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

**Zdeněk Kasner**  and  **Ondřej Dušek**
Charles University, Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics
Prague, Czech Republic
{kasner,odusek}@ufal.mff.cuni.cz

## 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 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 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 a model on the domain-specific data leads to overfitting to the particular benchmark, decreasing performance on out-of-domain
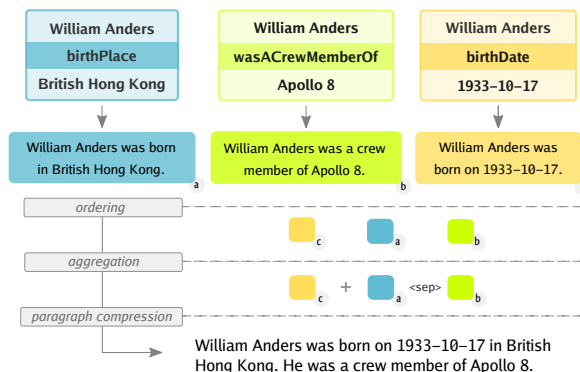


Figure 1: A scheme of our pipeline for zero-shot data-to-text generation from RDF triples: (1) ordering, (2) aggregation, (3) paragraph compression. Individual pipeline modules are trained on a large general-domain text corpus and operate over text in natural language. In-domain knowledge is included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

data (Laha et al., 2019). Gathering a large set of references for a particular domain is also costly and time-consuming as it usually requires collecting human-written references through crowdsourcing (Dušek et al., 2020). These problems can be partially mitigated using *few-shot* approaches (Chen et al., 2020b; Ke et al., 2021; Su et al., 2021a), which operate with only several dozens or hundreds of annotated examples, but the robustness of these approaches 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., 2022).

In this paper, we present a *zero-shot* 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

---

[1]Our code and data is available at https://github.com/kasnerz/zeroshot-d2t-pipeline.

we subsequently order, aggregate, and compress on the paragraph level to produce the resulting description of the data. In constrast to traditional pipeline systems, all our pipeline modules are built upon PLMs and operate over sentences in natural language. The modules are trained on our new WIKIFLUENT corpus, which contains 934k examples of first paragraphs from the English Wikipedia, each supplied with a synthesized set of simple template-like sentences which together convey the meaning of the original paragraph. Our approach allows generating natural language descriptions from RDF triples with a minimum amount of domain-specific rules or knowledge and without using training data from the D2T datasets. Although our approach is primarily a probe into the territory of zero-shot approaches and cannot yet match the quality of state-of-the-art models, we show that it can yield large improvements upon simple baselines and match older supervised systems on automatic metrics for text fluency. Moreover, the semantic accuracy metrics and our manual error analysis suggest that our approach offers a way to prevent omissions and hallucinations common in few-shot approaches.

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.[1]

## 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; Kasner and Dušek, 2020b). 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 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 any training data from target D2T datasets.

Laha et al. (2019) introduce a three-step pipeline for zero-shot D2T generation similar to ours. Unlike the approach we describe here, they use a semi-automatic template generation system,[2] their sentence fusion is rule-based, 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 in neural D2T pipelines. Unlike the aforementioned planners, which use predicates or keys from 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.

---

[2]As we describe in §5.1, we opted for a simpler way for generating the templates to showcase the results of our approach independently of the template generator quality.

(2021), who use a pointer network (Wang and Wan, 2019; Vinyals et al., 2015) on top of a PLM.

**Aggregating Input into Sentences** Typically, multiple pieces of input information need to be merged into 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 (see §3.1), into which we selectively insert delimiters to mark sentence boundaries.

**Paragraph Compression** We introduce *paragraph compression* (PC) as a new task and the final step in our D2T generation pipeline. This task combines several standard natural-language tasks including 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

In this section, we provide the formal description of our proposed approach. We focus on the task of producing a natural language description $Y$ for a set of $n$ RDF triples $X = \{x_1, \ldots, x_n\}$. Each triple $x_i = \{s_i, p_i, o_i\}$ consists of subject $s_i$, predicate $p_i$, and object $o_i$.

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

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

The individual steps are described in the following sections: 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 $T : X \to F$. We define a fact $f_i$ as 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. This step can be realized e.g. using a simple template for each predicate (cf. §5.1).

### 3.2 Ordering the Facts

We assume that the default order of triples $X$ is random and the same applies for the respective facts $F$. Note, however, that that $F$ is a indeed set of meaningful sentences. We can use this to our advantage and apply a sentence ordering model to maximize the coherency of the paragraph resulting from their concatenation. An example outcome of such operation may be grouping together facts mentioning *birth date* and *birth place* of a person, followed by their *occupation* (see Figure 1). The ordering module allows downstream modules to only focus on operations over neighboring sentences.

Formally, we apply the ordering model $O(F)$ to get an ordered sequence of facts: $F_o = \{f_{o_1}, \ldots, f_{o_n}\}$, where $o_{1:n}$ is a permutation of indices. We describe our ordering model in §5.2.

### 3.3 Aggregating the Facts

Some facts will be typically mentioned together in a single sentence. Considering the previous example, *occupation* is likely to be mentioned separately, while *birth date* and *birth place* are likely to be mentioned together. Using an ordered sequence of facts as input, we can apply an aggregation model to decide which facts should be merged into a single sentence.

Formally, the aggregation model takes a sequence of ordered facts $F_o$ as input and produces a sequence of sentence delimiters $A(F_o) = \{\delta_{o_1}, \delta_{o_2}, \ldots, \delta_{o_{n-1}}\}$; $\delta_i \in \{0, 1\}$. The output $\delta_i = 1$ means that the neighboring facts should be mentioned separately, i.e. the neighboring sentences should *not* be fused. Conversely, $\delta_i = 0$ means that the facts should be aggregated and their corresponding sentences should be fused. We describe our aggregation model in §5.3.

## 3.4 Paragraph Compression

The paragraph compression (PC) model is a generative model which outputs the final text description. It has two main objectives: (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 disfluencies in the templates, replacing noun phrases with refering expressions, etc. The goal of the task is to preserve the semantics of the text which is an already ordered sequence of sentences, so the edits will typically be minor. Formally, the model 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 the final text $Y$. We describe our PC model in §5.4.

## 4 WIKIFLUENT Corpus

Here we descibe 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).

Our goal is 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 using human-written paragraphs in English Wikipedia and applying *split-and-rephrase* and *coreference resolution* models to obtain synthetic source texts. 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 dump[3] using WikiExtractor (Attardi, 2015). Wikipedia is commonly used for large-scale pretraining of D2T generation models (Jin et al., 2020; Chen et al., 2020a). Although it is not bias-free, it provides more balanced sample of natural language use than typical D2T generation datasets. We used the first paragraphs of Wikipedia entries, which contain mostly concise, fact-based descriptions. We selected paragraphs with length between
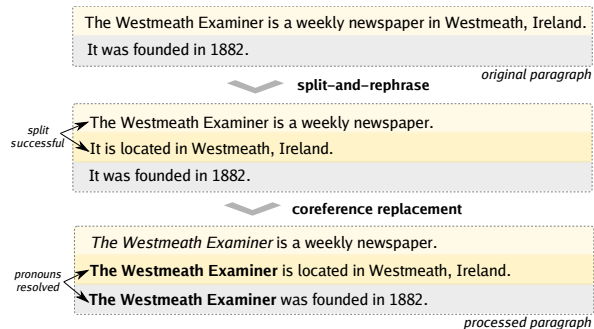
[3] `enwiki-20210401-pages-articles-multistream`



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. The result is a set of simple sentences which together convey the same meaning as the original paragraph. The synthesized sentences are used as *input* into our models, the original human-written texts are used as *ground truth*.

30-430 characters; filtering out lists, disambiguations, and repeated and malformed paragraphs. To balance the length of inputs, we selected 250k examples each from 4 equally sized length ranges (30-130 characters, etc.).

### 4.2 Split-and-Rephrase

To generate a set of simple 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). The process is illustrated in the upper part of Figure 2.

We train our split-and-rephrase model on the large-scale WikiSplit corpus by Botha et al. (2018), containing human-made sentence splits from Wikipedia edit history. Following the same setup as for a paragraph compression model (§3.4), we train BART-base (Lewis et al., 2020) on the WikiSplit dataset in a sequence-to-sequence setting. Next, we apply the trained split-and-rephrase model on each sentence in our Wikipedia-based corpus, 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

As the next step, we concatenate the split sentences and apply a coreference resolution model (Gardner et al., 2018; Lee et al., 2018) in order to replace referring expressions with their antecendents (e.g., pronouns with noun phrases). The motivation for this step is to match the style of the facts (see §3.1), which do not use pronouns since each fact describes a single triple only. Note that this procedure replaces the referring expressions only in the synthesized sentences (which are used as input) and keeps them in the original paragraphs (which are used as ground truth). As a consequence, the paragraph compression module is implicitly trained to generate referring expressions in the final description.

## 4.4 Filtering

To ensure that the generated sentences convey the same semantics as the original paragraph, we use a pretrained RoBERTa model[4] (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 75% 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).[5]

## 5.1 Templates

We transform triples into facts (§3.1) using a single-triple template $t_i$ for each predicate. For example, if $p_i = $ *instrument*, then $T(p_i) = $ "$s_i$ *plays* $o_i$" (cf. Table 1). We follow previous work in which simple hand-crafted templates have been used as an efficient way of introducing domain knowledge (Kale and Rastogi, 2020a; Kasner and Dušek, 2020a). Compared to more complex rule-based

| 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>*).

template generation engines (Laha et al., 2019; Heidari et al., 2021; Mehta et al., 2021), the approach may produce less fluent outputs, but it minimizes manual workload and makes it easier to control the quality of the input for the subsequent steps.

## 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 please refer to 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 ($d_2 = $<s>$f_{o_1}$</s>). 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 state $d_j$. The network is based on the scaled dot product attention, where $d_j$ is the query and encoder outputs $E_i$ are the keys:

$$Q = d_j W_Q$$
$$K = E W_K$$
$$P_j = \text{softmax}\left(\frac{QK^T}{\sqrt{b}}\right).$$

---

[4] https://huggingface.co/roberta-large-mnli
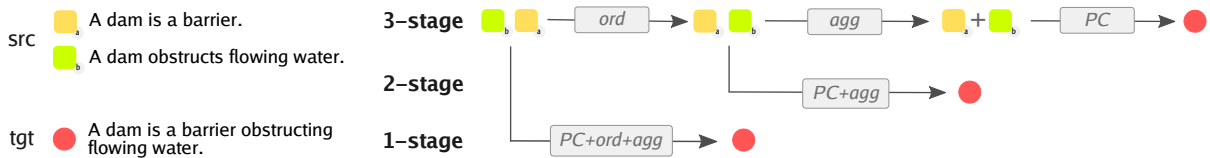[5] Our training setup details are included in Appendix C.

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 and §6 for the paragraph compression model (PC, PC+AGG, PC+ORD+AGG).

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 synthesized simple 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.[6] Similarly to the ordering model (§5.2), we input the sequence of (now ordered) 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). Subsequently, the token classification layer classifies each separator </s>$_i$ position into two classes $\{0, 1\}$ corresponding to the delimiter $\delta_i$. We ignore the outputs for the non-separator tokens while computing cross-entropy loss.

We create the training examples using the synthesized 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 finetune the model on the WIKIFLUENT corpus, concatenating the synthesized sentences on the input. We add delimiters between the sentences $i$ and $i+1$ where $\delta_i = 1$ using a special token <sep>, 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. Here we expect that the model will learn to fuse the sentences between which there are no delimiters

on the input. We evaluate how the model learns to respect the order and aggregation markers in §7.3.

## 6 Experiments

We train our pipeline modules on the WIKIFLUENT corpus as described in §5. Next, we use these modules *without finetuning* for generating descriptions for RDF triples on two English D2T datasets, WebNLG and E2E.

**Datasets** The datasets differ in domain, size, textual style, and number of predicates (see Appendix A for details):
- **WebNLG** (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 comparison to prior work. We hand-crafted templates for all 354 predicates, including unseen predicates in the test set.[7]
- **E2E** (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 transform the attribute-value pairs into RDF triples (using the restaurant name as a subject) and then apply the same setup as for WebNLG. We created a template for each of the 8 attributes manually.

**Pipeline versions** In order to evaluate individual components of our pipeline, we train three versions of the *paragraph compression* 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

---

[6]https://huggingface.co/transformers/model_doc/roberta.html#robertafortokenclassification

[7]See Appendix B for details on template creation.

(i.e., both ordering and aggregation are left implicitly to the model).

Correspondingly, 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 (ORD), the aggregation model (AGG) and the PC model (following the full pipeline from §3),
- 2-STAGE – a pipeline consisting of the ORD 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. 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). We also evaluate the performance of the content planning modules and the ability of the PC module to follow the content plan (§7.3). In §7.4, we include an intrinsic evaluation of our modules on the WIKIFLUENT test set.

### 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 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. For **WebNLG** (see Table 3), we compare our systems with the results of:

- UPF-FORGe and MELBOURNE – systems (grammar-based and supervised, respectively) from the first run of WebNLG Challenge (Gardent et al., 2017),
- Ke et al. (2021) – a state-of-the-art system with

a structure-aware encoder and task-specific pre-training,
- Laha et al. (2019) – the only other (to our knowledge) zero-shot D2T generation system applied to WebNLG.

For **E2E** (see Table 4), we compare our systems with the results of:

- TGEN (Dušek and Jurčíček, 2015) – the baseline system for the E2E Challenge (Dušek et al., 2020),
- Harkous et al. (2020) – a state-of-the-art supervised system on cleaned E2E data.

For both datasets, COPY denotes the baseline of copying the facts without further processing.

The automatic evaluation shows that our systems consistently outperform the COPY baseline (e.g., ~12 BLEU points for E2E), which is already strong thanks to our manually curated set of templates.[10] Automatic scores also suggest that our systems are comparable with some older supervised systems. Nevertheless, our systems still underperform the state-of-the-art supervised systems. For this reason, we further focus on manual *error analysis* in §7.2 to pinpoint the current shortcomings of our approach.

The 2-STAGE system is generally on par with the 3-STAGE system or better, which indicates that explicit aggregation using the AGG model may not be necessary. However, an advantage of having a separate aggregation module is the possibility to control the aggregation step explicitly. The models using the filtered version of the corpus generally produce better results, although they also bring in a larger number of omissions.

### 7.2 Manual Error Analysis

Since automatic performance metrics do not provide insights into specific weaknesses of the system (van Miltenburg et al., 2021), we manually examined 100 outputs of the models. We counted the number of errors: factual (hallucinations, omissions, incorrect fact merging, redundancies) and grammatical. The results are summarized 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 the 2-STAGE and 3-STAGE models, which produce

---

[8]We use the implementation from `https://github.com/tuetschek/e2e-metrics`.

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

[10]On WebNLG, our COPY baseline achieves 37.18 BLEU points, compared to 24.80 BLEU points of the *full system* of Laha et al. (2019), which uses automatic template generation.

| | 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 E 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. (2019)[†] | | 24.80 | 34.90 | - | - |
| COPY | | 37.18 | 38.77 | 0.000 | 0.000 |
| *full* | 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 |
| *filtered* | 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 taken from the respective works.

| | | **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 |
| *full* | 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 |
| *filtered* | 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 taken from their work.

| | | **WebNLG** | | | | | **E2E** | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | **H** | **I** | **O** | **R** | **G** | **H** | **I** | **O** | **R** | **G** |
| *full* | 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 |
| *filtered* | 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.

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 E for examples). This behavior is the main obstacle to ensure factual consistency of the output. As a possible remedy, we propose explicitly controlling the semantics of sentence fusion (Ben-David et al., 2020), e.g. using a variant of constrained decoding (Balakrishnan et al., 2019; Wang et al., 2021).

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. 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 and BLEU-2 score of our **ordering model** on WebNLG against the human-generated plans from Ferreira et al. (2018). The results are listed in Table 6 and compared against a RANDOM baseline (random ordering) and prior work. The results show that although our approach again lags behind state-of-the-art supervised ap-

| | 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 |
| Ours (BART+ptr) | 59.10 | 0.48 |

Table 6: Evaluation of our zero-shot ordering model based on Calizzano et al. (2021). B-2 = BLEU-2, Acc = accuracy. The results marked with † are copied from the respective papers.

| | | test (full) | test (filt.) |
|---|---|---|---|
| ORD | BLEU-2 | 64.8 | 71.9 |
| | Accuracy | 0.70 | 0.77 |
| AGG | Acc. per example | 0.68 | 0.68 |
| | Acc. per sent. bound. | 0.93 | 0.93 |
| PC | BLEU | 90.72 | 91.60 |
| | METEOR | 63.89 | 65.03 |

Table 7: Result of individual pipeline modules on the WIKIFLUENT test sets (full / filtered). The metrics correspond to the metrics used for evaluating the modules for D2T generation.

proaches, it can outperform both the random baseline and the Transformer-based approach from Ferreira et al. (2019) while not using any in-domain examples.

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, there is still room for improvement.

Finally, we manually evaluate how the **PC model** follows the content plan (i.e., keeping the predefined order and aggregating the sentences according to the delimiters) 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 or an extra boundary between sentences without a separator. We can thus conclude that the pretraining task successfully teaches the PC model to follow a given content plan.

### 7.4 Intrinsic Evaluation

Aside from the main D2T generation results, we also provide an intrinsic evaluation of our pipeline modules on the WIKIFLUENT test sets. We evaluated the ordering, aggregation, and paragraph compression modules trained on the *full* WIKIFLUENT corpus. The results for both *full* and *filtered* test sets are summarized in Table 7. The PC model achieves high scores, which follows from the fact that we provide it with ground truth content plans (i.e., the ordering and aggregation plan corresponding to the original paragraph). Accuracy of the ordering and aggregation modules is comparable to their performance on D2T datasets.

## 8 Future Work

Our experiments outline several possible future research directions. Automatic generation of facts without using hand-crafted templates (cf. §5.1) could allow applying zero-shot generation systems to datasets with a large number of predicates, such as ToTTo (Parikh et al., 2020). The task of paragraph compression could be used as a task-specific pretraining (Gururangan et al., 2020) for more efficient finetuning of D2T models, e.g., with a small amount of clean data. Consistency checks may be introduced in the pipeline to control the output from the modules, and individual modules may be improved by using more efficient model architectures.

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. As we suggested in §7.2, explicitly controlling the semantics of sentence fusion could help to mitigate this issue, while still keeping the advantages of a zero-shot approach.

## 9 Conclusion

We presented an approach for zero-shot D2T generation. The approach uses a pipeline of PLMs trained on general-domain lexical operations over natural language. The pipeline builds upon traditional approaches and consists of three interpretable intermediate steps. By avoiding noisy human-written references from the D2T datasets, our models produce more semantically consitent output. We believe that training models for zero-shot D2T generation using large cross-domain corpora will help to build D2T generation systems with good performance across various domains.

## 10 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, NAACL-HLT 2021*, pages 3554–3565, Online.

Giuseppppe Attardi. 2015. WikiExtractor. https://github.com/attardi/wikiextractor.

Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary G. Ives. 2007. DBpedia: A Nucleus for a Web of Open Data. In *The Semantic Web, 6th International Semantic Web Conference, 2nd Asian Semantic Web Conference, ISWC 2007 + ASWC 2007*, volume 4825 of *Lecture Notes in Computer Science*, pages 722–735, Busan, Korea.

Anusha Balakrishnan, Jinfeng Rao, Kartikeya Upasani, Michael White, and Rajen Subba. 2019. Constrained Decoding for Neural NLG from Compositional Representations in Task-Oriented Dialogue. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Volume 1: Long Papers*, pages 831–844, Florence, Italy.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Proceedings of the Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization@ACL 2005*, pages 65–72, Ann Arbor, MI, USA.

Regina Barzilay, Noémie Elhadad, and Kathleen R. McKeown. 2001. Sentence Ordering in Multidocument Summarization. In *Proceedings of the First International Conference on Human Language Technology Research, HLT 2001*, San Diego, CA, USA.

Regina Barzilay and Kathleen R. McKeown. 2005. Sentence Fusion for Multidocument News Summarization. *Comput. 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 *Findings of the Association for Computational Linguistics: EMNLP 2020*, Findings of ACL, pages 1491–1505, Online.

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 *FAccT '21: 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623, Online/Toronto, Canada.

Steven Bird. 2006. NLTK: The Natural Language Toolkit. In *ACL 2006, 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference*, Sydney, Australia.

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.

Rémi Calizzano, Malte Ostendorff, and Georg Rehm. 2021. Ordering sentences and paragraphs with pretrained encoder-decoder transformers and pointer ensembles. In *DocEng '21: ACM Symposium on Document Engineering 2021*, pages 10:1–10:9, Limerick, Ireland.

Ernie Chang, Xiaoyu Shen, Hui-Syuan Yeh, and Vera Demberg. 2021a. On Training Instance Selection for Few-Shot Neural Text Generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 2: Short Papers)*, pages 8–13, Online.

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, EACL 2021*, pages 758–768, Online.

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 2020*, pages 8635–8648, Online.

Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 2016. Neural Sentence Ordering. *CoRR*, abs/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, ACL 2020*, 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 2020*, pages 6310–6320, Online.

Robert Dale. 2020. Natural language generation: The commercial state of the art in 2020. *Nat. Lang. Eng.*, 26(4):481–487.

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, NAACL-HLT 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, MN, USA.

Ondřej Dušek, David M. Howcroft, and Verena Rieser. 2019. Semantic Noise Matters for Neural Natural Language Generation. In *Proceedings of the 12th International Conference on Natural Language Generation, INLG 2019*, pages 421–426, Tokyo, Japan.

Ondřej Dušek and Filip Jurčíček. 2015. Training a Natural Language Generator From Unaligned Data. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, Volume 1: Long Papers*, pages 451–461, Beijing, China.

Ondřej Dušek and Zdeněk Kasner. 2020. Evaluating Semantic Accuracy of Data-to-Text Generation with Natural Language Inference. In *Proceedings of the 13th International Conference on Natural Language Generation, INLG 2020*, pages 131–137, Dublin, Ireland.

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. *Comput. Speech Lang.*, 59:123–156.

William Falcon et al. 2019. Pytorch Lightning. https://github.com/PyTorchLightning/pytorch-lightning.

Thiago Ferreira, Claire Gardent, Nikolai Ilinykh, Chris van der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina. 2020. The 2020 Bilingual, Bi-Directional WebNLG+ Shared Task Overview and Evaluation Results (WebNLG+ 2020). In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*.

Thiago Castro Ferreira, Diego Moussallem, Emiel Krahmer, and Sander Wubben. 2018. Enriching the WebNLG corpus. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 171–176, Tilburg, The Netherlands.

Thiago Castro Ferreira, Chris van der Lee, Emiel van Miltenburg, and Emiel Krahmer. 2019. Neural data-to-text generation: A comparison between pipeline and end-to-end architectures. 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 2019*, pages 552–562, Hong Kong.

Katja Filippova and Yasemin Altun. 2013. Overcoming the Lack of Parallel Data in Sentence Compression. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1481–1491, Seattle, WA, USA.

Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. The WebNLG Challenge: Generating Text from RDF Data. In *Proceedings of the 10th International Conference on Natural Language Generation, INLG 2017*, pages 124–133, Santiago de Compostela, Spain.

Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew E. Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. AllenNLP: A Deep Semantic Natural Language Processing Platform. *CoRR*, abs/1803.07640.

Albert Gatt and Emiel Krahmer. 2018. Survey of the State of the Art in Natural Language Generation: Core tasks, applications and evaluation. *J. Artif. Intell. Res.*, 61:65–170.

Mor Geva, Eric Malmi, Idan Szpektor, and Jonathan Berant. 2019. DiscoFuse: A Large-Scale Dataset for Discourse-Based Sentence Fusion. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Volume 1 (Long and Short Papers)*, pages 3443–3455, Minneapolis, MN, USA.

Jingjing Gong, Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 2016. End-to-End Neural Sentence Ordering Using Pointer Network. *CoRR*, abs/1611.04953.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360.

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 *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020*, pages 2410–2424, Barcelona, Spain (Online).

Peyman Heidari, Arash Einolghozati, Shashank Jain, Soumya Batra, Lee Callender, Ankit Arun, Shawn Mei, Sonal Gupta, Pinar Donmez, Vikas Bhardwaj, Anuj Kumar, and Michael White. 2021. Getting to Production with Few-shot Natural Language Generation Models. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue, SIGdial 2021*, pages 66–76, Singapore/Online.

Chao Jiang, Mounica Maddela, Wuwei Lan, Yang Zhong, and Wei Xu. 2020. Neural CRF Model for Sentence Alignment in Text Simplification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020*, pages 7943–7960, Online.

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 *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020*, pages 2398–2409, Barcelona, Spain (Online).

Mihir Kale and Abhinav Rastogi. 2020a. Template Guided Text Generation for Task-Oriented Dialogue. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020*, pages 6505–6520, Online.

Mihir Kale and Abhinav Rastogi. 2020b. Text-to-Text Pre-Training for Data-to-Text Tasks. In *Proceedings of the 13th International Conference on Natural Language Generation, INLG 2020*, pages 97–102, Dublin, Ireland.

Zdeněk Kasner and Ondřej Dušek. 2020a. Data-to-Text Generation with Iterative Text Editing. In *Proceedings of the 13th International Conference on Natural Language Generation, INLG 2020*, pages 60–67, Dublin, Ireland.

Zdeněk Kasner and Ondřej Dušek. 2020b. Train Hard, Finetune Easy: Multilingual Denoising for RDF-to-Text Generation. In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*, pages 171–176, Online.

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. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pages 2526–2538, Online.

Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. CTRL: A Conditional Transformer Language Model for Controllable Generation. *CoRR*, abs/1909.05858.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015*, San Diego, CA, USA.

Anirban Laha, Parag Jain, Abhijit Mishra, and Karthik Sankaranarayanan. 2019. Scalable Micro-planned Generation of Discourse from Structured Data. *Comput. Linguistics*, 45(4):737–763.

Mirella Lapata. 2003. Probabilistic Text Structuring: Experiments with Sentence Ordering. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 545–552, Sapporo, Japan.

Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018. Higher-Order Coreference Resolution with Coarse-to-Fine Inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, Volume 2 (Short Papers)*, pages 687–692, New Orleans, LA, USA.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pretraining for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020*, pages 7871–7880, Online.

Jiwei Li and Dan Jurafsky. 2017. Neural Net Models of Open-domain Discourse Coherence. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017*, pages 198–209, Copenhagen, Denmark.

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. *CoRR*, abs/1907.11692.

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. *CoRR*, abs/2110.08467.

Amit Moryossef, Yoav Goldberg, and Ido Dagan. 2019a. Improving Quality and Efficiency in Plan-based Neural Data-to-text Generation. In *Proceedings of the 12th International Conference on Natural Language Generation, INLG 2019*, pages 377–382, Tokyo, Japan.

Amit Moryossef, Yoav Goldberg, and Ido Dagan. 2019b. 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, NAACL-HLT 2019, Volume 1 (Long and Short Papers)*, pages 2267–2277, Minneapolis, MN, USA.

Shashi Narayan, Claire Gardent, Shay B. Cohen, and Anastasia Shimorina. 2017. Split and Rephrase. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017*, pages 606–616, Copenhagen, Denmark.

Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2017. The E2E Dataset: New Challenges For End-to-End Generation. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, 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 *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, pages 311–318.

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. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020*, pages 1173–1186, Online.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019*, pages 8024–8035, Vancouver, BC, Canada.

Ratish Puduppully, Li Dong, and Mirella Lapata. 2019. Data-to-Text Generation with Content Selection and Planning. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019*, pages 6908–6915, Honolulu, HI, USA.

Clément Rebuffel, Marco Roberti, Laure Soulier, Geoffrey Scoutheeten, Rossella Cancelliere, and Patrick Gallinari. 2022. Controlling hallucinations at word level in data-to-text generation. *Data Min. Knowl. Discov.*, 36(1):318–354.

Ehud Reiter. 2007. An Architecture for Data-to-Text Systems. In *Proceedings of the Eleventh European Workshop on Natural Language Generation, ENLG 2007*, Schloss Dagstuhl, Germany.

Ehud Reiter and Robert Dale. 1997. Building applied natural language generation systems. *Nat. Lang. Eng.*, 3(1):57–87.

Leonardo F. R. Ribeiro, Martin Schmitt, Hinrich Schütze, and Iryna Gurevych. 2020. Investigating Pretrained Language Models for Graph-to-Text Generation. *CoRR*, abs/2007.08426.

Anna Rogers. 2021. Changing the World by Changing the Data. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers)*, pages 2182–2194, Online.

Zhihong Shao, Minlie Huang, Jiangtao Wen, Wenfei Xu, and Xiaoyan Zhu. 2019. Long and Diverse Text Generation with Planning-based Hierarchical Variational Model. 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 2019*, pages 3255–3266, Hong Kong.

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 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020*, pages 7155–7165, Online.

Yixuan Su, Zaiqiao Meng, Simon Baker, and Nigel Collier. 2021a. Few-Shot Table-to-Text Generation with Prototype Memory. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 910–917, Online/Punta Cana, Dominican Republic.

Yixuan Su, David Vandyke, Sihui Wang, Yimai Fang, and Nigel Collier. 2021b. Plan-then-Generate: Controlled Data-to-Text Generation via Planning. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 895–909, Online/Punta Cana, Dominican Republic.

Bayu Distiawan Trisedya, Jianzhong Qi, and Rui Zhang. 2020. Sentence Generation for Entity Description with Content-Plan Attention. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020*, pages 9057–9064, New York, NY, USA.

Emiel van Miltenburg, Miruna-Adriana Clinciu, Ondřej Dušek, Dimitra Gkatzia, Stephanie Inglis, Leo Leppänen, Saad Mahamood, Emma Manning, Stephanie Schoch, Craig Thomson, and Luou Wen. 2021. Underreporting of errors in NLG output, and what to do about it. In *Proceedings of the 14th International Conference on Natural Language Generation, INLG 2021*, pages 140–153, Aberdeen, Scotland, UK.

Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015*, pages 2692–2700, Montréal, QC, Canada.

Tianming Wang and Xiaojun Wan. 2019. Hierarchical Attention Networks for Sentence Ordering. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019*, pages 7184–7191, Honolulu, HI, USA.

Yufei Wang, Ian D. Wood, Stephen Wan, Mark Dras, and Mark Johnson. 2021. Mention Flags (MF): Constraining Transformer-based Text Generators. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers)*, pages 103–113, Online.

Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, LA, USA.

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, Brussels, Belgium.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. HuggingFace's Transformers: State-of-the-art Natural Language Processing. *CoRR*, abs/1910.03771.

Xinnuo Xu, Ondřej Dušek, Verena Rieser, and Ioannis Konstas. 2021. AggGen: Ordering and Aggregating while Generating. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers)*, pages 1419–1434, Online.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2020. PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020*, volume 119 of *Proceedings of Machine Learning Research*, pages 11328–11339, Online.

Chao Zhao, Marilyn A. Walker, and Snigdha Chaturvedi. 2020. Bridging the Structural Gap Between Encoding and Decoding for Data-To-Text Generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020*, pages 2481–2491, Online.

## A  Data Statistics

Statistics for the datasets described in the paper are listed in Table 9.

## B  Templates

The templates for our datasets are single-sentence and mostly clear-cut verbalizations of the predicates. The templates were created by one of the authors who had only the input data at their disposal, i.e. without using human references.

We have also considered extracting the templates for WebNLG from the training data by delexicalizing single-triple examples. However, the examples are noisy and such data would not be available in a zero-shot setup, which is why we decided not to use this option.

Although the templates were mostly unambiguous, we had to opt for the most general version in certain cases (e.g., using *country → "<s> is from <o>"*, even though *"<s> is a food from <o>."* would be possible in case the object is food).

Filling the templates also often results in minor disfluencies, e.g. *nationality → "<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 repository.

## C  Experimental Setup

We implemented the models for split-and-rephrase, aggregation, and paragraph compression in Py-Torch Lightning (Paszke et al., 2019), using the PyTorch (Falcon et al., 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 ($\beta_1 = 0.9, \beta_2 = 0.997, \varepsilon = 1^{-9}$) 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 will integrate the ordering model into our framework.

## D  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 |
| *full* | 3-STAGE | 0.0238 | 0.0698 | 0.0060 |
|  | 2-STAGE | 0.0054 | 0.0363 | 0.0000 |
|  | 1-STAGE | 0.0043 | 0.0330 | 0.0000 |
| *filtered* | 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 regex-based E2E slot error script.

However, please 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 problem is magnified 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"*.

## E  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 | #temp. |
|---|---|---|---|---|---|---|---|---|
| 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-*full* | 915,855 | 9,346 | 9,346 | 52.9 | 41.1 | 3.9 | 2.0 | - |
| WIKIFLUENT-*filtered* | 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), total number of templates.

| | |
|---|---|
| **Input** | *[1] (Andrews County Airport; elevationAboveTheSeaLevel (in metres); 973.0)* |
| | *[2] (Andrews County Airport; runwayLength; 896.0)* |
| | *[3] (Andrews County Airport; location; Texas)* |
| **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] (Akron, Ohio; country; United States)* |
| | *[2] (United States; language; English language)* |
| **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] (Juan Perón; party; Labour Party (Argentina))* |
|---|---|
| | *[2] (Alberto Teisaire; inOfficeWhilePresident; Juan Perón)* |
| | *[3] (Alberto Teisaire; nationality; Argentina)* |
| | *[4] (Argentina; language; Spanish language)* |
| **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] (Alfa Romeo 164; relatedMeanOfTransportation; Fiat Croma)* |
|---|---|
| | *[2] (Alfa Romeo 164; assembly; Italy)* |
| | *[3] (Italy; capital; Rome)* |
| **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** | *[0] name[The Cricketers]* |
| | *[1] eatType[restaurant]* |
| | *[2] near[All Bar One]* |
| | *[3] priceRange[cheap]* |
| | *[4] food[Chinese]* |
| | *[5] customerRating[average]* |
| | *[6] familyFriendly[yes]* |
| **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 | *[0] name[The Phoenix]* |
|---|---|
| | *[1] area[riverside]* |
| | *[2] priceRange[cheap]* |
| | *[3] familyFriendly[yes]* |
| | *[4] near[Raja Indian Cuisine]* |
| | *[5] eatType[restaurant]* |
| **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.