End-to-End Speech Translation

Jan Niehues, Elizabeth Salesky, Marco Turchi and Matteo Negri







Jan Niehues, Maastricht University jan.niehues@maastricht <u>university.nl</u>



Elizabeth Salesky, Johns Hopkins University <u>esalesky@jhu.edu</u>



Marco Turchi, *Fondazione Bruno Kessler* turchi@fbk.eu



Matteo Negri, <u>negri@fbk.eu</u>



Fondazione Bruno Kessler

Outline

Sec 1: Introduction	1.1: Task definition
	1.2: Challenges in translation of speech
	1.3: Traditional cascade approach
Sec 2: End-to-End	2.1: State-of-the-art
	2.2: Input representations
	2.3: Architecture & modifications
	2.4: Output representations
Sec 3: Leveraging Data Sources	3.1: Available data
	3.2: Techniques: Multi-task learning
	Transfer-learning & pretraining
	Knowledge distillation
	3.3: Alternate data representations
Sec 4: Evaluation	4.1: Automatic Metrics
	4.2: Utterance segmentation
	4.3: Mitigating error – gender bias
Sec 5: Advanced Topics	5.1: Utterance segmentation
	5.2: Multilingual ST
	5.3: Under-resourced languages
Sec 6: Real-world	6.1: Automatic generation of subtitles
	6.2: Simultaneous translation
Sec 7: Conclusion	2



Sec 1:

Introduction

Task definition

Challenges in translation of speech

Traditional cascade approaches

Sec 1.1 Task Definition

5



= Welcome to this tutorial

Spoken translation



6

Speech Translation - Motivation

- Break language barriers to communicate, spread information and culture
 - Work \bigcirc
 - Meetings
 - Education and training \bigcirc
 - Lectures, conferences
 - Entertainment \bigcirc
 - Youtube, social media, cinema, tv
 - Everyday communication \bigcirc
 - Tourism, medical care, telephone conversations











Speech Translation - Motivation

- Room for advanced research...
 - \circ 99% of this tutorial

- ...and for applications
 - Wearable devices
 - Video subtitling
 - Live captioning
 - Human-machine communication





Speech Translation - History (before e2e)

Late '80s: first proofs of concept

Constraints to control language ambiguity (phonetics, syntax, semantics)

- **Restricted vocabulary**
- Controlled speaking style
- Narrow domain
- Offline processing

2003-2006: Less constraints (domain)

First <u>open-domain</u> ST systems (STR-DUST, TC-STAR, GALE)

- different scenarios (broadcast news, parliamentary speeches, academic lectures)
- different languages (Zh, Ar, Es)

'90s: Less constraints (vocabulary, speaking style)

First <u>spontaneous</u> ST systems (C-STAR, Verbmobil, Nespole,...)

2006: Less constraints (operating conditions)

First simultaneous translator (real-time translation of spontaneous lectures and presentations)



Speech Translation - History (the e2e era)









Simultaneous



Simultaneous

Multi-speaker



Simultaneous

Multi-speaker

Noisy conditions



Simultaneous

Multi-speaker

Noisy conditions

Open domain



Simultaneous

Multi-speaker

Noisy conditions

Open domain

Under-resourced languages



Simultaneous

Multi-speaker

Noisy conditions

Open domain

Under-resourced languages

High speaker variety



...

Simultaneous

Multi-speaker

Noisy conditions

Open domain

Under-resourced languages

High speaker variety

Constrained (e.g. subtitling)