linguistics.

Regina Barzilay, Noemie Elhadad, and Kathleen McKeown. 2001. Sentence ordering in multidocument summarization. In *Proceedings of the first international conference on Human language technology research*.

Regina Barzilay and Kathleen R McKeown. 2005. Sentence fusion for multidocument news summarization. *Computational Linguistics*, 31(3):297–328.

Eyal Ben-David, Orgad Keller, Eric Malmi, Idan Szpektor, and Roi Reichart. 2020. Semantically driven sentence fusion: Modeling and evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 1491–1505.

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623.

Steven Bird. 2006. Nltk: the natural language toolkit. In *Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions*, pages 69–72.

Jan A. Botha, Manaal Faruqui, John Alex, Jason Baldridge, and Dipanjan Das. 2018. **Learning to split and rephrase from Wikipedia edit history**. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 732–737, Brussels, Belgium. Association for Computational Linguistics.

Rémi Calizzano, Malte Ostendorff, and Georg Rehm. 2021. Ordering sentences and paragraphs with pre-trained encoder-decoder transformers and pointer ensembles. In *Proceedings of the 21st ACM Symposium on Document Engineering*, pages 1–9.

Ernie Chang, Xiaoyu Shen, Hui-Syuan Yeh, and Vera Demberg. 2021a. On training instance selection for few-shot neural text generation. *arXiv preprint arXiv:2107.03176*.

Ernie Chang, Xiaoyu Shen, Dawei Zhu, Vera Demberg, and Hui Su. 2021b. Neural data-to-text generation with lm-based text augmentation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 758–768.

Wenhu Chen, Yu Su, Xifeng Yan, and William Yang Wang. 2020a. **KGPT: Knowledge-grounded pre-training for data-to-text generation**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8635–8648, Online. Association for Computational Linguistics.

Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 2016. Neural sentence ordering. *arXiv preprint arXiv:1607.06952*.

Zhiyu Chen, Harini Eavani, Wenhu Chen, Yinyin Liu, and William Yang Wang. 2020b. **Few-shot NLG with pre-trained language model**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 183–190, Online.

Baiyun Cui, Yingming Li, and Zhongfei Zhang. 2020. Bert-enhanced relational sentence ordering network. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6310–6320.

715	Robert Dale. 2020. Natural language generation: The commercial state of the art in 2020. <i>Natural Language Engineering</i> , 26(4):481–487.	769
716		770
717		771
718	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>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)</i> , pages 4171–4186.	772
719		773
720		774
721		775
722		776
723		777
724		778
725		779
726	Ondřej Dušek, David M Howcroft, and Verena Rieser. 2019. Semantic noise matters for neural natural language generation. In <i>Proceedings of the 12th International Conference on Natural Language Generation</i> , pages 421–426.	780
727		781
728		782
729		783
730		784
731		785
732	Ondřej Dušek and Filip Jurčiček. 2015. Training a natural language generator from unaligned data. In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 451–461.	786
733		787
734		788
735		789
736		790
737	Ondřej Dušek and Zdeněk Kasner. 2020. Evaluating semantic accuracy of data-to-text generation with natural language inference. In <i>Proceedings of the 13th International Conference on Natural Language Generation</i> , pages 131–137.	791
738		792
739		793
740		794
741		795
742		796
743	Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2020. Evaluating the state-of-the-art of end-to-end natural language generation: The e2e nlg challenge. <i>Computer Speech & Language</i> , 59:123–156.	797
744		798
745		799
746		800
747	et al. Falcon, WA. 2019. Pytorch lightning. <i>GitHub</i> . Note: https://github.com/PyTorchLightning/pytorch-lightning , 3.	801
748		802
749		803
750	Thiago Ferreira, Claire Gardent, Nikolai Ilinykh, Chris van der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina. 2020. The 2020 bilingual, bidirectional webnlg+ shared task overview and evaluation results (webnlg+ 2020). In <i>Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)</i> .	804
751		805
752		806
753		807
754		808
755		809
756	Thiago Castro Ferreira, Diego Moussallem, Emiel Kraemer, and Sander Wubben. 2018. Enriching the webnlg corpus. In <i>Proceedings of the 11th International Conference on Natural Language Generation</i> , pages 171–176.	810
757		811
758		812
759		813
760		814
761	Thiago Castro Ferreira, Chris van der Lee, Emiel van Miltenburg, and Emiel Kraemer. 2019. Neural data-to-text generation: A comparison between pipeline and end-to-end architectures. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 552–562.	815
762		816
763		817
764		818
765		819
766		820
767		821
768		822
		823
	Katja Filippova and Yasemin Altun. 2013. Overcoming the lack of parallel data in sentence compression. In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1481–1491.	
	Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. The WebNLG challenge: Generating text from RDF data . In <i>Proceedings of the 10th International Conference on Natural Language Generation</i> , pages 124–133, Santiago de Compostela, Spain.	
	Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. Allennlp: A deep semantic natural language processing platform. In <i>Proceedings of Workshop for NLP Open Source Software (NLP-OSS)</i> , pages 1–6.	
	Albert Gatt and Emiel Kraemer. 2018. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. <i>Journal of Artificial Intelligence Research</i> , 61:65–170.	
	Mor Geva, Eric Malmi, Idan Szpektor, and Jonathan Berant. 2019. Discofuse: A large-scale dataset for discourse-based sentence fusion. In <i>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)</i> , pages 3443–3455.	
	Jingjing Gong, Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 2016. End-to-end neural sentence ordering using pointer network. <i>arXiv preprint arXiv:1611.04953</i> .	
	Hamza Harkous, Isabel Groves, and Amir Saffari. 2020. Have your text and use it too! end-to-end neural data-to-text generation with semantic fidelity. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 2410–2424.	
	Peyman Heidari, Arash Einolghozati, Shashank Jain, Soumya Batra, Lee Callender, Ankit Arun, Shawn Mei, Sonal Gupta, Pinar Donmez, Vikas Bhardwaj, et al. 2021. Getting to production with few-shot natural language generation models. In <i>Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 66–76.	
	Chao Jiang, Mounica Maddela, Wuwei Lan, Yang Zhong, and Wei Xu. 2020. Neural crf model for sentence alignment in text simplification. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7943–7960.	
	Zhijing Jin, Qipeng Guo, Xipeng Qiu, and Zheng Zhang. 2020. Genwiki: A dataset of 1.3 million content-sharing text and graphs for unsupervised graph-to-text generation. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 2398–2409.	

824	Mihir Kale and Abhinav Rastogi. 2020a. Template guided text generation for task-oriented dialogue . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6505–6520, Online. Association for Computational Linguistics.	878
825		879
826		880
827		881
828		882
829		
830	Mihir Kale and Abhinav Rastogi. 2020b. Text-to-text pre-training for data-to-text tasks . In <i>Proceedings of the 13th International Conference on Natural Language Generation</i> , pages 97–102, Dublin, Ireland. Association for Computational Linguistics.	883
831		884
832		885
833		886
834		887
835	Zdeněk Kasner and Ondřej Dušek. 2020. Data-to-text generation with iterative text editing. In <i>Proceedings of the 13th International Conference on Natural Language Generation</i> , pages 60–67.	888
836		889
837		890
838		891
839	Pei Ke, Haozhe Ji, Yu Ran, Xin Cui, Liwei Wang, Linfeng Song, Xiaoyan Zhu, and Minlie Huang. 2021. Jointgt: Graph-text joint representation learning for text generation from knowledge graphs. <i>arXiv preprint arXiv:2106.10502</i> .	892
840		893
841		894
842		895
843		896
844	Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. <i>arXiv preprint arXiv:1909.05858</i> .	897
845		898
846		899
847		900
848	Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization . In <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> .	901
849		902
850		903
851		904
852		905
853	Anirban Laha, Parag Jain, Abhijit Mishra, and Karthik Sankaranarayanan. 2020. Scalable micro-planned generation of discourse from structured data. <i>Computational Linguistics</i> , 45(4):737–763.	906
854		907
855		908
856		909
857	Mirella Lapata. 2003. Probabilistic text structuring: Experiments with sentence ordering. In <i>Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics</i> , pages 545–552.	910
858		911
859		912
860		913
861	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880.	914
862		915
863		916
864		917
865		918
866		919
867		920
868		921
869	Jiwei Li and Dan Jurafsky. 2017. Neural net models of open-domain discourse coherence. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 198–209.	922
870		923
871		924
872		925
873	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	926
874		927
875		928
876		929
877		930
	Sanket Vaibhav Mehta, Jinfeng Rao, Yi Tay, Mihir Kale, Ankur Parikh, Hongtao Zhong, and Emma Strubell. 2021. Improving compositional generalization with self-training for data-to-text generation. <i>arXiv preprint arXiv:2110.08467</i> .	931
		932
		933
	Amit Moryossef, Yoav Goldberg, and Ido Dagan. 2019a. Improving quality and efficiency in plan-based neural data-to-text generation. In <i>Proceedings of the 12th International Conference on Natural Language Generation</i> , pages 377–382.	
	Amit Moryossef, Yoav Goldberg, and Ido Dagan. 2019b. Step-by-step: Separating planning from realization in neural data-to-text generation. <i>arXiv preprint arXiv:1904.03396</i> .	
	Shashi Narayan, Claire Gardent, Shay B. Cohen, and Anastasia Shimorina. 2017. Split and rephrase . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 606–616, Copenhagen, Denmark. Association for Computational Linguistics.	
	Jekaterina Novikova, Ondrej Dušek, and Verena Rieser. 2017. The E2E Dataset: New Challenges for End-to-End Generation . In <i>Proceedings of the 18th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 201–206, Saarbrücken, Germany.	
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	
	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. <i>Advances in Neural Information Processing Systems</i> , 32:8026–8037.	
	Ratish Puduppully, Li Dong, and Mirella Lapata. 2019. Data-to-text generation with content selection and planning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 33, pages 6908–6915.	
	Clément Rebuffel, Marco Roberti, Laure Soulier, Geoffrey Scoutheeten, Rossella Cancelliere, and Patrick Gallinari. 2021. Controlling hallucinations at word level in data-to-text generation. <i>arXiv preprint arXiv:2102.02810</i> .	
	Ehud Reiter. 2007. An architecture for data-to-text systems. In <i>proceedings of the eleventh European workshop on natural language generation (ENLG 07)</i> , pages 97–104.	
	Ehud Reiter and Robert Dale. 1997. Building applied natural language generation systems. <i>Natural Language Engineering</i> , 3(1):57–87.	

934	Leonardo FR Ribeiro, Martin Schmitt, Hinrich Schütze, and Iryna Gurevych. 2020. Investigating pretrained language models for graph-to-text generation. <i>arXiv preprint arXiv:2007.08426</i> .	989
935		990
936		991
937		992
938	Anna Rogers. 2021. Changing the world by changing the data . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2182–2194, Online. Association for Computational Linguistics.	993
939		994
940		995
941		996
942		997
943		998
944		999
945	Zhihong Shao, Minlie Huang, Jiangtao Wen, Wenfei Xu, and Xiaoyan Zhu. 2019. Long and diverse text generation with planning-based hierarchical variational model . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3257–3268, Hong Kong, China. Association for Computational Linguistics.	1000
946		1001
947		1002
948		1003
949		1004
950		1005
951		1006
952		1007
953		1008
954	Xiaoyu Shen, Ernie Chang, Hui Su, Cheng Niu, and Dietrich Klakow. 2020. Neural data-to-text generation via jointly learning the segmentation and correspondence . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7155–7165, Online. Association for Computational Linguistics.	1009
955		1010
956		1011
957		1012
958		1013
959		1014
960		1015
961	Yixuan Su, Zaiqiao Meng, Simon Baker, and Nigel Collier. 2021a. Few-shot table-to-text generation with prototype memory. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 910–917.	1016
962		1017
963		1018
964		1019
965		1020
966	Yixuan Su, David Vandyke, Sihui Wang, Yimai Fang, and Nigel Collier. 2021b. Plan-then-generate: Controlled data-to-text generation via planning. <i>arXiv preprint arXiv:2108.13740</i> .	1021
967		1022
968		1023
969		1024
970	Bayu Trisedya, Jianzhong Qi, and Rui Zhang. 2020. Sentence generation for entity description with content-plan attention. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 9057–9064.	1025
971		1026
972		1027
973		1028
974		1029
975	Emiel van Miltenburg, Miruna Clinciu, Ondřej Dušek, Dimitra Gkatzia, Stephanie Inglis, Leo Leppänen, Saad Mahamood, Emma Manning, Stephanie Schoch, Craig Thomson, et al. 2021. Underreporting of errors in nlg output, and what to do about it. In <i>Proceedings of the 14th International Conference on Natural Language Generation</i> , pages 140–153.	1030
976		1031
977		1032
978		1033
979		1034
980		1035
981		1036
982	Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. <i>Advances in Neural Information Processing Systems</i> , 28:2692–2700.	1037
983		1038
984		1039
985	Tianming Wang and Xiaojun Wan. 2019. Hierarchical attention networks for sentence ordering. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 7184–7191.	1040
986		1041
987		1042
988		1043
	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122.	1044
		1045
	Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2018. Learning neural templates for text generation. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3174–3187.	1046
		1047
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. <i>arXiv preprint arXiv:1910.03771</i> .	1048
		1049
	Xinnuo Xu, Ondřej Dušek, Verena Rieser, and Ioannis Konstas. 2021. Agggen: Ordering and aggregating while generating. <i>arXiv preprint arXiv:2106.05580</i> .	1050
		1051
	Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In <i>International Conference on Machine Learning</i> , pages 11328–11339. PMLR.	1052
		1053
	Chao Zhao, Marilyn Walker, and Snigdha Chaturvedi. 2020. Bridging the structural gap between encoding and decoding for data-to-text generation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2481–2491.	1054
		1055
		1056
		1057
		1058
		1059
		1060
		1061
		1062
		1063
		1064
		1065
		1066
		1067
		1068
		1069
		1070
		1071
		1072
		1073
		1074
		1075
		1076
		1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084
		1085
		1086
		1087
		1088
		1089
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200

Filling the templates also often results in minor disfluencies, e.g. *nationality* \rightarrow "*<s> is from <o>*" will produce a missing definite article for *<o>* = "*United States*" and ungrammatical sentence for *<o>* = "*French people*". In principle, the disfluencies may be fixed by rephrasing in the final step of the pipeline.

We provide all the templates we used in our experiments in our anonymized repository.

C Experimental Setup

We implemented the models for split-and-rephrase, aggregation, and paragraph compression in PyTorch Lightning (Paszke et al., 2019), using the PyTorch (Falcon, 2019) version of the BART and RoBERTa models from the Huggingface library (Wolf et al., 2019).

We use the Adam (Kingma and Ba, 2015) optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.997$, $\epsilon = 1^{-9}$) with learning rate 2^{-5} , linear scheduling and 0.1 warmup proportion; batches of size 8 and accumulating gradients with factor 4. We train the models for 1 epoch on a single GeForce RTX 3090 GPU with 24 GB RAM. Training times were approximately 24 hours for the ordering model and 3 hours for the aggregation and paragraph compression models. We use greedy decoding in all our experiments.

For training the ordering model, we used the implementation from Calizzano et al. (2021)¹¹ including their training parameters. We are on the way to integrating the ordering model into our framework.

D WIKIFLUENT Evaluation

Table 7 summarizes the results of the individual modules of our pipeline (§5) on the **WIKIFLUENT test sets**. The evaluation metrics correspond to the metrics for evaluation of D2T generation (§7). All the modules are trained on the full WIKIFLUENT corpus and evaluated on *full* and *filtered* test sets.

	test (full)	test (filtered)
Ordering – Ours (§5.2)		
BLEU-2	64.8	71.9
Accuracy	0.70	0.77
Aggregation – Ours (§5.3)		
Acc. – per example	0.68	0.68
Acc. – per sent. bound.	0.93	0.93
Paragraph Compression – Ours (§5.4)		
BLEU	90.72	91.60
METEOR	63.89	65.03

Table 7: Result of individual pipeline modules on the WIKIFLUENT test sets (full / filtered).

E Additional Results

We provide evaluation of semantic accuracy on the E2E dataset as evaluated with the slot-error script based on matching regular expressions in Table 8.¹²

		miss	add	miss+add
	TGEN	0.0060	0.0433	0.0016
	COPY	0.0000	0.0000	0.0000
<i>full</i>	3-STAGE	0.0238	0.0698	0.0060
	2-STAGE	0.0054	0.0363	0.0000
	1-STAGE	0.0043	0.0330	0.0000
<i>filtered</i>	3-STAGE	0.0444	0.0487	0.0076
	2-STAGE	0.0043	0.0368	0.0000
	1-STAGE	0.0043	0.0347	0.0000

Table 8: Proportion of output examples with missed only, added only, and both missed and added facts, according to the regular-expression-based E2E slot error script.

Note that our manual investigation of a sample of the data shows that the majority of the errors identified in our model outputs are false. For example, the following regular expression used in the slot-error script:

`prices?(?: range)?(?:w+)0,3 high`

matches "(...) price range and high customer rating (...)", incorrectly classifying the presence of the extra slot *priceRange[high]*. This importance of this problem is exacerbated by the consistent outputs of our models, which tend to repeat certain patterns. However, we also manually identified several cases in which an error was found correctly, e.g. the model hallucinating "3 out of 4 customer rating" instead of "3 out of 5 customer rating".

¹¹<https://github.com/airklizz/passage-ordering>

¹²https://github.com/tuetschek/e2e-cleaning/blob/master/slot_error.py

F Example Outputs

Tables 10, 11, 12, and 13 show examples of behavior of our models on the **WebNLG dataset**. Tables 14 and 15 show examples of behavior of our models on the **E2E dataset**.

The **green** color marks the model outputs which are completely correct, the **red** color marks the errors. For better readability of the input format, we add numeric order identifiers for the individual facts (bold, in squared brackets). These are subsequently used as references in the Order and Aggregation rows. Note that zero-th input in E2E is used as a subject in the RDF triples.

	# train	# dev	# test	tok/src	tok/tgt	sent/src	sent/tgt	# templates
WebNLG	18,102	870	1,862	26.8	22.6	3.0	1.4	354
Clean E2E	33,236	4,299	1,847	29.2	22.3	4.2	1.5	8
WIKIFLUENT- <i>full</i>	915,855	9,346	9,346	52.9	41.1	3.9	2.0	-
WIKIFLUENT- <i>filtered</i>	700,517	7,149	7,149	45.6	35.4	3.4	1.8	-

Table 9: Number of examples (train / dev / test), average number of tokens per source and target, average number of sentences per source and target (after filling the templates for the D2T datasets).

Input	[1] (<i>Andrews County Airport; elevationAboveTheSeaLevel (in metres); 973.0</i>) [2] (<i>Andrews County Airport; runwayLength; 896.0</i>) [3] (<i>Andrews County Airport; location; Texas</i>)
Facts	Andrews County Airport is 973.0 metres above the sea level. Andrews County Airport runway is 896.0 m long. Andrews County Airport is located in Texas.
Order	3 2 1
Aggregation	3 <sep> 2 1
Models (all)	Andrews County Airport is located in Andrews County, Texas. Its runway is 896.0 m long and 973.0 metres above the sea level.
Human	The runway length of Andrews County Airport (located in Texas and 973.0 above sea level) is 896.

Table 10: Correct behavior of the models on WebNLG. All the models ({1,2,3}-STAGE *full/filtered*) produce the same output.

Input	[1] (<i>Akron, Ohio; country; United States</i>) [2] (<i>United States; language; English language</i>)
Facts	Akron, Ohio is from United States. English language is spoken in United States.
Order	1 2
Aggregation	1 <sep> 2
3-stage full	Akron, Ohio is from United States. English language is spoken in Akron, Ohio, United States, which is located in the U.S. state of Ohio.
2-stage full	Akron, Ohio is from United States. English language is spoken in Akron, Ohio.
1-stage full	English language is spoken in the United States. Akron, Ohio is from the state of Ohio in the U.S. state of Illinois.
3-stage filtered	Akron, Ohio is from United States. English language is spoken in the United States of America.
2-stage filtered	Akron, Ohio is from United States. English language is spoken in Akron, Ohio.
1-stage filtered	English language is spoken in the United States. Akron, Ohio is from the United Kingdom.
Human	The country of Akron, Ohio, is United States, where English is the language spoken.

Table 11: Incorrect behavior on WebNLG: besides the minor disfluencies caused by the templates ("Akron, Ohio is from..."), the models (except for 3-STAGE *filtered*) tend to hallucinate and merge the facts incorrectly.

Input	[1] (<i>Juan Perón; party; Labour Party (Argentina)</i>) [2] (<i>Alberto Teisaire; inOfficeWhilePresident; Juan Perón</i>) [3] (<i>Alberto Teisaire; nationality; Argentina</i>) [4] (<i>Argentina; language; Spanish language</i>)
Facts	Juan Perón belongs to the Labour Party Argentina. Alberto Teisaire was in office while Juan Perón was a president. Alberto Teisaire is from Argentina. Spanish language is spoken in Argentina.
Order	3 2 1 4
Aggregation	3 2 <sep> 1 <sep> 4
3-stage full	Alberto Teisaire is from Argentina, who was in office while Juan Perón was a president. He belongs to the Labour Party Argentina. Spanish language is spoken in Argentina.
2-stage full	Alberto Teisaire is from Argentina. He was in office while Juan Perón was a president. He belongs to the Labour Party Argentina. Spanish language is spoken in Argentina.
1-stage full	Alberto Teisaire is from Argentina. He was in office while Juan Perón was a president. Spanish language is spoken in Argentina. Perón belongs to the Labour Party Argentina.
3-stage filtered	Alberto Teisaire is from Argentina, in office while Juan Perón was a president. He belongs to the Labour Party Argentina. Spanish language is spoken in Argentina.
2-stage filtered	Alberto Teisaire is from Argentina. He was in office while Juan Perón was a president. He belongs to the Labour Party Argentina. Spanish language is spoken in Argentina.
1-stage filtered	Alberto Teisaire is from Argentina. He was in office while Juan Perón was a president. Spanish language is spoken in Argentina. Juan belongs to the Labour Party Argentina.
Human	Alberto Teisaire, is a national of Argentina, where Spanish is spoken. Alberto Teisaire, was in office while Juan Perón was president. Juan Perón was a member of the Labour Party in Argentina.

Table 12: Incorrect behavior on WebNLG: by using a pronoun "He", the models tend to imply that *Alberto Teisaire* (instead of *Juan Perón*) belongs to the Labour Party.

Input	[1] (<i>Alfa Romeo 164; relatedMeanOfTransportation; Fiat Croma</i>) [2] (<i>Alfa Romeo 164; assembly; Italy</i>) [3] (<i>Italy; capital; Rome</i>)
Facts	Alfa Romeo 164 is related to Fiat Croma. Alfa Romeo 164 was assembled in Italy. Italy's capital is Rome.
Order	2 3 1
Aggregation	2 3 <sep> 1
3-stage full	Alfa Romeo 164 was assembled in Italy's capital, Rome . It is related to Fiat Croma.
2-stage full	Alfa Romeo 164 was assembled in Italy's capital, Rome . It is related to Fiat Croma.
1-stage full	Alfa Romeo 164 was assembled in Italy and related to Fiat Croma. Italy's capital is Rome.
3-stage filtered	Alfa Romeo 164 was assembled in Italy's capital, Rome . It is related to Fiat Croma.
2-stage filtered	Alfa Romeo 164 was assembled in Italy's capital, Rome . It is related to Fiat Croma.
1-stage filtered	Alfa Romeo 164 was assembled in Italy. It is related to Fiat Croma. It was assembled at the Italian capital is Rome.
Human	The Fiat Croma and the Alfa Romeo 164 are both cars. The latter was built in Italy where the capital is Rome.

Table 13: Incorrect behavior on WebNLG: the models tend to merge the country of vehicle's assembly and its capital.

Input	<i>[0] name[The Cricketers]</i> <i>[1] eatType[restaurant]</i> <i>[2] near[All Bar One]</i> <i>[3] priceRange[cheap]</i> <i>[4] food[Chinese]</i> <i>[5] customerRating[average]</i> <i>[6] familyFriendly[yes]</i>
Facts	The Cricketers is a restaurant. The Cricketers is near All Bar One. The Cricketers has cheap price range. The Cricketers serves Chinese food. The Cricketers has average customer rating. The Cricketers is family-friendly.
Order	1 4 2 6 3 5
Aggregation	1 4 2 <sep> 6 3 5
3-stage full	The Cricketers is a restaurant serving Chinese food near All Bar One. It is family-friendly, has cheap price range and average customer rating.
2-stage full	The Cricketers is a restaurant serving Chinese food near All Bar One. It is family-friendly with cheap price range and average customer rating.
1-stage full	The Cricketers is a restaurant serving Chinese food. It is family-friendly and has cheap price range. It serves Chinese food near All Bar One. It has average customer rating.
3-stage filtered	The Cricketers is a restaurant serving Chinese food near All Bar One. It is family-friendly with cheap price range and average customer rating.
2-stage filtered	The Cricketers is a restaurant serving Chinese food near All Bar One. It is family-friendly with cheap price range and average customer rating.
1-stage filtered	The Cricketers is a restaurant serving Chinese food. It is family-friendly and has average customer rating. It serves Chinese food near All Bar One. It has cheap price range.
Human	For a central family friendly venue serving cheap Chinese food, try The Cricketers restaurant. It has average customer ratings and is near All Bar One.

Table 14: Correct behavior of the models on E2E: the models can deal even with a higher amount of input values.

Input	<i>[0] name[The Phoenix]</i> <i>[1] area[riverside]</i> <i>[2] priceRange[cheap]</i> <i>[3] familyFriendly[yes]</i> <i>[4] near[Raja Indian Cuisine]</i> <i>[5] eatType[restaurant]</i>
Facts	The Phoenix is in the riverside. The Phoenix has cheap price range. The Phoenix is family-friendly. The Phoenix is near Raja Indian Cuisine. The Phoenix is a restaurant.
Order	5 1 4 <sep> 3 2
Aggregation	5 1 4 3 2
3-stage full	The Phoenix is a restaurant in the riverside near Raja Indian Cuisine. It is family-friendly with cheap price range.
2-stage full	The Phoenix is a restaurant in the riverside near Raja Indian Cuisine. It is family-friendly with cheap price range.
1-stage full	The Phoenix is a restaurant near Raja Indian Cuisine in the riverside. It is family-friendly and has cheap price range. It has cheap food range.
3-stage filtered	The Phoenix is a restaurant in the riverside near Raja Indian Cuisine. It is family-friendly with cheap price range.
2-stage filtered	The Phoenix is a restaurant in the riverside near Raja Indian Cuisine. It is family-friendly with cheap price range.
1-stage filtered	The Phoenix is a restaurant near Raja Indian Cuisine in the riverside. It is family-friendly and has cheap price range. It has cheap food.
Human	Cheap food and a family friendly atmosphere at The Phoenix restaurant. Situated riverside near the Raja Indian Cuisine.

Table 15: Incorrect behavior on E2E: the 1-STAGE models add redundant information to the output.