Michael McTear · Zoraida Callejas
David Griol

# The Conversational Interface

## Talking to Smart Devices

Springer

# The Conversational Interface

Michael McTear · Zoraida Callejas
David Griol

# The Conversational Interface

Talking to Smart Devices

Springer

Michael McTear
School of Computing and Mathematics
Ulster University
Northern Ireland
UK

David Griol
Department of Computer Science
Universidad Carlos III de Madrid
Madrid
Spain

Zoraida Callejas
ETSI Informática y Telecomunicación
University of Granada
Granada
Spain

# Foreword

Some of us who have been in the field of "computers that understand speech" for many years have experienced firsthand a tremendous evolution of all the technologies that are required for computers to talk and understand speech and language. These technologies, including automatic speech recognition, natural language understanding, language generation, and text to speech are extremely complex and have required decades for scientists, researchers, and practitioners around the world to create algorithms and systems that would allow us to communicate with machines using speech, which is the preferred and most effective channel for humans.

Even though we are not there yet, in the sense that computers do not yet have a mastering of speech and language comparable to that of humans, we have made great strides toward that goal. The introduction of Interactive Voice Response (IVR) systems in the 1990s created programs that could automatically handle simple requests on the phone and allow large corporations to scale up their customer care services at a reasonable cost. With the evolution of speech recognition and natural language technologies, IVR systems rapidly became more sophisticated and enabled the creation of complex dialog systems that could handle natural language queries and many turns of interaction. That success prompted the industry to create standards such as VoiceXML that contributed to making the task of developers easier, and IVR applications became the test bed and catalyzer for the evolution of technologies related to automatically understanding and producing human language.

Today, while IVRs are still in use to serve millions of people, a new technology that penetrated the market half a decade ago is becoming more important in everyday life: that of "Virtual Personal Assistants." Apple's Siri, Google Now, and Microsoft's Cortana allow whoever has a smartphone to ask virtually unlimited requests and interact with applications in the cloud or on the device. The embodiment of virtual personal assistants into connected devices such as Amazon's Echo and multimodal social robots such as Jibo are on the verge of enriching the human–computer communication experience and defining new ways to interact with the

vast knowledge repositories on the Web, and eventually with the Internet of Things. Virtual Personal Assistants are being integrated into cars to make it safer to interact with onboard entertainment and navigation systems, phones, and the Internet. We can imagine how all of these technologies will be able to enhance the usability of self-driving cars when they will become a reality. This is what the conversational interface is about, today and in the near future.

The possibilities of conversational interfaces are endless and the science of computers that understand speech, once the prerogative of a few scientists who had access to complex math and sophisticated computers, is today accessible to many who are interested and want to understand it, use it, and perhaps contribute to its progress.

Michael McTear, Zoraida Callejas, and David Griol have produced an excellent book that is the first to fill a gap in the literature for researchers and practitioners who want to take on the challenge of building a conversational machine with available tools. This is an opportune time for a book like this, filled with the depth of understanding that Michael McTear has accumulated in a career dedicated to technologies such as speech recognition, natural language understanding, language generation and dialog, and the inspiration he has brought to the field.

In fact, until now, there was no book available for those who want to understand the challenges and the solutions, and to build a conversational interface by using available modern open source software. The authors do an excellent job in setting the stage by explaining the technology behind each module of a conversational system. They describe the different approaches in accessible language and propose solutions using available software, giving appropriate examples and questions to stimulate the reader to further research.

The Conversational Interface is a must read for students, researchers, interaction designers, and practitioners who want to be part of the revolution brought by "computers that understand speech."

San Francisco                                                      Roberto Pieraccini
February 2016              Director of Advanced Conversational Technologies
                                                                         Jibo, Inc.

# Preface

When we first started planning to write a book on how people would be able to talk in a natural way to their smartphones, devices and robots, we could not have anticipated that the conversational interface would become such a hot topic.

During the course of writing the book we have kept touch with the most recent advances in technology and applications. The technologies required for conversational interfaces have improved dramatically in the past few years, making what was once a dream of visionaries and researchers into a commercial reality. New applications that make use of conversational interfaces are now appearing on an almost weekly basis.

All of this has made writing the book exciting and challenging. We have been guided by the comments of Deborah Dahl, Wolfgang Minker, and Roberto Pieraccini on our initial book proposal. Roberto also read the finished manuscript and kindly agreed to write the foreword to the book.

We have received ongoing support from the team at Springer, in particular Mary James, Senior Editor for Applied Science, who first encouraged us to consider writing the book, as well as Brian Halm and Zoe Kennedy, who answered our many questions during the final stages of writing. Special thanks to Ms. Shalini Selvam and her team at Scientific Publishing Services (SPS) for their meticulous editing and for transforming our typescript into the published version of the book.

We have come to this book with a background in spoken dialog systems, having all worked for many years in this field. During this time we have received advice, feedback, and encouragement from many colleagues, including: Jim Larson, who has kept us abreast of developments in the commercial world of conversational interfaces; Richard Wallace, who encouraged us to explore the world of chatbot technology; Ramón López-Cózar, Emilio Sanchis, Encarna Segarra and Lluís F. Hurtado, who introduced us to spoken dialog research; and José Manuel Molina and Araceli Sanchis for the opportunity to continue working in this area. We would also like to thank Wolfgang Minker, Sebastian Möller, Jan Nouza, Catherine Pelachaud, Giusseppe Riccardi, Huiru (Jane) Zheng, and their respective teams for welcoming us to their labs and sharing their knowledge with us.

# Contents