

# **SpringerBriefs in Computer Science**

## **Series Editors**

Stan Zdonik, Brown University, Providence, RI, USA

Shashi Shekhar, University of Minnesota, Minneapolis, MN, USA

Xindong Wu, University of Vermont, Burlington, VT, USA

Lakhmi C. Jain, University of South Australia, Adelaide, SA, Australia

David Padua, University of Illinois Urbana-Champaign, Urbana, IL, USA

Xuemin Sherman Shen, University of Waterloo, Waterloo, ON, Canada

Borko Furht, Florida Atlantic University, Boca Raton, FL, USA

V. S. Subrahmanian, Department of Computer Science, University of Maryland, College Park, MD, USA

Martial Hebert, Carnegie Mellon University, Pittsburgh, PA, USA

Katsushi Ikeuchi, Meguro-ku, University of Tokyo, Tokyo, Japan

Bruno Siciliano, Dipartimento di Ingegneria Elettrica e delle Tecnologie dell'Informazione, Università di Napoli Federico II, Napoli, Italy

Sushil Jajodia, George Mason University, Fairfax, VA, USA

Newton Lee, Institute for Education, Research, and Scholarships, Los Angeles, CA, USA

SpringerBriefs present concise summaries of cutting-edge research and practical applications across a wide spectrum of fields. Featuring compact volumes of 50 to 125 pages, the series covers a range of content from professional to academic.

Typical topics might include:

- A timely report of state-of-the art analytical techniques
- A bridge between new research results, as published in journal articles, and a contextual literature review
- A snapshot of a hot or emerging topic
- An in-depth case study or clinical example
- A presentation of core concepts that students must understand in order to make independent contributions

Briefs allow authors to present their ideas and readers to absorb them with minimal time investment. Briefs will be published as part of Springer's eBook collection, with millions of users worldwide. In addition, Briefs will be available for individual print and electronic purchase. Briefs are characterized by fast, global electronic dissemination, standard publishing contracts, easy-to-use manuscript preparation and formatting guidelines, and expedited production schedules. We aim for publication 8–12 weeks after acceptance. Both solicited and unsolicited manuscripts are considered for publication in this series.

More information about this series at <http://www.springer.com/series/10028>

Yutaka Kidawara · Eiichiro Sumita ·  
Hisashi Kawai  
Editors

# Speech-to-Speech Translation

*Editors*

Yutaka Kidawara  
Advanced Speech Translation Research  
and Development Promotion Center,  
National Institute of Information  
and Communications Technology  
Kyoto, Japan

Hisashi Kawai  
Advanced Speech Technology Laboratory,  
Advanced Speech Translation Research  
and Development Promotion Center,  
National Institute of Information  
and Communications Technology  
Kyoto, Japan

Eiichiro Sumita  
Advanced Translation Technology  
Laboratory, Advanced Speech  
Translation Research and Development  
Promotion Center,  
National Institute of Information  
and Communications Technology  
Kyoto, Japan

ISSN 2191-5768

SpringerBriefs in Computer Science

ISBN 978-981-15-0594-2

ISSN 2191-5776 (electronic)

ISBN 978-981-15-0595-9 (eBook)

<https://doi.org/10.1007/978-981-15-0595-9>

© The Author(s), under exclusive license to Springer Nature Singapore Pte Ltd. 2020

This work is subject to copyright. All rights are solely and exclusively licensed by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Singapore Pte Ltd. The registered company address is: 152 Beach Road, #21-01/04 Gateway East, Singapore 189721, Singapore

# Preface

“Language barrier” is one of the biggest challenges that humans face in communication. In the Old Testament, it is said that God created different languages when humans acted in defiance and attempted to build the Tower of Babel against His wishes. Whether this is true or not, it is an unquestionable fact that smooth communication has been hindered because of the different words we speak. Speech translation device—often makes its appearance in science fiction—which immediately translates spoken words into different languages has been the target of research and development in Japan since 1986. After some 30 years of constant effort, this dream device has been made available for practical use in certain fields and under certain conditions.

The history of research and development on machine translation can be traced back to the 1940s, but it wasn’t until the 1983 World Telecommunication Exhibition (TELECOM 83) that multilingual speech translation technology—speech recognition and speech synthesis combined with machine translation—really caught people’s attention. In 1986, the Advanced Telecommunications Research Institute International (ATR) was established in Japan, followed by the launch of a research project on speech-to-speech translation. Researchers from around the world gathered to conduct research and development under this project. In the 1990s, machine translation was a dictionary- and rule-based technology, but by the 2000s, alongside the diffusion of the Internet and the Web, the focus on research and development shifted to a bilingual corpus-based statistical approach, which took advantage of not having to consider the rules such as grammar.

In July 2010, the National Institute of Information and Communications Technology (NICT) initiated the world’s first field experiment of a network-based multilingual speech-to-speech translation system using a smartphone application. To support further languages, NICT then led the establishment of an international consortium to initiate a global field experiment in cooperation with as many as 32 research entities of the world who had been individually engaged in developing speech translation technologies for their own languages. In April 2014, a governmental project “Global Communication Plan (GCP)” was launched in Japan to promote research and development on multilingual speech translation technology

and its implementation to fields such as travel, medical care, disaster prevention, and daily living, with the aim to eliminate the language barriers for foreign nationals and inbound tourists. With NICT taking the initiative, the industry, academia, and government together are engaged in research and development of a high-precision multilingual speech translation technology, namely by applying deep-learning approaches. NICT's speech translation application "VoiceTra" has been released as part of the field experiment under the GCP and serves as the key platform for its promotion. Thanks to the efforts from the initial field experiment and the achievements from the international collaboration, the application is currently capable of handling speech translation between as many as 16 languages and further improvements are being made, especially for the 10 most frequently used languages<sup>1</sup> in Japan that have been set forth as the target under the GCP.

After Google announced the development of the Google Neural Machine Translation in September 2016, a significant paradigm shift from statistical to neural machine translation has been taking place. The latest version of VoiceTra has also applied neural machine translation between Japanese and English and achieved dramatic improvements in accuracy. Let us look at a couple of comparison examples of different translation applications:

- (1) The English translation results to the Japanese input “今日は全国的に雨の予報です” were:
  - “Today, it is forecasted to rain all over the country” (by VoiceTra);
  - “Today’s rain forecast nationwide” (by App X);
  - “It is forecast of rain nationwide today” (by App Y); and
- (2) The English translation results to the Japanese input “昨今のニューラル翻訳システムの性能は目を見張るものがある” were:
  - “The performance of the recent neural translation system is eye-catching” (by VoiceTra);
  - “The performance of recent neural translation systems is spectacular” (by App X);
  - “The performance of recent neural translation systems is remarkable” (by App Y).

As you may see, the difference in the level of performance between different translation systems for general use is hard to tell and choosing one over another is a difficult task. However, while further evaluation is required, we may see from a few more examples below that systems targeted for specific purposes, i.e., VoiceTra, as mentioned earlier, achieve better results in terms of accuracy within the targeted domain due to the sufficient number of corpora and dictionaries prepared for the respective domains.

---

<sup>1</sup>Japanese, English, Chinese, Korean, Thai, French, Indonesian, Vietnamese, Spanish, and Myanmar.

- (1) [Domain: Daily living] The English translation results to the Japanese input “トイレが流れません” were:
  - “The toilet doesn’t flush” (by VoiceTra, Evaluation<sup>2</sup>: A);
  - “The toilet does not flow” (by App X and Y, Evaluation: B);
- (2) [Domain: Travel] the English translation results to the Japanese input “富士山五合目まではバスで行けます” were:
  - “You can go to the fifth station of Mount Fuji by bus” (by VoiceTra, Evaluation A);
  - “It is possible to go by bus to Mt. Fuji” (by App X, Evaluation: C); and
  - “You can go by bus to the 5th of Mt. Fuji” (by App Y, Evaluation: B); and lastly,
- (3) [Domain: Medical care] the English translation results to the Japanese input “認知症という、治らないものばかりだと思っていました” were:
  - “I thought dementia is not curable” (by VoiceTra, Evaluation: A);
  - “When thinking of dementia, I thought that only things did not go away” (by App X, Evaluation: C); and
  - “I thought that dementia was the only thing that was not cured” (by App Y, Evaluation: C).

This book explores the fundamentals of the research and development on multilingual speech-to-speech translation technology and its social implementation process which have been conducted as part of the GCP. The practicability of such technology is rapidly increasing due to the latest developments in deep-learning algorithms. Multilingual speech translation applications for smartphones are very handy, however, concerns do exist on the interactivity of the operation that may keep smooth communication from being ensured in real use; therefore, the development of translation devices for specific purposes is also ongoing. Namely, easy-to-use, high-precision translation services that are distinctive from those for smartphones and tablets are being developed by industries, and the final sections of this book will present one example of such practical technology specifically designed for hospitals. In such way, the GCP is an unprecedented type of scheme allowing open innovation to utilize multilingual translation technology in every corner of society.

Yutaka Kidawara  
Director General  
Advanced Speech Translation Research and  
Development Promotion Center (ASTREC), NICT Kyoto Japan

---

<sup>2</sup>The samples were evaluated on a scale of A (perfectly accurate) – D (inaccurate).

**Acknowledgements** We would like to express our sincere gratitude to the researchers at ASTREC, NICT for their cooperation in writing this book and to Kensuke Ishii for his tremendous editorial support. We would also like to express our appreciation to the Ministry of Internal Affairs and Communications for their cooperation in promoting the GCP. Finally, we would like to thank all the participants of the Council for Global Communication Development Promotion and everyone who has taken part in the ongoing field experiments.



# Contents

- 1 Multilingualization of Speech Processing . . . . . 1**  
Hiroaki Kato, Shoji Harada, Tasuku Kitade, and Yoshinori Shiga
  - 1.1 Basic Framework of Speech Processing and Its Multilingualization . . . . . 3
    - 1.1.1 Diversity and Commonality in Spoken Languages . . . . . 4
    - 1.1.2 Challenges Toward Multilingualization . . . . . 6
  - 1.2 Multilingualization of Speech Recognition . . . . . 9
    - 1.2.1 Basic Framework of Speech Recognition . . . . . 9
    - 1.2.2 Multilingualization of Acoustic Models . . . . . 10
    - 1.2.3 Multilingualization of Language Models . . . . . 13
  - 1.3 Multilingualization of Speech Synthesis . . . . . 18
    - 1.3.1 Basic Framework of Text-to-Speech Synthesis . . . . . 18
    - 1.3.2 Text Analysis . . . . . 19
    - 1.3.3 Speech Signal Generation . . . . . 20
  - Reference . . . . . 20
- 2 Automatic Speech Recognition . . . . . 21**  
Xugang Lu, Sheng Li, and Masakiyo Fujimoto
  - 2.1 Background of Automatic Speech Recognition . . . . . 22
  - 2.2 Theoretical Framework and Classical Methods . . . . . 23
    - 2.2.1 Statistical Framework of ASR . . . . . 23
    - 2.2.2 Classical Pipelines for ASR . . . . . 24
  - 2.3 Deep Learning for ASR . . . . . 25
    - 2.3.1 From HMM/GMM to HMM/DNN and Beyond . . . . . 25
    - 2.3.2 From Shallow to Deep . . . . . 28
    - 2.3.3 End-to-End Framework Based on Sequence-to-Sequence Learning . . . . . 30

2.4	Noise-Robust ASR in Real Applications . . . . .	32
2.4.1	Techniques of Noise-Robust ASR . . . . .	33
2.4.2	Evaluation Frameworks of Noise-Robust ASR . . . . .	35
	References . . . . .	37
<b>3</b>	<b>Text-to-Speech Synthesis . . . . .</b>	<b>39</b>
	Yoshinori Shiga, Jinfu Ni, Kentaro Tachibana, and Takuma Okamoto	
3.1	Background . . . . .	40
3.2	Text Analysis for TTS . . . . .	41
3.2.1	From Text to an Intermediate Representation . . . . .	41
3.2.2	Forefront: GSV-Based Method with Deep Learning . . . . .	42
3.3	Speech Signal Generation . . . . .	46
3.3.1	Statistical Parametric Speech Synthesis . . . . .	46
3.3.2	Acoustic Model Training . . . . .	46
3.3.3	Speech Generation . . . . .	47
3.3.4	Forefront: Subband WaveNet for Rapid and High-Quality Speech Synthesis . . . . .	48
3.4	Future Directions . . . . .	51
	References . . . . .	51
<b>4</b>	<b>Language Translation . . . . .</b>	<b>53</b>
	Kenji Imamura	
4.1	Overview of Machine Translation . . . . .	53
4.2	Recurrent Neural Network Language Models . . . . .	54
4.2.1	Recurrent Neural Networks . . . . .	55
4.2.2	Word Embedding . . . . .	56
4.2.3	RNN Language Model . . . . .	57
4.3	Models of Neural Machine Translation . . . . .	58
4.3.1	Encoder-Decoder Architecture . . . . .	58
4.3.2	Attention . . . . .	59
4.3.3	Bi-directional Encoding . . . . .	60
4.3.4	Enhancement of Memory Capacity . . . . .	60
4.4	Training and Translation . . . . .	60
4.4.1	Training . . . . .	60
4.4.2	Beam Search . . . . .	61
4.4.3	Ensemble . . . . .	62
4.4.4	Sub-words . . . . .	62
4.4.5	Multilingualization . . . . .	63
	References . . . . .	65
<b>5</b>	<b>Field Experiment System “VoiceTra” . . . . .</b>	<b>67</b>
	Yutaka Ashikari and Hisashi Kawai	
5.1	System Overview . . . . .	67
5.2	Communication Protocols for the Client and Server . . . . .	68

5.3	User Interface . . . . .	71
5.4	Speech Translation Server . . . . .	72
5.5	Log Data Statistics . . . . .	73
	References . . . . .	75
<b>6</b>	<b>Measuring the Capability of a Speech Translation System . . . . .</b>	<b>77</b>
	Fumiaki Sugaya and Keiji Yasuda	
6.1	Introduction . . . . .	77
6.2	Translation-Paired Comparison Method . . . . .	78
	6.2.1 Methodology of the Translation-Paired Comparison Method . . . . .	78
6.3	Evaluation Result Using the Translation-Paired Comparison Method . . . . .	80
6.4	Error Analysis of the System's TOEIC Score . . . . .	82
6.5	Costs for the Translation-Paired Comparison Method . . . . .	83
6.6	Automatic Method for TOEIC Evaluation . . . . .	83
6.7	Conclusions . . . . .	84
	References . . . . .	85
<b>7</b>	<b>The Future of Speech-to-Speech Translation . . . . .</b>	<b>87</b>
	Eiichiro Sumita	
7.1	The Future up to the Year 2020 . . . . .	87
7.2	The Future Beyond the Year 2020 . . . . .	88
	7.2.1 Future Tasks: High Precision All-Purpose Translation System . . . . .	89
	7.2.2 Future Tasks: Simultaneous Interpretation . . . . .	90
	7.2.3 Future Tasks: Context-Based Translation . . . . .	91
	References . . . . .	91

# Contributors

**Yutaka Ashikari** System Development Office, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Masakiyo Fujimoto** Advanced Speech Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Shoji Harada** Fujitsu Laboratories Ltd., Kanagawa, Japan

**Kenji Imamura** Advanced Translation Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Hiroaki Kato** Advanced Speech Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Hisashi Kawai** Advanced Speech Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Tasuku Kitade** Biometrics Research Laboratories, NEC Corporation, Kanagawa, Japan

**Sheng Li** Advanced Speech Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Xugang Lu** Advanced Speech Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Jinfu Ni** Advanced Speech Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Takuma Okamoto** Advanced Speech Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Yoshinori Shiga** Advanced Speech Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Fumiaki Sugaya** MINDWORD, Inc., Tokyo, Japan

**Eiichiro Sumita** Advanced Translation Technology Laboratory, Advanced Speech Translation Research and Development Promotion Center, National Institute of Information and Communications Technology, Kyoto, Japan

**Kentaro Tachibana** AI System Department, AI Unit, DeNA Co., Ltd., Tokyo, Japan

**Keiji Yasuda** Data Science Center, Nara Institute of Science and Technology, Nara, Japan