

# E2E NLG Challenge: Neural Models vs. Templates

Yevgeniy Puzikov and Iryna Gurevych

Ubiquitous Knowledge Processing Lab (UKP-TUDA)  
Department of Computer Science, Technische Universität Darmstadt  
Research Training Group AIPHES  
[www.ukp.tu-darmstadt.de](http://www.ukp.tu-darmstadt.de)

## Abstract

E2E NLG Challenge is a shared task on generating restaurant descriptions from sets of key-value pairs. This paper describes the results of our participation in the challenge. We develop a simple, yet effective neural encoder-decoder model<sup>1</sup> which produces fluent restaurant descriptions and outperforms a strong baseline. We further analyze the data provided by the organizers and conclude that the task can also be approached with a template-based model developed in just a few hours.

## 1 Introduction

Natural Language Generation (NLG) is the task of generating natural language utterances from structured data representations. The E2E NLG Challenge<sup>2</sup> is a shared task which focuses on end-to-end data-driven NLG methods. These approaches attract a lot of attention, because they perform joint learning of textual structure and surface realization patterns from non-aligned data, which allows for a significant reduction of the amount of human annotation effort needed for NLG corpus creation (Wen et al., 2015; Mei et al., 2016; Dušek and Jurcicek, 2016; Lampouras and Vlachos, 2016).

The contribution of our submission to the challenge can be summarized as follows: (1) we show how exploiting data properties allows us to design more accurate neural architectures; (2) we develop a simple template-based system which achieves performance comparable to neural approaches.

<sup>1</sup><https://github.com/UKPLab/e2e-nlg-challenge-2017>

<sup>2</sup><http://www.macs.hw.ac.uk/InteractionLab/E2E>

## MR:

<i>name</i> [The Eagle]	<i>eatType</i> [coffee shop]
<i>food</i> [French]	<i>priceRange</i> [moderate]
<i>customerRating</i> [3/5]	<i>area</i> [riverside]
<i>kidsFriendly</i> [yes]	<i>near</i> [Burger King]

## Human Natural Language Reference:

*“The three star coffee shop, The Eagle, gives families a mid-priced dining experience featuring a variety of wines and cheeses. Find The Eagle near Burger King.”*

Figure 1: E2E NLG Challenge data specification.

## 1.1 Task Definition

The organizers of the shared task provided a crowd-sourced data set of 50k instances in the restaurant domain (Novikova et al., 2017b). Each training instance consists of a dialogue act-based meaning representation (MR) and up to 16 references in natural language (Figure 1).

The data was collected using pictorial representations as stimuli, with the intention of creating more natural, informative and diverse human references compared to the ones one might generate from textual inputs.

The task is to generate an utterance from a given MR, which is both similar to human-generated reference texts and highly rated by humans. Similarity is assessed using standard evaluation metrics: BLEU (Papineni et al., 2002), NIST (Dodington, 2002), METEOR (Lavie and Agarwal, 2007), ROUGE-L (Lin, 2004), CIDEr (Vedantam et al., 2015). However, the final assessment is done via human ratings obtained using a mixture of crowd-sourcing and expert annotations.

To facilitate a better assessment of the proposed approaches, the organizing team used TGen (Dušek and Jurcicek, 2016), one of the recent E2E data-driven systems, as a baseline. It is a sequence-to-sequence neural system with attention (Bahdanau et al., 2014). TGen uses beam search for decod-

ing, incorporates a reranker over the top  $k$  outputs, penalizing the candidates that do not verbalize all attributes from the input MR. TGen also includes a delexicalization module which deals with sparsely occurring MR attributes (*name*, *near*) by mapping such values to placeholder tokens when preprocessing the input data, and substituting the placeholders with actual values as a post-processing step.

## 2 Our Approach

This section describes two different approaches we developed for the shared task.

The first one (Model-D, for “data-driven”) is an encoder-decoder neural system which is similar to TGen, but uses a more efficient encoder module. The second approach is a simple template-based model (Model-T, for “template-based”) which we developed based on the results of the data analysis.

### 2.1 Model-D

Model-D was motivated by two important properties of the E2E NLG Challenge data:

- fixed number of unique MR attributes
- low diversity of the lexical instantiations of the MR attribute values

Each input MR contains a fixed number of unique attributes (between three and eight), which allows us to associate a positional id with each attribute and omit the corresponding attribute names (or keys) from the encoding procedure. This shortens the encoded sequence, presumably making the learning procedure easier for the encoder. This also unifies the lengths of input MRs and thus allows us to use simpler and more efficient neural networks which are not sequential and process input sequences in one step (e.g. multilayer perceptron (MLP) networks).

One might argue that using an MLP would be complicated by the fact that neither the number of active (non-null value) input MR keys nor the number of tokens constituting the corresponding values is fixed. For example, an MR key *price* may have a one-token value of “low” or a more lengthy “less than £10”. However, realizations of the MR attribute values exhibit low variability: six out of eight keys have less than seven unique values, while the remaining two keys (*name*, *near*) denote named entities and thus are easy to delexicalize. This allows us to treat each value as a single token,

<i>posID</i>	<i>Key</i>	<i>Value</i>
1	<i>area</i>	<i>PAD</i>
2	<i>customerRating</i>	<i>high</i>
3	<i>eatType</i>	<i>PAD</i>
4	<i>familyFriendly</i>	<i>yes</i>
5	<i>food</i>	<i>PAD</i>
6	<i>name</i>	<i>Wrestlers</i>
7	<i>near</i>	<i>PAD</i>
8	<i>priceRange</i>	<i>PAD</i>

Table 1: Input representation of the running example using positional ids.

even if it consists of multiple words (e.g. “more than £30”, “Fast food”).

Each predicted output is a textual description of a restaurant. As reported by Novikova et al. (2017b), the average number of words per reference is 20.1. We used the value of 50 as a cut-off threshold, filtering out training instances with long restaurant descriptions.

The overall architecture of our model is shown in Figure 2. The system is an encoder-decoder model (Cho et al., 2014b; Sutskever et al., 2014) consisting of three main modules: an embedding matrix, one dense hidden layer as an encoder and a RNN-based decoder with gated recurrent units (GRU) (Cho et al., 2014a).

Let us first describe the input specifications of the model. We will use the following MR instance as a running example:

*name*[*Wrestlers*]    *customerRating*[*high*]  
*familyFriendly*[*yes*]

Considering the alphabetic ordering of the MR key names, we can assign positional ids to the keys as shown in Table 1. The remaining five keys are assigned dummy *PAD* values.

Given an instance of a (*MR*, *text*) pair, we decompose the MR into eight components ( $mr_j$  in Figure 2), each corresponding to a value for a unique MR key, and add an end-of-sentence symbol (*EOS*) to denote the end of the encoded sequence. Each component is represented as a high-dimensional embedding vector. Each embedding vector is further mapped to a dense hidden representation via an affine transformation followed by a ReLu (Nair and Hinton, 2010) function. These hidden representations are further used by the decoder network, which in our case is a unidirectional GRU-based RNN with an attention module (Bahdanau et al., 2014). The decoder is initialized with an average of the encoder outputs.

The decoder generates a sequence of tokens, one

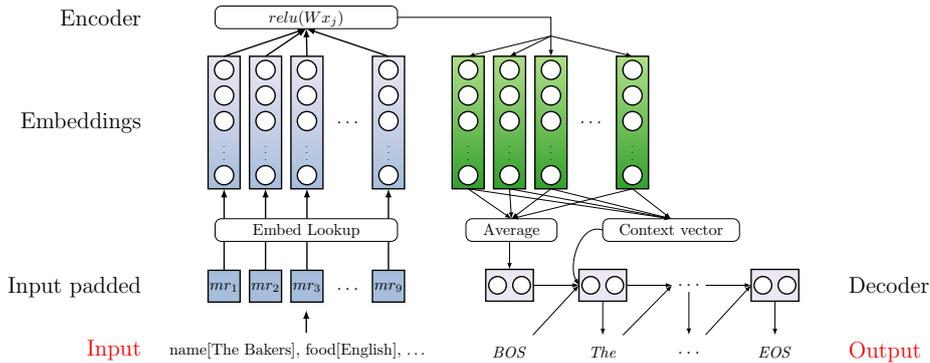


Figure 2: Schematic view of the neural network architecture (Model-D).

token at a time, until it predicts the *EOS* token. Our model employs the greedy search decoding strategy and does not use any reranker module.

## 2.2 Model-T

Taking into consideration low lexical variation of the MR attribute values, one might be interested in whether it is possible to design a deterministic NLG system to tackle the task. We examined the ways MR attribute keys and values are verbalized in the training data and discovered that the majority of textual descriptions follow a similar ordering of MR attribute verbalizations:

```
[name] is a [familyFriendly] [eatType]
which serves [food] food in the [price]
price range. It has a [customerRating]
customer rating. It is located in the
[area] area, near [near].
```

Here  $[X]$  denotes the value of the MR key  $X$ . This pattern became a central template of Model-T. Not all MR attribute verbalizations fit into this schema. For example, a key-value pair *customerRating*[3 out of 5] would be verbalized as "... has a 3 out of 5 customer rating", which is not the best phrasing one can come up with. A better way to describe it is "... has a customer rating of 3 out of 5". We incorporate such variations into Model-T with a set of simple rules which modify the general template depending on a specific value of an MR attribute.

As mentioned in Section 2.1, each instance’s input can have up to eight MR attributes. In order to account for this fact, we decomposed the general template into smaller components, each corresponding to a specific MR attribute mentioned in the input. We further developed a set of rules which activate each component depending on whether an MR attribute is part of the input. For example, if

Metric	TGen	Model-D	Model-T
BLEU	0.6925	<b>0.7128</b> $\pm$ 0.013	0.6051
NIST	8.4781	<b>8.5020</b> $\pm$ 0.092	7.5257
CIDEr	2.3987	<b>2.4432</b> $\pm$ 0.088	1.6997
ROUGE-L	0.7257	<b>0.7378</b> $\pm$ 0.015	0.6890
METEOR	0.4703	<b>0.4770</b> $\pm$ 0.012	0.4678

Table 2: Evaluation results according to automatic metrics (development set).

*price* is not in the set of input MR attributes, then the general template becomes:

```
[name] is a [familyFriendly] [eatType]
which serves [food] food. It has a
[customerRating] customer rating.
It is located in the [area] area,
near [near].
```

Finally, we also add a simple post-processing step to handle specific punctuation and article choices.

## 3 Metric Evaluation

Table 2 shows the results of metric evaluation of the systems. Since we were provided with only one TGen prediction file and a single performance score, comparing score distributions is not possible and statistical significance tests are not meaningful due to the non-deterministic nature of the approaches based on neural networks and randomized training procedures (Reimers and Gurevych, 2017). In order to facilitate a fair comparison with other competing systems, we report the mean development score of Model-D (averaged across twenty runs with different random seeds) and performance variance for each automatic metric. Model-T is a deterministic system, so it is sufficient to report the results of a single run.

The results show that Model-D outperforms

Error type	TGen	Model-D	Model-T
dropped contents	9	<b>49</b>	0
punctuation errors	1	<b>12</b>	0
modified contents	4	4	0
bad grammar	4	1	0

Table 3: Common errors made by the compared models (100 randomly sampled development instances).

TGen as measured by all five metrics, albeit the performance variance is quite large. Model-T clearly scores below both TGen and Model-D. This is expected, since Model-T is not data-driven, and hence the texts it generates might be different from the reference outputs.

Previous studies have shown that widely used automatic metrics (including the ones used in our competition) lack strong correlation with human judgments (Scott and Moore, 2007; Reiter and Belz, 2009; Novikova et al., 2017a). We decided to examine the predictions made by the compared systems on one hundred randomly sampled input instances, focusing on generic errors, which make sense to look out for in many NLG scenarios. Table 3 shows the error types and the number of mistakes found in each of the prediction files. The error types should be self-explanatory (sample predictions are given in Appendix A.2).

As far as the (subjective) manual analysis goes, Model-T outputs descriptions with the best linguistic quality. Table 3 shows that the predictions of the template-based system contain no errors – this is because we incorporated our notion of grammaticality into the templates’ definition, which allowed Model-T to avoid the errors found in predictions of the other two approaches.

The majority of errors made by Model-D are either wrong verbalizations of the input MR values or punctuation mistakes. The latter ones are limited to the cases of missing a comma between clauses or not finishing a sentence with a full stop. An easy solution to this problem is adding a post-processing step which fixes punctuation mistakes before outputting the text.

Crucially, Model-D often drops or modifies some MR attribute values. According to the organizers, 40% of the data by design contain either additional or omitted information on the output side (Novikova et al., 2017b): crowd workers were allowed to not lexicalize attribute values which they deemed unimportant. We decided to examine the

training data and find out if the discrepancies of Model-D were learned from the data.

## 4 Training Data Analysis

The E2E NLG Challenge is based on noisy data, but the organizers provided multiple instances to account for this noise. In order to better understand the behaviour of Model-D and determine if it took advantage of having multiple references per training instance, we have randomly sampled a hundred training instances and manually checked their linguistic quality. Table 4 shows the most common errors we encountered.

Most mistakes come from ungrammatical constructions, e.g. incorrect phrase attachment decisions (“The price of the food is high and is located ...”), incorrect usage of articles (“located in riverside”), repetitive constructions (“Cotto, an Indian coffee shop located in ..., is an Indian coffee shop ...”). Some restaurant descriptions follow a tweet-style narration pattern which is understandable, but ungrammatical (“The Golden Palace Italian riverside coffee shop price range moderate and customer rating 1 out of 5”).

A considerable number of instances have restaurant descriptions which contain information that does not entirely follow from the given input MR. These are cases in which input content elements are modified or dropped, which goes in line with what we observed in the outputs of Model-D.

Some instances (10%) contained descriptions which we marked as questionable due to pragmatic and/or stylistic considerations. For example, restaurants which have *familyFriendly[no]* as part of the input MR are often described by crowd workers as “adults-only” establishments, which has an undesirable connotation. Finally, it is necessary to mention that some crowd workers followed inconsistent spelling and punctuation rules when hyphenating compound modifiers (“family friendly restaurant”, “the restaurant is family friendly”) or capitalizing MR attributes (“Riverside”, “Fast food”). Punctuation errors were mainly restricted to missing a full stop at the end of a restaurant description or failing to delimit sentence clauses with commas.

The results of manual data analysis show that Model-D indeed generates texts that are similar to the restaurant descriptions in the provided data set. Unfortunately, our data-driven approach is not flexible enough to make use of multiple references; it cannot cancel out the noise present in some train-

Error type	Example	%
bad grammar	“it’s French food falls within a high price range”	15
modified contents	<i>area[riverside]</i> → “city centre”	12
dropped contents	<i>priceRange[high]</i> → $\emptyset$	10
questionable lexicalization	“Adult-only Chinese restaurant, The Waterman, offers top-rated food in the city centre”	9
punctuation errors	“X is a coffee shop and also a Japanese restaurant great for family and close to Crowne Plaza Hotel”	6

Table 4: Data annotation discrepancies (100 randomly sampled training instances).

	Model-T	Best result
<b>Metric evaluation</b>		
BLEU	0.5657	0.6805
NIST	7.4544	8.7777
METEOR	0.4529	0.4571
ROUGE-L	0.6614	0.7084
CIDEr	1.8206	2.3371
<b>Human evaluation</b>		
Quality	0.228/(2.0, 4.0)/2	0.300/(1.0, 1.0)/1
Naturalness	0.077/(5.0, 10.0)/2	0.211/(1.0, 1.0)/1

Table 5: Final evaluation results on the test set. Human evaluation results have the following format: *score/(range)/cluster*.

ing instances. One way of alleviating this problem could be reformulating the loss function to inform the system about the existence of multiple ways of generating a good restaurant description. Given a training instance, Model-D would generate a corresponding candidate text which could be compared to all human references. Each comparison results in computing a certain cost; the gradients could be then computed on the minimal cost among all comparisons.

#### 4.1 Final Evaluation

For the final submission we have chosen Model-T’s predictions – despite lower metric scores, they contained most grammatical outputs and kept all input information in the generated text.

The results of the final evaluation on the test data are presented in Table 5. They were produced by the TrueSkill algorithm (Sakaguchi et al., 2014), which performs pairwise system comparisons and clusters them into groups. For completeness, we include the highest reported scores among all the participants (rightmost column). Note, however, that the numerical scores are not directly interpretable, but the relative ranking of a system in terms of its range and cluster is important – systems within one cluster are considered tied.

Model-T was assigned to the second best cluster both in terms of quality and naturalness, despite the much lower metric scores. Retrospectively, this justifies our decision to choose Model-T instead of Model-D for the final submission. The E2E NLG Challenge focuses on end-to-end data-driven NLG methods, which is why systems like Model-T might not exactly fit into the task setup. Nevertheless, we view such a system as a necessary candidate for comparison, since the E2E NLG Challenge data was designed to learn models that produce “more natural, varied and less template-like system utterances” (Novikova et al., 2017b).

## 5 Conclusion

In this paper we have presented the results of our participation in the E2E NLG Challenge. We have developed two conceptually different approaches and analyzed their performance, both in quantity and in quality. We have shown that sometimes the costs of developing complex data-driven models are not justified and one is better off approaching the problem with simpler techniques. We hope that our observations and conclusions shed some light on the limitations of modern NLG approaches and possible ways of overcoming them.

## Acknowledgments

This work was supported by the German Research Foundation (DFG) under grant No. GU 798/17-1 and the DFG-funded research training group “Adaptive Preparation of Information from Heterogeneous Sources” (AIPHES, GRK 1994/1). The first author of the paper is supported by the FAZIT Foundation scholarship. We thank the anonymous reviewers and our colleagues Michael Bugert, Tristan Miller, Maxime Peyrard, Nils Reimers and Markus Zopf who provided insightful comments and suggestions that greatly assisted our research.

## References

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. [Neural Machine Translation by Jointly Learning to Align and Translate](#). *CoRR*, abs/1409.0473.
- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014a. [On the Properties of Neural Machine Translation: Encoder-Decoder Approaches](#). In *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, pages 103–111, Doha, Qatar. Association for Computational Linguistics.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014b. [Learning Phrase Representations Using RNN Encoder-decoder for Statistical Machine Translation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.
- George Doddington. 2002. [Automatic Evaluation of Machine Translation Quality Using N-gram Co-occurrence Statistics](#). In *Proceedings of the Second International Conference on Human Language Technology Research*, pages 138–145, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Ondřej Dušek and Filip Jurcicek. 2016. [Sequence-to-sequence Generation for Spoken Dialogue via Deep Syntax Trees and Strings](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 45–51, Berlin, Germany. Association for Computational Linguistics.
- Gerasimos Lampouras and Andreas Vlachos. 2016. [Imitation Learning for Language Generation from Unaligned Data](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1101–1112, Osaka, Japan. The COLING 2016 Organizing Committee.
- Alon Lavie and Abhaya Agarwal. 2007. [METEOR: An Automatic Metric for MT Evaluation with High Levels of Correlation with Human Judgments](#). In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 228–231, Prague, Czech Republic. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. [ROUGE: a Package for Automatic Evaluation of Summaries](#). In *Text Summarization Branches Out: Proceedings of the ACL-04 Workshop*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Hongyuan Mei, Mohit Bansal, and Matthew R. Walter. 2016. [What to Talk About and How? Selective Generation Using LSTMs with Coarse-to-fine Alignment](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 720–730, San Diego, California. Association for Computational Linguistics.
- Vinod Nair and Geoffrey E. Hinton. 2010. [Rectified Linear Units Improve Restricted Boltzmann Machines](#). In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pages 807–814, USA. Omnipress.
- Jekaterina Novikova, Ondřej Dušek, Amanda Cercas Curry, and Verena Rieser. 2017a. [Why We Need New Evaluation Metrics for NLG](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2231–2242, Copenhagen, Denmark. Association for Computational Linguistics.
- Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2017b. [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. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [BLEU: a Method for Automatic Evaluation of Machine Translation](#). In *Proceedings of 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2017. [Reporting Score Distributions Makes a Difference: Performance Study of LSTM-networks for Sequence Tagging](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 338–348, Copenhagen, Denmark. Association for Computational Linguistics.
- Ehud Reiter and Anja Belz. 2009. [An Investigation into the Validity of Some Metrics for Automatically Evaluating Natural Language Generation Systems](#). *Computational Linguistics*, 35(4):529–558.
- Keisuke Sakaguchi, Matt Post, and Benjamin Van Durme. 2014. [Efficient Elicitation of Annotations for Human Evaluation of Machine Translation](#). In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 1–11, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Donia Scott and Johanna Moore. 2007. [An NLG Evaluation Competition? Eight Reasons to Be Cautious](#). In *Proceedings of the NSF Workshop on Shared Tasks and Comparative Evaluation in Natural Language Generation*, pages 22–23. National Science Foundation.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. [Sequence to Sequence Learning with Neural Networks](#). In Z. Ghahramani, M. Welling, C. Cortes,

N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 3104–3112. Curran Associates, Inc.

Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. [Cider: Consensus-based image description evaluation](#). In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4566–4575. IEEE Computer Society.

Tsung-Hsien Wen, Milica Gasic, Nikola Mrkšić, Pei-Hao Su, David Vandyke, and Steve Young. 2015. [Semantically Conditioned LSTM-based Natural Language Generation for Spoken Dialogue Systems](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1711–1721, Lisbon, Portugal. Association for Computational Linguistics.

## A Supplemental Material

### A.1 Manual Data Analysis Results

As mentioned in Section 4, manual analysis of the training data revealed certain annotation problems. Below we provide sentence IDs of the instances which we considered as containing errors:

- modified contents: 4136, 34141, 32915, 35936, 6152, 2005, 1463, 14529, 14339, 21804, 25779, 11049;
- dropped contents: 4136 (price and food), 17455 (familyFriendly), 38742 (area), 1463 (customerRating), 27664 (food, priceRange), 19882 (priceRange), 25779 (customerRating, area), 40293 (food), 807 (familyFriendly), 15216 (familyFriendly);
- ungrammatical: 39810, 1904, 36233, 33671, 28698;
  - tweet-style descriptions: 19425, 38875, 14873;
  - word attachment problems, bad phrasing: 15276, 30950, 31973, 36468, 12442, 18623, 25773;
- punctuation errors: 6152, 14993, 38742, 14339, 35967, 25795;
- questionable realization: customerRating (8806, 24983, 29189), familyFriendly (20220, 10729), priceRange (29189), style (14993, 28362, 32494).

### A.2 Example Predictions

The majority of predictions from all three systems are grammatically correct. Here is a random sample:

- **Input:** name[Cocum], eatType[coffee shop], food[English], priceRange[moderate], customerRating[3 out of 5], familyFriendly[no]
- **TGen:** Cocum is a moderately priced coffee shop serving english food. It has a customer rating of 3 out of 5 and is not kid friendly.
- **Model-D:** Cocum is a moderately priced coffee shop serving English food . It is not kid friendly and has a customer rating of 3 out of 5 .
- **Model-T:** Cocum is a coffee shop which serves English food in the moderate price range. It has a customer rating of 3 out of 5. It is not family friendly.

To illustrate common prediction mistakes (regardless of the system), we have hand-picked some samples. Full predictions by all models can be found in the GitHub repository.

#### A.2.1 Bad Grammar

- **Input:** name[The Wrestlers], eatType[coffee shop], food[Chinese], priceRange[less than £20], area[riverside], familyFriendly[no], near[Raja Indian Cuisine]
  - **TGen:** The Wrestlers is a coffee shop providing chinese food in the less than £20 price range. It is located in the riverside. It is near Raja Indian Cuisine. Its no children feel allowed.
  - **Model-D:** The Wrestlers is a coffee shop providing Chinese food in the less than £20 price range . It is located in the riverside . It is near Raja Indian Cuisine .
  - **Model-T:** The Wrestlers is a coffee shop which serves Chinese food in the price range of less than £20. It is located in the riverside area, near Raja Indian Cuisine. It is not family friendly.
- 
- **Input:** name[Wildwood], eatType[coffee shop], food[English], priceRange[high], customerRating[1 out of 5], near[Ranch]
  - **TGen:** Wildwood is a high priced coffee shop that serves english food. It is located near Ranch and has a customer rating of 1 out of 5.
  - **Model-D:** Wildwood is a coffee shop near Ranch , serving English food , has a high price range and a customer rating of 1 out of 5 .
  - **Model-T:** Wildwood is a coffee shop which serves English food in the high price range. It has a customer rating of 1 out of 5 and is located near Ranch.

### A.2.2 Modified and/or Dropped Contents

- **Input:** name[The Golden Palace], eatType[coffee shop], food[English], priceRange[high], customerRating[average], area[riverside]
- **TGen:** The Golden Palace is a high priced coffee shop serving english food in the riverside area. It has an average customer rating.
- **Model-D:** The Golden Palace is a high priced coffee shop located in the city centre .
- **Model-T:** The Golden Palace is a coffee shop which serves English food in the high price range. It has an average customer rating and is located in the riverside area.

### A.2.3 Dropped Contents

- **Input:** name[The Eagle], eatType[coffee shop], food[English], priceRange[moderate], customerRating[1 out of 5], area[riverside], familyFriendly[yes], near[Burger King]
- **TGen:** The Eagle is a moderately priced coffee shop that serves english food. It is located in the riverside area near Burger King. It is kid friendly and has a customer rating of 1 out of 5.
- **Model-D:** The Eagle is a moderately priced coffee shop located near Burger King .
- **Model-T:** The Eagle is a family-friendly coffee shop which serves English food in the moderate price range. It has a customer rating of 1 out of 5 and is located in the riverside area, near Burger King.