# Survey on Evaluation Methods for Dialogue Systems

**Authors:** Jan Deriu[1], Alvaro Rodrigo[2], Arantxa Otegi[3], Guillermo Echegoyen[2], Sophie Rosset[4], Eneko Agirre[3], Mark Cieliebak[1]

**Affiliation:** (1) Zurich University of Applied Sciences , (2) Universidad Nacional de Educación a Distancia , (3) IXA NLP Group, University of the Basque Country , (4) LIMSI, CNRS, Université Paris Saclay

# Abstract

In this paper we survey the methods and concepts developed for the evaluation of dialogue systems. Evaluation is a crucial part during the development process. Often, dialogue systems are evaluated by means of human evaluations and questionnaires. However, this tends to be very cost and time intensive. Thus, much work has been put into finding methods, which allow to reduce the involvement of human labour. In this survey, we present the main concepts and methods. For this, we differentiate between the various classes of dialogue systems (task-oriented dialogue systems, conversational dialogue systems, and question-answering dialogue systems). We cover each class by introducing the main technologies developed for the dialogue systems and then by presenting the evaluation methods regarding this class.

# Contents

# 1 Introduction

As the amount of digital data grows continuously, users demand technologies that offer a quick access to such data. In fact, users are relying on systems able to support an interaction for searching information such as SIRI[1], Google Assistant[2], Amazon Alexa[3] or Microsoft XiaoIce [Zhou *et al.*, 2018], etc. These technologies, called Dialogue systems (DS), allow the user to converse with a computer system using natural language. Dialogue Systems are applied to a variety of tasks, e.g.:

- Virtual Assistants, aid users in every-day tasks, such as scheduling appointments. They usually operate on predefined actions, which can be triggered by voice command.

- Information-seeking systems, provide users with information about a question (e.g. the most suitable hotel in town). These questions also include factoid questions as well as more complex questions.

- E-learning, where the dialogue systems train students for various situations. For instance, train the interaction with medical patients or train military personnel in questioning a witness.

One crucial step in the development of DS is evaluation. That is, to measure how well the DS is performing. However, evaluating a dialogue system is tricky because there are two important factors to be considered. First, the definition of what constitutes a high quality dialogue is not always clear and often depends on the application. Even if a definition is assumed, it is not always clear how to measure it. For instance, if we assume that a high quality dialogue system is defined by its ability to respond with an appropriate utterance, it is not clear how to measure appropriateness or what appropriateness means for a particular system. Moreover, one might ask the users if the responses were appropriate, but as we will discuss below user feedback might not always be reliable for various reasons.
The second factor is that the evaluation of dialogue systems is very cost and time intensive. This is especially true when the evaluation is carried out by a user study, which requires careful preparation, inviting and paying users to participate.

Over the past decades, many different evaluation methods have been proposed. The evaluation methods are closely tied to the characteristics of the dialogue system they aim to evaluate. Thus, quality is defined in the context of the functionality of the dialogue system. For instance, a system designed to answer questions will be evaluated on the basis of correctness, which is not necessarily a suitable metric for evaluating a conversational agent.

Most methods aim to automate the evaluation or at least automate certain aspects of

---

[1]https://www.apple.com/es/siri/
[2]https://assistant.google.com/
[3]https://www.amazon.com

the evaluation. The goal of an evaluation method is to have automated and repeatable evaluation procedures, which allow to efficiently compare the quality of different dialogue strategies.

The survey is structured as follows. In the next section we give a general overview over the different classes of dialogue systems and their characteristics. Then we introduce the evaluation task in more detail, with focus on the goals of an evaluation and the requirements on an evaluation metric. In Sections 3, 4, and 5, we introduce each dialogue system class (i.e. task-oriented systems, conversational agents and question answering dialogue systems). We give an overview of the characteristics, the dialogue behaviour, the ideas behind the methods used to implement the various dialogue systems and finally we present the evaluation methods and the ideas behind them. Here, we set a focus on how these methods are tied to the dialogue system class and what aspects of the evaluation are automated. In Section 6, we give a short overview of the relevant datasets and evaluation campaigns in the domain of dialogue systems. In Section 7, we discuss the issues and challenges of devising automated evaluation methods and discuss the state of automation achieved.

# 2   A general Overview

## 2.1   Dialogue Systems

Dialogue Systems (DS) usually structure dialogues in *turns*, each turn is defined by one or more *utterances* from one speaker. Two consecutive turns between two different speakers is called an *exchange*. Multiple exchanges constitute a *dialogue*. Another correlated view, is to interpret each turn or each utterance as an action (more on this later).

The main component of a dialogue system is the dialogue manager, which defines the content of the next utterance and thus the behaviour of the dialogue system. There are many different approaches to design a dialogue manager, which are partly dictated by the application of the dialogue system. However, there are three broad classes of dialogue systems, which we encounter in the literature: task-oriented systems, conversational agents and interactive question answering systems[4].

We identified the following characteristic features, which help differentiate between the three different classes: is the system developed to solve a task, does the dialogue follow a structure, is the domain restricted or is it open domain, does the dialogue span over multiple turns, are the dialogues rather long or efficient, who takes the initiative, and what is the interface used (text, speech, multi-modal). In Table 1 the characteristics for each

---

[4]In recent literature, the distinction is made only between the first two classes of dialogue systems [Serban *et al.*, 2017d; Chen *et al.*, 2017; Jurafsky and Martin, 2017]. However, interactive question answering systems cannot be completely placed in either of the two categories.

of the dialogue system classes is depicted. In the Table, we can see the following main features for each class:

- Task-oriented systems are developed to help the user solve a specific task as efficiently as possible. The dialogues are characterized by following a clearly defined structure, which is derived from the domain. The dialogues follow mixed initiative: both, the user and the system can take the lead. Usually, the systems found in the literature are built for speech input and output. However, task-oriented systems in the domain of assistance are built on multi-modal input and output.

- Conversational agents display a more unstructured conversation, as their purpose is to have open-domain dialogues with no specific task to solve. Most of these systems are built to emulate social interactions and thus longer dialogues are desired.

- Question Answering (QA) systems are built for the specific task of answering questions. The dialogues are not defined by a structure as with task-oriented systems, however they mostly follow the question and answer style pattern. QA systems may be built for a specific domain, but also be tilted towards more open domain questions. Usually, the domain is dictated by the underlying data, e.g. knowledge bases or text snippets from forums. Traditional QA systems work on a singe-turn interaction, however, there exist systems that allow multiple turns to cover follow-up questions. The initiative is mostly done by the user who asks questions.

|  | Task-oriented DS | Conversational Agents | Interactive QA |
|---|---|---|---|
| Task | Yes - clear defined | No | Yes - answer questions |
| Dial. Structure | Very structured | Not structured | No |
| Domain | Restricted | Mostly open domain | Mixed |
| Turns | Multi | Multi | Single/Multi |
| Length | Short | Long | - |
| Initiative | Mixed/ system init | mixed/user init | user init |
| Interface | multi-modal | multi-modal | mostly text |

Table 1: Characterizations of the different dialogue system types.

## 2.2 Evaluation

Evaluating dialogue systems is a challenging task and subject of much research. We define the goal of an evaluation method as having an automated, repeatable evaluation procedure with high correlation to human judgments, which is able to differentiate between various dialogue strategies and is able to explain which features of the dialogue systems are important. Thus, the following requirements can be stated:

- Automatic: in order to reduce the dependency on human labour, which is time and cost intensive as well as not necessarily repeatable, the evaluation method should be

automated.

- Repeatable: the evaluation method should yield the same result if applied various times to the same dialogue system under the same circumstances.

- Correlated to human judgments: the procedure should yield ratings, which correlate to human judgments.

- Differentiate between different dialogue systems: the evaluation procedure should be able to differentiate between different strategies. For instance, if one wants to test the effect of a *barge-in* feature, the evaluation procedure should be able to highlight the effects.

- Explainable: the method should give insights into which features of the dialogue system are correlated to the quality. For instance, the methods should reveal that the *word-error rate* of the automatic speech recognition system has a high influence on the dialogue quality.

In this survey, we focus on the efforts of automating the evaluation process. This is a very hard but crucial task, as human evaluations are cost and time intensive. Although much progress has been made in automating the evaluations of dialogue systems, the reliance on human evaluation is still present. Here, we give a condensed overview on the human based evaluations used in the literature.

**Human Evaluation.** There are various approaches to a human evaluation. The test subjects can take on two roles: they interact with the system, they rate a dialogue or utterance or they do both. In the following, we differentiate among different types of user populations, among each of the populations the subjects can take on any of the two roles.

- Lab-experiments: Before crowd sourcing was popular, the dialogue systems were evaluated in a lab environment. Users were invited to participate in the lab where they interacted with the dialogue system and subsequently filled a questionnaire. For instance, in [Young *et al.*, 2010] the authors recruited 36 subjects, which were instructed and presented with various scenarios. The subjects were asked to solve a task using a spoken dialogue system. Furthermore, a supervisor was present to guide the users. The lab environment is very controlled, which is not necessarily comparable to the real-world [Black *et al.*, 2011; Schmitt and Ultes, 2015].

- In-field experiments: Here, the evaluation is performed by collecting feedback from real users of the dialogue systems [Lamel *et al.*, 2000]. For instance, for the Spoken Dialogue Challenge [Black *et al.*, 2011], the systems were developed to provide schedule information in Pittsburgh. The evaluation was performed by redirecting the evening calls to the dialogue systems and getting the user feedback at the end of the

conversation. The Alexa Prize [5] also followed the same strategy, i.e. let real users interact with operational systems and gather the user feedback over a span of several months.

- Crowd-sourcing: Recently, the human evaluation has shifted from a lab environment to using crowd-sourcing platforms such as Amazon Mechanical Turk (AMT). These platforms provide large amounts of recruited users. In [Jurcícek *et al.*, 2011] the authors evaluate the validity of using crowd-sourcing for evaluating dialogue systems, their experiments suggest that using enough crowd-sourced users, the quality of the evaluation is comparable to the lab conditions. Current research relies on crowd-sourcing for human evaluation [Serban *et al.*, 2017c; Wen *et al.*, 2017].

Human based evaluation is difficult to set-up and to carry out. Much care has to be taken to setup the experiments: the users need to be properly instructed, the tasks need to be prepared so that the experiment is close to real-world conditions. Furthermore, one needs to take into account the high variability of user behaviour, which is especially present in crowd-sourced environments.

**Automated Evaluation Procedures**  A procedure which satisfies the aforementioned requirements has not yet been developed. Most evaluation procedures either require a degree of involvement of humans in order to be somewhat correlated to human judgement or they require a significant engineering effort. The methods for evaluation, which we cover in this survey can be categorized into: model the human judges, model the user behaviour or finer-grained methods, which evaluate a specific aspect of the dialogue system (e.g. its ability to stick to a topic). Methods modelling human judges rely on human judgements to be collected beforehand so as to fit a model which predicts the human rating. User behaviour models involve a significant engineering step in order to build a model which emulates the human behaviour. The finer-grained methods need a certain degree of engineering as well, which depends on the feature under evaluation. All the evaluation methods have in common that they depend on the characteristics of the dialogue system under consideration.

## 2.3   Modular Structure of this Article

Different evaluation procedures have been proposed based on the characteristics of the dialogue system class. For instance, the evaluation of task-oriented systems exploits the highly structured dialogues. The goal can be precisely defined and measured to compute the task-success rate. On the other hand conversational agents generate more unstructured dialogues, which can be evaluated on the basis of appropriateness of the responses, which has been shown to be hard to automate. We introduce each type of dialogue system to

---

[5]https://developer.amazon.com/alexaprize

highlight the respective characteristics and methods used to implement dialogue system. With this knowledge, we introduce the most important concepts and methods developed to evaluate the respective class of dialogue system. In the following survey, we discuss each of the three classes of dialogue systems separately. Thus Section 3: *Task Oriented Dialogue Systems*, Section 4: *Conversational Agents*, and Section 5: *Interactive Question Answering* can be read independently from each-other.

# 3 Task Oriented Dialogue System

## 3.1 Characteristics

As the name suggests, a task-oriented dialogue system is developed to perform a clearly defined task. These dialogue systems are usually characterized by a clearly defined and measurable goal, a structured dialogue behaviour, a closed domain to work on and a focus on efficiency. Usually, the task involves finding information within a database and returning it to the user, performing an action or retrieving information from user. For instance, a restaurant information dialogue system helps the user to find a restaurant which satisfies the users constraints. An in-car entertainment dialogue system can be ordered to start playing music via voice commands or querying the agenda (See Figure 1 for an example). The commonality is that the dialogue system infers the task constraints

| Event | Time | Date | Party |
|-------|------|------|-------|
| trau. appt. | 9am | March 15th | brother |
| lunch | 1pm | March 16th | Cathy |
| trau. appt. | 5pm | April 1st | Mike |
| trau. appt. | 10am | April 4th | Lisa |
| ... | ... | ... | ... |

**DRIVER:** I need to find out the time and parties attending my traumatologist appointment.

**CAR:** I have 3 appointments scheduled, with your brother, Mike and Lisa. Which are you referring to?

**DRIVER:** I want to know about the one that my brother is joining me at.

**CAR:** That appointment is at 9am

Figure 1: Example Dialogue where the driver can query the agenda via a voice command. The dialogue system guides the driver through the various options.

through the dialogue and retrieves the information requested by the user. For a ticket-reservation system, the dialogue system needs to know the origin station, the destination, the departure time and date. In most cases the dialogue system is designed for a specific

domain, such as restaurant information. The nature of these dialogue systems makes the dialogues very structured and tailored. The ideal dialogue satisfies the user goal with as little interactions as possible. The dialogues are characterized by mixed initiatives, the user states its goal but the dialogue system pro-actively asks questions to retrieve the required constraints.

## 3.2   Dialogue Structure

The dialogue structure for task-oriented systems is defined by two aspects: the content of the conversation and the strategy used within the conversation.

**Content.**   The content of the conversation is derived from the domain ontology. The domain ontology is usually defined as a list of slot-value pairs. For instance, in Table 2, the domain ontology for the restaurant domain is shown [Novikova *et al.*, 2017]. Each slot has a type and a list of values, which the slot can be filled with.

| Slot | Type | Example Values |
|------|------|----------------|
| name | verbatim string | Alimentum, .. |
| eatType | dictionary | restaurant, pub, coffee shop |
| familyFriendly | boolean | yes, no |
| food | dictionary | Italian, French, English, ... |
| near | verbatim string | Burger King |
| area | dictionary | riverside, city center |
| customerRating | dictionary | 1 of 5, 3 of 5, 5 of 5, low, average, high |
| priceRange | dictionary | <£20, £20-25, >£30 cheap, moderate, high |

Table 2:   Domain ontology of the E2E dataset [Novikova *et al.*, 2017]. There are eight different slots (or attributes), each has a type and a set of values it can take.

**Strategy.**   While the domain ontology defines the content of the dialogue, the strategy to fill the required slots during the conversation is modelled as a sequence of actions [Austin, 1962]. These actions are so-called *dialogue acts*. A dialogue act is defined by its type (e.g. inform, query, confirm, and housekeeping) and a list of arguments it can take. Each utterance corresponds to an action performed by an interlocutor.

Table 3 shows the dialogue acts proposed by [Young *et al.*, 2010].

For instance, the *inform* act is used to inform the user about its arguments, i.e. inform(food="French", area="riverside") informs the user that there is a French restaurant

| Dialogue Act | Description |
| --- | --- |
| hello($a = x, b = y, ..$) | Open a dialogue and give info $a = x, b = y, ..$ |
| inform($a = x, b = y, ..$) | Give information $a = x, b = y, ..$ |
| request($a, b = x, ..$) | Request value for $a$ given $b = x, ...$ |
| reqalts($a = x, ..$) | Request alternative with $a = x, ..$ |
| confirm($a = x, b = y, ..$) | Explicitly confirm $a = x, b = y, ..$ |
| confreq($a = x, ..., d$) | Implicitly confirm $a = x, ..$ and request value of $d$ |
| select($a = x, a = y$) | Select either $a = x$ or $a = y$ |
| affirm($a = x, b = y$) | Affirm and give further info $a = x, b = y, ..$ |
| negate($a = x$) | Negate and give corrected value $a = x$ |
| deny($a = x$) | Deny that $a = x$ |
| bye() | Close a dialogue |

Table 3: Dialogue acts proposed by [Young *et al.*, 2010].

at the riverside area. On the other hand, the request-act is used to request a value for a given list of slot-value pairs.

| Utterance | Dialogue Act |
| --- | --- |
| U: Hi, I am looking for somewhere to eat | hello(task = find,type=restaurant) |
| S: You are looking for a restaurant. What type of food? | confreq(type = restaurant,food) |
| U: I'd like an Italian somewhere near the museum. | inform(food = Italian,near=museum) |
| S: Roma is a nice Italian restaurant near the museum. | inform(name = "Roma", type = restaurant, food = Italian, near = museum) |
| U: Is it reasonably priced? | confirm(pricerange = moderate) |
| S: Yes, Roma is in the moderate price range. | affirm(name = "Roma", pricerange = moderate) |
| U: What is the phone number? | request(phone) |
| S: The number of Roma is 385456. | inform(name = "Roma", phone = "385456") |
| U: Ok, thank you goodbye. | bye() |

Table 4: Sample Dialogue and corresponding dialogue acts.

Table 4 shows an example dialogue with the corresponding dialogue acts. Each user utterance is translated into a dialogue act, and each dialogue act of the dialogue system is translated into an utterance in natural language. For instance, the utterance "Hi, I am looking for somewhere to eat" corresponds to the act of "hello". The parameters describe the task that the user intends to solve, i.e. find a restaurant. For a formal description of dialogue acts, refer to [Traum, 1999; Young, 2007].

## 3.3   Technologies

We have just seen that content and strategy are the two main aspects driving the structure of a dialogue, but their influence reaches down to the different functionalities making a classic dialogue system architecture. It is composed of several parts which are built around the idea of modelling the dialogue as a sequence of actions.

The central component is the so-called *dialogue manager*. It defines the dialogue policy, which consists in deciding which action to take at each dialogue turn. The input to
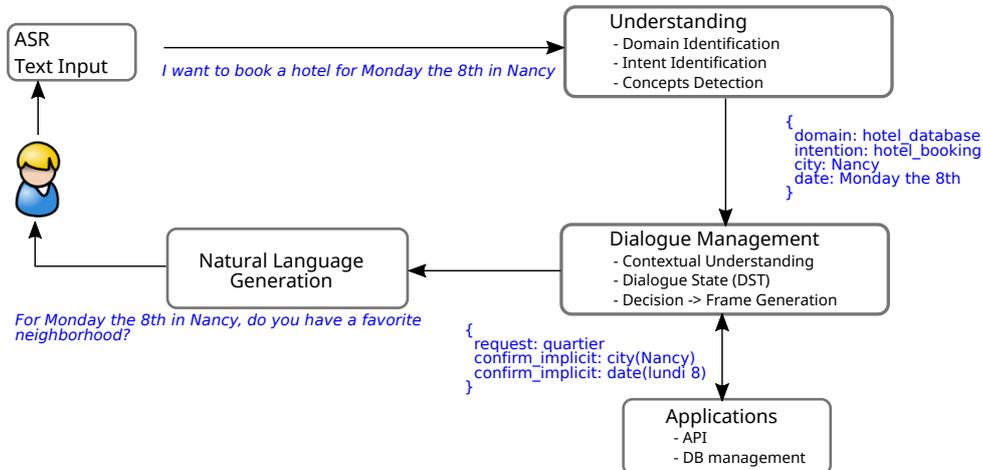
Figure 2: General overview of a task-oriented dialogue system.

the dialogue manager is the current state of the conversation. The output of the dialogue manager is a dialogue act, which represents the system's action. Other components convert the user's input into a dialogue act and the dialogue manager's output into a natural language utterance.

Usually, the user's input is processed by a natural language understanding (NLU) unit, which extracts the slots and their values from the utterance and identify corresponding the dialogue act. This information is passed to the dialogue state tracker (DST), which infers the current state of the dialogue.Finally the output of the dialogue manager is passed to a natural language generation (NLG) component.

Traditionally, these components were assembled into a pipelined architecture, but recent approaches based on end-to-end trainable neural networks offer a promising alternative. In the following, we briefly introduce the modules of the pipelined architecture and the deep neural network based approach.

### 3.3.1   Pipelined Systems

Usually, these four components are put into a pipelined architecture, where the output of one component is fed as the input of the next component (see Figure 2). The input of the dialogue system is either a chat-interface or an automatic speech recognition (ASR) system. The input to the NLU unit is the utterance of the user in text format or, in the case of automatic speech recognition (ASR) a list of the N-best last user utterance transcriptions.

**NLU**   The goal of the natural language understanding (NLU) unit is to detect the slot-value pairs expressed in the current user utterance. Since the early 2000s, the natural

language understanding task is often seen as a set of sub tasks [Tur and Mori, 2011]: (i) identification of domain (if multiple domains), (ii) identification of intents (that is the question type, the dialogue act, etc.) and (iii) identification of the slot.

In an utterance such as *I want to book a hotel room for Monday 8th*, the domain is *hotel*, the intent *hotel booking* and the slot-value pair is *date(Monday, 8th)*. The first two tasks are formalized as a classification task and all classification methods may be used. For concept detection one makes use of sequence labelling methods such as *Conditional Random Field* (CRF) [Hahn *et al.*, 2010] or recurrent neural network, typically bi-LSTM with CRF layer [Yao *et al.*, 2014; Mesnil *et al.*, 2015].

Before this split, all these tasks were done in one step, whatever the methods (rule-based, Hidden Markov Model, SVM, or CRF, etc.), also one can notice that some available data still used for NLU task do not include this split into three subtasks (see for example [Dinarelli *et al.*, 2017]). Recently, methods to learn joint model for intent detection and slot filling tasks have been proposed [Guo *et al.*, 2014; Zhang and Wang, 2016].

**Dialogue State Tracking**   The Dialogue State Tracker (DST) infers the current *belief state* of the conversation, given the dialogue history up to the current point $t$ [Williams *et al.*, 2016]. The current belief state encodes the user's goal (e.g. which price range the user prefers) and the relevant dialogue history, i.e. it is an internal representation of the state of the conversation. It is important to take the previous belief states into account in order to handle misunderstandings. For instance, in Figure 3 the confidence that the user wants an Italian restaurant is low. In the successive turn, the ASR system still gives low confidence to the Italian restaurant. However, since the state tracker takes into account that the Italian restaurant could have been mentioned in the previous turn, it assigns a higher overall probability to it.

The main challenge for the DST module is to handle the uncertainty, which stems from the errors made by the ASR module and the NLU unit. Typically, the output of the DST unit is represented as a probability distribution over multiple possible dialogue states $b(s)$, which provides a representation of the uncertainty. Generative methods have been widely used to manage this task, for example, dynamic Bayesian network (DBN) along with a beam search [Young *et al.*, 2007]. Those methods present some limits which are widely discussed in [Metallinou *et al.*, 2013], the most important being that all the correlations in the input features have to be modeled (even the unseen cases).

Discriminative models were then proposed to overcome these limits. [Metallinou *et al.*, 2013] proposed to use linear classifier where the dialogue history in the input features and [Henderson *et al.*, 2013] proposed to map directly the ASR hypotheses onto a dialogue state by means of recurrent neural networks which integrated both NLU and DST into a single function. Nowdays, neural approaches are becoming more and more popular [Mrkšić *et al.*, 2017].

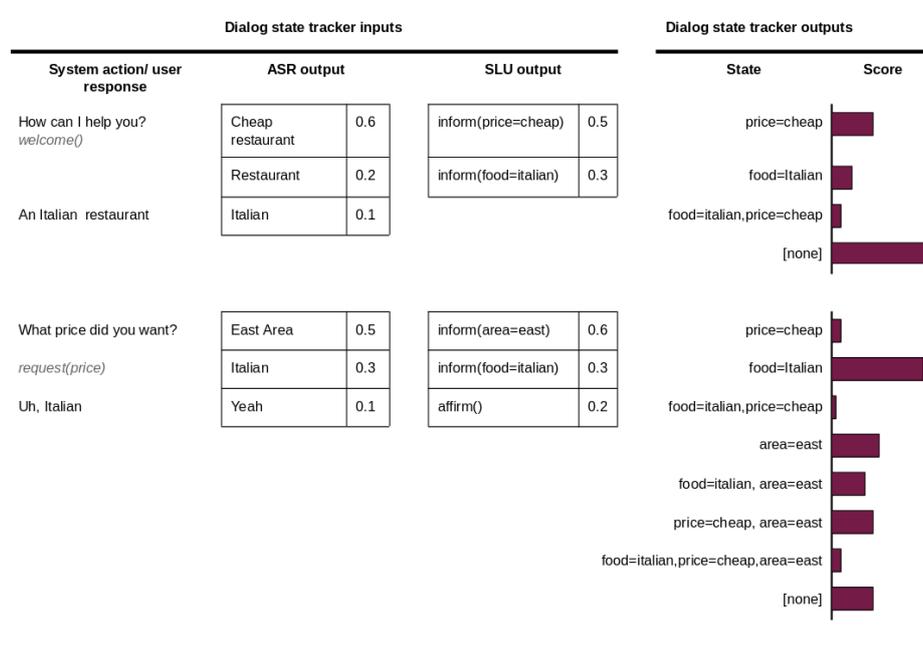| Dialog state tracker inputs | | | | | | Dialog state tracker outputs | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| System action/ user response | ASR output | | SLU output | | | State | Score |
| How can I help you? *welcome()* | Cheap restaurant | 0.6 | inform(price=cheap) | 0.5 | | price=cheap | |
| | Restaurant | 0.2 | inform(food=italian) | 0.3 | | food=Italian | |
| An Italian restaurant | Italian | 0.1 | | | | food=italian,price=cheap | |
| | | | | | | [none] | |
| What price did you want? | East Area | 0.5 | inform(area=east) | 0.6 | | price=cheap | |
| *request(price)* | Italian | 0.3 | inform(food=italian) | 0.3 | | food=Italian | |
| Uh, Italian | Yeah | 0.1 | affirm() | 0.2 | | food=italian,price=cheap | |
| | | | | | | area=east | |
| | | | | | | food=italian, area=east | |
| | | | | | | price=cheap, area=east | |
| | | | | | | food=italian,price=cheap,area=east | |
| | | | | | | [none] | |

Figure 3: Overview of a DST module. The input to the DST module is the combined output of the ASR and the NLU model.

**Strategy**   The strategy is learned by the dialogue manager. The input is the current belief state $b(s)$ computed by the DST module. The DM generates the next action of the system, which is represented as a dialogue act. In other words, based on the current turn values and on the value history the system performs an action (e.g. retrieve data from a database, ask for a missing information, etc.). Deciding which action to take is part of the dialogue control.

In earlier systems, the dialogue control was based on a finite state automate in which the nodes represent the questions of the system and the transitions the possible user' answers. This method, while being rigid, is efficient when the domain and the task are simple. It has been widely used to design dialogue systems and many toolkits are available such as the one from the Center for Spoken Language Understanding [Cole, 1999] or VoiceXML.[6] The main issue is the rigid dialogue structure as well as the tendency to be error-prone. In fact, such a system does not model discourse phenomena like ellipsis (a part of the sentence structure that can be inferred from the context is omitted) or anaphoric references (which can be resolved only in a given context).

To overcome these inefficiencies, a dialogue manager is designed which keeps track of the interaction history and controls the dialogue strategy. This is called a frame-based dialogue control and management. Frame-based techniques rely on schemas specifying what the system has to solve instead of representing what the system has to do and when. This

---

[6]See `https://www.w3.org/TR/voicexml20/`

allows for a more flexible dialogue and the possibility to handle errors [McTear *et al.*, 2005; van Schooten *et al.*, 2007].

Initially dialogue managers were implemented using rule-based approaches. When data had become available in sufficient amount, data-driven methods were proposed for learning dialogue strategies from data. The dialogue is represented as a Markov decision problem (MDPs), following the intuition that a dialogue can be represented as a sequence of actions [Levin *et al.*, 1998; Singh *et al.*, 2000]. These actions are referred to as *speech acts* or *dialogue acts*[Austin, 1962; Searle, 1969; Searle, 1975]. However, MDPs cannot handle uncertainty coming from speech recognition errors [Young *et al.*, 2013].

Thus, partially observable MDPs (POMDP) were adopted, as they introduce the belief state, which models the uncertainty of the current state [Paek, 2006; Lemon and Pietquin, 2012; Young *et al.*, 2013]. Although this alleviated the problem of hand-crafting the dialogue policy, the domain ontology still needs to be manually created. Furthermore, these dialogue systems are trained on a static and well-defined domain, once trained the policy works only on this domain. Finally, the dialogue systems need large amounts of data to be trained efficiently, mostly using a user simulation for training [Schatzmann *et al.*, 2006].

To mitigate the issues arising from the lack of data, [Gašić *et al.*, 2011] applied Gaussian process POMDP optimization [Engel *et al.*, 2005], which exploits the correlation between different belief states and speeds up the learning process. The authors showed, that a reasonable policy can be learned with on-line user feedback after a few hundreds of dialogues. In [Gasic *et al.*, 2013; Gasic *et al.*, 2014] the authors showed that it is possible to adapt the policy if the domain is extended dynamically. Note also the work of [Wang *et al.*, 2015] which aims at enabling domain-transfer by introducing a domain-independent ontology parametrisation framework.

**NLG**   The natural language generation module translates the dialogue act represented in a semantic frame into an utterance in natural language [Bangalore *et al.*, 2001]. The task of NLG is usually divided into separate sub tasks such as content selection, sentence planning, and surface realization [Stent *et al.*, 2004]. Traditionally, the task has been solved by relying on rule-based methods and canned texts. Statistical methods were also proposed and used, such as phrase-based NLG with statistical language model [Mairesse *et al.*, 2010] or CRF based on semantic trees [Dethlefs *et al.*, 2013]. Recently, deep learning techniques have become more prominent for NLG. With these techniques, there now exists a large variety of different network architectures, each addressing a different aspect of NLG: [Wen *et al.*, 2015] propose an extension to the vanilla LSTM [Hochreiter and Schmidhuber, 1997] to control the semantic properties of an utterance, whereas [Hu *et al.*, 2017] use variational autoencoder (VAE) and generative adversarial networks to control the generation of texts by manipulating the latent space; [Mei *et al.*, 2016] employ an encoder-decoder architecture extended by a coarse-to-fine aligner to solve the problem of content selection; [Wen *et al.*,

2016] apply data counter-fitting to generate out-of-domain training data for pretraining a model where there is little in-domain data available; [Semeniuta *et al.*, 2017; Bowman *et al.*, 2015] use a VAE trained in an unsupervised fashion on large amounts of data to sample texts from the latent space; and [Dušek and Jurcicek, 2016] use a sequence-to-sequence model with attention to generate natural language strings as well as deep syntax dependency trees from dialogue acts.

### 3.3.2   End-to-end trainable Systems

Traditionally, task-oriented dialogue systems were designed along the pipelined architecture, where each module has to be designed and trained separately. There are several drawbacks to this approach. As the architecture is modular, each component needs to be designed separately, which involves lots of hand-crafting, the costly generation of annotated data for each module, and training each component [Wen *et al.*, 2017]. Furthermore, the pipelined architecture leads to the propagation and amplification of errors through the pipeline as each module depends on the output of the previous module [Li *et al.*, 2017b; Liu *et al.*, 2018].

Related to the architecture there is a credit assignment problem, as the dialogue system is evaluated as a whole, it is hard to determine which module is responsible for which reward. Furthermore, this architecture leads to interdependence among the modules, i.e. when one module is changed, all the subsequent modules need to be adapted as well [Zhao and Eskenazi, 2016].

Finally, the slot-filling architecture, which is often used, makes these systems inherently hard to scale to new domains, since there is a need to handcraft the representation of the state and action space [Bordes *et al.*, 2017].

To overcome these limitations, current research focuses on end-to-end trainable architectures, where the dialogue system is trained as a single module. In [Wen *et al.*, 2017] the authors model the dialogue as a sequence to sequence mapping, where the traditional pipeline elements are modelled as interacting neural networks. The policy network takes as input the results form the intent network, the belief tracker network, the database operator and selects the next action, based on the selected action, the generation network produces the output utterance.

[Bordes *et al.*, 2017] propose a set of synthetic tasks to evaluate the feasibility of end-to-end models in the task-oriented setting, for which they use a memory network to model the conversation. These approaches learn the dialogue policy in a supervised fashion from the data. In contrast the work by [Li *et al.*, 2017b; Zhao and Eskenazi, 2016] train the system using reinforcement learning. Note that all these approaches rely on huge amounts of dialogue corpus.

## 3.4    Evaluation

The evaluation of task-oriented dialogue systems is built around the structured nature of the interaction. Two main aspects are evaluated, which have been shown to define the quality of the dialogue: task-success and dialogue efficiency. Two main flavours of evaluation methods have been proposed:

- User Satisfaction Modelling: here, the assumption is that the usability of the system can be approximated by the satisfaction of its users, which can be measured by questionnaires. These approaches aim to model the human judgements, i.e. creating models which give the same ratings as the human judges. First, a human evaluation is performed where subjects interact with the dialogue system. Afterwards, the dialogue system is rated via questionnaires. Finally, the ratings are used as target labels to fit a model based on objectively measurable features (e.g. task success rate, word error rate of the ASR system).

- User Simulation: Here, the idea is to simulate the behaviour of the users. There are two applications of user simulation: first, to evaluate a functioning system with the goal of finding weaknesses, second, the user simulation is used as an environment to train a reinforcement learning based system. The evaluation in the latter is based on the reward achieved by the dialogue manager under the user simulation.

Both these approaches rely on measuring task-success rate and dialogue efficiency. Before we introduce the methods themselves, we will go over the ways to measure performance along these two dimensions.

**Task-Success Rate.** The goal or the task of the dialogue can be split into two parts [Schatzmann *et al.*, 2007] (see Figure 4):

- Set of Constraints, which define the target information to be retrieved. For instance, the specifications of the venue (e.g. a bar in the central area, which serves beer) or the travel route (e.g. ticket from Torino to Milan at 8pm).

- Set of Requests, which define what information the user wants. For instance the name, address and the phone number of the venue.

The task-success rate measures how well the dialogue system fulfills the information requirements dictated by the user goals. For instance, this includes if the correct type of venue has been found by the dialogue system and if the dialogue system returned all the requested information. One possibility to measure this is via a confusion matrix (see Table 5), which represents the errors made over several dialogues. Based on this representation the Kappa coefficient [Carletta, 1996] can be applied to measure the success (see [Powers, 2012] for Kappa shortcomings).

$$C_0 = \begin{bmatrix} type = bar \\ drinks = beer \\ area = central \end{bmatrix}$$

$$R_0 = \begin{bmatrix} name = \\ addr = \\ phone = \end{bmatrix}$$

| attribute | actual value |
|---|---|
| depart-city | Torino |
| arrival-city | Milano |
| depart-range | evening |
| depart-time | 8pm |

Figure 4:   Examples of goals from [Schatzmann *et al.*, 2007] and [Walker *et al.*, 1997]

| | KEY | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DEPART-CITY | | | | ARRIVAL-CITY | | | | DEPART-RANGE | | DEPART-TIME | | | |
| DATA | v1 | v2 | v3 | v4 | v5 | v6 | v7 | v8 | v9 | v10 | v11 | v12 | v13 | v14 |
| v1 | **22** | | 1 | | 3 | | | | | | | | | |
| v2 | | **29** | | | | | | | | | | | | |
| v3 | 4 | | **16** | 4 | | | 1 | | | | | | | |
| v4 | 1 | 1 | 5 | **11** | | | 1 | | | | | | | |
| v5 | | | | | **20** | | | | | | | | | |
| v6 | | | | | | **22** | | | | | | | | |
| v7 | | | | | 1 | 1 | **20** | 5 | | | | | | |
| v8 | | | | | 1 | 2 | 8 | **15** | | | | | | |
| v9 | | | | | | | | | **45** | 10 | | | | |
| v10 | | | | | | | | | 5 | **40** | | | | |
| v11 | | | | | | | | | | | **20** | | 2 | |
| v12 | | | | | | | | | | | 1 | **19** | 2 | 4 |
| v13 | | | | | | | | | | | 2 | | **18** | |
| v14 | | | | | | | | | | | 2 | 6 | 3 | **21** |
| sum | 30 | 30 | 25 | 15 | 25 | 25 | 30 | 20 | 50 | 50 | 25 | 25 | 25 | 25 |

Table 5:   Confusion matrix from [Walker *et al.*, 1997]

**Dialogue Efficiency.**   Dialogue efficiency or dialogue costs are measures which are related to the length of the dialogue[Walker *et al.*, 1997] . For instance, the number of turns or the elapsed time are such measures. More intricate measures could include the number of inappropriate repair utterances or the number of turns required for a sub-dialogue to fill a single slot.

In the following, we introduce the most important research for both of the aforementioned evaluation procedures.

### 3.4.1   User Satisfaction Modelling

User satisfaction modelling is based on the idea that the usability of a system can be approximated by the satisfaction of its users. The research in this area is concerned with three goals: measure the impact of the properties of the dialogue system on the user satisfaction (explainability requirement), then automate the evaluation process based on these properties (automation requirement) and use the models to evaluate different dialogue strategies (differentiability requirement). Usually, a predictive model is fit, which takes the properties as input and uses the human judgements as target variable. Thus, modelling the user satisfaction as either a regression or a classification task. There are different approaches to measure the user satisfaction, which are based on two questions: who evaluates the dialogue and at which granularity is the dialogue evaluated? The first question allows for two groups: either the dialogue is evaluated by the users themselves or by objective judges. The second question allows for different points on a spectrum: on one end, the evaluation takes place on the dialogue level, on the other end the evaluation takes place at the exchange level. Especially, the question of who evaluates the dialogue is often at the centre of discussion. Here, we shortly summarize the main points.

**User or Expert ratings**    There are three main criticisms regarding the judgments made by users:

- Reliability: [Evanini *et al.*, 2008] state as main argument that users tend to interpret the questions on the questionnaires differently, thus making the evaluation unreliable. [Gašić *et al.*, 2011] noted that also in the lab setting, where users are given a predefined goal, users tend to forget the task requirements, thus, incorrectly assessing the task success. Furthermore, in the in-field setting, where the feedback is given optionally, the judgements are likely to be skewed towards the positive interactions.

- Cognitive demand: [Schmitt and Ultes, 2015] note that rating the dialogue puts more cognitive demand on users. This is especially true if the evaluation has to be done at the exchange level. This would falsify the judgments about the interaction.

- Impracticability: [Ultes *et al.*, 2013] note the impracticability of having a user rate the live dialogue, as he would have to press a button on the phone, or have a special installation to give feedback.

[Ultes *et al.*, 2013] analyzed the relation between the user ratings and ratings given by objective judges (called *experts*). Especially, they investigated if the ratings from the experts could be used to predict the ratings of the users. Their results showed that the user ratings and the expert ratings are highly correlated with a score Spearman's $\rho = 0.66(p < 0.01)$. Thus, expert ratings can be used as replacement for user judgments. Furthermore, they trained classifiers using the expert rating as targets and evaluated on the user ratings as targets. The best performing classifier achieved an unweighed average recall (UAR) of 0.34 compared to the best classifier trained on user satisfaction, which

achieved $UAR = 0.5$. These results indicate that it is not possible to precisely predict the user satisfaction. However the correlation scores show that the predicted scores of both models correlate equally to the user satisfaction $p = 0.6$. Although the models cannot be used to exactly predict the user satisfaction, the authors showed that the expert ratings are strongly related to user ratings.

In the following, we present different approaches to user satisfaction modelling. We cover the most important research for each of the various categories.

**PARADISE framework**    PARADISE (PARAdigm for DIalog System Evaluation) [Walker *et al.*, 1997] is the most known evaluation framework proposed for task-oriented systems. It is a general framework, which can be applied to any task-oriented system, since it is domain independent. It belongs to the evaluation methods which are based on user ratings on the dialogue level, although it allows for evaluations of sub dialogues.

Originally, the motivation was to produce an evaluation procedure, which can distinguish between different dialogue strategies. At that time the most widely used automatic approach was based on the comparison of utterances with a reference answer [Hirschman *et al.*, 1990]. Methods based on comparisons to reference answers suffer from various drawbacks: they cannot discriminate between different strategies, they are not capable to attribute the performance on system specific properties, and the approach is not generalizable to other tasks.

The main idea of PARADISE is to combine different measures of performance into a single metric, and in turn assess the contribution of each of these measures to the final user satisfaction. PARADISE originally uses two objective measures for performance: task-success and measures that define the dialogue cost (as explained above).

In Figure 5, an overview of the PARADISE framework is depicted. The user interacts with the dialogue system and completes a questionnaire after the dialogue ends. From the questionnaire, a user satisfaction score is computed, which is used as the target variable. The input variables to the linear regression models are extracted from the logged conversation data. The extraction can be done automatically (e.g. for task-success as discussed above) or manually by an expert (e.g. for inappropriate repair utterances). Finally, a linear regression model is fitted to predict the user satisfaction for a given set of input variables.

Thus, PARADISE models the (subjective) performance of the system with a linear combination of objective measures (task-success and dialogue costs). Applying multiple linear regressions showed that only the task-success measure and the number of repetitions are significant. In a follow-up study [Walker *et al.*, 2000] the authors further investigated PARADISEs ability to generalize to other systems and user populations and its predictive power. For this, they applied PARADISE on three different dialogue systems. In a large-scale user study they collected 544 dialogues over 42 hours of speech. For these experiments the authors worked with an extended number of quality measures: e.g. num-
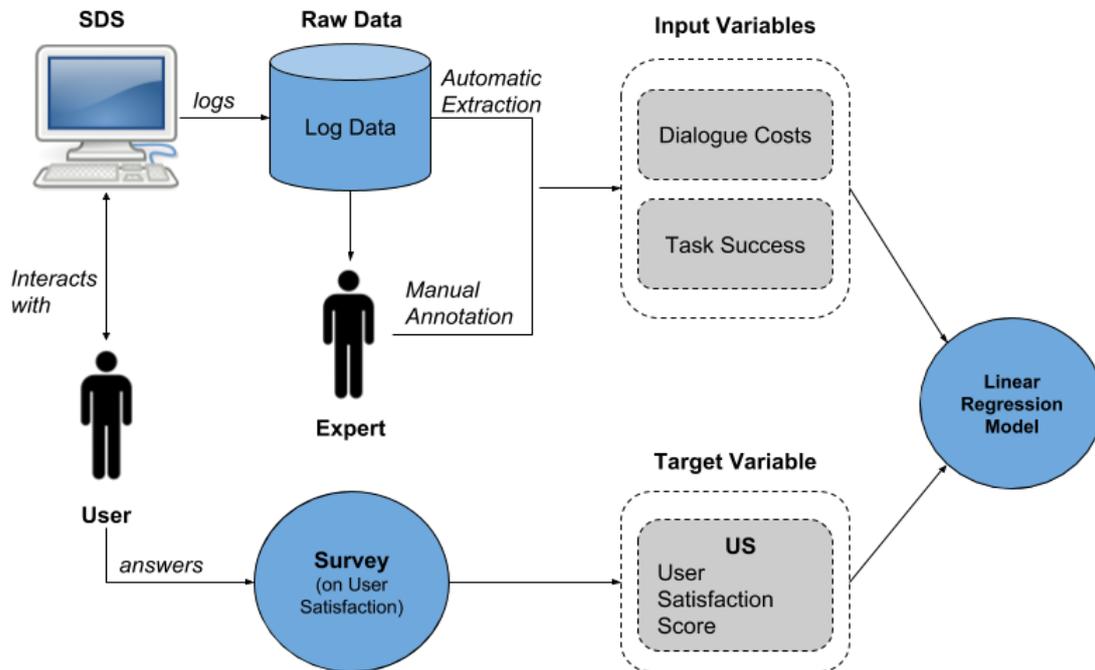
Figure 5: PARADISE Overview

ber of barge-ins, number of cancel operations, number of help requests. A survey at the end of the dialogue was used to measure the user satisfaction. The survey asked about various aspects: e.g. speech recognition performance, ease of the task, if the user would use the system again. Table 6 shows the generalization scores of PARADISE for different scenarios. According to these scores, we obtain the following observations:

| Training Set [a] | $R^2$ Training (SE) [b] | Test Set [c] | $R^2$ Test (SE) [d] |
|---|---|---|---|
| ALL 90% | 0.47 (0.004) | ALL 10% | 0.50 (0.035) |
| ELVIS 90% | 0.42 | TOOT | 0.55 |
| ELVIS 90% | 0.42 | ANNIE | 0.36 |
| NOVICES | 0.47 | ANNIE EXPERTS | 0.04 |

Table 6: Predictive power of PARADISE .

- A linear regression model is fitted on 90% of the data and evaluated on the remaining 10%. The results show that the model is able to explain $R^2 = 50\%$ of the variance, which is considered to be a good predictor by the authors.

- Training the regression model on the data for one system and evaluating the model on the data for another dialogue system (e.g. train on the ELVIS data and evaluate on the TOOT data) show high variability as well. The evaluation on the TOOT

system data yields much higher scores than evaluating on the ANNIE data. These results show that the model is able to generalize to data of other dialogue systems to a certain degree.

- The evaluation of the generalizability of the model across different populations of users yields a negative result. When trained on dialogue data from conversation by novice users, the linear model is not capable of predicting the scores by experienced users of the dialogue system.

The PARADISE framework is not only able to find the factors, which have the most impact on the rating, it is also capable of predicting the ratings. However, the experiments also revealed that the framework is not capable of distinguishing between different user groups. This result was confirmed by [Engelbrecht *et al.*, 2008], which tested the predictive power of PARADISE for individual users.

**User satisfaction at the exchange level**   In contrast to rate the dialogue as a whole, in some cases it is important to know the rating at each point in time. This is especially useful for online dialogue breakdown detection. There are two approaches to modelling the user satisfaction at the exchange level: annotate dialogues at the exchange level either by users [Engelbrecht *et al.*, 2009a] or by experts [Higashinaka *et al.*, 2010; Schmitt and Ultes, 2015]. Different models can be fitted with the sequential data: Hidden Markov Models (HMM), Conditional Random Fields or Recurrent Neural Networks are the most obvious choice but also SVM based approaches are possible.

In [Engelbrecht *et al.*, 2009a] the authors model user satisfaction as a continuous process evolving over time, where the current judgment depends on the current dialogue events and the previous judgments. Users interacted with the dialogue system and judged the dialogue after each turn on a 5-point scale using a number pad. Based on these target values and annotated dialogue features a HMM was trained. Some input features were manually annotated, which is not a reasonable setting for online break-down detection.

[Higashinaka *et al.*, 2010] modelled the evaluation similarly as in [Engelbrecht *et al.*, 2009a]. In their study they evaluated different models (HMM and CRF), different measures to evaluate the trained model, and addressed the question of subjectivity of the annotators. The input features to the model were the dialogue acts and the target variables were the annotations by experts, which listened to the dialogue. The low inter-rater agreement and the fact of only using dialogue acts as inputs made the model perform only marginally better than the random baseline.

A different approach was taken by [Hara, 2010] who relied on dialogue-level ratings but trained the model on n-grams of dialogue-acts. More precisely, they used as input features $n$ consecutive dialogue acts and used the dialogue-level rating as target variable (on a 5-point scale and an extra class to denote unsuccessful task). The model achieved an accuracy of only 34.4% using a 3-gram model. Further testing yielded that the model is

able to predict the task-success with an accuracy of 94.7%.

These approaches suffer from the following problems: they either rely on manual feature extraction, which is not useful for online breakdown detection or they used only dialogue acts as input features, which does not cover the whole dialogue complexity. Furthermore, the approaches had issues with data annotation, either having low inter-rater agreement or using dialogue-level annotation. [Schmitt and Ultes, 2015] addressed these issues by proposing *Interaction Quality* as approximation to user ratings at the exchange level.

**Interaction Quality**   Interaction Quality is a metric proposed by [Schmitt and Ultes, 2015] with the goal to allow the automatic detection of problematic dialogue situations. The approach is based on letting experts rate the quality of the dialogue at each point in time - the median rating of several expert ratings at the exchange level is called Interaction Quality.

Figure 6 shows the overview of the Interaction Quality procedure. The user interacts with the dialogue system and the conversation relevant data is logged. From the logs, the input variables are automatically extracted. The target variables are manually annotated by experts, from which the target variable is derived. Based on the input and target variables a support vector machine (SVM) is fitted.
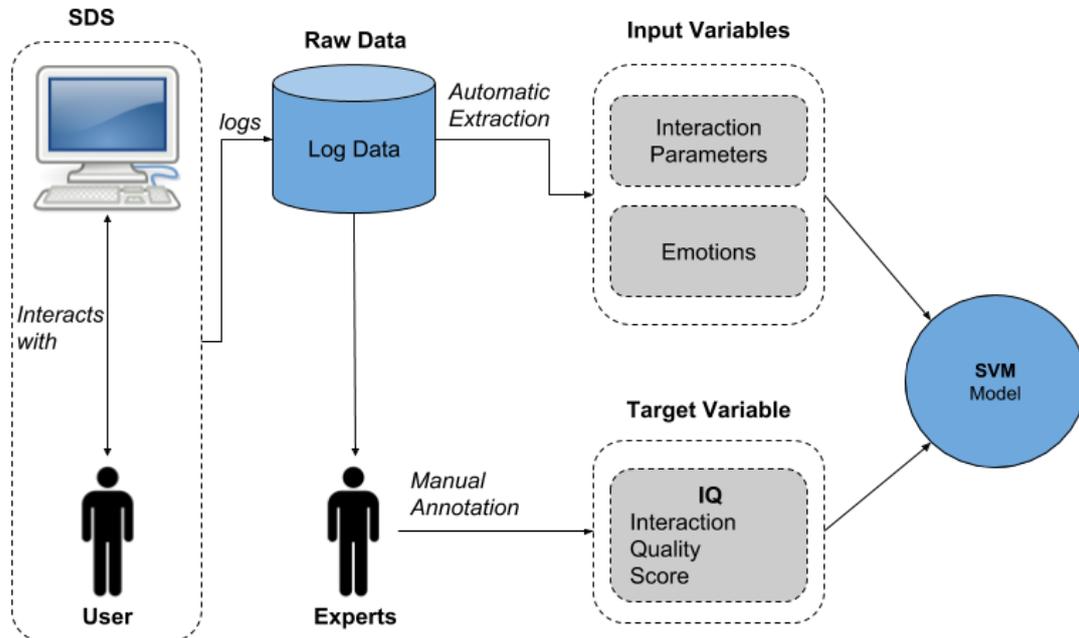


Figure 6:   Overview of the Interaction Quality procedure.

Interaction Quality is meant to approximate the user satisfaction. In this study the authors showed that Interaction Quality is an objective and valid approximation to user satisfac-

tion, which is easier to obtain. This is especially important for in-field evaluations of dialogue system, which are practically infeasible to be rated by users at the exchange level. Thus, it is important that in-field dialogues can be rated by experts at the exchange level. The challenge is to make sure that the ratings are objective, i.e. eliminate the subjectivity of the experts as much as possible.

Based on these Interaction Quality scores a predictive model is trained to automatically judge the dialogue at any point in time. Since there is no possibility to gather user satisfaction scores at the exchange level from in-field conditions, the authors relied on user satisfaction scores from lab experiments and Interaction Quality scores over dialogues from both in-field and lab conditions. The authors found a strong correlation (Spearman's $\rho = 0.66$) between Interaction Quality and user satisfaction in the lab environment, which means that Interaction Quality is a valid substitute for user satisfaction. A similarly strong correlation (Spearman's $\rho = 0.72$) exists between the Interaction Quality of the in-filed and the lab conditions.

In order to automatically predict Interaction Quality the input variable need to be automatically extractable from the dialogue system. From each subsystem of a task-oriented dialogue system (Figure 2) various values are extracted. Additionally, the authors experimented with hand-annotated features such as emotions and user specific features (e.g. age, gender, etc.) as well as semi-automatically annotated data such as the dialogue acts (similar to [Higashinaka *et al.*, 2010]). Based on these input variables the authors trained various SVM, one for each target variable, namely Interaction Quality for both in-field and the lab data as well as the user satisfaction label for the lab data. Table 7 shows the scores achieved for the various target variables and input feature groups.

|                    | $IQ_{field}$ | $IQ_{lab}$ | $US_{lab}$ |
|--------------------|--------------|------------|------------|
| ASR                | 0.753        | 0.811      | 0.625      |
| AUTO               | 0.776        | 0.856      | 0.668      |
| AUTO + EMO         | 0.785        | 0.856      | 0.669      |
| AUTO + EMO + USER  | -            | 0.888      | 0.741      |

Table 7: Model performance (in terms of $\rho$) on the test set. [Schmitt and Ultes, 2015].

The in-field Interaction Quality model ($IQ_{field}$) achieves a score of $\rho = 0.776$ on the automatically extracted features, with the ASR features alone the score lies at $\rho = 0.753$. The addition of the emotional and user -specific features do not increase the scores significantly. A similar behaviour is measured for the lab Interaction Quality model ($IQ_{lab}$), which achieves high scores with ASR features alone ($\rho = 0.856$) and profits only marginally from the inclusion of the emotional features. However, the model improves when including user specific features ($\rho = 0.894$). The lab based user satisfaction model ($US_{lab}$) achieves lower scores with $\rho = 0.668$ for the automatic features.

Table 8 shows the cross model evaluation. The $IQ_{field}$ can be used to predict $IQ_{lab}$ labels and vice versa ($\rho \sim 0.66$). Furthermore, the $IQ_{lab}$ model is able to predict the $US_{lab}$

| Feature set | Test | Train | $\rho$ |
|---|---|---|---|
| Auto | $US_{lab}$ | $IQ_{lab}$ | 0.667 |
| Auto | $IQ_{lab}$ | $IQ_{field}$ | 0.647 |
| Auto | $IQ_{field}$ | $IQ_{lab}$ | 0.696 |

Table 8: Model performance (in terms of $\rho$, $\kappa$ and UAR) on the test set. [Schmitt and Ultes, 2015]

variable. These results show that Interaction Quality is a good substitute to user satisfaction and that the models based on Interaction Quality yield high predictive performance when trained on the automatically extracted features. This allows to evaluate an ongoing dialogue in real time at the exchange level and ensures high correlation to the actual user satisfaction.

### 3.4.2 User Simulation

User Simulators (US) are tools, which are designed to simulate the users behaviour. There are two main applications for US: i) for training the dialogue manager in an off-line environment, and ii) to evaluate the dialogue policy.

**Training Environment**   User Simulations are used as a learning environment to train reinforcement learning based dialogue managers. They mitigate the problem of recruiting humans to interact with the systems, which is both time and cost intensive. There is a vast amount of literature on designing User Simulations as training environment, for a comprehensive survey refer to [Schatzmann *et al.*, 2006]. There are several considerations to be made when building a User Simulation.

- Interaction level: does the interaction take place at the semantic level (i.e. on the level of dialogue acts) or at the surface level (i.e. using natural language understanding and generation)?

- User goal: does the simulator update the goal during the conversation or not. The dialogues in the DSTC2 data [Henderson *et al.*, 2014] contain a large amount of examples where the user change their goal during the interaction. Thus, it is more realistic to model these changes as well.

- Error model: if and how to realistically model the errors made by the components of the dialogue system?

- Evaluation of the user simulation: for a discussion on this topic refer to [Pietquin and Hastie, 2013]. There are two main evaluation strategies: direct and indirect evaluation. The direct evaluation of the simulation are based on metrics (e.g. precision

and recall on dialogue acts, perplexity). The indirect evaluation measure the utility of the user simulation (e.g. by evaluating the trained dialogue manager).

The most popular approach to user simulation is based on the agenda based user simulation (ABUS) [Schatzmann *et al.*, 2007]. The simulations takes place at the semantic level, the user goal stays fixed throughout the interaction, and the user behaviour is represented as a priority ordered stack of necessary user actions. The ABUS was evaluated using indirect methods, by performing a human study on a dialogue system trained with the ABUS. The results show that the DS achieved an average task success rate of 90.6% based on 160 dialogues. The ABUS system works by randomly generating a hidden user goal (i.e. the goal is unknown to the dialogue system), which consists of constraints and request slots. From this goal, the ABUS system generates a stack of dialogue acts in order to reach the goal, which is the agenda. During the interaction with the dialogue system, the ABUS adapts the stack after each turn, e.g. if the dialogue system misunderstood something, the ABUS system pushes a negation act onto the stack.

Similar to other aspects of dialogue systems, more recent work is based on neural network based approaches. The Neural User Simulator (NUS) by [Kreyssig *et al.*, 2018] proposes an end-to-end trainable architecture based on neural networks. The system performs the interaction on the surface instead of the semantic level, during the training it considers variable user goals, and the evaluation is performed indirectly. The indirect evaluation is performed from two different perspectives. First, the dialogue system, which is trained with the NUS is compared to a dialogue system trained with ABUS in the context of a human evaluation. Here, the authors report the average reward and the success rate. In both cases the NUS-trained system performs significantly better. The second evaluation is performed in a cross-model evaluation [Schatztnann *et al.*, 2005], i.e. the NUS-trained dialogue system is evaluated using the ABUS system and vice-versa. Here, the NUS system performed significantly better as well. This indicates that the NUS system is diverse and realistic.

**Model Based Evaluation**  The idea of model based evaluation is to model the user behaviour but to put more emphasis on modelling a large variety of behavioural aspects. Here, the focus does not lie in the shaping of rewards for reinforcement learning, rather the focus lies on understanding the effects of different types of behaviour on the quality of the interaction. Furthermore, the goal is to gain insights on the effects of adapting a dialogue strategy, i.e. evaluate the changes made to the dialogue system. In [Engelbrecht *et al.*, 2009b] the MeMo workbench is introduced, which allows to model user simulations. The main focus is to model different types of users and typical errors the users make. In [Möller *et al.*, 2006] the authors introduced various types of conceptual errors, which users tend to make. There errors arise from the discrepancy between how the user expects the system to behave and the actual system behaviour. For instance:

- State errors arise when the user input cannot be interpreted in the current state, but

might be interpretable in a different state.

- Capability errors arise when the system cannot execute the users commands due to missing capability.

- Modelling errors arise due to discrepancies in how the user and the system model the world. For instance, when presented with a list of options and the system allows to address the elements in the list by their positions but the user addresses them by their name.

On the other hand, the workbench allows the definition of various user groups based on different characteristics of a user. The characteristics used in [Engelbrecht *et al.*, 2009b] include: affinity to technology, anxiety, problem solving strategy, domain expertise, age and deficits (e.g. hearing impairment). Behavioural rules are associated to each of the characteristics. For instance, a user with high domain expertise might use a more specific vocabulary. The rules are manually curated and are engineered to influence the probabilities of user actions. During the interaction, the user model selects a task to solve similar to the aforementioned approaches for reinforcement learning environments. In order to evaluate the user simulation, the authors compared the results of an experiment conducted with real users to the experiments conducted with the MeMo workbench. This evaluation procedure is aimed at finding whether the simulation yields the same insights as a user study. For this, they invited user from two user groups, namely older and younger users. The participants interacted with two version of a smart-home device control system: the versions differed in the way they provide help to the users. The comparison between the user simulation and the user study results was done at various levels:

- High-level features such as concept error rates (CER) or average number of semantic concepts per user turn (#) AVP. Here, the results show that the simulation was not always able to recreate the absolute values, it was able to replicate the relative results. Which is helpful, as it would lead to the same conclusions for the same questions.

- User judgment prediction based on a predictive model trained using the PARADISE framework. Here, the authors compared the real user judgments to the predicted judgments (where the linear model predicted the judgments of the simulated dialogue).Again, the results show that the user model would yield the same conclusions as the user study, namely that young users rated the system higher than the older users and that old user judged the dynamic help system worse than the other.

- Precision and Recall of predicted actions. Here, the simulation is used to predict the next user-action for a given context from a dialogue corpus. The predicted user action is compared to the real user action and based on this precision and recall is computed. The results show that precision and recall are relatively low.

The model based user simulations are designed with the idea of allowing the evaluation of dialogue system early in the development. Furthermore, they emphasize the need of interpretability, i.e. being able to understand how a certain change in the dialogue system

influences the quality of the dialogue. This lies in contrast to the user simulations for reinforcement learning, which are aimed at training a dialogue system and use the reward as a measure of quality. However, the reward is often only based on the task success and the number of turns.

# 4   Conversational Dialogue Systems

## 4.1   Characteristics

Conversational dialoge systems (also refferred to as chatbots, social bots) are usually developed for unstructured, open-domain conversations with its users. They are often not developed with a specific goal in mind, other than to maintain an engaging conversation with the user [Zhou *et al.*, 2018]. These systems are usually built with the intention to mimic human behaviour, which is traditionally assessed by the Turing Test (more on this later). However, Conversational dialogue systems might also be developed for practical applications. "Virtual Humans", for instance, are a class of conversational agents developed for training or entertainment purposes. They mimic certain human behaviours for specific situations. For instance, a Virtual Patient mimics the behaviour of a patient, which is then used to train medical students [Kenny *et al.*, 2009; Mazza *et al.*, 2018]. Early versions of conversational agents stem from the psychology community with *ELIZA* [Weizenbaum, 1966] and *PARRY* [Colby, 1981]. ELIZA was developed to mimic a Rogerian psychologist, whereas PARRY was developed to mimic a paranoid mind.

**Modelling Approaches**     Generally, there are two main approaches for modelling a Conversational dialogue system: *rule-based systems* and *corpus-based systems*.

Early systems, such as *ELIZA* [Weizenbaum, 1966] and *PARRY* [Colby, 1981] are based on a set of rules which determine their behaviour. ELIZA works on pattern recognition and transformation rules, which take the users input and apply transformations to it in order to generate responses.

Recently, Conversational dialogue systems are recently gaining a renewed attention in the research community, as shown by the recent effort to generate and collect data for the (RE-)WOCHAT workshops[7]. This renewed attention is motivated by the opportunity of exploiting large amounts of dialogue data (see [Serban *et al.*, 2018] for an extensive study and Section 6) to automatically author a dialogue strategy that can be used in conversational systems such as chatbots [Banchs and Li, 2012; Charras *et al.*, 2016]. Most recent approaches train Conversational agents in and end-to-end fashion using deep neural

---

[7]See `http://workshop.colips.org/re-wochat/` and `http://workshop.colips.org/wochat/`

networks, which mostly rely on the sequence-to-sequence architecture [Sutskever *et al.*, 2014].

In the following, we focus on the corpus-based approaches used to model conversational agents. First we describe the general concepts, and then the technologies used to implement conversational agents. Finally, we cover the various evaluation methods which have been developed in the research community.

## 4.2   Modelling Conversational Dialogue Systems

Generally, there are two different strategies to exploit large amounts of data:

- Utterance Selection: Here, the dialogue is modelled as an information retrieval task. A set of candidate utterances is ranked by relevance. The dialogue structure is thus defined by the utterances in a dialogue database [Lee *et al.*, 2009]. The idea is to retrieve the most relevant answer to a given utterance, thus, learning to map multiple semantically equivalent user-utterances to an appropriate answer.

- Generative Models: Here, the dialogue systems are based on deep neural networks, which are trained to generate the most likely response to a given conversation history. Usually the dialogue structure is learned from a large corpus of dialogues. Thus, the corpus defines the dialogue behaviour of the conversational agent.

Utterance selection methods can be interpreted as an approximation to generative methods. This approach is often used for modelling the dialogue system of Virtual Humans. Usually, the dialogue database is manually curated and the dialogue system is trained to map different utterances of the same meaning to the same response utterance. Another application of utterance selection is applied to integrate different systems [Serban *et al.*, 2017a; Zhou *et al.*, 2018]. Here the utterance selection system selects from a candidate list, which is comprised of outputs of different subsystems. Thus, given a set of dialogue systems, the utterance selection module is trained to select for the given context, the most suitable output from the various dialogue systems. This approach is especially interesting for dialouge systems, which work on a large number of domains and incorporate a large amount of skills (e.g. set alarm clock, report the news, return the current weather forecast). Here, we present the technologies for corpus-based approaches, namely the neural generative models and the utterance selection models.

### 4.2.1   Neural Generative Models

The architectures are inspired by the machine translation literature [Ritter *et al.*, 2011], especially neural machine translation. Neural machine translation models are based on the Sequence to Sequence (seq2seq) architecture [Sutskever *et al.*, 2014], which is composed of an encoder and a decoder. They are usually based on a Recurrent Neural Network (RNN).

The encoder maps the input into a latent representation which is used to condition the decoder on. Usually, the latent representation of the encoder is used as initial state of the recurrent cell in the decoder. The earliest approaches were proposed by [Shang *et al.*, 2015; Vinyals and Le, 2015], which trained a seq2seq model on a large amount of dialogue data. There are two fundamental weaknesses with the neural conversational agents: First, they do not take into account the context of the conversation. Since the encoder only reads the current user input, all previous states are ignored. This leads to dialogues, where the dialogue system does not refer to previous information, which might lead to nonsensical dialogues. And second, the models tend to generate generic answers, that follow the most common pattern in the corpus. This renders the dialogue monotonous and in the worst case leads to repeating the same answer, regardless of the current input. We briefly discuss these two aspects in the following.

**Context.**   The context of the conversation is usually defined as the previous turns in the conversations. It is important to take these into account as they contain information relevant to the current conversation. [Sordoni *et al.*, 2015] propose to model the context by adding the dialogue history as a bag of words representation. The decoder is then conditioned on the encoded user utterance and the context representation. An alternative approach was proposed by [Serban *et al.*, 2016] who proposed the hierarchical encoder decoder architecture (HRED), shown in Figure 7, which works in three steps:

1. A turn-encoder (usually a recurrent neural network) encodes each of the previous utterances in the dialogue history, including the last user utterance. Thus, for each of the preceding turns a latent representation is created.

2. A context-encoder (a recurrent neural network) takes the latent turn representations as input and generates a context representation.

3. The decoder is conditioned on the latent context representation and generates the final output.

**Variability.**   There are two main approaches on dealing with the issue of repetitive and universal responses:

- Adapt the loss functions. The main idea is to adapt the loss function in order to penalize generic responses and promote more diverse responses. [Li *et al.*, 2016a] propose two loss functions based on maximum mutual information: One is based on an anti-language model, which penalizes high-frequency words, the other is based on the probability of the source given the target. [Li *et al.*, 2016b] propose to train the neural conversational agent using the reinforcement learning framework. This allows to learn a policy which can plan in advance and generate more meaningful responses. The major focus is the reward function, which encapsulates various aspects: ease of answering (reduce the likelihood of producing a dull response), information flow
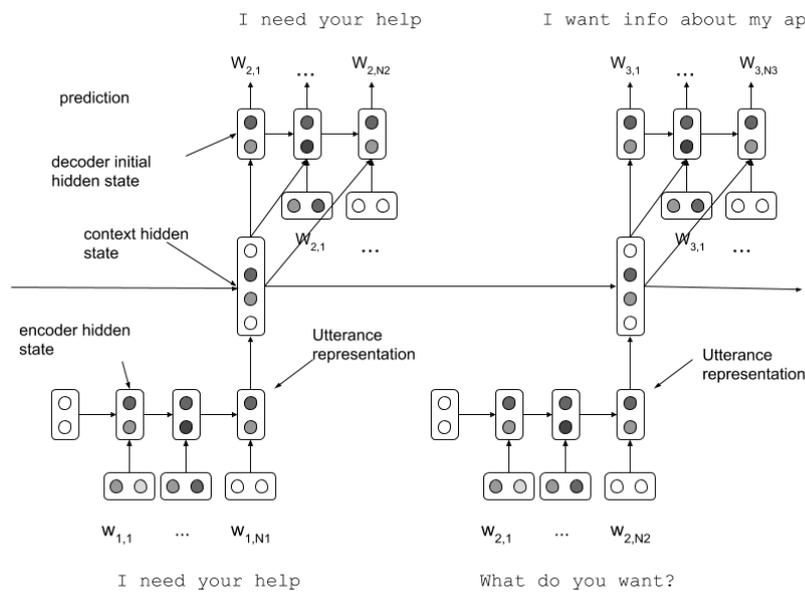
Figure 7:   Overview of the HRED architecture. There are two levels of encoding: (i) the utterance encoder, which encodes a single utterance and (ii) the context encoder, which encodes the sequence of utterance encodings. The decoder is conditioned on the context encoding.

(penalize answers which are semantically similar to a previous answer given), and semantic coherence (based on the mutual information).

- Condition the decoder. The seq2seq models perform a *shallow* generation process. This means that each sampled word is only conditioned on the previously sampled words. There are two methods for conditioning the generation process: condition on stochastic latent variables or on topics. [Serban *et al.*, 2017b] enhance the HRED model with stochastic latent variables at the utterance level and on the word level. At the decoding stage, first the latent variable is sampled from a multivariate normal distribution and then the output sequence is generated. In [Xing *et al.*, 2017] the authors add a topic-attention mechanism in their generation architecture, which takes as inputs *topic words* which are extracted using the Twitter LDA model [Zhao *et al.*, 2011]. The work by [Ghazvininejad *et al.*, 2017] extends the seq2seq model with a *Facts Encoder*. The "facts" are represented as a large collection of raw texts (Wikipedia, Amazon reviews,...), which are indexed by named entities.

### 4.2.2   Utterance Selection Methods

Utterance selection methods generally try to devise a similarity measure, which measures the similarity between the dialogue history and the candidate utterances. There are roughly

three different types of such measures:

- *Surface form similarity.* This measures the similarity on token level. This include measures such as: Levenshtein distance, METEOR [Lavie and Denkowski, 2009], or Term Frequency-Inverse Document Frequency (TF-IDF) retrieval models [Charras *et al.*, 2016; Dubuisson Duplessis *et al.*, 2016]. For instance, in [Dubuisson Duplessis *et al.*, 2017], the authors propose an approach that exploits recurrent surface text patterns to represent dialogue utterances.

- *Multi-class classification task.* These methods model the selection task as a multi-class classification problem, where each candidate response is a single class. For instance, [Gandhe and Traum, 2013] model each utterance as a separate class, and the training data consists of utterance-context pairs on which features are extracted. Then a perceptron model is trained to select the most appropriate response utterance. This approach is suitable for applications with a small amount ($\sim 100$) of candidate answers.

- *Neural network based approaches.* To leverage large amounts of training data, neural network architectures were introduced. Usually, they are based on a siamese architecture, where both the current utterance and a candidate response are encoded. Based on this representation a binary classifier is trained to distinguish between relevant responses and irrelevant. One well-known example is the dual encoder architecture by [Lowe *et al.*, 2017b]. Dual Encoders transform the user input and a candidate response into a distributed representation. Based on the two representations a logistic regression layer is trained to classify the pair of utterance and candidate response as either relevant or not. The softmax score of the relevant class is used to sort the candidate responses. The authors experimented with different neural network architectures for modelling the encoder, such as recurrent neural networks or long short-term memory networks [Hochreiter and Schmidhuber, 1997].

## 4.3   Evaluation Methods

Automatically evaluating conversational dialogue systems is an open problem. The difficulty to automate this step can be attributed to the characteristics of the conversational dialogue system. Without a clearly defined goal or task to solve and a lack of structure in the dialogues, it is not clear which attributes of the conversation are relevant to measure the quality. Two common approaches to assess the quality of a conversational dialogue system is to measure the appropriateness of its responses, or to measure the human likeness. Both these approaches are very coarse-grained and might not reveal the complete picture. Nevertheless, most approaches to evaluation follow these principles. Depending on the characteristics of a specific dialogue system, more fine-grained approaches to evaluation can be applied, which measure the capability of the specific characteristic. For instance, a system built to increase the variability of its answers might be evaluated based on lexical

complexity measures. In the following, we introduce the automated approaches for evaluating conversational dialogue systems. In the first part, we discuss the general metrics which can be applied to both the generative models as well as for the selection based models. We then survey the approaches specifically designed for the utterance selection approaches, as they can exploit various metrics from information retrieval.

### 4.3.1  General Metrics for Conversational Dialogue Systems

To evaluate a conversational dialogue system there are generally two levels: coarse-grained evaluations and fine-grained evaluations. The coarse-grained evaluations focus on the adequacy of the responses generated or selected by the dialogue system. On the other hand, fine-grained evaluations focus on specific aspects of its behaviour. Coarse-grained evaluations are based on two concepts: adequacy (or appropriateness) of a response, and the human likeness. Fine-grained evaluations focus on specific behaviours that a dialogue system should manifest. Here, we focus on the methods devised for coherence and the ability of maintaining the topic of a conversation. In the following, we give an overview over the methods, which have been designed to automatically evaluate the above dimensions.

**Appropriateness**    is a coarse-grained concept to evaluate a dialogue, as it encapsulates many finer-grained concepts, e.g. coherence, relevance, or correctness among others. There are two main approaches in the literature: word-overlap based metrics and methods based on predictive models inspired by the PARADISE framework (see Section 3.4.1).

- *Word-overlap metrics* were originally proposed by the machine translation (MT) and the summarization community. They were initially a popular choice of metrics for evaluating dialogue systems, as they are easily applicable. Popular metrics such as BLEU score [Papineni *et al.*, 2002] and ROUGE [Lin, 2004] were used as approximation for the appropriateness of an utterance. However, the authors of [Liu *et al.*, 2016] showed that neither of the word-overlap based scores have any correlation to human judgments.
  Based on the criticism of the word-overlap metrics, several new metrics have been proposed. In [Galley *et al.*, 2015] the authors propose to include human judgments into the BLEU score, which they call $\Delta BLEU$. The human judges rated the reference responses of the test set according to the relevance to the context. The ratings are used to weight the BLEU score to reward high-rated responses and penalize low-rated responses. The correlation to human judgments was measured by means of Spearman's $\rho$. $\Delta BLEU$ has a correlation of $\rho = 0.484$, which is significantly higher than the correlation of the BLEU score, which lies at $\rho = 0.318$. Although this increases the correlation of the metric to the human judgments, this procedure involves human judgments to label the reference sentences.

- Model based methods: In [Lowe *et al.*, 2017a] the authors present ADEM, a recurrent

neural network trained to predict appropriateness ratings by human judges. The human ratings were collected via Amazon Mechanical Turk (AMT), where the judges were presented with a dialogue context and a candidate response, which they rated on appropriateness on a scale from 1 to 5. Based on the ratings, a recurrent neural network (RNN) was trained to score the model response, given the context and the reference response. The Pearson's correlation between ADEM and the human judgments is computed on two levels: the utterance level and at the system level, where the system level rating is computed as the average score at the utterance-level achieved by the system.

The Pearson's correlation for ADEM lies at 0.41 on the utterance level and at 0.954 on the system level. For comparison, the correlation to human judgments for the ROUGE score only lies at 0.062 on the utterance level and at 0.268 at the system level.

**Human Likeness**   The classic approach to measure the quality of a conversational agent is the Turing Test devised by [Turing, 1950]. The idea is to measure if the conversational dialogue system is capable of fooling a human into thinking that it is a human as well. Thus, according to this test, the main measure is the ability to imitate human behaviour. Inspired by this idea, the use of *adversarial learning* [Goodfellow *et al.*, 2014] can be applied to evaluate a dialogue system. The framework of a generative adversarial model is composed of two parts: the generator, which generates data, and the discriminator, which tries to distinguish whether the data is real, or artificially generated. The two components are trained in an adversarial manner: the generator tries to fool the discriminator, and the discriminator learns at the same time to identify if the data is real or artificial. Adversarial Evaluation of dialogue systems was first studied by [Kannan and Vinyals, 2016], where the authors trained a generative adversarial network (GAN) on dialogue data, and used the performance of the discriminator as indicator for the quality of the dialogue. The discriminator achieved an accuracy of 62.5% which indicates a weak generator. However, the authors did not evaluate whether the discriminator score is a viable metric for evaluating a dialogue system.

A study on the viability of adversarial evaluation was conducted by [Bruni and Fernandez, 2017]. For this, they compared the performance of discriminators to the performance of humans on the task of discriminating between real and artificially generated dialogue excerpts. Three different domains were used, namely: MovieTriples (46k dialogue passages) [Serban *et al.*, 2016], SubTle (3.2M dialogue passages) [Banchs, 2012] and Switchboard (77k dialogue passages) [Godfrey *et al.*, 1992]. The GAN was trained on the concatenation of the three datasets. The evaluation was conducted on 900 dialogue passages, 300 per dataset, which were rated by humans as real or artificially generated. The results show that the annotator agreement among humans was low, with a Fleiss [Fleiss, 1971] $\pi = 0.3$, which shows that the task is difficult. The agreement between the discriminator and the humans

is on pair with the agreement among the humans, except for the Switchboard corpus, where $\pi = 0.07$. Human annotators achieve an accuracy score w.r.t. the ground-truth of $64\% - 67.7\%$ depending on the domain. The discriminator achieves lower accuracy scores on the Switchboard dataset but higher scores than humans on the other two datasets.

In order to evaluate the discriminators ability on different models, a Seq2Seq model was trained on the OpenSubtitles dataset [Tiedemann, 2009] (80M dialogue passages). The discriminator and the human performance on the dialogues generated by the Seq2Seq model was evaluated. The results show that the discriminator performs better than the humans, which the authors attribute to the fact that the discriminators may pick up on patterns which are not apparent to humans. The agreement between humans and the discriminator is very low.

**Fine-grained Metrics**   The above methods for evaluating conversational dialogue systems work on a coarse-grained level. The dialogue is evaluated on the basis of producing adequate responses or its ability to emulate human behaviour. These concepts encompass more finer-grained concepts. In this section, we look at topic-based evaluation.

*Topic-based evaluation* measures the ability of a conversational agent to talk about different topics in a cohesive manner. In [Guo *et al.*, 2018] the authors propose two dimensions of topic-based evaluation: topic breadth (can the system talk about a large variety of topics?) and topic depth (can the system sustain a long and cohesive conversation about one topic?). For topic classification, a Deep Averaging Network (DAN) was trained on a large amount of question data. Deep Averaging Networks do topic classification and the detection of topic-specific keywords. The conversational data used to evaluate the topic-based metrics stems from the Alexa-Prize challenge [8], which consists of millions of dialogues and hundred of thousands of live user ratings (on a scale from 1 to 5). Using the DAN, the authors classified the dialogue utterances according to the topics.

Conversational *topic depth* is measured by the average length of a sub-conversation on a specific topic, i.e. multiple consecutive turns where the utterances are classified as being the same topic. The conversational breadth is measured on a coarse grained and fine grained level. Coarse grained topic breadth is measured as the average number of topics a bot converses about during a conversation. On the other hand, *topic breadth* measures looks at the total number of distinct topic keywords across all conversations.

To measure the validity of the proposed metrics, the correlations between the metric and the human judgments is computed. The conversational topic depth metric has a correlation of $\rho = 0.707$ with the human judgments. The topic breadth metric has a correlation of $\rho = 0.512$ with the human judgments. The lower correlation of the topic breadth is

---

[8]https://developer.amazon.com/alexaprize

attributed to the fact that the users may not have noticed a bot repeating itself as they only conversed with a bot a few times.

### 4.3.2  Utterance Selection Metrics

The evaluation of dialogue systems based on utterance selection differs from the evaluation of generation-based dialogue systems. Here, the evaluation is based on metrics used in information retrieval, especially Recall@k (R@k). R@k measures the percentage of relevant utterances among the top-k selected utterances. One major drawback of this approach is that potentially correct utterances among the candidates could be regarded as incorrect.

**Next Utterance Selection.**  In [Lowe *et al.*, 2016] the authors evaluate the impact of this limitation and evaluate whether the Next Utterance Classification (NUC) task is suitable to evaluate dialogue systems. For this, they invited 145 participants from Amazon Mechanical Turk (AMT) and 8 experts from their lab. The task was to select the correct response given a dialogue context (of at most six turns) and five candidate utterances, of which exactly one is correct. Note that the other four utterances could also be relevant but are regarded as incorrect in this experiment. The study was performed on dialogues of three different domains: the SubTle Corpus [Banchs, 2012] consisting of movie dialogues, the Twitter Corpus [Ritter *et al.*, 2010] consisting of user dialogues, and the Ubuntu Dialogue Corpus [Lowe *et al.*, 2015], which consists of conversations about Ubuntu related topics.

The human performance was compared to the performance of an artificial neural network (ANN), which is trained to solve the same task. The performance was measured by means of R@1 score. The results show that for all domains, the human performance was significantly above random, which indicates that the task is feasible. Furthermore, the results show that the human performance varies depending on the domain and the expertise level. In fact, the lab participants performed significantly better on the Ubuntu domain, which is regarded as harder as it requires expert knowledge. This shows that there is a range of performance which can be achieved. Finally, the results showed that the ANN achieved similar performance to the human non-experts and performed worse than the experts. This shows that this task is not trivial and by far not solved. However, the authors did not take into account the fact that multiple candidates responses could be regarded as correct. This is possible since the selection of the candidate response is performed by sampling at random from the corpus. On the other hand, it is not clear if their evaluation suffered from this potential limitation, as their results showed the feasibility and relevance of the NUC task.

The authors of [DeVault *et al.*, 2011] and [Gandhe and Traum, 2016] tackle the problem of having multiple relevant candidate utterances and propose a metric which takes

this into account. Their metrics are both dependent on human judges and measure the appropriateness of an utterance.

**Weak Agreement**   The authors of [DeVault *et al.*, 2011] propose the *weak agreement* metric. This metric is based on the observation that human judges only agree in about 50% of the cases on the same utterance for a given context. The authors attribute this to the fact that multiple utterances could be regarded as acceptable choices. Thus, the weak agreement metric regards an utterance as appropriate if at least one annotator chose this utterance to be appropriate.

The authors apply the weak agreement metric on the evaluation of a virtual human which simulates a witness in a war-zone and is designed to train military personnel in Tactical Questioning [Gandhe *et al.*, 2009]. They gathered 19 dialogues and 296 utterances in a Wizard-of-Oz experiment. To allow for more diversity, they let human experts write paraphrases of the commander-role to ensure that the virtual character understands a larger variety of inputs. Furthermore, the experts expanded the set of possible answers by the virtual character by annotating other candidate utterances as appropriate.

The weak agreement metric was able to measure the improvement of the system when the extended dataset was applied: The simple system based on the raw Wizard-of-Oz data achieved a weak agreement of 43%; augmented with the paraphrases, the system achieved a score of 56%; and, finally, adding the manual annotation increases the score to 67%. Thus, the metric is able to measure the improvements made by the variety in the data.

**Voted Appropriateness**   One major drawback of the weak agreement is that it depends on human annotations and is not applicable to large amounts of data. The authors of [Gandhe and Traum, 2016] improve upon the idea of weak agreement by introducing the *Voted Appropriateness* metric. Voted Appropriateness takes the number of judges into account which selected an utterance for a given context. In contrast to weak agreement, which regarded each adequate utterance equally, Voted Appropriateness weights each utterance.

Similarly to the PARADISE approach, the authors of Voted Appropriateness fit a linear regression model on the pairs of utterances and contexts labelled with the amount of judges that selected the utterance. The fitted model only explains 23.8% of the variance. The authors compared the correlation of the Voted Appropriateness and the weak agreement metric to human judgments. The correlation was computed on the individual utterance level and the system level. For the system level, the authors used data from 7 different dialogue systems and averaged the ratings over all dialogues of one system. On the interaction level, the Voted Appropriateness achieved a correlation score of $0.479$ ($p < 0.001, n = 397$), and the weak agreement achieved $0.485$ ($p < 0.001, n = 397$). On the system level, Voted Appropriateness achieved $0.893$ ($p < 0.01, n = 7$) and weak agreement achieved $0.803$ ($p < 0.001, n = 397$). Thus, on the system level Voted Appropriateness performs closer

to human judgments. Both metrics rely heavily on human annotations, which makes the metrics hardly suitable for large-scale data driven approaches.

# 5   Question Answering Dialogue Systems

A different form of task-oriented systems are Question Answering systems. Here, the task is defined as finding the correct answer to a question. This setting differs from the aforementioned task-oriented systems in following ways:

- Task-oriented systems are developed for a multitude of tasks (e.g. restaurant reservation, travel information system, virtual assistant, ..),whereas the QA systems are developed to find answers to specific questions.

- Task-oriented systems are usually domain-specific, i.e. the domain is defined in advance through an ontology and remains fixed. In contrast, QA systems usually work on broader domains (e.g. factoid QA can be done over different domains at once), although there are also some QA systems focused only on a specific domain [Sarrouti and Ouatik El Alaoui, 2017; Do *et al.*, 2017].

- The dialogue aspect for QA systems is not tailored to sound human-like, rather the focus is set on the completion of the task. That is, to provide a correct answer to the input question.

Generally, QA systems allow the users to search for information using a natural language interface, and return short answers to the users' question [Voorhees, 2006]. QA systems can be broadly categorized into three categories [Bernardi and Kirschner, 2010]:

- Single-turn QA: this is the most common type of system. Here, the system is developed to return a single answer to the users' question without any further interaction. These systems work very well for factoid questions [Voorhees, 2006]. However, they have difficulties handling complex questions, which require several inference steps [Iyyer *et al.*, 2017] or situations where systems need additional information from the user [Li *et al.*, 2017a].

- Context QA: These are systems which allow for follow-up questions to resolve ambiguities or keeping track of a sequence of inference steps. Thus, new questions can refered to entities in previous questions [Peñas *et al.*, 2012].

- Interactive QA (IQA): These systems combine context QA systems and task-oriented dialogue systems. The main purpose of the conversation module is to handle under-or-over constrained questions [Qu and Green, 2002]. E.g. if a question does not yield any results, the system might propose to relax some constraints. In contrast, if a question yields to many results, the interaction can be used to introduce new

constraints to filter a list of results [Rieser and Lemon, 2009]. For a more in-depth discussion on IQA systems, refer to [Konstantinova and Orasan, 2013].

There is a large amount of research in the area of single-turn QA and there are several survey, we refer the reader to [Kolomiyets and Moens, 2011; Diefenbach *et al.*, 2018; Mishra and Jain, 2016]. In this survey, we focus on the evaluation of multi-tun QA systems, which is a much less researched area.

## 5.1   Evaluation of QA Dialogue Systems

The nature of multi-turn QA systems makes quite hard to design accurate evaluation frameworks. In fact, a proper evaluation of multi-turn QA systems requires humans to interact with the systems. However, the metrics used to assess the quality are often based on metrics used in Information Retrieval (IR). For instance, [Li *et al.*, 2017a] report the error rate of the system. Also [Kelly and Lin, 2007] base the valuation on F-scores, which is computed over "information nuggets". These nuggets are retrieved by the assessors of the system, and thus, this evaluation method is dependent on heavy human involvement. Both methods do not take into consideration the dialogue aspect of the interaction. The first evaluation framework designed specifically for IQA systems is based on a series of questionnaires [Kelly *et al.*, 2009] to capture different aspects of the system. The authors argue that metrics based on the relevance of the answers are not sufficient to evaluate an IQA system (e.g. it does not take the user feedback into account). Thus, they evaluate the usage of different questionnaires in order to assess the different systems. The questionnaires they propose are:

- NASA TLX (cognitive workload questionnaire): used to measure the cognitive workloads as subjects completed different scenarios.

- Task Questionnaire: after each task the questionnaire is filled out, which focuses on the experiences of using a system for a specific task.

- System Questionnaire: compiled after using a system for multiple tasks. This measures the overall experiences of the subjects.

Their evaluation showed that the Task Questionnaire is the most effective at distinguishing among different systems. To the best of our knowledge, there are no other evaluation frameworks available yet.

## 6   Evaluation Datasets and Challenges

Datasets play an important role for the evaluation of dialogue systems, together with challenges open to public participation. A large number of datasets have been used and made publicly available for the evaluation of dialogue systems in the last decades, but the

coverage across dialogue components and evaluation methods (cf. Sections 3 and 4) is uneven. Note also that datasets are not restricted to specific evaluation methods, as they can be used to feed more than one evaluation method or metric interchangeably. In this section, we cover the most relevant datasets and challenges, starting with some selected datasets. For further references, see a wide survey of publicly available datasets that have already been used to build and evaluate dialogue systems carried out by [Serban *et al.*, 2018].

## 6.1 Datasets for Task-oriented systems

Datasets are usually designed to evaluate some specific dialogue components, and very few public datasets are able to evaluate a complete **task-oriented dialogue system** (cf. Section 3). The evaluation of these kind of systems is very system specific, and it is therefore hard to reuse the dataset with other systems. They also require high human effort, as the involvement of individual users or external evaluators is usually needed. For example, in [Gasic *et al.*, 2013], which is a POMDP-based dialogue system mentioned in Section 3.3.1 for the restaurants domain, the evaluation of policies is done by crowd-sourcers via the Amazon Mechanical Turk service. Mechanical Turk users were asked first to find some specific restaurants, and after each dialogue was finished, they had to fill in a feedback form to indicate if the dialogue had been successful or not. Similarly, for the end-to-end dialogue system by [Wen *et al.*, 2017] (cf. Section 3.3.2), also for the restaurants domain, human evaluation was conducted by users recruited via Amazon Mechanical Turk. Each evaluator had to follow a given task and to rate the system's performance. More specifically, they had to grade the subjective success rate, the perceived comprehension ability and naturalness of the responses.

Most of the task-oriented datasets are designed to evaluate components of dialogue systems. For example, several datasets have been released through different editions of the Dialog State Tracking Challenge[9], focused on the development and evaluation of the dialogue state tracker component. However, even if these datasets were designed to test state tracking, [Bordes *et al.*, 2017] used them to build and evaluate a whole dialogue system, readjusting the dataset by ignoring the state annotation and reusing only the transcripts of dialogues.

PyDial[10] partially addresses these shortage of evaluation datasets for task-oriented systems, as it offers the opportunity to develop a Reinforcement Learning based Dialogue Management benchmarking environment [Ultes *et al.*, 2017]. Thus, it makes possible the evaluation and comparison of different task-oriented dialogue systems in the same conditions. This toolkit provides not only domain independent implementations of different modules of a dialogue system, but also simulated users (cf. Section 3.4.2). It uses two

---

[9]https://www.microsoft.com/en-us/research/event/dialog-state-tracking-challenge/
[10]http://www.camdial.org/pydial/

metrics for the evaluation: (1) the average success rate and (2) the average reward for each evaluated policy model of reinforcement learning algorithms. Success rate is defined as the percentage of dialogues which are completed successfully, and thus, it is closely related to the task completion metric used by PARADISE framework (see Section 3.4.1).

MultiWOZ (Multi-Domain Wizard-of-Oz) [Budzianowski *et al.*, 2018] dataset represents a significant breakthrough in the scarcity of dialogues as it contains around 10K dialogues, which is at least one order of magnitude larger than any structured corpus available up to date. It is annotated with dialogue belief states and dialogue actions, so it can be used for the development of the individual components of a dialogue system. But its considerable size makes it very appropriate for the training of the end-to-end based dialogue systems. The main topic of the dialogues is tourism, and it contains seven domains, such as attraction, hospital, police, hotel, restaurant, taxi and train. Each dialogue can contain more than one of these domains.

## 6.2   Data for Question Answering Dialogue Systems

With respect to **QA dialogue systems**, two datasets have been created based on the human interactions from technical chats or forums. The first one is the Ubuntu Dialogue Corpus, containing almost one million multi-turn dialogues extracted from the Ubuntu chat logs, which was used to receive technical support for various Ubuntu-related problems [Lowe *et al.*, 2015]. Similarly, MSDialog contains dialogues from a forum dedicated to Microsoft products. MSDialog also contains the user intent of each interaction [Qu *et al.*, 2018].

ibAbI represents another approach for creating multi-turn QA datasets [Li *et al.*, 2017a]. ibAbI adds interactivity to the bAbI dataset that are going to be presented later on, adding sentences and ambiguous questions with the corresponding disambiguation question that should be asked by an automatic system. The authors evaluate their system regarding the successful tasks. However, it is unclear how to evaluate a system if it produces a modified version of the disambiguation question.

Recently, several datasets which are very appropriate for the context of QA dialogue systems have been released. QuAC (Question Answering in Context) consists of 14K information-seeking QA dialogues (100K total QA pairs) over sections from Wikipedia articles (all selected articles are about people) [Choi *et al.*, 2018]. What makes different from other datasets so far is that some of the questions are unanswerable and that context is needed in order to answer some of the questions. CoQA (Conversational Question Answering) dataset contains 8K dialogues and 127K conversation turns [Reddy *et al.*, 2018]. As opposed to QuAC, the answers from CoQA are free-form text with their corresponding evidence highlighted in the passage. It is multi domain, as the passages are selected from several sources, covering seven different domains: children's stories, literature, middle and high school English exams, news, articles from Wikipedia, science and Reddit. Amazon

Mechanical Turk was used to collect the dialogues for both datasets.

## 6.3  Data for Conversational Dialogue Systems

Regarding the evaluation of **Conversational dialogue systems** presented in Section 4, datasets derived from conversations on micro-blogging or social media websites. Twitter or Reddit are good candidates, as they contain general-purpose or non-task-oriented conversations that are orders of magnitude larger than other dialogue datasets used before. For instance, Switchboard [Godfrey *et al.*, 1992] (telephone conversations on pre-specified topics), British National Corpus [Leech, 1992] (British dialogues many contexts, from formal business or government meetings to radio shows and phone-ins) and SubTle Corpus [Ameixa and Coheur, 2013] (aligned interaction-response pairs from movie subtitles) are three datasets released earlier which have 2,400, 854 and 3.35M dialogues and 3M, 10M and 20M words, respectively. These sizes are relatively small if we compared to the huge Reddit Corpus[11] which contains over 1.7 billion of comments[12], or the Twitter Corpus described below.

Because of the limit on the number of characters permitted in each message on Twitter, the utterances are quite short, very colloquial and chat-like. Moreover, as the conversations happen almost in real-time, the conversations of this micro-blogging website are very similar to spoken dialogues between humans. There are two publicly available large corpus extracted from Twitter. The former one is the Twitter Corpus presented in [Ritter *et al.*, 2010], which contains roughly 1.3 million conversations and 125M words drawn from Twitter. The latter one is a collection of 4,232 three-step (context-message-response) conversational snippets extracted from Twitter logs[13]. It is labeled by crowdsourced annotators measuring quality of the response in the context [Sordoni *et al.*, 2015].

Alternatively, [Lowe *et al.*, 2015] hypothesized that chat-room style messaging is more closely correlated to human-to-human dialogues than micro-blogging websites like Twitter, or forum-based sites such as Reddit. Thus, they presented the Ubuntu Dialogue Corpus, just mentioned above. This large-scale corpus targets a specific domain, so it could be be used as a task-oriented dataset for research and evaluate the dialogue state trackers accordingly. But it also has the property of the unstructured nature of interactions from microblog services, which makes it appropriate for the evaluation of non-task-oriented dialogue systems, for example, dialogue managers based on neural language models that make use of large amounts of unlabeled data.

---

[11]https://www.reddit.com/r/datasets/comments/3bxlg7/i_have_every_publicly_available_reddit_comment/

[12]As far as we know, this dataset has not been used in any research work. Researchers have used smaller and more curated versions of the Reddit dataset like Reddit Domestic Abuse Corpus [Schrading, 2015], which contains 21,133 dialogues.

[13]https://www.microsoft.com/en-us/download/details.aspx?id=52375

These two large datasets are adequate for the three subtypes of non-task-oriented dialogue systems, unsupervised, trained and utterance selection metrics. Notice that, additionally, some human judgments could be needed in some cases, like in [Lowe *et al.*, 2017a] for the system ADEM, where they use some human judgments collected via Amazon Mechanical Turk in addition to the evaluation using the Twitter dataset.

Apart from the afore-mentioned two datasets, the five datasets generated recently for bAbI tasks [Bordes *et al.*, 2017] are appropriate for evaluation using the next utterance classification method (see Section 4.3.2). These tasks were designed for testing end-to-end dialogue systems in the restaurant domain, but they check whether the systems can predict the appropriate utterances among a fixed set of candidates, and are not useful for systems that generate the utterance directly. The ibAbI dataset mentioned before it has been created based on bAbI to cover several representative multi-turn QA tasks.

Another interesting resource is the ParlAI framework[14] for dialogue research, as it contains many popular datasets available all in one place with the goal of sharing, training and evaluating dialogue models across many tasks [Miller *et al.*, 2017]. Some of the dialogue datasets that are included have been already mentioned, bAbI Dialog tasks and Ubuntu Dialog Corpus, but it also contains conversation mined from OpenSubtitles[15] and Cornell Movie[16].

## 6.4 Evaluation Challenges

We complete this section by summarizing some of the recent **evaluation challenges** that are popular for benchmarking state-of-the-art dialogue system. They have an important role in the evaluation of dialogue systems, but not only because they offer a good benchmark scenario to test and compare the systems on a common platform, but also because they often release the dialogue datasets for later evaluation.

Perhaps one of the most popular challenges is the Dialog State Tracking Challenge (DSTC)[17] mentioned before in this section. It started to provide a common testbed for the task of dialogue state tracking in 2013, and continued yearly with a remarkable success. For its sixth edition it renamed itself as Dialog System Technology Challenges due to the interest of the research community in a wider variety of dialogue related problems. Different well-known datasets have been produced and released for every edition: DSTC1 has human-computer dialogues in the bus timetable domain; DSTC2 and DSTC3 used human-computer dialogues in the restaurant information domain; DSTC4 dialogues were human-human and in the tourist information domain; DSTC5 also is from the tourist information domain, but training dialogues are provided in one language and test dialogs are in a different

---

[14]http://parl.ai/

[15]http://opus.lingfil.uu.se/OpenSubtitles.php

[16]https://www.cs.cornell.edu/~cristian/Cornell_Movie-Dialogs_Corpus.html

[17]https://www.microsoft.com/en-us/research/event/dialog-state-tracking-challenge/

language; and finally, as DSTC6 edition consisted of 3 parallel tracks, different datasets were released for each track, such as, a transaction dialogue dataset for the restaurant domain, two datasets that are part of OpenSubtitles and Twitter datasets, and different chat-oriented dialogue datasets with dialogue breakdown annotations in Japanese and English.

A more recent challenge that started last year and continued this year with its second edition is the Conversational Intelligence Challenge (ConvAI)[18]. This challenge conducted under the scope of NIPS has the aim to unify the community around the task of building systems capable of intelligent conversations. In its first edition teams were expected to submit dialogue systems able to carry out intelligent and natural conversations about specific news articles with humans. The aim of the task of the second edition has been to model normal conversation when two interlocutors meet for the first time, and get to know each other. The dataset of this task consists of 10,981 dialogues with 164,356 utterances, and it is available in the ParlAI framework mentioned above.

Finally, the Alexa Prize[19] has attracted mass media and researcher attention alike. This annual competition for university teams is dedicated at accelerating the field of conversational AI in the framework of the Alexa technology. The participants have to create socialbots that can converse coherently and engagingly with humans on news events and popular topics such as entertainment, sports, politics, technology and fashion. Unfortunately no datasets have been released.

# 7   Challenges and Future Trends

We stated in the introduction that the goal of dialogue evaluation is to find methods that are automated, repeatable, correlate to human judgements, are capable of differentiating among various dialogue strategies and explain which features of the dialogue system contribute to the quality. The main motivation behind this is the need to reduce the human effort as much as possible, since human involvement produces high costs and takes a long time. In this survey, we presented the main concepts regarding evaluation of dialogue systems and showcased the most important methods. However, evaluation of dialogue systems is still an open area of research. In this section we summarize the current challenges and future trends that we deem most important.

**Automation.**   The evaluation methods covered in this survey all achieve a certain degree of automation. The automation is however achieved with a significant engineering effort or by a loss of correlation to human judgements. Word-overlap metrics (see Section 4.3.1), which are borrowed from the machine translation and summarization community, are fully

---

[18]http://convai.io/
[19]https://developer.amazon.com/alexaprize

automated. However, they do not correlate with human judgements on the turn level. On the other hand, BLEU becomes more competitive when applied on the corpus-level or system-level [Galley *et al.*, 2015; Lowe *et al.*, 2017a]. More recent metrics such as ΔBLEU and ADEM (see Section 4.3.1) have significantly higher correlations to human judgements while requiring a significant amount of human annotated data as well as thorough engineering.

Task-oriented dialogue systems can be evaluated semi-automatically or even fully automatic. These systems benefit from having a well defined task, whose success can be measured. Thus, user satisfaction modelling (see Section 3.4.1) as well as user simulations (see Section 3.4.2) exploit this to automate their evaluation. However, both approaches need a significant amount of engineering and human annotation: user satisfaction modelling usually requires prior annotation effort, which is followed by fitting a model which predicts the judgements. In addition to this effort, the process potentially has to be repeated for each new domain or new functionality the dialogue system incorporates. Although in some cases the model fitted on the data for one dialogue system can be reused to predict another dialogue system,this not always possible.

On the other hand, user simulations require two steps: gather data to develop a first version of the simulation, and then build the actual user simulation. The first step is only required for user simulations which are based on corpora to train (e.g. the neural user simulation). A significant drawback is that the user simulation is only capable of simulating the behaviour which is represented in the corpus or the rules. This means, that it cannot cover unseen behaviour well. Furthermore, the user simulation can hardly be used to train or evaluate dialogue systems for other tasks or domains.

Automation is, thus, achieved to a certain degree but with significant drawbacks. Hence, finding ways to facilitate the automation of evaluation methods is clearly an open challenge.

**High Quality Dialogues.**   One major objective for a dialogue system is to deliver high quality interactions with its users. However, it is often not clear how "high quality" is defined in this context or how to measure it. For task oriented dialogue systems, the mostly used definition of quality is often measured by means of task success and number of dialogue turns (e.g. a reward of 20 for task-success minus the number of turns needed to achieve the goal). But, this definition is not applicable to conversational dialogue systems and it might ignore other aspects of the interaction (e.g. frustration of the user). Thus, the current trend is to let humans judge the *appropriateness* of the system utterances. However, the notion of appropriateness is highly subjective and entails several finer-grained concepts (e.g. ability to maintain the topic, the coherence of the utterance, the grammatical correctness of the utterance itself, etc.). Currently, appropriateness is modelled by means of latent representations (e.g. ADEM), which are derived again from annotated data.

Other aspects of quality concern the purpose of the dialogue system in conjunction with

the functionality of the system. For instance, [Zhou *et al.*, 2018] define the purpose of their conversational dialogue system to build an emotional bond between the dialogue system and the user. This goal differs significantly from the task of training a medical student in the interaction with patients. Both systems need to be evaluated with respect to their particular goal. The ability to build an emotional bond can be evaluated by means of the interaction length (longer interactions are an indicator of a higher user engagement), whereas training (or e-learning) systems are usually evaluated regarding their ability of selecting an appropriate utterance for the given context.

The target audience plays an important role as well. Since quality is mainly a subjective measure, different user groups prefer different types of interactions. For instance, depending on the level of domain knowledge, novice users prefer instructions which use less specialized wordings, whereas domain experts might prefer a more specialized vocabulary.

The notion of quality is, thus, dependent on a large amount of factors. Depending on the purpose, the target audience, and the dialogue system implementation itself, evaluation needs to be adapted to take all these factors into account.

**Lifelong Learning.** The notion of lifelong learning for machine learning systems has gained traction recently. The main concept of lifelong learning is that a deployed machine learning system continues to improve by interaction with its environment [Chen and Liu, 2016]. Lifelong learning for dialogue systems is motivated by the fact that it is not possible to encounter all possible situations during training, thus, a component which allows the dialogue system to retrain itself and adapt its strategy during deployment seems the most logical solution.

To achieve lifelong learning, the evaluation step is critical. Since the dialogue system relies on the ability to automatically find critical dialogue states where it needs assistance a module is needed which is able to evaluate the ongoing dialogue. One step in this direction is done by [Hancock *et al.*, 2019], who present a solution that relies on a satisfaction module that is able of to classify the current dialogue state as either satisfactory or not. If this module finds an unsatisfactory dialogue state, a feedback module asks the user for feedback. The feedback data is then used to improve the dialogue system.

The aspect of lifelong learning brings a large variety of novel challenges. First, the lifelong learning system requires a module which self-monitors its behaviour and notices when a dialogue is going wrong. For this, the module needs to rely on evaluation methods which work automatically or at least-semi automated. The second challenge lies in the evaluation of the lifelong learning system itself. The self-monitoring module as well as the adaptive behaviour need to be evaluated. This brings a new dimension of complexity into the evaluation procedure.

**Conclusion**   Evaluation is a critical task when developing and researching dialogue systems. Over the past decades lots of methods and concepts have been proposed. These methods and concepts are related to the different requirements and functionalities of the dialogue systems. These are in turn dependent on the current development stage of the dialogue system technology. Currently, the trend is going towards building end-to-end trainable dialogue system based on large amounts of data. These systems have different different requirements for evaluation than a finite state machine based system. Thus, the problem of evaluation is evolving in parallel to the progress in the dialogue system technology itself. This survey presents the current state-of-the-art research in evaluation.

# Acknowledgements

# References

[Ameixa and Coheur, 2013] David Ameixa and Luisa Coheur. From subtitles to human interactions: introducing the SubTle Corpus. In *Technical report*, 2013.

[Austin, 1962] John Langshaw Austin. *How to do things with words*. William James Lectures. Oxford University Press, 1962.

[Banchs and Li, 2012] Rafael E. Banchs and Haizhou Li. IRIS: a chat-oriented dialogue system based on the vector space model. In *Proceedings of the ACL 2012 Demonstrations*, pages 37–42, 2012.

[Banchs, 2012] Rafael E. Banchs. Movie-DiC: a Movie Dialogue Corpus for Research and Development. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 203–207. Association for Computational Linguistics, 2012.

[Bangalore *et al.*, 2001] S Bangalore, O Rambow, and M Walker. Natural language generation in dialog systems. In *Proceedings of Human Language Technology Conference*, pages 67–73, 2001.

[Bernardi and Kirschner, 2010] Raffaella Bernardi and Manuel Kirschner. From artificial questions to real user interaction logs: Real challenges for interactive question answering systems. In *Proc. of Workshop on Web Logs and Question Answering (WLQA'10)*, Malta, 2010.

[Black *et al.*, 2011] Alan W Black, Susanne Burger, Alistair Conkie, Helen Hastie, Simon Keizer, Oliver Lemon, Nicolas Merigaud, Gabriel Parent, Gabriel Schubiner, Blaise Thomson, Jason D. Williams, Kai Yu, Steve Young, and Maxine Eskenazi. Spoken dialog challenge 2010: Comparison of live and control test results. In *Proceedings of the SIGDIAL 2011 Conference*, pages 2–7, Portland, Oregon, June 2011. Association for Computational Linguistics.

[Bordes *et al.*, 2017] Antoine Bordes, Y-Lan Boureau, and Jason Weston. Learning End-to-End Goal-Oriented Dialog. In *ICLR 2017*, 2017. arXiv:1605.07683.

[Bowman *et al.*, 2015] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349.*, 2015.

[Bruni and Fernandez, 2017] Elia Bruni and Raquel Fernandez. Adversarial evaluation for open-domain dialogue generation. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 284–288. Association for Computational Linguistics, 2017.

[Budzianowski *et al.*, 2018] Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Ultes Stefan, Ramadan Osman, and Milica Gašić. MultiWOZ - A

Large-Scale Multi-Domain Wizard-of-Oz Dataset for Task-Oriented Dialogue Modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.

[Carletta, 1996] Jean Carletta. Squibs and discussions: Assessing agreement on classification tasks: The kappa statistic. *Computational Linguistics*, 22(2):249–254, June 1996.

[Charras *et al.*, 2016] F. Charras, G. Dubuisson Duplessis, V. Letard, A.-L. Ligozat, and S. Rosset. Comparing system-response retrieval models for open-domain and casual conversational agent. In *Workshop on Chatbots and Conversational Agent Technologies*, 2016.

[Chen and Liu, 2016] Zhiyuan Chen and Bing Liu. Lifelong machine learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 10(3):1–145, 2016.

[Chen *et al.*, 2017] Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. A Survey on Dialogue Systems: Recent Advances and New Frontiers. *SIGKDD Explor. Newsl.*, 19(2):25–35, November 2017.

[Choi *et al.*, 2018] E. Choi, H. He, M. Iyyer, M. Yatskar, W.-t. Yih, Y. Choi, P. Liang, and L. Zettlemoyer. QuAC: Question Answering in Context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018. arXiv:1808.07036.

[Colby, 1981] Kenneth Mark Colby. Modeling a paranoid mind. *Behavioral and Brain Sciences*, 4(4):515–534, 1981.

[Cole, 1999] Ron Cole. Tools for research and education in speech science. In *Proceedings of the International Conference of Phonetic Sciences*, volume 1, pages 277–1. Citeseer, 1999.

[Dethlefs *et al.*, 2013] Nina Dethlefs, Helen Hastie, Heriberto Cuayáhuitl, and Oliver Lemon. Conditional random fields for responsive surface realisation using global features. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 1254–1263, 2013.

[DeVault *et al.*, 2011] David DeVault, Anton Leuski, and Kenji Sagae. Toward Learning and Evaluation of Dialogue Policies with Text Examples. In *Proceedings of the SIGDIAL 2011 Conference*, SIGDIAL '11, pages 39–48, Stroudsburg, PA, USA, 2011. Association for Computational Linguistics.

[Diefenbach *et al.*, 2018] Dennis Diefenbach, Vanessa Lopez, Kamal Singh, and Pierre Maret. Core techniques of question answering systems over knowledge bases: A survey. *Knowl. Inf. Syst.*, 55(3):529–569, June 2018.

[Dinarelli *et al.*, 2017] Marco Dinarelli, Vedran Vukotic, and Christian Raymond. Label-dependency coding in Simple Recurrent Networks for Spoken Language Understanding. In *Interspeech*, Stockholm, Sweden, August 2017.

[Do *et al.*, 2017] Phong-Khac Do, Huy-Tien Nguyen, Chien-Xuan Tran, Minh-Tien Nguyen, and Minh-Le Nguyen. Legal question answering using ranking SVM and deep convolutional neural network. *CoRR*, abs/1703.05320, 2017.

[Dubuisson Duplessis *et al.*, 2016] Guillaume Dubuisson Duplessis, Vincent Letard, Anne-Laure Ligozat, and Sophie Rosset. Purely corpus-based automatic conversation authoring. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, Paris, France, may 2016. European Language Resources Association (ELRA).

[Dubuisson Duplessis *et al.*, 2017] Guillaume Dubuisson Duplessis, Franck Charras, Vincent Letard, Anne-Laure Ligozat, and Sophie Rosset. Utterance retrieval based on recurrent surface text patterns. In *European Conference on Information Retrieval (ECIR 2017)*, page 12p, Aberdeen, Scotland UK, 09/04 au 13/04 2017.

[Dušek and Jurcicek, 2016] Ondřej Dušek and Filip Jurcicek. 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. Association for Computational Linguistics, 2016.

[Engel *et al.*, 2005] Yaakov Engel, Shie Mannor, and Ron Meir. Reinforcement Learning with Gaussian Processes. In *Proceedings of the 22Nd International Conference on Machine Learning*, ICML '05, pages 201–208, New York, NY, USA, 2005. ACM.

[Engelbrecht *et al.*, 2008] Klaus-Peter Engelbrecht, Sebastian Möller, Robert Schleicher, and Ina Wechsung. Analysis of paradise models for individual users of a spoken dialog system. *Proc. of ESSV*, pages 86–93, 2008.

[Engelbrecht *et al.*, 2009a] Klaus-Peter Engelbrecht, Florian Gödde, Felix Hartard, Hamed Ketabdar, and Sebastian Möller. Modeling User Satisfaction with Hidden Markov Model. In *Proceedings of the SIGDIAL 2009 Conference: The 10th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, SIGDIAL '09, pages 170–177, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics.

[Engelbrecht *et al.*, 2009b] Klaus-Peter Engelbrecht, Michael Quade, and Sebastian Möller. Analysis of a New Simulation Approach to Dialog System Evaluation. *Speech Commun.*, 51(12):1234–1252, December 2009.

[Evanini *et al.*, 2008] K. Evanini, P. Hunter, J. Liscombe, D. Suendermann, K. Dayanidhi, and R. Pieraccini. Caller Experience: A method for evaluating dialog systems and its automatic prediction. In *2008 IEEE Spoken Language Technology Workshop*, pages 129–132, Dec 2008.

[Fleiss, 1971] Joseph L Fleiss. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378, 1971.

[Gašić *et al.*, 2011] M. Gašić, F. Jurčíček, B. Thomson, K. Yu, and S. Young. On-line policy optimisation of spoken dialogue systems via live interaction with human subjects.

In *2011 IEEE Workshop on Automatic Speech Recognition Understanding*, pages 312–317, Dec 2011.

[Galley *et al.*, 2015] Michel Galley, Chris Brockett, Alessandro Sordoni, Yangfeng Ji, Michael Auli, Chris Quirk, Margaret Mitchell, Jianfeng Gao, and Bill Dolan. deltaBLEU: A Discriminative Metric for Generation Tasks with Intrinsically Diverse Targets. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 445–450. Association for Computational Linguistics, 2015.

[Gandhe and Traum, 2013] Sudeep Gandhe and David R. Traum. Surface text based dialogue models for virtual humans. In *Proceedings of the SIGDIAL*, 2013.

[Gandhe and Traum, 2016] Sudeep Gandhe and David Traum. *A Semi-automated Evaluation Metric for Dialogue Model Coherence*, pages 217–225. Springer International Publishing, Cham, 2016.

[Gandhe *et al.*, 2009] Sudeep Gandhe, Nicolle Whitman, David Traum, and Ron Artstein. An integrated authoring tool for tactical questioning dialogue systems. In *6th IJCAI Workshop on Knowledge and Reasoning in Practical Dialogue Systems*, page 10, 2009.

[Gasic *et al.*, 2013] Milica Gasic, Catherine Breslin, Matthew Henderson, Dongho Kim, Martin Szummer, Blaise Thomson, Pirros Tsiakoulis, and Steve Young. POMDP-based dialogue manager adaptation to extended domains. In *Proceedings of the SIGDIAL 2013 Conference*, pages 214–222. Association for Computational Linguistics, 2013.

[Gasic *et al.*, 2014] Milica Gasic, Dongho Kim, Pirros Tsiakoulis, Catherine Breslin, Matthew Henderson, Martin Szummer, Blaise Thomson, and Steve J. Young. Incremental on-line adaptation of POMDP-based dialogue managers to extended domains. In *INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication Association, Singapore, September 14-18, 2014*, pages 140–144, 2014.

[Ghazvininejad *et al.*, 2017] Marjan Ghazvininejad, Chris Brockett, Ming-Wei Chang, Bill Dolan, Jianfeng Gao, Wen-tau Yih, and Michel Galley. A knowledge-grounded neural conversation model. *arXiv preprint arXiv:1702.01932*, 2017.

[Godfrey *et al.*, 1992] J. J. Godfrey, E. C. Holliman, and J. McDaniel. SWITCHBOARD: telephone speech corpus for research and development. In *[Proceedings] ICASSP-92: 1992 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 1, pages 517–520 vol.1, Mar 1992.

[Goodfellow *et al.*, 2014] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014.

[Guo *et al.*, 2014] Daniel Guo, Gokhan Tur, Wen-tau Yih, and Geoffrey Zweig. Joint semantic utterance classification and slot filling with recursive neural networks. In *Spoken Language Technology Workshop (SLT), 2014 IEEE*, pages 554–559. IEEE, 2014.

[Guo *et al.*, 2018] Fenfei Guo, Angeliki Metallinou, Chandra Khatri, Anirudh Raju, Anu Venkatesh, and Ashwin Ram. Topic-based evaluation for conversational bots. *arXiv preprint arXiv:1801.03622*, 2018.

[Hahn *et al.*, 2010] Stefan Hahn, Marco Dinarelli, Christian Raymond, Fabrice Lefèvre, Patrick Lehen, Renato De Mori, Alessandro Moschitti, Hermann Ney, and Giuseppe Riccardi. Comparing stochastic approaches to spoken language understanding in multiple languages. *IEEE Transactions on Audio, Speech and Language Processing (TASLP)*, 16:1569–1583, 2010.

[Hancock *et al.*, 2019] Braden Hancock, Antoine Bordes, Pierre-Emmanuel Mazare, and Jason Weston. Learning from dialogue after deployment: Feed yourself, chatbot! *arXiv preprint arXiv:1901.05415*, 2019.

[Hara, 2010] Sunao Hara. Estimation method of user satisfaction using N-gram-based dialog history model for spoken dialog system. *Proceedings of LREC2010*, pages 78–83, 2010.

[Henderson *et al.*, 2013] Matthew Henderson, Blaise Thomson, and Steve Young. Deep neural network approach for the dialog state tracking challenge. In *Proceedings of the SIGDIAL 2013 Conference*, pages 467–471, 2013.

[Henderson *et al.*, 2014] Matthew Henderson, Blaise Thomson, and Jason Williams. The second dialog state tracking challenge. In *Proceedings of SIGDIAL*. ACL – Association for Computational Linguistics, June 2014.

[Higashinaka *et al.*, 2010] Ryuichiro Higashinaka, Yasuhiro Minami, Kohji Dohsaka, and Toyomi Meguro. Issues in Predicting User Satisfaction Transitions in Dialogues: Individual Differences, Evaluation Criteria, and Prediction Models. In Gary Geunbae Lee, Joseph Mariani, Wolfgang Minker, and Satoshi Nakamura, editors, *Spoken Dialogue Systems for Ambient Environments*, pages 48–60, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg.

[Hirschman *et al.*, 1990] Lynette Hirschman, Deborah A Dahl, Donald P McKay, Lewis M Norton, and Marcia C Linebarger. Beyond class A: A proposal for automatic evaluation of discourse. In *In Proceedings of the Speech and Natural Language Workshop*, pages 109–113, 1990.

[Hochreiter and Schmidhuber, 1997] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, pages 1735–1780, 1997.

[Hu *et al.*, 2017] Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing. Toward Controlled Generation of Tex. *International Conference on Machine Learning*, pages 1587–1596, 2017.

[Iyyer *et al.*, 2017] Mohit Iyyer, Wen-tau Yih, and Ming-Wei Chang. Search-based neural structured learning for sequential question answering. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1821–1831. Association for Computational Linguistics, 2017.

[Jurafsky and Martin, 2017] Daniel Jurafsky and James H. Martin. *Speech and Language Processing*, chapter Dialog Systems and Chatbots. Draft of 3rd edition edition, 2017.

[Jurcícek *et al.*, 2011] Filip Jurcícek, Simon Keizer, Milica Gasic, François Mairesse, Blaise Thomson, Kai Yu, and Steve J. Young. Real user evaluation of spoken dialogue systems using amazon mechanical turk. In *INTERSPEECH 2011, 12th Annual Conference of the International Speech Communication Association, Florence, Italy, August 27-31, 2011*, pages 3061–3064, 2011.

[Kannan and Vinyals, 2016] Anjuli Kannan and Oriol Vinyals. Adversarial evaluation of dialogue models. *Workshop on Adversarial Training at Neural Information Processing Systems 2016*, 2016.

[Kelly and Lin, 2007] Diane Kelly and Jimmy Lin. Overview of the trec 2006 ciqa task. *SIGIR Forum*, 41(1):107–116, June 2007.

[Kelly *et al.*, 2009] Diane Kelly, Paul B Kantor, Emile L Morse, Jean Scholtz, and Ying Sun. Questionnaires for eliciting evaluation data from users of interactive question answering systems. *Natural Language Engineering*, 15(1):119–141, 2009.

[Kenny *et al.*, 2009] Patrick G. Kenny, Thomas D. Parsons, and Albert A. Rizzo. Human computer interaction in virtual standardized patient systems. In *Proceedings of the 13th International Conference on Human-Computer Interaction. Part IV: Interacting in Various Application Domains*, pages 514–523, Berlin, Heidelberg, 2009. Springer-Verlag.

[Kolomiyets and Moens, 2011] Oleksandr Kolomiyets and Marie-Francine Moens. A survey on question answering technology from an information retrieval perspective. *Inf. Sci.*, 181(24):5412–5434, December 2011.

[Konstantinova and Orasan, 2013] Natalia Konstantinova and Constantin Orasan. Interactive question answering. In *Emerging Applications of Natural Language Processing: Concepts and New Research*, pages 149 –. 10 2013.

[Kreyssig *et al.*, 2018] Florian Kreyssig, Inigo Casanueva, Pawel Budzianowski, and Milica Gasic. Neural user simulation for corpus-based policy optimisation for spoken dialogue systems. *arXiv preprint arXiv:1805.06966*, 2018.

[Lamel *et al.*, 2000] Lori Lamel, Sophie Rosset, Jean-Luc Gauvain, Samir Bennacef, Martine Garnier-Rizet, and Bernard Prouts. The limsi arise system. *Speech Communication*, 31(4):339–353, 2000.

[Lavie and Denkowski, 2009] Alon Lavie and Michael J. Denkowski. The meteor metric for automatic evaluation of machine translation. *Machine Translation*, 23(2-3):105–115, September 2009.

[Lee *et al.*, 2009] Cheongjae Lee, Sangkeun Jung, Seokhwan Kim, and Gary Geunbae Lee. Example-based dialog modeling for practical multi-domain dialog system. *Speech Communication*, 51(5):466–484, 2009.

[Leech, 1992] G.N. Leech. 100 million words of english: the british national corpus (BNC). *Language Research*, 28:1–13, 01 1992.

[Lemon and Pietquin, 2012] Oliver Lemon and Olivier Pietquin. *Data-Driven Methods for Adaptive Spoken Dialogue Systems: Computational Learning for Conversational Interfaces*. Springer Publishing Company, Incorporated, 2012.

[Levin *et al.*, 1998] E. Levin, R. Pieraccini, and W. Eckert. Using Markov decision process for learning dialogue strategies. In *Acoustics, Speech and Signal Processing, 1998. Proceedings of the 1998 IEEE International Conference on*, volume 1, pages 201–204 vol.1, May 1998.

[Li *et al.*, 2016a] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A Diversity-Promoting Objective Function for Neural Conversation Models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119. Association for Computational Linguistics, 2016.

[Li *et al.*, 2016b] Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. Deep reinforcement learning for dialogue generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1192–1202. Association for Computational Linguistics, 2016.

[Li *et al.*, 2017a] Huayu Li, Martin Renqiang Min, Yong Ge, and Asim Kadav. A context-aware attention network for interactive question answering. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, pages 927–935, New York, NY, USA, 2017. ACM.

[Li *et al.*, 2017b] Xiujun Li, Yun-Nung Chen, Lihong Li, Jianfeng Gao, and Asli Celikyilmaz. End-to-End Task-Completion Neural Dialogue Systems. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 733–743. Asian Federation of Natural Language Processing, 2017.

[Lin, 2004] Chin-Yew Lin. ROUGE: A Package for Automatic Evaluation of Summaries. In Stan Szpakowicz Marie-Francine Moens, editor, *Text Summarization Branches Out: Proceedings of the ACL-04 Workshop*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.

[Liu *et al.*, 2016] Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. How NOT To Evaluate Your Dialogue System: An Em-

pirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2122–2132. Association for Computational Linguistics, 2016.

[Liu *et al.*, 2018] Bing Liu, Gokhan Tür, Dilek Hakkani-Tür, Pararth Shah, and Larry Heck. Dialogue Learning with Human Teaching and Feedback in End-to-End Trainable Task-Oriented Dialogue Systems. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2060–2069. Association for Computational Linguistics, 2018.

[Lowe *et al.*, 2015] Ryan Joseph Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. The Ubuntu Dialogue Corpus: A Large Dataset for Research in Unstructured Multi-Turn Dialogue Systems. In *Proceedings of the SIGDIAL 2015 Conference*, pages 285–294. Association for Computational Linguistics, 2015.

[Lowe *et al.*, 2016] Ryan Lowe, Iulian Vlad Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. On the Evaluation of Dialogue Systems with Next Utterance Classification. In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 264–269. Association for Computational Linguistics, 2016.

[Lowe *et al.*, 2017a] Ryan Lowe, Michael Noseworthy, Iulian Vlad Serban, Nicolas Angelard-Gontier, Yoshua Bengio, and Joelle Pineau. Towards an Automatic Turing Test: Learning to Evaluate Dialogue Responses. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1116–1126. Association for Computational Linguistics, 2017.

[Lowe *et al.*, 2017b] Ryan Lowe, Nissan Pow, Iulian V. Serban, Laurent Charlin, Chia-Wei Liu, and Joelle Pineau. Training End-to-End Dialogue Systems with the Ubuntu Dialogue Corpus. *Dialogue & Discourse*, 8(1):31–65, 2017.

[Mairesse *et al.*, 2010] François Mairesse, Milica Gašić, Filip Jurčíček, Simon Keizer, Blaise Thomson, Kai Yu, and Steve Young. Phrase-based statistical language generation using graphical models and active learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1552–1561. Association for Computational Linguistics, 2010.

[Mazza *et al.*, 2018] Riccardo Mazza, L. Ambrosini, N. Catenazzi, S. Vanini, Don Tuggener, and G. Tavarnesi. Behavioural simulator for professional training based on natural language interaction. In *10th International Conference on Education and New Learning Technologies*, pages 3204–3214, 2018.

[McTear *et al.*, 2005] Michael McTear, Ian O'Neill, Philip Hanna, and Xingkun Liu. Handling errors and determining confirmation strategies—an object-based approach. *Speech Communication*, 45(3):249–269, 2005.

[Mei *et al.*, 2016] Hongyuan Mei, TTI UChicago, Mohit Bansal, and Matthew R Walter. What to talk about and how? Selective Generation using LSTMs with Coarse-to-Fine Alignment. In *Proceedings of NAACL-HLT*, pages 720–730, 2016.

[Mesnil *et al.*, 2015] Grégoire Mesnil, Yann Dauphin, Kaisheng Yao, Yoshua Bengio, Li Deng, Dilek Hakkani-Tur, Xiaodong He, Larry Heck, Gokhan Tur, Dong Yu, and Geoffrey Zweig. Using recurrent neural networks for slot filling in spoken language understanding. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 23(3):530–539, March 2015.

[Metallinou *et al.*, 2013] Angeliki Metallinou, Dan Bohus, and Jason Williams. Discriminative state tracking for spoken dialog systems. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 466–475, Sofia, Bulgaria, August 2013. Association for Computational Linguistics.

[Miller *et al.*, 2017] A. H. Miller, W. Feng, A. Fisch, J. Lu, D. Batra, A. Bordes, D. Parikh, and J. Weston. ParlAI: A Dialog Research Software Platform. *arXiv preprint arXiv:1705.06476*, 2017.

[Mishra and Jain, 2016] Amit Mishra and Sanjay Kumar Jain. A survey on question answering systems with classification. *J. King Saud Univ. Comput. Inf. Sci.*, 28(3):345–361, July 2016.

[Möller *et al.*, 2006] Sebastian Möller, Roman Englert, Klaus Engelbrecht, Verena Hafner, Anthony Jameson, Antti Oulasvirta, Alexander Raake, and Norbert Reithinger. MeMo: towards automatic usability evaluation of spoken dialogue services by user error simulations. In *Ninth International Conference on Spoken Language Processing*, 2006.

[Mrkšić *et al.*, 2017] Nikola Mrkšić, Diarmuid Ó Séaghdha, Tsung-Hsien Wen, Blaise Thomson, and Steve Young. Neural belief tracker: Data-driven dialogue state tracking. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1777–1788. Association for Computational Linguistics, 2017.

[Novikova *et al.*, 2017] Jekaterina Novikova, Ondrej Dušek, and Verena Rieser. The E2E Dataset: New Challenges for End-to-End Generation. In *Proceedings of the 18th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, Saarbrücken, Germany, 2017. arXiv:1706.09254.

[Paek, 2006] Tim Paek. Reinforcement learning for spoken dialogue systems: Comparing strengths and weaknesses for practical deployment. Technical report, 2006.

[Papineni *et al.*, 2002] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, 2002.

[Peñas *et al.*, 2012] Anselmo Peñas, Bernardo Magnini, Pamela Forner, Richard F. E. Sutcliffe, Álvaro Rodrigo, and Danilo Giampiccolo. Question answering at the cross-

language evaluation forum 2003-2010. *Language Resources and Evaluation*, 46(2):177–217, 2012.

[Pietquin and Hastie, 2013] Olivier Pietquin and Helen Hastie. A survey on metrics for the evaluation of user simulations. *The Knowledge Engineering Review*, 28(1):59–73, 2013.

[Powers, 2012] David Martin Ward Powers. The problem with kappa. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 345–355, Avignon, France, April 2012. Association for Computational Linguistics.

[Qu and Green, 2002] Yan Qu and Nancy Green. A constraint-based approach for cooperative information-seeking dialogue. In *Proceedings of the International Natural Language Generation Conference*, pages 136–143, 2002.

[Qu *et al.*, 2018] Chen Qu, Liu Yang, W. Bruce Croft, Johanne R. Trippas, Yongfeng Zhang, and Minghui Qiu. Analyzing and characterizing user intent in information-seeking conversations. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018*, pages 989–992, 2018.

[Reddy *et al.*, 2018] S. Reddy, D. Chen, and C. D. Manning. CoQA: A Conversational Question Answering Challenge. 2018. arXiv:1808.07042.

[Rieser and Lemon, 2009] Verena Rieser and Oliver Lemon. Does this list contain what you were searching for? learning adaptive dialogue strategies for interactive question answering. *Natural Language Engineering*, 15(1):55–72, 2009.

[Ritter *et al.*, 2010] Alan Ritter, Colin Cherry, and Bill Dolan. Unsupervised Modeling of Twitter Conversations. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, HLT '10, pages 172–180, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.

[Ritter *et al.*, 2011] Alan Ritter, Colin Cherry, and William B. Dolan. Data-driven Response Generation in Social Media. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, EMNLP '11, pages 583–593, Stroudsburg, PA, USA, 2011. Association for Computational Linguistics.

[Sarrouti and Ouatik El Alaoui, 2017] Mourad Sarrouti and Said Ouatik El Alaoui. A passage retrieval method based on probabilistic information retrieval model and umls concepts in biomedical question answering. *J. of Biomedical Informatics*, 68(C):96–103, April 2017.

[Schatzmann *et al.*, 2006] Jost Schatzmann, Kark Weilhammer, Matt Stuttle, and Steve Young. A survey of statistical user simulation techniques for reinforcement-learning of

dialogue management strategies. *The Knowledge Engineering Review*, 21(2):97–126, 06 2006. Copyright - 2006 Cambridge University Press; Zuletzt aktualisiert - 2015-08-15.

[Schatzmann *et al.*, 2007] Jost Schatzmann, Blaise Thomson, Karl Weilhammer, Hui Ye, and Steve Young. Agenda-based user simulation for bootstrapping a pomdp dialogue system. In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Companion Volume, Short Papers*, NAACL-Short '07, pages 149–152, Stroudsburg, PA, USA, 2007. Association for Computational Linguistics.

[Schatztnann *et al.*, 2005] Jost Schatztnann, Matthew N Stuttle, Karl Weilhammer, and Steve Young. Effects of the user model on simulation-based learning of dialogue strategies. In *Automatic Speech Recognition and Understanding, 2005 IEEE Workshop on*, pages 220–225. IEEE, 2005.

[Schmitt and Ultes, 2015] Alexander Schmitt and Stefan Ultes. Interaction Quality: Assessing the quality of ongoing spoken dialog interaction by experts—And how it relates to user satisfaction. *Speech Communication*, 74:12 – 36, 2015.

[Schrading, 2015] J. N. Schrading. Analyzing domestic abuse using natural language processing on social media data. In *Master's thesis, Rochester Institute of Technology*, 2015.

[Searle, 1969] John R. Searle. *Speech Acts: An Essay in the Philosophy of Language*. Cambridge University Press, Cambridge, London, 1969.

[Searle, 1975] John R. Searle. Indirect speech acts. In P. Cole and J. Morgan, editors, *Syntax and Semantics 3: Speech Acts*, pages 59–82. Academic Press, New York, 1975.

[Semeniuta *et al.*, 2017] Stanislau Semeniuta, Aliaksei Severyn, and Erhardt Barth. A Hybrid Convolutional Variational Autoencoder for Text Generation. *arXiv preprint arXiv:1702.02390*, 2017.

[Serban *et al.*, 2016] Iulian V. Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. Building End-to-end Dialogue Systems Using Generative Hierarchical Neural Network Models. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI'16, pages 3776–3783. AAAI Press, 2016.

[Serban *et al.*, 2017a] Iulian V Serban, Chinnadhurai Sankar, Mathieu Germain, Saizheng Zhang, Zhouhan Lin, Sandeep Subramanian, Taesup Kim, Michael Pieper, Sarath Chandar, Nan Rosemary Ke, et al. A deep reinforcement learning chatbot. *arXiv preprint arXiv:1709.02349*, 2017.

[Serban *et al.*, 2017b] Iulian V. Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. A Hierarchical Latent Variable Encoder-Decoder Model for Generating Dialogues. In *Thirty-First AAAI Conference on Artificial Intelligence*, page 1583, 2017.

[Serban *et al.*, 2017c] Iulian Vlad Serban, Tim Klinger, Gerald Tesauro, Kartik Tala-madupula, Bowen Zhou, Yoshua Bengio, and Aaron C. Courville. Multiresolution Recurrent Neural Networks: An Application to Dialogue Response Generation. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pages 3288–3294, 2017.

[Serban *et al.*, 2017d] Iulian Vlad Serban, Ryan Lowe, Peter Henderson, Laurent Charlin, and Joelle Pineau. A survey of available corpora for building data-driven dialogue systems. *arXiv preprint arXiv:1512.05742*, 2017.

[Serban *et al.*, 2018] Iulian Vlad Serban, Ryan Lowe, Peter Henderson, Laurent Charlin, and Joelle Pineau. A Survey of Available Corpora for Building Data-Driven Dialogue Systems: The Journal Version. *Dialogue & Discourse*, 1(9), 2018.

[Shang *et al.*, 2015] Lifeng Shang, Zhengdong Lu, and Hang Li. Neural Responding Machine for Short-Text Conversation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1577–1586. Association for Computational Linguistics, 2015.

[Singh *et al.*, 2000] Satinder P. Singh, Michael J. Kearns, Diane J. Litman, and Marilyn A. Walker. Reinforcement Learning for Spoken Dialogue Systems. In S. A. Solla, T. K. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems 12*, pages 956–962. MIT Press, 2000.

[Sordoni *et al.*, 2015] Alessandro Sordoni, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. A Neural Network Approach to Context-Sensitive Generation of Conversational Responses. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 196–205. Association for Computational Linguistics, 2015.

[Stent *et al.*, 2004] Amanda Stent, Rashmi Prasad, and Marilyn Walker. Trainable Sentence Planning for Complex Information Presentation in Spoken Dialog Systems. In *Proceedings of the 42Nd Annual Meeting on Association for Computational Linguistics*, ACL '04, Stroudsburg, PA, USA, 2004. Association for Computational Linguistics.

[Sutskever *et al.*, 2014] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'14, pages 3104–3112, Cambridge, MA, USA, 2014. MIT Press.

[Tiedemann, 2009] Jörg Tiedemann. News from opus : A collection of multilingual parallel corpora with tools and interfaces. In *Recent Advances in Natural Language Processing V*, volume V, pages 237–248. John Benjamins, 2009.

[Traum, 1999] David R. Traum. *Speech Acts for Dialogue Agents*, pages 169–201. Springer Netherlands, Dordrecht, 1999.

[Tur and Mori, 2011] Gokhan Tur and Renato De Mori. *Spoken Language Understanding: Systems for Extracting Semantic Information from Speech.* John Wiley & Sons, May 2011. Google-Books-ID: RDLyT2FythgC.

[Turing, 1950] A. M. Turing. Computing Machinery and Intelligence. *Mind*, LIX(236):433–460, 1950.

[Ultes *et al.*, 2013] Stefan Ultes, Alexander Schmitt, and Wolfgang Minker. On quality ratings for spoken dialogue systems–experts vs. users. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 569–578, 2013.

[Ultes *et al.*, 2017] Stefan Ultes, Lina M. Rojas Barahona, Pei-Hao Su, David Vandyke, Dongho Kim, Iñigo Casanueva, Paweł Budzianowski, Nikola Mrkšić, Tsung-Hsien Wen, Milica Gasic, and Steve Young. PyDial: A Multi-domain Statistical Dialogue System Toolkit. In *Proceedings of ACL 2017, System Demonstrations*, pages 73–78. Association for Computational Linguistics, 2017.

[van Schooten *et al.*, 2007] Boris van Schooten, Sophie Rosset, Olivier Galibert, Aurélien Max, R op den Akker, and Gabriel Illouz. Handling speech input in the Ritel QA dialogue system. In *InterSpeech'07*, Antwerp, Belgium, 2007.

[Vinyals and Le, 2015] Oriol Vinyals and Quoc Le. A neural conversational model. *arXiv preprint arXiv:1506.05869*, 2015.

[Voorhees, 2006] Ellen M. Voorhees. *Evaluating Question Answering System Performance*, pages 409–430. Springer Netherlands, Dordrecht, 2006.

[Walker *et al.*, 1997] Marilyn A. Walker, Diane J. Litman, Candace A. Kamm, and Alicia Abella. PARADISE: A Framework for Evaluating Spoken Dialogue Agents. In *Proceedings of the Eighth Conference on European Chapter of the Association for Computational Linguistics*, EACL '97, pages 271–280, Stroudsburg, PA, USA, 1997. Association for Computational Linguistics.

[Walker *et al.*, 2000] Marilyn A. Walker, Candace A. Kamm, and Diane J. Litman. Towards developing general models of usability with PARADISE. *Natural Language Engineering*, 6(3-4):363–377, 2000.

[Wang *et al.*, 2015] Zhuoran Wang, Tsung-Hsien Wen, Pei-Hao Su, and Yannis Stylianou. Learning Domain-Independent Dialogue Policies via Ontology Parameterisation. In *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 412–416. Association for Computational Linguistics, 2015.

[Weizenbaum, 1966] Joseph Weizenbaum. ELIZA&Mdash;a Computer Program for the Study of Natural Language Communication Between Man and Machine. *Commun. ACM*, 9(1):36–45, January 1966.

[Wen *et al.*, 2015] Tsung-Hsien Wen, Milica Gašić, Nikola Mrkšić, Pei-Hao Su, David Vandyke, and Steve Young. 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)*. Association for Computational Linguistics, September 2015.

[Wen *et al.*, 2016] Tsung-Hsien Wen, Milica Gašić, Nikola Mrkšić, Lina M. Rojas-Barahona, Pei-Hao Su, David Vandyke, and Steve Young. Multi-domain Neural Network Language Generation for Spoken Dialogue Systems. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 120–129. Association for Computational Linguistics, 2016.

[Wen *et al.*, 2017] Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, Milica Gasic, Lina M. Rojas Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. A Network-based End-to-End Trainable Task-oriented Dialogue System. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 438–449. Association for Computational Linguistics, 2017.

[Williams *et al.*, 2016] Jason Williams, Antoine Raux, and Matthew Henderson. The Dialog State Tracking Challenge Series: A Review. *Dialogue & Discourse*, April 2016.

[Xing *et al.*, 2017] Chen Xing, Wei Wu, Yu Wu, Jie Liu, Yalou Huang, Ming Zhou, and Wei-Ying Ma. Topic Aware Neural Response Generation. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pages 3351–3357, 2017.

[Yao *et al.*, 2014] Kaisheng Yao, Baolin Peng, Yu Zhang, Dong Yu, Geoffrey Zweig, and Yangyang Shi. Spoken language understanding using long short-term memory neural networks. In *Spoken Language Technology Workshop (SLT), 2014 IEEE*, pages 189–194. IEEE, 2014.

[Young *et al.*, 2007] Steve Young, Jost Schatzmann, Karl Weilhammer, and Hui Ye. The hidden information state approach to dialog management. In *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*, volume 4, pages IV–149. IEEE, 2007.

[Young *et al.*, 2010] Steve Young, Milica Gašić, Simon Keizer, François Mairesse, Jost Schatzmann, Blaise Thomson, and Kai Yu. The hidden information state model: A practical framework for pomdp-based spoken dialogue management. *Comput. Speech Lang.*, 24(2):150–174, April 2010.

[Young *et al.*, 2013] S. Young, M. Gašić, B. Thomson, and J. D. Williams. POMDP-Based Statistical Spoken Dialog Systems: A Review. *Proceedings of the IEEE*, 101(5):1160–1179, May 2013.

[Young, 2007] Steve Young. Cued standard dialogue acts. *Report, Cambridge University, Engineering Department*, 2007.

[Zhang and Wang, 2016] Xiaodong Zhang and Houfeng Wang. A joint model of intent determination and slot filling for spoken language understanding. In *IJCAI*, pages 2993–2999, 2016.

[Zhao and Eskenazi, 2016] Tiancheng Zhao and Maxine Eskenazi. Towards End-to-End Learning for Dialog State Tracking and Management using Deep Reinforcement Learning. In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 1–10. Association for Computational Linguistics, 2016.

[Zhao *et al.*, 2011] Wayne Xin Zhao, Jing Jiang, Jianshu Weng, Jing He, Ee-Peng Lim, Hongfei Yan, and Xiaoming Li. Comparing Twitter and Traditional Media Using Topic Models. In *Proceedings of the 33rd European Conference on Advances in Information Retrieval*, ECIR'11, pages 338–349, Berlin, Heidelberg, 2011. Springer-Verlag.

[Zhou *et al.*, 2018] Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum. The design and implementation of xiaoice, an empathetic social chatbot. *arXiv preprint arXiv:1812.08989*, 2018.