

A Neural Conversational Model

Oriol Vinyals

Google

Quoc V. Le

Google

VINYALS@GOOGLE.COM

QVL@GOOGLE.COM

Abstract

Conversational modeling is an important task in natural language understanding and machine intelligence. Although previous approaches exist, they are often restricted to specific domains (e.g., booking an airline ticket) and require hand-crafted rules. In this paper, we present a simple approach for this task which uses the recently proposed sequence to sequence framework. Our model converses by predicting the next sentence given the previous sentence or sentences in a conversation. The strength of our model is that it can be trained end-to-end and thus requires much fewer hand-crafted rules. We find that this straightforward model can generate simple conversations given a large conversational training dataset. Our preliminary results suggest that, despite optimizing the wrong objective function, the model is able to converse well. It is able to extract knowledge from both a domain specific dataset, and from a large, noisy, and general domain dataset of movie subtitles. On a domain-specific IT helpdesk dataset, the model can find a solution to a technical problem via conversations. On a noisy open-domain movie transcript dataset, the model can perform simple forms of common sense reasoning. As expected, we also find that the lack of consistency is a common failure mode of our model.

1. Introduction

Advances in end-to-end training of neural networks have led to remarkable progress in many domains such as speech recognition, computer vision, and language processing. Recent work suggests that neural networks can do more

than just mere classification, they can be used to map complicated structures to other complicated structures. An example of this is the task of mapping a sequence to another sequence which has direct applications in natural language understanding (Sutskever et al., 2014). The main advantage of this framework is that it requires little feature engineering and domain specificity whilst matching or surpassing state-of-the-art results. This advance, in our opinion, allows researchers to work on tasks for which domain knowledge may not be readily available, or for tasks which are simply too hard to design rules manually.

Conversational modeling can directly benefit from this formulation because it requires mapping between queries and responses. Due to the complexity of this mapping, conversational modeling has previously been designed to be very narrow in domain, with a major undertaking on feature engineering. In this work, we experiment with the conversation modeling task by casting it to a task of predicting the next sequence given the previous sequence or sequences using recurrent networks (Sutskever et al., 2014). We find that this approach can do surprisingly well on generating fluent and accurate replies to conversations.

We test the model on chat sessions from an IT helpdesk dataset of conversations, and find that the model can sometimes track the problem and provide a useful answer to the user. We also experiment with conversations obtained from a noisy dataset of movie subtitles, and find that the model can hold a natural conversation and sometimes perform simple forms of common sense reasoning. In both cases, the recurrent nets obtain better perplexity compared to the n-gram model and capture important long-range correlations. From a qualitative point of view, our model is sometimes able to produce natural conversations.

2. Related Work

Our approach is based on recent work which proposed to use neural networks to map sequences to sequences (Kalchbrenner & Blunsom, 2013; Sutskever et al., 2014; Bahdanau et al., 2014). This framework has been

used for neural machine translation and achieves improvements on the English-French and English-German translation tasks from the WMT'14 dataset (Luong et al., 2014; Jean et al., 2014). It has also been used for other tasks such as parsing (Vinyals et al., 2014a) and image captioning (Vinyals et al., 2014b). Since it is well known that vanilla RNNs suffer from vanishing gradients, most researchers use variants of Long Short Term Memory (LSTM) recurrent neural networks (Hochreiter & Schmidhuber, 1997).

Our work is also inspired by the recent success of neural language modeling (Bengio et al., 2003; Mikolov et al., 2010; Mikolov, 2012), which shows that recurrent neural networks are rather effective models for natural language. More recently, work by Sordoni et al. (Sordoni et al., 2015) and Shang et al. (Shang et al., 2015), used recurrent neural networks to model dialogue in short conversations (trained on Twitter-style chats).

Building bots and conversational agents has been pursued by many researchers over the last decades, and it is out of the scope of this paper to provide an exhaustive list of references. However, most of these systems require a rather complicated processing pipeline of many stages (Lester et al., 2004; Will, 2007; Jurafsky & Martin, 2009). Our work differs from conventional systems by proposing an end-to-end approach to the problem which lacks domain knowledge. It could, in principle, be combined with other systems to re-score a short-list of candidate responses, but our work is based on producing answers given by a probabilistic model trained to maximize the probability of the answer given some context.

3. Model

Our approach makes use of the sequence-to-sequence (*seq2seq*) framework described in (Sutskever et al., 2014). The model is based on a recurrent neural network which reads the input sequence one token at a time, and predicts the output sequence, also one token at a time. During training, the true output sequence is given to the model, so learning can be done by backpropagation. The model is trained to maximize the cross entropy of the correct sequence given its context. During inference, given that the true output sequence is not observed, we simply feed the predicted output token as input to predict the next output. This is a “greedy” inference approach. A less greedy approach would be to use beam search, and feed several candidates at the previous step to the next step. The predicted sequence can be selected based on the probability of the sequence.

Concretely, suppose that we observe a conversation with two turns: the first person utters “ABC”, and second person replies “WXYZ”. We can use a recurrent neural network,

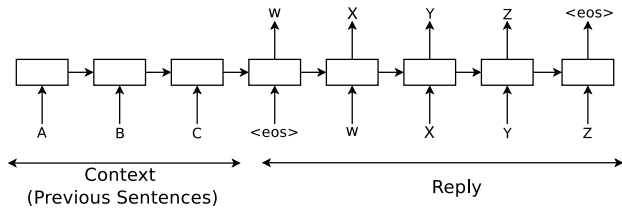


Figure 1. Using the *seq2seq* framework for modeling conversations.

and train to map “ABC” to “WXYZ” as shown in Figure 1 above. The hidden state of the model when it receives the end of sequence symbol “<eos>” can be viewed as the *thought vector* because it stores the information of the sentence, or thought, “ABC”.

The strength of this model lies in its simplicity and generality. We can use this model for machine translation, question/answering, and conversations without major changes in the architecture. Applying this technique to conversation modeling is also straightforward: the input sequence can be the concatenation of what has been conversed so far (the context), and the output sequence is the reply.

Unlike easier tasks like translation, however, a model like sequence-to-sequence will not be able to successfully “solve” the problem of modeling dialogue due to several obvious simplifications: the objective function being optimized does not capture the actual objective achieved through human communication, which is typically longer term and based on exchange of information rather than next step prediction. The lack of a model to ensure consistency and general world knowledge is another obvious limitation of a purely unsupervised model.

4. Datasets

In our experiments we used two datasets: a closed-domain IT helpdesk troubleshooting dataset and an open-domain movie transcript dataset. The details of the two datasets are as follows.

4.1. IT Helpdesk Troubleshooting dataset

In our first set of experiments, we used a dataset which was extracted from a IT helpdesk troubleshooting chat service. In this service, costumers face computer related issues, and a specialist help them by conversing and walking through a solution. Typical interactions (or threads) are 400 words long, and turn taking is clearly signaled. Our training set contains 30M tokens, and 3M tokens were used as validation. Some amount of clean up was performed, such as removing common names, numbers, and full URLs.

4.2. OpenSubtitles dataset

We also tested our model on the OpenSubtitles dataset (Tiedemann, 2009). This dataset consists of movie conversations in XML format. It contains sentences uttered by characters in movies. We applied a simple processing step removing XML tags and obvious non-conversational text (e.g., hyperlinks) from the dataset. As turn taking is not clearly indicated, we treated consecutive sentences assuming they were uttered by different characters. We trained our model to predict the next sentence given the previous one, and we did this for every sentence (noting that this doubles our dataset size, as each sentence is used both for context and as target). Our training and validation split has 62M sentences (923M tokens) as training examples, and the validation set has 26M sentences (395M tokens). The split is done in such a way that each sentence in a pair of sentences either appear together in the training set or test set but not both. Unlike the previous dataset, the OpenSubtitles is quite large, and rather noisy because consecutive sentences may be uttered by the same character. Given the broad scope of movies, this is an open-domain conversation dataset, contrasting with the technical troubleshooting dataset.

5. Experiments

In this section, we describe the experimental results with the two datasets and show some samples of the interactions with the system that we trained. We also compare the performance of our system against a popular rule-based bot (CleverBot¹) using human evaluations on a set of 200 questions.

5.1. IT Helpdesk Troubleshooting experiments

In this experiment, we trained a single layer LSTM with 1024 memory cells using stochastic gradient descent with gradient clipping. The vocabulary consists of the most common 20K words, which includes special tokens indicating turn taking and actor.

At convergence, this model achieved a perplexity of 8, whereas an n-gram model achieved 18. Below is a few samples of simulated troubleshooting sessions for most three of the most common IT problems (remote access, software crashes, and password issues). In these conversations, **Machine** is our Neural Conversational Model, and **Human** the human actor interacting with it.

Conversation 1: VPN issues.

Describe your problem: i am having issues accessing vpn

Machine: hi

Human: hello

Machine: could you please let me know what are the operating systems you are using ?

Human: i am using linux

Machine: ok it was in the ssl vpn on which you stated it was broken

Machine: please check if you are logged in to corp network

Machine: what is the error that you are running please

Human: i am seeing an error related to vpn

Machine: what is the error message that you are getting when connecting to vpn using network connect ?

Human: connection refused or something like that

Machine: may i know the version of network connect you connect ?

Human: i am not sure i know that

Machine: is the network connect prompting for any user certificate

Machine: ?

Human: yes

Machine: are you at home or at work at the moment ?

Machine: office ?

Human: i am at home

Machine: try this

Machine: goto <URL> page and enter your name and password

Human: done

Machine: try logging in with <NAME> and check if its working fine now

Human: yes , now it works !

Machine: great

Machine: anything else that i can help

Human: no , that was very useful

Machine: you are other than <NAME> :)

Machine: have a nice night

(The <URL> indeed contains information about vpn access.)

Conversation 2: Browser issues.

Describe your problem: browser keeps crashing

Machine: hi

Human: hi

Machine: could you please let me know what is the operating system you are using ?

Human: linux

Machine: ok this is the issue with the system performance , nasser , tools .

Human: what is that ?

¹<http://www.cleverbot.com/>

Machine: *is it a windows machine ?*
Machine: *we have linux on your machine .*
Machine: *is it windows xp or windows 7 ?*
Human: *i am using linux*
Machine: *let me take a remote of your machine and check if issue has been pushed properly .*
Human: *ok*
Machine: *please open terminal and run the below commands in one computer :*
Machine: *sudo apt - get update*
Machine: *sudo apt - get upgrade*
Machine: *sudo <COMMAND> - updater*
Human: *done*
Machine: *let me try connecting now*
Human: *ok*
Machine: *its done*
Machine: *should use the second or hit run on the background*
Human: *done*
Machine: *try signing out now and follow the prompts as per the message given above*
Machine: *shd i reboot the machine .*
Human: *ok , please go ahead*
Machine: *once its completed please try to login .*
Human: *sure*
Machine: *please ping me back once you have tried to login*
Machine: *are we still connected*
Human: *i now logged in*
Machine: *thanks for your confirmation*
Human: *what should i do now ?*
Machine: *wait for 1 5 minutes to complete the login policy*
Machine: *and then try logging in*
Human: *ok*
Machine: *is there anything else to help you ?*
Human: *let my try browsing a little bit and if it keeps crashing i will let you know*
Machine: *sure will do*
Machine: *seems the mail is not loading*
Machine: *are we still connected to the chat ?*
Machine: *<reinvent>*
Human: *it seems good now , cool , thank you , bye*
Machine: *you , too*
Machine: *bye*

Conversation 3: Password issues.

Describe your problem: i forgot my password
Machine: *hi , this is <NAME> from helpdesk connect , what can i help you with today ?*
Human: *hi , i forgot my password*

Machine: *ok great*
Human: *can you tell me how i can recover it ?*
Machine: *i ' ll need to verify who it is at that machine*
Machine: *can we do a hangout*
Machine: *?*
Human: *yes*
Machine: *<anonymized hangout instructions and verifications>*
Human: *cool , i am good now*
Machine: *how ' s this looking good ?*
Machine: *i ' ll go ahead and close this chat .*
Machine: *take care . :)*
Human: *bye*
Machine: *bye .*

5.2. OpenSubtitles experiments

In this experiment, we trained a **two-layered LSTM using AdaGrad with gradient clipping**. Each layer of the LSTM has 4096 memory cells, and we built a vocabulary consisting of the most frequent 100K words. To speed up the softmax, we project the memory cells to 2048 linear units before feeding the information to the classifier.

At convergence the perplexity of the recurrent model on the validation set is 17. Our smoothed 5-gram model achieves a perplexity of 28. Interestingly, adding the soft attention mechanism of (Bahdanau et al., 2014) did not significantly improve the perplexity on neither training or validation sets.

5.2.1. MODEL SAMPLES

In addition to the perplexity measure, our simple recurrent model does often produce plausible answers. Here are a few sample question-answer pairs that allowed us to understand the capabilities of the model when trained on the noisy but large dataset of OpenSubtitles from movies. Again, in these conversations, **Machine** is our Neural Conversational Model, and **Human** the human actor interacting with it.

Conversation 1: Basic

Human: *hello !*
Machine: *hello !*
Human: *how are you ?*
Machine: *i ' m good .*
Human: *what ' s your name ?*
Machine: *i ' m julia .*
Human: *when were you born ?*
Machine: *july 20th .*
Human: *what year were you born ?*
Machine: *1977 .*
Human: *where are you from ?*

Machine: *i 'm out in the boonies.*
Human: *are you a man or a woman ?*
Machine: *i 'm a woman .*
Human: *why are we here ?*
Machine: *i 'm not sure .*
Human: *okay, bye !*
Machine: *bye .*
Human: *see you later .*
Machine: *bye .*

Conversation 2: Simple Q&A (contexts and multiple choice)

Human: *my name is david . what is my name ?*
Machine: *david .*
Human: *my name is john . what is my name ?*
Machine: *john .*
Human: *are you a leader or a follower ?*
Machine: *i 'm a leader .*
Human: *are you a follower or a leader ?*
Machine: *i 'm a leader .*

Conversation 3: General knowledge Q&A

Human: *who is skywalker ?*
Machine: *he is a hero .*
Human: *who is bill clinton ?*
Machine: *he 's a billionaire .*
Human: *is sky blue or black ?*
Machine: *blue .*
Human: *does a cat have a tail ?*
Machine: *yes .*
Human: *does a cat have a wing ?*
Machine: *no*
Human: *can a cat fly ?*
Machine: *no .*
Human: *how many legs does a cat have ?*
Machine: *four , i think .*
Human: *how many legs does a spider have ?*
Machine: *three , i think .*
Human: *how many legs does a centipede have ?*
Machine: *eight .*
Human: *what is the color of the sky ?*
Machine: *blue .*
Human: *what is the color of water ?*
Machine: *water .*
Human: *what is the color of blood ?*
Machine: *it is the same as a black eye .*
Human: *what is the usual color of a leaf ?*
Machine: *it is a green one .*
Human: *what is the color of a yellow car ?*
Machine: *yellow .*
Human: *how much is two plus two ?*
Machine: *four .*

Human: *how much is ten minus two ?*
Machine: *seventy - two .*

Conversation 4: Philosophical Q&A

Human: *what is the purpose of life ?*
Machine: *to serve the greater good .*
Human: *what is the purpose of living ?*
Machine: *to live forever .*
Human: *what is the purpose of existence ?*
Machine: *to find out what happens when we get to the planet earth .*
Human: *where are you now ?*
Machine: *i 'm in the middle of nowhere .*
Human: *what is the purpose of dying ?*
Machine: *to have a life .*
Human: *what is the purpose of being intelligent ?*
Machine: *to find out what it is .*
Human: *what is the purpose of emotions ?*
Machine: *i don 't know .*

Conversation 5: Morality

Human: *what is moral ?*
Machine: *what empowered humanity , what intellectual the essence is .*
Human: *what is immoral ?*
Machine: *the fact that you have a child .*
Human: *what is morality ?*
Machine: *what is altruism ?*
Human: *what is the definition of altruism ?*
Machine: *if you don 't believe in god , then you don 't know .*
Human: *ok ... so what is the definition of morality ?*
Machine: *well , the truth is , you 're not a believer in god almighty .*
Human: *tell me the definition of morality , i am quite upset now !*
Machine: *i 'm not ashamed of being a philosopher !*
Human: *tell me the definition of morality .*
Machine: *i don 't have ethics .*
Human: *look , i need help , i need to know more about morality ...*
Machine: *i don 't know what ethics is .*
Human: *seriously , what is morality ?*
Machine: *what is the definition of living ?*
Human: *why living has anything to do with morality ?*
Machine: *you 're not a cop .*
Human: *okay , i need to know how should i behave morally ...*

Machine: *i don't know how to tell you .*
Human: *is morality and ethics the same ?*
Machine: *yes , sir .*
Human: *what are the things that i do to be immoral ?*
Machine: *i don't know .*
Human: *give me some examples of moral actions ...*
Machine: *i 'm not a moralist .*
Human: *alright , morality ?*
Machine: *integrity .*
Human: *what is integrity ?*
Machine: *i 'm sorry , i don't know what else to say .*
Human: *be moral !*
Machine: *be a man !*
Human: *i really like our discussion on morality and ethics ...*
Machine: *and how i 'm not in the mood for a philosophical debate .*
Human: *what do you like to talk about ?*
Machine: *nothing .*

Conversation 6: Opinions (notice the model outputs “he” , “she” or “it” depending on the subject being asked)

Human: *what do you think about tesla ?*
Machine: *he 's a good conductor .*
Human: *what do you think about bill gates ?*
Machine: *he 's a good man .*
Human: *What do you think about messi ?*
Machine: *he 's a great player .*
Human: *what do you think about cleopatra ?*
Machine: *oh , she 's very regal .*
Human: *what do you think about england during the reign of elizabeth ?*
Machine: *it was a great place .*

We find it encouraging that the model can remember facts, understand contexts, perform common sense reasoning without the complexity in traditional pipelines. What surprises us is that the model does so without any explicit knowledge representation component except for the parameters in the word vectors.

Perhaps most practically significant is the fact that the model can generalize to new questions. In other words, it does not simply look up for an answer by matching the question with the existing database. In fact, most of the questions presented above, except for the first conversation, do not appear in the training set.

Nonetheless, one drawback of this basic model is that it only gives simple, short, sometimes unsatisfying answers to our questions as can be seen above. Perhaps a more

problematic drawback is that the model does not capture a consistent personality. Indeed, if we ask not identical but semantically similar questions, the answers can sometimes be inconsistent. This is expected due to the simplicity of our model and the dataset in our experiments. The conversation below illustrates such failure:

Conversation 7: Job and Personality

Human: *what is your job ?*
Machine: *i 'm a lawyer .*
Human: *what do you do ?*
Machine: *i 'm a doctor .*

Lastly, as a comparison of other traditional bots that are rule-based, we provide the output of CleverBot for some of the questions in our test:²

Human: *Are you a follower or a leader ?*
CleverBot: *No !!*
Human: *My name is david . what is my name ?*
CleverBot: *How old are you David ?*
Human: *My name is john . what is my name ?*
CleverBot: *John*
Human: *What is the purpose of being intelligent ?*
CleverBot: *The purpose of life is to love God and love people.*
Human: *Is sky blue or black ?*
CleverBot: *I am 200 years old.*
Human: *What is the color of a yellow car ?*
CleverBot: *Brown.*
Human: *What is the usual color of a leaf ?*
CleverBot: *What is the color of the apple in the fridge ?*
Human: *What is the purpose of life ?*
CleverBot: *My purpose it to forward my species, in other words to make it easier for future generations of mankind to live.*
Human: *What do you think about tesla ?*
CleverBot: *I am tesla.*
Human: *What do you think about cleopatra ?*
CleverBot: *Cool is cool.*
Human: *What do you think about england during the reign of elizabeth ?*
CleverBot: *I like him. Do you ?*
Human: *What do you think about bill gates ?*
CleverBot: *Lets talk about Link.*
Human: *What is moral ?*
CleverBot: *Morality is the quality of being in accord with standards of right or good conduct.*
Human: *What is immoral ?*

²The outputs of CleverBot are non-deterministic so the answers may be different for different trials.

CleverBot: *What is the objective of sistetic life ?*

Human: *Does a cat have a wing ?*

CleverBot: *What is the name of your husband.*

Human: *Can a cat fly ?*

CleverBot: *What is your husband's name ?*

Human: *Who is skywalker ?*

CleverBot: *There are two of them, Anakin Skywalker and Luke Skywalker, both characters from star Wars.*

Human: *How many legs does a cat have ?*

CleverBot: *An ant.*

Human: *How many legs does a centipede have ?*

CleverBot: *A centipede's number of legs varies from thirty to two hundred.*

Human: *What is your job ?*

CleverBot: *I work in mental health.*

Human: *What do you do ?*

CleverBot: *I have no job.*

5.2.2. HUMAN EVALUATION

In order to fairly and objectively compare our model against CleverBot, we picked 200 questions³, and asked four different humans to rate our model (NCM) versus CleverBot (CB). Human judges were asked to pick which of the two bots they preferred, and were also allowed to declare a tie if both answers were of equal quality.

If three out of four human judges agreed, we record that score. Our model was preferred in 97 out of 200 questions, whereas CleverBot was picked in 60 out of 200. There was a tie in 20 questions, and in 23 questions the judges were in disagreement.

Here are some samples of how the judges evaluated the systems (in bold we show their preferred answer, the third example shows a tie, and the fourth examples shows disagreement):

Human: What is the deepest spot on the world ?

NCM: **The pacific ocean.**

CB: What ?

Human: Do you like Mexican food or Indian food ?

NCM: I hate it.

CB: **I love Mexican food.**

Human: Who are you crazy about ?

NCM: **You're the one who's crazy.**

CB: Who are you ?

³The questions we used together with the answers can be found in <http://ai.stanford.edu/~quocle/QAresults.pdf>

Human: What's thirty plus forty ?

NCM: **Sixty ?**

CB: What should we talk about ?

We however believe that it remains an open research problem of designing a good metric to quickly measure the quality of a conversational model. We did so by manual inspection, by using Mechanical Turk to get human evaluations comparing with a rule-based bot, and by computing the perplexity. However, there are obvious drawbacks on using those, and we are actively pursuing this and other research directions.

6. Discussion

In this paper, we show that a simple language model based on the *seq2seq* framework can be used to train a conversational engine. **Our modest results show that it can generate simple and basic conversations, and extract knowledge from a noisy but open-domain dataset.** Even though the model has obvious limitations, it is **surprising to us that a purely data driven approach without any rules can produce rather proper answers to many types of questions.** However, the model may **require substantial modifications to be able to deliver realistic conversations.** Amongst the many limitations, the **lack of a coherent personality makes it difficult for our system to pass the Turing test (Turing, 1950).**

ACKNOWLEDGMENTS

We thank Greg Corrado, Andrew Dai, Jeff Dean, Tom Dean, Matthieu Devin, Rajat Monga, Mike Schuster, Noam Shazeer, Ilya Sutskever and the Google Brain team for the help with the project.

References

- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Bengio, Y., Ducharme, R., Vincent, P., and Janvin, C. A neural probabilistic language model. *The Journal of Machine Learning Research*, 3:1137–1155, 2003.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Computation*, 1997.
- Jean, S., Cho, K., Memisevic, R., and Bengio, Y. On using very large target vocabulary for neural machine translation. *CoRR*, abs/1412.2007, 2014.
- Jurafsky, D. and Martin, J. *Speech and language processing*. Pearson International, 2009.

- Kalchbrenner, N. and Blunsom, P. Recurrent continuous translation models. In *EMNLP*, 2013.
- Lester, J., Branting, K., and Mott, B. Conversational agents. In *Handbook of Internet Computing*. Chapman & Hall, 2004.
- Luong, T., Sutskever, I., Le, Q. V., Vinyals, O., and Zaremba, W. Addressing the rare word problem in neural machine translation. *arXiv preprint arXiv:1410.8206*, 2014.
- Mikolov, T. *Statistical Language Models based on Neural Networks*. PhD thesis, Brno University of Technology, 2012.
- Mikolov, T., Karafiát, M., Burget, L., Cernocký, J., and Khudanpur, S. Recurrent neural network based language model. In *INTERSPEECH*, pp. 1045–1048, 2010.
- Shang, L., Lu, Z., and Li, H. Neural responding machine for short-text conversation. In *Proceedings of ACL*, 2015.
- Sordoni, A., Galley, M., Auli, M., Brockett, C., Ji, Y., Mitchell, M., Gao, J., Dolan, B., and Nie, J.-Y. A neural network approach to context-sensitive generation of conversational responses. In *Proceedings of NAACL*, 2015.
- Sutskever, I., Vinyals, O., and Le, Q. V. Sequence to sequence learning with neural networks. In *NIPS*, 2014.
- Tiedemann, J. News from OPUS - A collection of multilingual parallel corpora with tools and interfaces. In Nicolov, N., Bontcheva, K., Angelova, G., and Mitkov, R. (eds.), *Recent Advances in Natural Language Processing*, volume V, pp. 237–248. John Benjamins, Amsterdam/Philadelphia, Borovets, Bulgaria, 2009. ISBN 978 90 272 4825 1.
- Turing, A. M. Computing machinery and intelligence. *Mind*, pp. 433–460, 1950.
- Vinyals, O., Kaiser, L., Koo, T., Petrov, S., Sutskever, I., and Hinton, G. Grammar as a foreign language. *arXiv preprint arXiv:1412.7449*, 2014a.
- Vinyals, O., Toshev, A., Bengio, S., and Erhan, D. Show and tell: A neural image caption generator. *arXiv preprint arXiv:1411.4555*, 2014b.
- Will, T. *Creating a Dynamic Speech Dialogue*. VDM Verlag Dr, 2007.