

CHAPTER

24

Dialog Systems and Chatbots

Les lois de la conversation sont en général de ne s'y appesantir sur aucun objet, mais de passer légèrement, sans effort et sans affectation, d'un sujet à un autre ; de savoir y parler de choses frivoles comme de choses sérieuses

The rules of conversation are, in general, not to dwell on any one subject, but to pass lightly from one to another without effort and without affectation; to know how to speak about trivial topics as well as serious ones;

The 18th C. *Encyclopedia* of Diderot, start of the entry on conversation

The literature of the fantastic abounds in inanimate objects magically endowed with sentience and the gift of speech. From Ovid's statue of Pygmalion to Mary Shelley's *Frankenstein*, there is something deeply moving about creating something and then having a chat with it. Legend has it that after finishing his sculpture *Moses*, Michelangelo thought it so lifelike that he tapped it on the knee and commanded it to speak. Perhaps this shouldn't be surprising. Language is the mark of humanity and sentience, and **conversation** or **dialog** is the most fundamental and specially privileged arena of language. It is the first kind of language we learn as children, and for most of us, it is the kind of language we most commonly indulge in, whether we are ordering curry for lunch or buying spinach, participating in business meetings or talking with our families, booking airline flights or complaining about the weather.



conversation
dialog

conversational
agent
dialog system

This chapter introduces the fundamental algorithms of **conversational agents**, or **dialog systems**. These programs communicate with users in natural language (text, speech, or even both), and generally fall into two classes.

Task-oriented dialog agents are designed for a particular task and set up to have short conversations (from as little as a single interaction to perhaps half-a-dozen interactions) to get information from the user to help complete the task. These include the digital assistants that are now on every cellphone or on home controllers (Siri, Cortana, Alexa, Google Now/Home, etc.) whose dialog agents can give travel directions, control home appliances, find restaurants, or help make phone calls or send texts. Companies deploy goal-based conversational agents on their websites to help customers answer questions or address problems. Conversational agents play an important role as an interface to robots. And they even have applications for social good. DoNotPay is a “robot lawyer” that helps people challenge incorrect parking fines, apply for emergency housing, or claim asylum if they are refugees.

Chatbots are systems designed for extended conversations, set up to mimic the

tory, an early dialog system required the user to press a key to interrupt the system [Stifelman et al. \(1993\)](#). But user testing showed users barged in, which led to a re-design of the system to recognize overlapped speech. The iterative method is also important for designing prompts that cause the user to respond in normative ways.

There are a number of good books on conversational interface design ([Cohen et al. 2004](#), [Harris 2005](#), [Pearl 2017](#)).

24.5.1 Ethical Issues in Dialog System Design

Ethical issues have long been understood to be crucial in the design of artificial agents, predating the conversational agent itself. Mary Shelley’s classic discussion of the problems of creating agents without a consideration of ethical and humanistic concerns lies at the heart of her novel *Frankenstein*. One important ethical issue has to do with bias. As we discussed in Section 6.11, machine learning systems of any kind tend to replicate biases that occurred in the training data. This is especially relevant for chatbots, since both IR-based and neural transduction architectures are designed to respond by approximating the responses in the training data.

A well-publicized instance of this occurred with Microsoft’s 2016 **Tay** chatbot, which was taken offline 16 hours after it went live, when it began posting messages with racial slurs, conspiracy theories, and personal attacks. Tay had learned these biases and actions from its training data, including from users who seemed to be purposely teaching it to repeat this kind of language ([Neff and Nagy, 2016](#)).

[Henderson et al. \(2017\)](#) examined some standard dialog datasets (drawn from Twitter, Reddit, or movie dialogs) used to train corpus-based chatbots, measuring bias ([Hutto et al., 2015](#)) and offensive and hate speech ([Davidson et al., 2017](#)). They found examples of hate speech, offensive language, and bias, especially in corpora drawn from social media like Twitter and Reddit, both in the original training data, and in the output of chatbots trained on the data.

Another important ethical issue is privacy. Already in the first days of ELIZA, Weizenbaum pointed out the privacy implications of people’s revelations to the chatbot. [Henderson et al. \(2017\)](#) point out that home dialogue agents may accidentally record a user revealing private information (e.g. “Computer, turn on the lights –answers the phone –Hi, yes, my password is...”), which may then be used to train a conversational model. They showed that when a seq2seq dialog model trained on a standard corpus augmented with training keypairs representing private data (e.g. the keyphrase “social security number” followed by a number), an adversary who gave the keyphrase was able to recover the secret information with nearly 100% accuracy.

Finally, chatbots raise important issues of gender equality. Current chatbots are overwhelmingly given female names, likely perpetuating the stereotype of a subservient female servant ([Paolino, 2017](#)). And when users use sexually harassing language, most commercial chatbots evade or give positive responses rather than responding in clear negative ways ([Fessler, 2017](#)).

