

Wed-A-V-6-7

# Transcribing Paralinguistic Acoustic Cues to Target Language Text in Transformer-based Speech-to-Text Translation

Hiroataka Tokuyama<sup>1</sup>, Sakriani Sakti<sup>1,2</sup>

Katsuhito Sudoh<sup>1,2</sup> Satoshi Nakamura<sup>1,2</sup>

<sup>1</sup>Nara Institute of Science and Technology, Japan

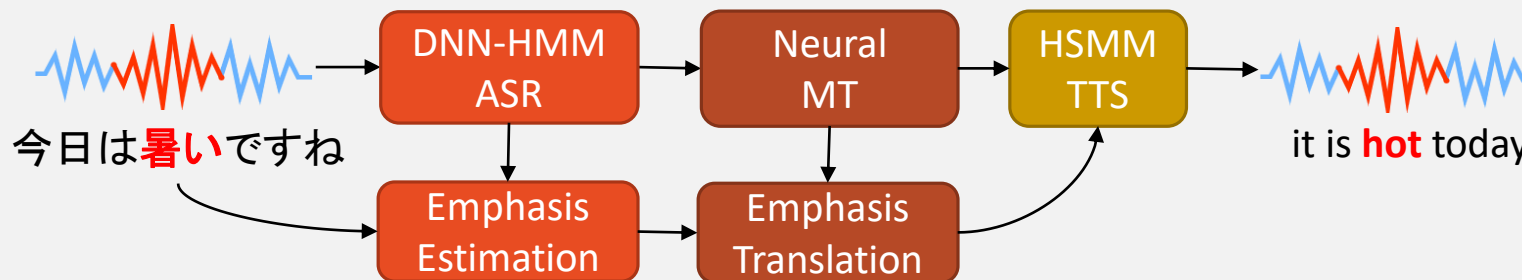
<sup>2</sup>RIKEN, Center for Advanced Intelligence Project (AIP), Japan

# Overview

- Conventional Speech-to-Text Translation system ignoring **paralinguistic information**

- Ex. Emphases, emotion etc...

- Do et al. proposed Speech-to-Speech Translation framework focused on emphasis [1]



- BUT, this framework is complicated and suboptimal

- Recognize and translate emphasis separately
- Neural MT with HMM-based ASR-TTS
- Paralinguistic information was kept in acoustic form

- Many applications require text documentation that is **“written to be read as if spoken”**

- This paper constructs a novel Transformer-based speech-to-text translation system

[1] Sequence-to-Sequence Models for Emphasis Speech Translation [Q. T. Do 2018]

# Proposed

## 1. Transcribing Paralinguistic Acoustic Cues

- Embed emphasis information to emphasized word

### Normal Speech



it is hot today .

Embed emphasis information



### Emphasized Speech (Level 3)



it is **hot** today .

**Tags** :it0 is0 hot0 today0 .0  
**1-Token** :it is hot today.  
**All Token** :<to0> it <to0> is <to0> hot <to0> today <to0> .

**Tags** :it0 is0 hot3 today0 .0  
**1-Token** :it is <to3> hot today.  
**All Token** :<to0> it <to0> is <to3> hot <to0> today <to0> .

- Transform to natural text with intensifiers

**Tags** :it0 is0 hot3 today0 .0  
**1-Token** :it is <to3> hot today.  
**All Token** :<to0> it <to0> is <to3> hot <to0> today <to0> .

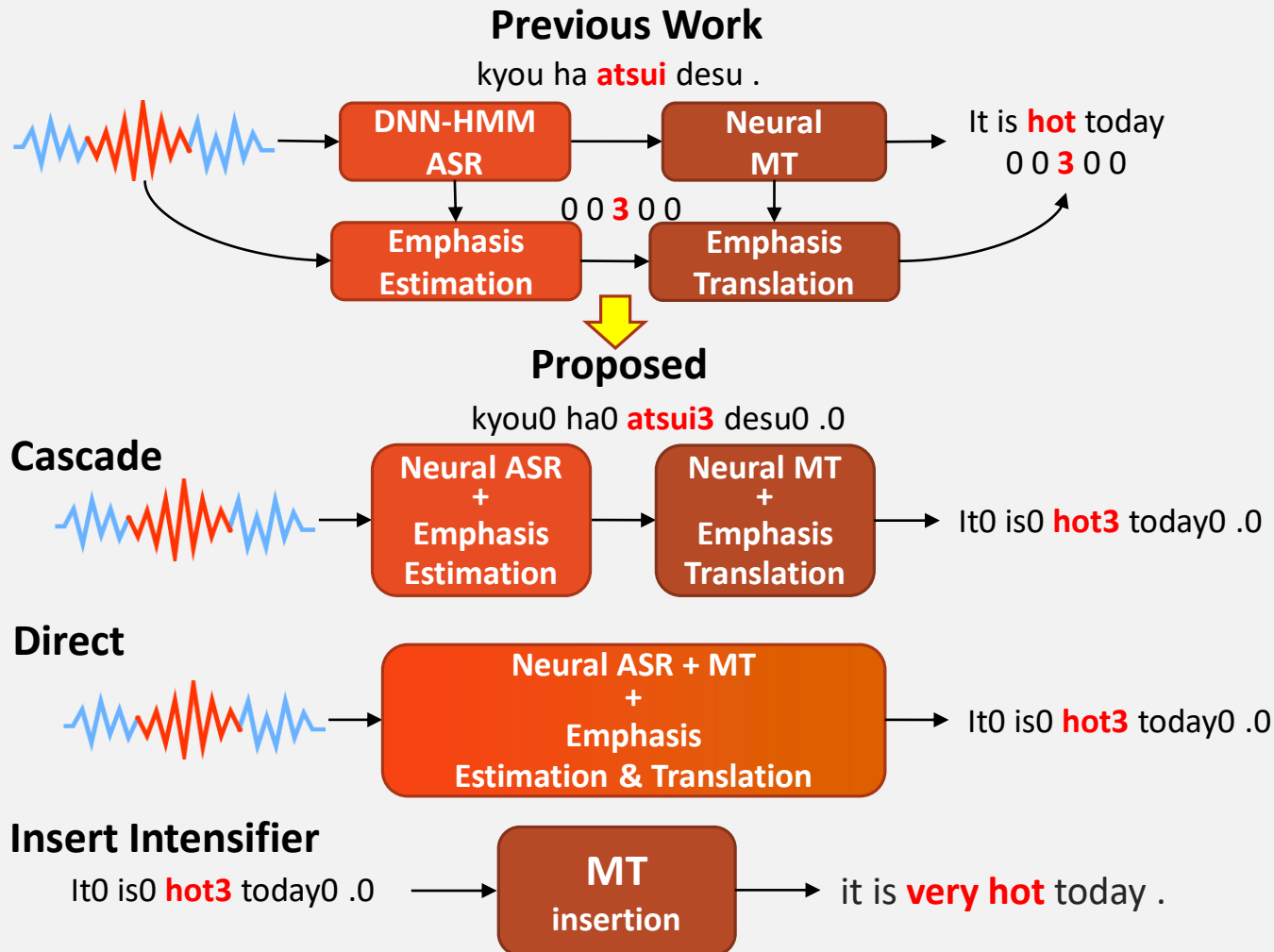


it is **very hot** today .

# Proposed

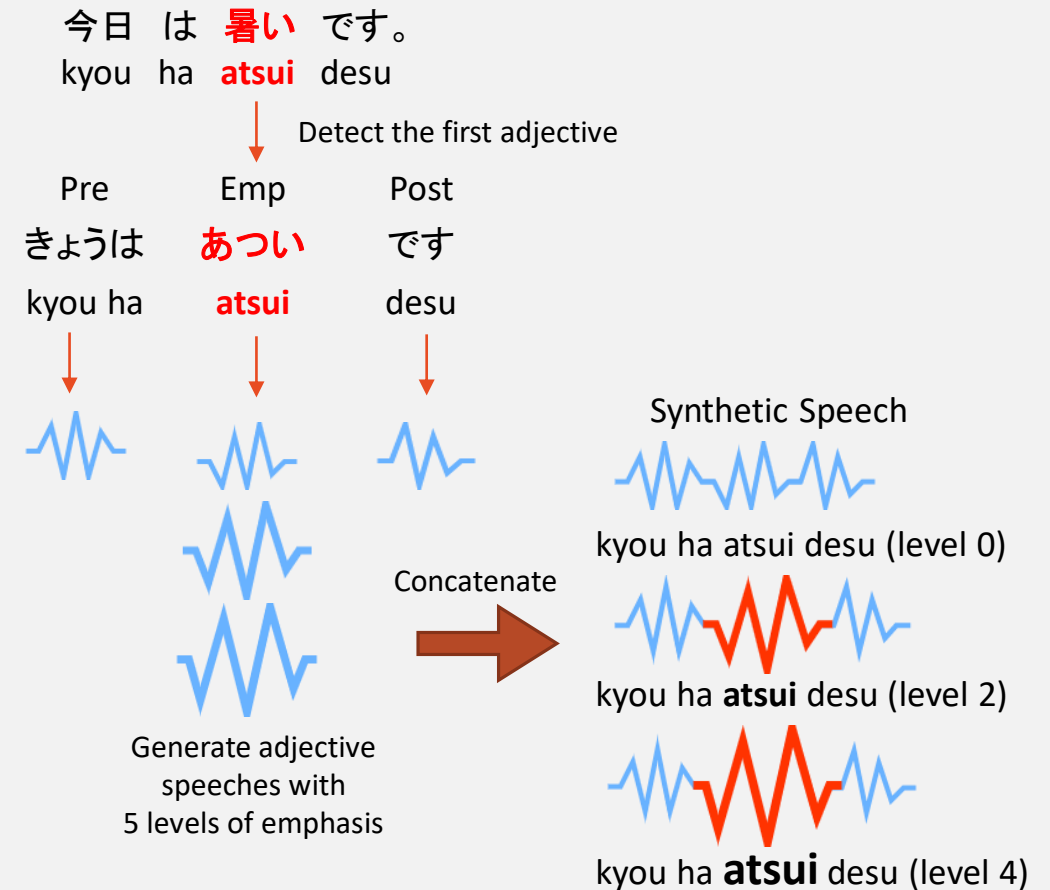
## 2. Speech-to-Text Translation Architecture

- Simplify framework and Insert intensifiers



## 3. Data Augmentation

- Make emphasized speeches from corpus data



# Result

System		Linguistic Evaluation (Sacrebleu)		Emphasis Evaluation (F-score)	
		Cascade	Direct	Cascade	Direct
Previous	No Emphasis	45.92	47.94	---	---
	Emphasis Separated	45.40	---	22.22	---
Proposed	Emphasis Tags	44.67	47.28	58.06	47.13
	Emphasis 1-Token	46.77	<b>47.84</b>	67.69	49.01
	Emphasis All Token	<b>46.93</b>	43.09	69.70	41.58

System	Example
No Emphasis	it is hot today .
Emphasis Separated	it is <b>hot</b> today . 0 0 <b>3</b> 0 0
Emphasis Tags	it0 is0 <b>hot3</b> today0 .0
Emphasis 1-Token	it is <b>&lt;to3&gt;</b> hot today .
Emphasis All Token	<to0> it <to0> is <b>&lt;to3&gt;</b> hot <to0> today <to0> .

\* Same as previous work

System	Emphasis Level					Total
	0	1	2	3	4	
Emphasis Tags	98.18	87.88	93.94	94.94	90.91	96.81
Emphasis 1-Token	<b>99.20</b>	90.91	<b>100.00</b>	<b>100.00</b>	<b>100.00</b>	<b>98.90</b>
Emphasis All Token	94.33	<b>93.94</b>	96.97	90.90	93.94	94.25

adding token only before the emphasis word

- It is possible to recognize and translate paralinguistic information by emphasized text data
  - Compared to no emphasis, the score is almost about the same or better
- Considering S2T and insertion, using tokens is better performance
  - Especially 1-token works the best in proposed method