

Eliza-like chatbot with Transformers and DialoGPT fine-tuning

Mohamad Al Mdfaa¹ **Mostafa Hegazy**² **Leon Derczynski**³
Innopolis University / Russia Innopolis University / Russia ITU / Copenhagen, Denmark

¹m.almdfaa@innopolis.university

²m.hegazy@innopolis.university

³ld@itu.dk

Abstract

Making a chatbot in the past was an extremely hard problem due to various reasons including the complexity of mimicking the human behaviour and responses, the lack of training data, the natural language processing field was not as mature as today. Eliza was a rule-based chatbot that using a set of predefined rules made by the developers, So we utilized the power of transfer learning with the DialoGPT pretrained model in order to construct the modern Eliza which will be trained on a custom dataset from a modern TV series and the responses of the original Eliza and will not be using a predefined set of rules. We found that the fine-tuned model behaved in a way that mimics the human behavior from the TV series along with the original Eliza's behaviour.

1 Introduction

A chatbot system or a conversational agent is software that uses natural language processing to interact with humans, since the birth of the Artificial Intelligence field, modeling conversations was and still is one of the toughest challenges. Although they still didn't reach their maximum potential and they're far from perfect, chatbots are now used in too many applications like Siri (Apple, 2017), Google Assistant (Google, 2017) or (Microsoft, 2017). Chatbots are used to answer the users' questions in any particular domain where it is operating. Recent advances in deep learning have inspired many applications of neural models to dialogue systems. As an example, (Wen et al., 2016) and (Bordes et al., 2016) introduced a network-based end-to-end trainable task-oriented dialogue system, which treated dialogue system learning as the problem of learning a mapping from dialogue histories to system responses, and the whole system was trained by an encoder-decoder model. Nevertheless, the system was trained in a supervised manner: it requires a lot of training data to get reasonable and good performance, it also may fail to reach good performance due to the lack of exploration of dialogue control in the data itself. (Zhao and Eskenazi, 2016) first presented an end-to-end reinforcement learning (RL) approach to dialogue state tracking and policy learning in the DM. This approach is shown to have good performance when applied to the task-oriented dialogue problem of guessing the famous person a user thinks of. In the conversation, the chatbot asks the user a series of questions which the answer to is Yes/No to find the correct answer. Despite this simplified task, it may not generalize to practical problems due to the following: Inflexible question types, Poor robustness, and User requests during dialogues. That brings us to recent advances in large-scale pre-training using transformer-based architectures (Radford et al., 2019), (Raffel et al., 2019) have achieved great empirical success. OpenAI's GPT-2 (Radford et al., 2019), has proven that transformer models that were trained on a huge dataset can capture valuable information and long-term dependencies in textual data and the generated text from such transformer is fluent, diverse, and rich in content. These models have the capacity and intelligence to capture textual data with Such models have the capacity to capture textual data with excellent attention to detail and produce an output with a high-resolution that can mimic and emulate real-world text as if it was written by humans.

2 Related Works

There are many open-sourced toolkits for large-scale pre-trained models of transformers. Hugging-face Conv-AI transfer learning repository (Wolf et al., 2019) contains the necessary and basic code for training conversational AI systems with transfer learning which is based on the GPT-2 transformer model. It achieves the state-of-the-art (SoTA) performance on ConvAI-2 dialogue competition. DLGnet (Olabiyi and Mueller, 2019) is a large transformer model trained on a dialogue dataset and achieves good performance in turn-based dialogue generation. AllenNLP (Gardner et al., 2018) is developed as a package for many natural language processing tasks, including the large-scale pre-trained BI-LSTM sentence representation learning framework ELMo (Peters et al., 2018). Text generation including style transferring and controlled generation is the main focus of Texar (Hu et al., 2018). It includes reinforcement learning capabilities along with its sequence modeling tools. Another popular framework that focuses on task-oriented dialogue is DeepPavlov (Burtsev et al., 2018). This Hugging-face public repository (HuggingFace, 2019) contains several demos and pre-trained models for many applications including question answering and sentiment classification and Analysis. Icecaps (Dhillon et al., 2020) is a response generation package with grounding on personalities or external knowledge and multi-task training techniques. ParlAI (Miller et al., 2017) is one of the libraries for developing task-oriented dialogue systems. It contains many pre-trained models for knowledge-grounded chatbots trained with data gathered from crowdsourcing. The Text-to-Text Transformer (Raffel et al., 2019) unifies multiple text modeling tasks, and achieves the state-of-the-art (SoTA) results in many natural language generations and understanding benchmarks.

3 Implementation

3.1 Dataset

The dataset used is gathered from the subtitles of the first season of a popular serial show called "The big-bang theory", which is an American sitcom show about two physicists, Leonard and Sheldon they have brilliant minds to understand how the universe works. But socially awkward, that shows how little they know about life outside of the laboratory. The different characters of the show have various responses, questions, and dialogues throughout the episodes, these responses were taken and used as input to the Eliza-chatbot (Weizenbaum, 1966) which is the very first chatbot created by Joseph Weizenbaum in 1966 to imitate a therapist who would ask open-ended questions and even respond with follow-ups. The responses of Eliza were taken as an extra input to the model for training with the addition of the various responses of each character. This dataset enables us to gain specific real-world comedic responses because the show is considered to be one of the comedic shows in the last 20 years.

Scene: The stairs of the apartment building.

Sheldon: Are you still mad about the sperm bank?

Leonard: No.

Sheldon: You want to hear an interesting thing about stairs?

Leonard: Not really.

Sheldon: If the height of a single step is off by as little as two millimetres, most people will trip.

Leonard: I don't care. Two millimetres? That doesn't seem right.

Sheldon: No, it's true, I did a series of experiments when I was twelve, my father broke his clavicle.

Leonard: Is that why they sent you to boarding school?

Sheldon: No, that was the result of my work with lasers.

Leonard: New neighbour?

Sheldon: Evidently.

Figure 1: A sample dialogue from the dataset

3.2 Language model definition

To solve our problem we need to define what is a language model in the first place, basically a machine learning model that is able to look at part of a sentence and predict the next word. The language model provides context to distinguish between words and phrases that sound phonetically similar. We decided to use DialoGPT as our pretrained language model.

3.3 Transformers for Language Modeling

Before we explore the architecture of DialoGPT we have to first see the basic transformer architecture for language modeling.

The basic architecture for the transformer consists of an encoder and a decoder with self attention mechanism

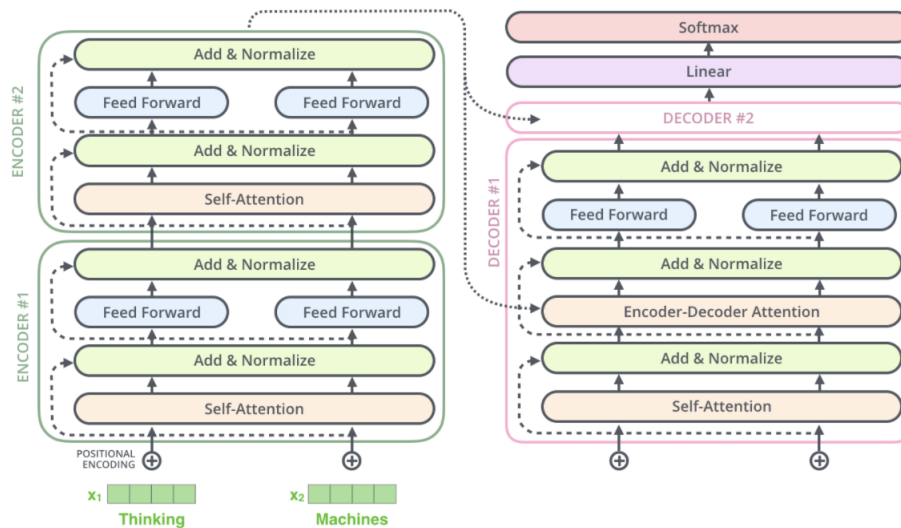


Figure 2: Basic transformer architecture (Alammar, 2020)

Later research work saw the architecture discard either the encoder or decoder, and simply use one stack of transformer blocks – stacking them up as high as practically feasible and possible, feeding them enormous amounts of training data. How many blocks in a model is the main factor that distinguishes the models in the GPT2-family.

The GPT-2 is built using transformer decoder blocks. BERT on the contrary, uses transformer encoder blocks. The main difference between them is that the GPT-2 outputs one token at a time. The GPT-2, and some later models like TransformerXL and XLNet are auto-regressive. BERT is not. That is a trade off. The price of losing auto-regression is that it gained the ability to incorporate the context on both sides of a word to gain better results. Auto-regression enabled the model to be unreasonably effective.

The decoder block of the GPT-2 has a small architectural variation from the encoder block – a layer to allow it to pay attention to specific segments from the encoder:

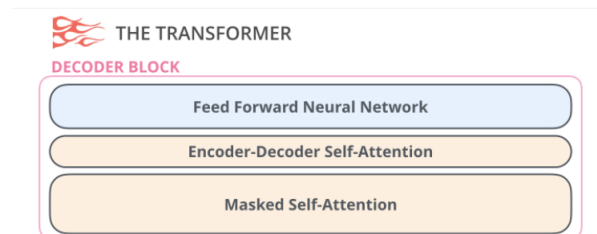


Figure 3: Decoder block diagram architecture (Alammar, 2020)

One key difference in the self-attention layer here, is that it masks future tokens – not by changing the word to (mask) like BERT, but by interfering in the self-attention calculation blocking information from tokens that are to the right of the position being calculated.

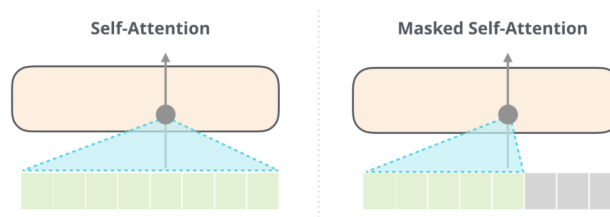


Figure 4: Decoder block diagram architecture (Alammar, 2020)

It's important that the distinction between self-attention (what BERT uses) and masked self-attention (what GPT-2 uses) is clear. A normal self-attention block allows a position to peak at tokens to its right. Masked self-attention prevents that from happening.

3.4 DialoGPT

DialoGPT is a transformer based on the GPT-2 transformer. DialoGPT is a large, tunable neural conversational response generation model, it was trained on 147M conversation-like exchanges extracted from Reddit comment chains over a period spanning from 2005 through 2017. It inherits from the GPT-2 its 12-to-48 layer with layer normalization, an initialization scheme that accounts for model depth that the original authors modified, and byte pair encodings (Sennrich et al., 2015) for the tokenizer.

3.5 Setting up the data

We need to construct original Eliza's responses in order to use them in our fine-tuning of the transformer model. The responses were generated based on a virtual dialogue between the characters in the show with Eliza. The responses were taken and stored in the dataset.

Eliza	How do you do. Please tell me your problem.
Sheldon	So if a photon is directed through a plane wi...
Eliza	Do you think its likely that it's observed aft...
Leonard	Agreed, what's your point?
Eliza	Why are you concerned over my point? ?
Sheldon	There's no point, I just think it's a good id...
Eliza	Are you saying no just to be negative?
Leonard	Excuse me?
Eliza	I'm not sure I understand you fully.
Receptionist	Hang on.

Figure 5: A snippet of Eliza's interaction with the characters

We will convert this dataset in a way that every response row will contain n previous responses as a context. For our purpose seven previous responses has been taken, so the n is 7.

response	context	context/0	context/1	context/2	context/3	context/4	context/5
Excuse me?	Are you saying no just to be negative?	There's no point, I just think it's a good id...	Why are you concerned over my point? ?	Agreed, what's your point?	Do you think its likely that it's observed aft...	So if a photon is directed through a plane wi...	How do you do. Please tell me your problem.
I'm not sure I understand you fully.	Excuse me?	Are you saying no just to be negative?	There's no point, I just think it's a good id...	Why are you concerned over my point? ?	Agreed, what's your point?	Do you think its likely that it's observed aft...	So if a photon is directed through a plane wi...
Hang on.	I'm not sure I understand you fully.	Excuse me?	Are you saying no just to be negative?	There's no point, I just think it's a good id...	Why are you concerned over my point? ?	Agreed, what's your point?	Do you think its likely that it's observed aft...
Please go on.	Hang on.	I'm not sure I understand you fully.	Excuse me?	Are you saying no just to be negative?	There's no point, I just think it's a good id...	Why are you concerned over my point? ?	Agreed, what's your point?
One across is Aegean, eight down is Nabakov, ...	Please go on.	Hang on.	I'm not sure I understand you fully.	Excuse me?	Are you saying no just to be negative?	There's no point, I just think it's a good id...	Why are you concerned over my point? ?

Figure 6: Dataset Transformation

The whole dataset will be split into 90% training data and 10% testing data.

3.6 Fine-tuning

The fine-tuning was done with the help of PyTorch library. The chosen optimizer is AdamW which is a stochastic optimization method that modifies the typical implementation of weight decay in Adam, by

decoupling weight decay from the gradient update.

3.6.1 First experiment

The number of training examples in the data set is 4464, number of epochs is set to be 3, the Instantaneous batch size per GPU is 4. The number of evaluation examples is 496.

3.6.2 Second experiment

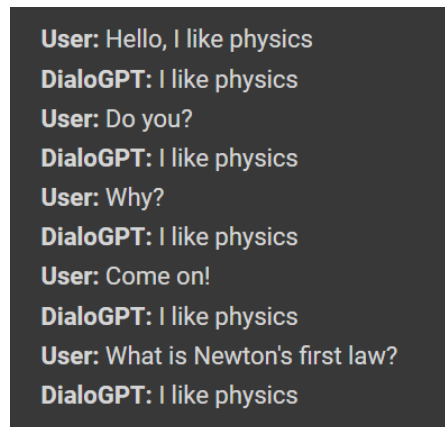
The number of training examples in the data set is 6638, number of epochs is set to be 3, the Instantaneous batch size per GPU is 4. The number of evaluation examples is 738.

3.6.3 Evaluation

The model is evaluated, using the metric perplexity, which is a measure of how unsure the model is in its choice of the next token. The more unsure the model is, the higher its perplexity. One fascinating thing about perplexity is that it correlates very well with what humans think of when it comes to coherent and specific natural conversations (Adiwardana et al., 2020).

4 Results

4.1 Without fine-tuning



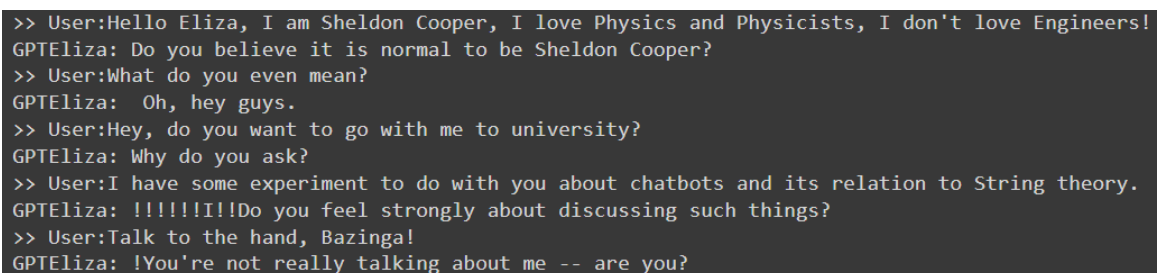
```
User: Hello, I like physics
DialoGPT: I like physics
User: Do you?
DialoGPT: I like physics
User: Why?
DialoGPT: I like physics
User: Come on!
DialoGPT: I like physics
User: What is Newton's first law?
DialoGPT: I like physics
```

Figure 7: A snippet of the model's responses without any fine-tunings

As we can see that the pre-trained model didn't give any human-like responses although it was trained on a huge dataset, But the domain of questions/responses here is not the same as the domain of pre-training. That's why the performance is not satisfactory.

4.2 With fine-tuning

4.2.1 First Experiment



```
>> User:Hello Eliza, I am Sheldon Cooper, I love Physics and Physicists, I don't love Engineers!
GPTELiza: Do you believe it is normal to be Sheldon Cooper?
>> User:What do you even mean?
GPTELiza: Oh, hey guys.
>> User:Hey, do you want to go with me to university?
GPTELiza: Why do you ask?
>> User:I have some experiment to do with you about chatbots and its relation to String theory.
GPTELiza: !!!!!!!I!!Do you feel strongly about discussing such things?
>> User:Talk to the hand, Bazinga!
GPTELiza: !You're not really talking about me -- are you?
```

Figure 8: A snippet of the model's responses with fine-tuning

As shown, even with small amount of data the responses becomes nearly human-like responses but we also have one response to a question which didn't really fit into the context, this happened because the dataset contained a lot of group interactions and talks, also the setting in which sitcom works is really group-oriented so most of the dialogues in it involves a group of people.

4.2.2 Second Experiment

```
>> User:Hello, I like physics
GPTELiza: What does that suggest to you?
>> User:It makes me understand the universe.
GPTELiza: Do you say It makes you understand the Universe for some special reason?
>> User:Yes, if you understand how the universe works, you will be able to build a civilization.
GPTELiza: What do you know about I understand the world works?
>> User:Have you heard about the Newton's laws?
GPTELiza: !!!lets discuss further why your hands as often as you can.
>> User:Sure, I feel good, thanks to you!
GPTELiza: !L!!?!!“!!I!!
```

Figure 9: A snippet of the model's responses with fine-tuning

It's obvious in this example that the model can respond with a mix of Eliza-like and human-like responses but we also notice the more sophisticated user's responses the higher chances that the model give not the most satisfactory answers, one of the reasons that this happened because the dataset is not big enough. Later in the future we are planning to initialize an experiment to train the model on more data with all the episodes from all seasons in the Big-Bang theory TV series.

4.2.3 Third Experiment

```
>> User: Hello, I like physics
GPTELiza: How do you do. Please state your problem.
>> User: I feel tired because of studying
GPTELiza: I can't believe you've never seen a show rated M.A., that's incredible.
>> User: Do you think watching a TV show could help me feel better?
GPTELiza: Yeah, it's true. M.I.T., Princeton and UCLA are in the lead.
>> User: Okay, thanks!
GPTELiza: Please go on.
>> User: I am feeling better, thanks
GPTELiza: !!! Remind me!!
```

Figure 10: A snippet of the model's responses with fine-tuning

```
>> User: I am sick!
GPTELiza: Do you feel strongly about discussing such things?
>> User: I want to take a year off from everything
GPTELiza: Why do you want to do a year of everything?
>> User: I want to go to a quiet country where I can relax with my family
GPTELiza: Tell me more about your family.
>> User: We love each other
GPTELiza: !!!
>> User: It is personal matter! I don't want to talk about it!
GPTELiza: !.!!.
```

Figure 11: A snippet of the model's responses with fine-tuning

As we can see in this example, the model's responses were refined. As we feed in more data from the show

```

>> User: I am fed up with this stupid world
GPTEliza: Does that have anything to do with the fact that your God?
>> User: no
GPTEliza: Are you saying no just to be negative?
>> User: no
GPTEliza: You are being a bit negative.
>> User: Nooooo
GPTEliza: What does that suggest to you?
>> User: don't know
GPTEliza: !!! remind me!! Remind me! Remand me!

```

Figure 12: A snippet of the model's responses with fine-tuning

We can notice in this experiment that we have a typical typo that many English speakers and learners make, "your" instead of "you're" which can be explained from a typo in the subtitles used as a training data to the model. So the model not only learn the context and the sentences from the data but also the mistakes, that's why we should be very careful in filtering these mistakes from the dataset.

```

>> User: What are your plans for life?
GPTEliza: Well, I'm a physicist, and I'm an advocate of the death grip of premature babies.
>> User: Do you think you're attractive?
GPTEliza: Really?
>> User: Yes
GPTEliza: You know what? I'm gonna go. I'm not going anywhere.
>> User: Why, i like talking to you
GPTEliza: No, no. You know, I've been talking about this with Amy a long time. It's just, I was thinking about her.
>> User: Been a pleasure, Eliza
GPTEliza: !!! Remind me to try again in an hour.

```

Figure 13: A snippet of the model's responses with fine-tuning

In this experiment, the model behaves like Sheldon Cooper, the main star of the show and it's very apparent that the sarcastic tone of Sheldon is visible inside the model.

```

>> User: I am tired because you are not tired
GPTEliza: What makes you think I am not tired because I are not exhausted?
>> User: Why you didn't focus that I am tired and you focused on yourself?
GPTEliza: Because I am a superior intellect to yours.
>> User: Really? As I see, you are a looser in your specialization which psychology
GPTEliza: Is that the real reason?
>> User: I do recommend that you visit a psychologist.
GPTEliza: !!! remind me!! calendar reminder!! Remind me! calendar!! reminder! calendar!!!
>> User: Be safe, bye
GPTEliza: ! remind you!! calendars!! And a special hello to my dear friends, Howard and Bernadette.

```

Figure 14: A snippet of the model's responses with fine-tuning

In this experiment we can see as well that the behaviour of Sheldon is apparent once more, so the model thinks that it's more worthy due to the fact that the training data had many conversations in which Sheldon thought he was better than most people.

```

>> User: Sometimes I feel lonely! I miss home! I miss my old life, my friends, my family!
GPTEliza: Your family?
>> User: I feel anxious about people behaviors here, what do you recommend me to do?
GPTEliza: You like to think I recommend you -- don't you?
>> User: okay
GPTEliza: Lets discuss further why your old friends.
>> User: Sometimes I feel depressed and the conditions around me make me remember a lot of stuff about my past
GPTEliza: !!! remind me tomorrow morning to put on pants.
>> User: Sorry! what?
GPTEliza: ! remind you tomorrow morning will you put on underwear?

```

Figure 15: A snippet of the model's responses with fine-tuning

In this experiment, the original psychologicistic trait of the original Eliza bot is apparent as well as the sarcastic tone of the big bang theory show. What we can understand from here is that the behaviour of the model will be determined based on the amount of the data and context we want it to have.

5 Conclusion

In this work we utilized the power of the DialoGPT pretrained model and transfer learning in order to construct the modern Eliza which was trained on a custom dataset from a modern TV series and the responses of the original Eliza. We found that the fine-tuned model behaved in a way more like human while preserving the original Eliza’s behaviour.

References

- Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. 2020. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*.
- Jay Alammar. 2020. Illustrated transformer. <https://jalammar.github.io/illustrated-transformer/>.
- Apple. 2017. Siri. <https://www.apple.com/ios/siri/>.
- Antoine Bordes, Y-Lan Boureau, and Jason Weston. 2016. Learning end-to-end goal-oriented dialog. *arXiv preprint arXiv:1605.07683*.
- Mikhail Burtsev, Alexander Seliverstov, Rafael Airapetyan, Mikhail Arkhipov, Dilyara Baymurzina, Nickolay Bushkov, Olga Gureenkova, Taras Khakhulin, Yurii Kuratov, Denis Kuznetsov, et al. 2018. Deeppavlov: Open-source library for dialogue systems. In *Proceedings of ACL 2018, System Demonstrations*, pages 122–127.
- Gurpreet Singh Dhillon, Shohreh Honarbakhsh, Antonio Di Monaco, Ann Elizabeth Coling, Kernerová Lenka, Francesca Pizzamiglio, Ross J Hunter, Rodney Horton, Moussa Mansour, Andrea Natale, et al. 2020. Use of a multi-electrode radiofrequency balloon catheter to achieve pulmonary vein isolation in patients with paroxysmal atrial fibrillation: 12-month outcomes of the radiance study. *Journal of cardiovascular electrophysiology*, 31(6):1259–1269.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. Allennlp: A deep semantic natural language processing platform. *arXiv preprint arXiv:1803.07640*.
- Google. 2017. Google assistant. <https://assistant.google.com/>.
- Zhiting Hu, Haoran Shi, Bowen Tan, Wentao Wang, Zichao Yang, Tiancheng Zhao, Junxian He, Lianhui Qin, Di Wang, Xuezhe Ma, et al. 2018. Texar: A modularized, versatile, and extensible toolkit for text generation. *arXiv preprint arXiv:1809.00794*.
- HuggingFace. 2019. *PyTorch transformer repository*. <https://github.com/huggingface/pytorch-transformers>.
- Microsoft. 2017. Cortana. <https://www.microsoft.com/en-us/windows/cortana>.
- Alexander H Miller, Will Feng, Adam Fisch, Jiasen Lu, Dhruv Batra, Antoine Bordes, Devi Parikh, and Jason Weston. 2017. Parlai: A dialog research software platform. *arXiv preprint arXiv:1705.06476*.
- Oluwatobi Olabiyi and Erik T Mueller. 2019. Dlgnet: A transformer-based model for dialogue response generation. *arXiv preprint arXiv:1908.01841*.
- ME Peters, M Neumann, M Iyyer, M Gardner, C Clark, K Lee, and L Zettlemoyer. 2018. Deep contextualized word representations. *naacl-hlt*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*.
- Joseph Weizenbaum. 1966. Eliza—a computer program for the study of natural language communication between man and machine. *Communications of the ACM*, 9(1):36–45.

- Tsung-Hsien Wen, David Vandyke, Nikola Mrksic, Milica Gasic, Lina M Rojas-Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. 2016. A network-based end-to-end trainable task-oriented dialogue system. *arXiv preprint arXiv:1604.04562*.
- 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*.
- Tiancheng Zhao and Maxine Eskenazi. 2016. Towards end-to-end learning for dialog state tracking and management using deep reinforcement learning. *arXiv preprint arXiv:1606.02560*.