

Hello, It's GPT-2 - How Can I Help You?

Towards the Use of Pretrained Language Models for Task-Oriented Dialogue Systems

Paweł Budzianowski^{1,2,3} and Ivan Vulić^{2,3}

¹Engineering Department, Cambridge University, UK

²Language Technology Lab, Cambridge University, UK

³PolyAI Limited, London, UK

pfb30@cam.ac.uk, iv250@cam.ac.uk

Abstract

Data scarcity is a long-standing and crucial challenge that hinders quick development of task-oriented dialogue systems across multiple domains: task-oriented dialogue models are expected to learn grammar, syntax, dialogue reasoning, decision making, and language generation from absurdly small amounts of task-specific data. In this paper, we demonstrate that recent progress in language modeling pre-training and transfer learning shows promise to overcome this problem. We propose a task-oriented dialogue model that operates solely on text input: it effectively bypasses explicit policy and language generation modules. Building on top of the TransferTransfo framework (Wolf et al., 2019) and generative model pre-training (Radford et al., 2019), we validate the approach on complex multi-domain task-oriented dialogues from the MultiWOZ dataset. Our automatic and human evaluations show that the proposed model is on par with a strong task-specific neural baseline. In the long run, our approach holds promise to mitigate the data scarcity problem, and to support the construction of more engaging and more eloquent task-oriented conversational agents.

1 Introduction

Statistical conversational systems can be roughly clustered into two main categories: 1) task-oriented modular systems and 2) open-domain chit-chat neural models. The former typically consist of independently trained constituent modules such as language understanding, dialogue management, and response generation. The main goal of such systems is to provide meaningful system responses which are invaluable in building conversational agents of practical value for restricted domains and tasks. However, data collection and annotation for such systems is complex, time-intensive, expensive, and not easily transferable

(Young et al., 2013). On the other hand, open-domain conversational bots (Li et al., 2017; Serban et al., 2017) can leverage large amounts of freely available unannotated data (Ritter et al., 2010; Henderson et al., 2019a). Large corpora allow for training end-to-end neural models, which typically rely on sequence-to-sequence architectures (Sutskever et al., 2014). Although highly data-driven, such systems are prone to producing unreliable and meaningless responses, which impedes their deployment in the actual conversational applications (Li et al., 2017).

Due to the unresolved issues with the end-to-end architectures, the focus has been extended to retrieval-based models. Here, the massive datasets can be leveraged to aid task-specific applications (Kannan et al., 2016; Henderson et al., 2017, 2019b). The retrieval systems allow for the full control over system responses, but the behaviour of the system is often highly predictable. It also depends on the pre-existing set of responses, and the coverage is typically insufficient for a multitude of domains and tasks. However, recent progress in training high-capacity language models (e.g., GPT, GPT-2) (Radford et al., 2018, 2019) on large datasets reopens the question of whether such generative models can support task-oriented dialogue applications. Recently, Wolf et al. (2019) and Golovanov et al. (2019) showed that the GPT model, once fine-tuned, can be useful in the domain of personal conversations. In short, their approach led to substantial improvements on the Persona-Chat dataset (Zhang et al., 2018), showcasing the potential of exploiting large pretrained generative models in the conversational domain.¹

In this paper, we demonstrate that large generative models pretrained on large general-domain

¹E.g., TransferTransfo (Wolf et al., 2019) yields gains in all crucial dialogue evaluation measures such as fluency, consistency and engagingness on the Persona-Chat dataset.

<p>Dialogue Context: Visitor: Hi I need a help with finding a hotel Belief state: hotel() Database state: hotel(many)</p> <p style="text-align: right;">System: Hello, any requirements?</p> <p>Visitor: I need to find a luxury hotel in the centre please. Belief state: hotel(price-expensive;area-centre) Database state: hotel(three)</p> <p>Dialogue-Context-to-Text: System: <restaurant_name> should fit you.</p>

Figure 1: Dialogue-context-to-text task.

corpora can support *task-oriented dialogue applications*. We first discuss how to combine a set of diverse components such as word tokenization, multi-task learning, and probabilistic sampling to support task-oriented applications. We then show how to adapt the task-oriented dialogue framework to operate entirely on text input, effectively bypassing an explicit dialogue management module and a domain-specific natural language generation module. The proposed model operates entirely in the sequence-to-sequence fashion, consuming only simple text as input. The entire dialogue context, which includes the belief state, the database state and previous turns, is provided to the decoder as raw text. The proposed model follows the recently proposed TransferTransfo framework (Wolf et al., 2019), and relies on pretrained models from the GPT family (Radford et al., 2018, 2019).

Our results in the standard Dialogue-Context-to-Text task (see Figure 1) on the multi-domain MultiWOZ dataset (Budzianowski et al., 2018b) suggest that our GPT-based task-oriented dialogue model learns to generate and understand domain-specific tokens, which in turn leads to a seamless adaptation to particular focused domains. While automatic evaluation indicates that our framework still falls slightly short of a strong task-specific neural baseline, it also hints at the main advantage of our framework: it is widely portable and easily adaptable to a large number of domains, bypassing the intricate modular design only at a small cost in performance. Furthermore, user-centered evaluations suggest that there is no significant difference between the two models.

2 From Unsupervised Pretraining to Dialogue Modeling

Task-oriented dialogue modeling requires substantial amounts of domain-specific manually labeled

data. A natural question to ask is: Can we leverage transfer learning through generative pretraining on large unlabelled corpora to enable task-oriented dialogue modeling. In this work, we rely on the standard language modeling (LM) pretraining, where the task is to predict the next word given the preceding word sequence (Bengio et al., 2003). The objective maximizes the likelihood over the word sequence $S = \{w_1, \dots, w_{|S|}\}$:

$$\mathcal{L}_1(S) = \sum_{i=1}^{|S|} \log P(w_i | w_0, w_1, \dots, w_{i-1}). \quad (1)$$

Transfer learning based on such LM pretraining combined with the Transformer decoder model (Vaswani et al., 2017) resulted in significant progress across many downstream tasks (Rei, 2017; Howard and Ruder, 2018; Radford et al., 2018, 2019).

2.1 TransferTransfo Framework

Golovanov et al. (2019) and Wolf et al. (2019) achieved a first successful transfer of a generative pretrained GPT model to an open-domain dialogue task. The pretrained GPT model is fine-tuned in a multi-task learning fashion following the original work (Radford et al., 2018). The LM objective from Eq. (1) is combined with the next utterance classification task:

$$p(c, a) = \text{softmax}(h_l * W_h). \quad (2)$$

c and a represent the context of the conversation (c) and a proposed answer (a), h_l is the last hidden state of the transformer decoder, and W_h is learnt during the fine-tuning phase. The model significantly improves upon previous baselines over all automatic dialogue evaluation metrics as well as in evaluation with human subjects when evaluated on the Persona-Chat dataset (Zhang et al., 2018).

The GPT input consists of token embeddings and positional embeddings. In order to move from a single-speaker setting to a setting with two interlocutors, Wolf et al. (2019) introduced *dialogue-state* embeddings. These embeddings inform the model whether the current token comes from an utterance of the first speaker or an utterance of the second speaker. The dialogue-state embeddings are learned during the fine-tuning phase.

3 Domain Transfer for (Task-Oriented) Dialogue Modeling

We now briefly discuss several advances in modeling of natural language that facilitate applicability

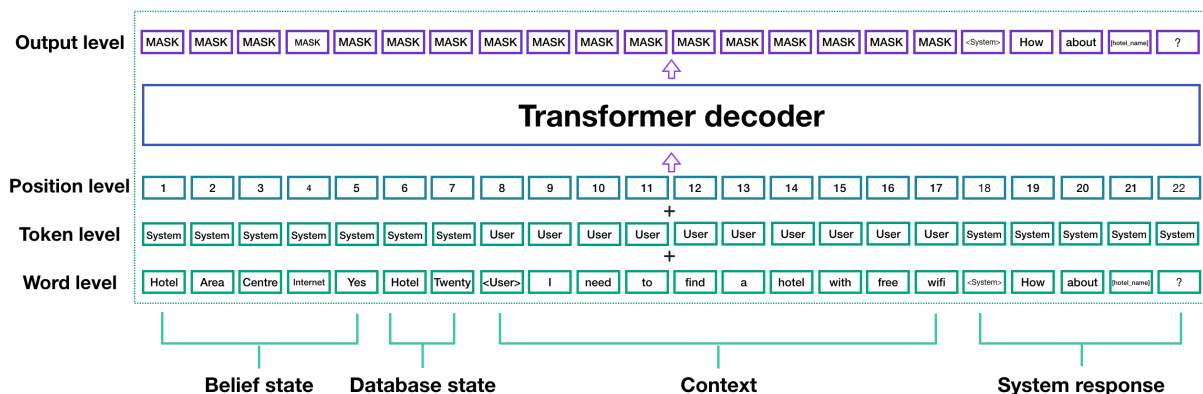


Figure 2: The framework for modeling task-oriented conversations based on a pretrained GPT model which uses only unstructured simple text as input. The context, belief state, and database state are joined together without explicit standalone dialogue policy and generation modules. The token-level (i.e., dialogue-state) embeddings are learned following Wolf et al. (2019).

of pretrained generative models in task-oriented dialogue modeling. To the best of our knowledge, this work is first to combine these existing components to enable task-oriented dialogue modeling.

3.1 Domain Adaptation and Delexicalization

Dealing with out-of-vocabulary (OOV) words has been a long-standing challenge in dialogue modeling, e.g., it is crucial for task-oriented generation where the generated output is often *delexicalized* (Wen et al., 2015). Delexicalization replaces slot values by their corresponding (generic) slot tokens and it allows learning value-independent parameters. Recently, owing to subword-level tokenisation (Sennrich et al., 2016), language models are now able to deal with OOVs and domain-specific vocabularies more effectively (Radford et al., 2018).

3.2 Simple Text-Only Input

There have been some empirical validations recently which suggest that posing NLP tasks in the form of simple text can yield improvements with unsupervised architectures (Wolf et al., 2019; Radford et al., 2019). For instance, in task-oriented dialogue modeling the Sequicity model (Lei et al., 2018) sees the classification over the belief state as a generation problem. That way, the entire dialogue model pipeline is based on the sequence-to-sequence architecture: the output from one model is the input to the subsequent recurrent model. We follow this approach by providing both the belief state and the knowledge base state in a simple text format to the generator. This significantly simplifies the paradigm of building task-oriented models: any new source of information can be simply

added to as another part of the text-only input provided in “natural language”.

3.3 Transferring Language Generation Capabilities

Transformer architecture shows ability to learn new (i.e., domain-specific) token embeddings in the fine-tuning phase (Radford et al., 2018; Wolf et al., 2019). This means that the GPT models can adapt through special tokens to particular tasks. By providing the input representation as text with domain-specific tokens, we can use off-the-shelf architectures and adapt to the domain-specific input without the need of training new dialogue sub-modules. As mentioned in §2.1, the token level layer (Figure 2) informs the transformer decoder what part of the input comes from the system side or from the user side. In our framework, we create two task-oriented specific tokens (*System* and *User* tokens) that are learned during fine-tuning.

3.4 Generation Quality

Finally, the long-standing problem of dull and repetitive response generation (Li et al., 2017) has been in the focus of recent work (Kulikov et al., 2018; Holtzman et al., 2019). Owing to new sampling strategies, generative models are now able to create longer and more coherent sequence outputs. This has been validated also for open-domain dialogue modeling (Wolf et al., 2019; Golovanov et al., 2019). We experiment with standard decoding strategies as well as with the recently proposed *nucleus* sampling procedure (Holtzman et al., 2019). A standard *greedy sampling* strategy

chooses the most probable word as :

$$\arg \max_{w_i} = \log P(w_i | w_0, w_1, \dots, w_{i-1}).$$

On the other hand, nucleus sampling is restricted only to words from the p -th percentile of the distribution during generation. The probabilities of words for which the cumulative sum exceeds the percentile are rescaled and the sequence is sampled from this subset. We probe the ability of such large pretrained models to generate more varied and semantically richer responses relying on nucleus sampling in lieu of greedy sampling without hurting the actual performance.

4 Fine-Tuning GPT on MultiWOZ

To evaluate the ability of transferring the GPT generation capability to constrained/focused dialogue tasks and domains, we rely on the multi-domain MultiWOZ dataset (Budzianowski et al., 2018b). MultiWOZ consists of 7 domains and 10,438 dialogues and it is substantially larger than previous available datasets (Wen et al., 2017; El Asri et al., 2017). The conversations are natural as they were gathered through human-human interactions. However, the dialogues are based on domain-specific vocabulary such as booking IDs or telephone numbers that need to be delexicalized as they are entirely database-dependent.

Natural Language as (the Only) Input. GPT operates solely on the text input. This is in opposition to the standard task-oriented dialogue architectures (Wen et al., 2017; Zhao et al., 2017) where the belief state and the database state are encoded in a numerical form. For example, the database state is typically defined as n -bin encodings representing a number of available entities at the current state of the conversation (Wen et al., 2017). Therefore, we transform the belief state and the knowledge base representation to a simple text representation. The belief state takes the following form:

```
Domain1 Slot1 Value1 Slot2 Value2
Domain2 Slot1 ...
```

and the database representation is provided as:

```
Domain1 # of entities
Domain2 # of entities ...
```

This is also similar in spirit to the Sequicity architecture (Lei et al., 2018) where the second recurrent model takes as input the belief state in the natural language (i.e., simple text-only) form. In this

work, we also transform the knowledge base state to a similar natural language format. These two pieces of information are then concatenated with the history of the conversation forming the full dialogue context, see Figure 2. Following Wolf et al. (2019), we add new token embeddings for two parties involved in the conversation to inform the attention layers what part of the context comes from the user, and what part is related to the system. Figure 2 presents the final architecture.

Training Details. We use the open-source implementation of the GPT architecture that provides both GPT and GPT-2 fine-tunable checkpoints.² Following previous work (Radford et al., 2018; Wolf et al., 2019), we set the weight on the language model loss to be two times higher than the one for the response prediction. The parameters for the batch size (24), learning rate (1e-5) and the number of candidates per sequence (2) were chosen based on the grid search.³

5 Results and Analysis

Following prior work (Budzianowski et al., 2018b; Zhao et al., 2019; Chen et al., 2019), our evaluation task is the dialogue-context-to-text generation task (see Figure 1). Given a dialogue history, the oracle belief state and the database state, the model needs to output the adequate response. By relying on the oracle belief state, prior work has bypassed the possible errors originating from natural language understanding (Budzianowski et al., 2018b).

The main evaluation is based on the comparison between the following two models: 1) the baseline is a neural response generation model with an oracle belief state obtained from the wizard annotations as in (Budzianowski et al., 2018a); 2) the model proposed in §4 and shown in Figure 2 that works entirely with text-only format as input (see §4). We test all three available pretrained GPT models - the original GPT model (Radford et al., 2018). and two GPT-2 models referred to as small (GPT2) and medium (GPT2-M) (Radford et al., 2019).

²<https://github.com/huggingface/transfer-learning-conv-ai>

³We searched over the following values: learning rates \in {1-e4, 1-e5, 5-e6, 1-e6}, batch sizes \in {8, 12, 16, 20, 24} and candidate set sizes \in {1, 2, 4, 6}.

	Baseline	GPT	GPT2-S	GPT2-M
Inform (%)	76.7	71.53	66.43	70.96
Success (%)	64.63	55.36	55.16	61.36
BLEU (%)	18.05	17.80	18.02	19.05

Table 1: Evaluation on MultiWOZ with the greedy sampling procedure.

	Baseline	GPT	GPT2-S	GPT2-M
Inform (%)	72.57	70.43	69.3	73.96
Success (%)	57.63	51.0	54.93	61.20
BLEU (%)	15.75	15.65	15.64	16.55

Table 2: Evaluation on MultiWOZ with the nucleus sampling procedure.

5.1 Evaluation with Automatic Measures

We report scores with three standard automatic evaluation measures. Two of them relate to the dialogue task completion: whether the system has provided an appropriate entity (*Inform*) and then answered all requested attributes (*Success* rate). Finally, fluency is measured by the BLEU score (Papineni et al., 2002).

First, three versions of GPT were fine-tuned on MultiWOZ and evaluated with greedy sampling. The results are summarized in Table 1). They show that the baseline obtains the highest score on task-related metrics while the highest BLUE score was achieved by GPT2-M. Although the results are lower for the GPT-based methods, we note the design simplicity of the GPT-based task-oriented dialogue models. Further, the gap in performance might be partially attributed to the chosen greedy sampling procedure which puts too much focus on the properties of the original pretraining phase (Holtzman et al., 2019).

Therefore, we also report the results with the nucleus sampling method in Table 2. The scores confirm the importance of choosing the correct sampling method. The GPT2 models improve the score on *Inform* and *Success* metrics. It is worth noting the consistent drop in BLUE scores across all models. This comes from the fact that nucleus sampling allows for increased variability: this might reduce the probability of generating domain-specific tokens.

We have also qualitatively analyzed a sample of successful dialogues. Only around 50% of dialogues are successful both with the baseline and with the GPT-based models. Moreover, there are no clearly observed distinct patterns between successful dialogues for the two model types. This

Model 1	vs		Model 2
GPT	59 %	41 %	Baseline
GPT	46 %	54 %	Target
GPT2	46 %	54 %	Target
GPT2	45 %	55 %	Baseline
Baseline	43 %	57 %	Target
GPT2	51 %	49 %	GPT

Table 3: Human ranking of responses between all pairs of four analyzed models and the original responses.

suggests that they might be effectively ensembled using a ranking model to evaluate the score of each response (Henderson et al., 2019b). We will investigate the complementarity of the two approaches along with ensemble methods in future work.

5.2 Human Evaluation

In another, now user-centered experiment, the goal was to analyze the generation quality. Turkers, native speakers of English, were asked to rate their binary preference when presented with one-turn responses from the baseline, GPT, GPT2-M and the original dialogues (*Target*). The turkers were required to choose what response they prefer when presented with two responses from two different models, resulting in more than 300 scores per each model pair.

The results are summarized in Table 3, while some example dialogues with responses are provided in Figure 3. As expected, the original responses are ranked higher than all neural models with the largest difference observed between the oracle and the baseline model. Although the generated output from the GPT is strongly preferred against the neural baseline, interestingly the opposite is observed with the GPT2 model. These inconclusive results call for further analyses in future work, and also show that there are no substantial differences in the quality of generated responses when comparing the strong neural baseline and the GPT-based models.

6 Conclusion

In this paper, we have made a first step towards leveraging large pretrained generative models for modeling task-oriented dialogue in multiple domains. The simplicity of the fine-tuning procedure where all necessary information can be encoded as simple text enables a quick adaptation to constrained domains and domain-specific vo-

cabularies. We hope that this framework will inform and guide future research in hope of simultaneously improving and simplifying the design of task-oriented conversational systems.

References

- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2003. [A neural probabilistic language model](#). *Journal of machine learning research*, 3(Feb):1137–1155.
- Paweł Budzianowski, Iñigo Casanueva, Bo-Hsiang Tseng, and Milica Gašić. 2018a. [Towards end-to-end multi-domain dialogue modelling](#). *Tech. Rep. CUED/F-INFENG/TR.706, University of Cambridge, Engineering Department*.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gasic. 2018b. [MultiWOZ-A Large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling](#). In *Proceedings of EMNLP*, pages 5016–5026.
- Wenhu Chen, Jianshu Chen, Pengda Qin, Xifeng Yan, and William Yang Wang. 2019. [Semantically conditioned dialog response generation via hierarchical disentangled self-attention](#). *arXiv preprint arXiv:1905.12866*.
- Layla El Asri, Hannes Schulz, Shikhar Sharma, Jeremie Zumer, Justin Harris, Emery Fine, Rahul Mehrotra, and Kaheer Suleman. 2017. [Frames: A corpus for adding memory to goal-oriented dialogue systems](#). In *Proceedings of SIGDIAL*, pages 207–219.
- Sergey Golovanov, Rauf Kurbanov, Sergey Nikolenko, Kyril Truskovskiy, Alexander Tselousov, and Thomas Wolf. 2019. [Large-scale transfer learning for natural language generation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6053–6058, Florence, Italy. Association for Computational Linguistics.
- Matthew Henderson, Rami Al-Rfou, Brian Strope, Yunhsuan Sung, Laszlo Lukacs, Ruiqi Guo, Sanjiv Kumar, Balint Miklos, and Ray Kurzweil. 2017. [Efficient natural language response suggestion for smart reply](#). *arXiv preprint arXiv:1705.00652*.
- Matthew Henderson, Paweł Budzianowski, Iñigo Casanueva, Sam Coope, Daniela Gerz, Girish Kumar, Nikola Mrkšić, Georgios Spithourakis, Pei-Hao Su, Ivan Vulić, and Tsung-Hsien Wen. 2019a. [A repository of conversational datasets](#). In *Proceedings of the Workshop on NLP for Conversational AI*. Data available at github.com/PolyAI-LDN/conversational-datasets.
- Matthew Henderson, Ivan Vulić, Daniela Gerz, Iñigo Casanueva, Paweł Budzianowski, Sam Coope, Georgios Spithourakis, Tsung-Hsien Wen, Nikola Mrkšić, and Pei-Hao Su. 2019b. [Training neural response selection for task-oriented dialogue systems](#). In *Proceedings of ACL*.
- Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin Choi. 2019. [The curious case of neural text degeneration](#). *arXiv preprint arXiv:1904.09751*.
- Jeremy Howard and Sebastian Ruder. 2018. [Universal language model fine-tuning for text classification](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339.
- Anjali Kannan, Karol Kurach, Sujith Ravi, Tobias Kaufmann, Andrew Tomkins, Balint Miklos, Greg Corrado, Laszlo Lukacs, Marina Ganea, Peter Young, et al. 2016. [Smart reply: Automated response suggestion for email](#). In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 955–964. ACM.
- Ilya Kulikov, Alexander H Miller, Kyunghyun Cho, and Jason Weston. 2018. [Importance of a search strategy in neural dialogue modelling](#). *arXiv preprint arXiv:1811.00907*.
- Wenqiang Lei, Xisen Jin, Min-Yen Kan, Zhaochun Ren, Xiangnan He, and Dawei Yin. 2018. [Sequicity: Simplifying task-oriented dialogue systems with single sequence-to-sequence architectures](#). In *Proceedings of ACL*, pages 1437–1447.
- Jiwei Li, Will Monroe, Tianlin Shi, Sébastien Jean, Alan Ritter, and Dan Jurafsky. 2017. [Adversarial learning for neural dialogue generation](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2157–2169.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [BLEU: A method for automatic evaluation of machine translation](#). In *Proceedings of ACL*, pages 311–318.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. [Improving language understanding by generative pre-training](#).
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#). *OpenAI Blog*, 1(8).
- Marek Rei. 2017. [Semi-supervised multitask learning for sequence labeling](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2121–2130.
- Alan Ritter, Colin Cherry, and Bill Dolan. 2010. [Un-supervised modeling of twitter conversations](#). In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 172–180.

- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 1715–1725.
- Iulian Vlad Serban, Alessandro Sordani, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2017. [A hierarchical latent variable encoder-decoder model for generating dialogues](#). In *Proceedings of AAAI*, pages 3295–3301.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. [Sequence to sequence learning with neural networks](#). In *Proceedings of NeurIPS*, pages 3104–3112.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in neural information processing systems*, pages 5998–6008.
- Tsung-Hsien Wen, Milica Gašić, 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 (EMNLP)*.
- Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, Milica Gašić, Lina M Rojas-Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. 2017. [A network-based end-to-end trainable task-oriented dialogue system](#). In *Proceedings of EACL*, pages 438–449.
- Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. 2019. [Transfertransfo: A transfer learning approach for neural network based conversational agents](#). *arXiv preprint arXiv:1901.08149*.
- Steve Young, Milica Gašić, Blaise Thomson, and Jason D. Williams. 2013. [POMDP-based statistical spoken dialog systems: A review](#). *Proceedings of the IEEE*, 101(5):1160–1179.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213.
- Tiancheng Zhao, Kaige Xie, and Maxine Eskenazi. 2019. [Rethinking action spaces for reinforcement learning in end-to-end dialog agents with latent variable models](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1208–1218.
- Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017. [Learning discourse-level diversity for neural dialog models using conditional variational autoencoders](#). In *Proceedings of ACL*, pages 654–664.

Dialogue 1:
Visitor: I need to book a taxi to arrive by [value_time] .

Baseline: Where will you be departing from ?
GPT: I would be happy to help with your request, first i'll need to know where you are going and what time you want to leave or arrive by.

Visitor: I need to be picked up from [restaurant_name] and i am headed to the [attraction_name] .

Baseline: Your taxi has been booked . the car will be a [taxi_type] and the contact number is [taxi_phone].
GPT: Your taxi is booked. You will be picked up at [restaurant_name] and the contact number is [taxi_phone]. Is there anything else i can help you with?
...

Dialogue 2:
Visitor: I am looking for some places to go, do you have any attractions in the [value_area] ?

Baseline: There are [value_count] attractions in the [value_area]. Do you have a particular type of attraction in mind?
GPT: Sure! There are many colleges and museums in the area. Do you have a preference for type or price range?

Visitor: i really do not have a preference . can you recommend a few ?

Baseline: I recommend [attraction_name] It is located in the [value_area].
GPT: How about the [attraction_name]?

Visitor: Sure. Also could you help me find a restaurant serving unusual food in the same part of town ?

Baseline: There are no restaurants in the [value_area] of town. Would you like to try a different type of food?
GPT: Unfortunately, there are no restaurants that match that description. Would you like to try a different area or type of food?

Figure 3: The comparison of generated responses from the baseline model and GPT2-M.