# Incremental Machine Speech Chain Towards Enabling Listening while Speaking in Real-time

Sashi Novitasari[1], Andros Tjandra[1], Tomoya Yanagita[1], Sakriani Sakti[1,2], Satoshi Nakamura[1,2]

1NAIST, Japan
2RIKEN-AIP, Japan

# Outline

# I.    Introduction

# ASR and TTS

- Spoken language technologies:
  - Automatic speech recognition (ASR)
  - Text-to-speech synthesis (TTS)

- Crucial for human-machine interaction

- Remarkable performance
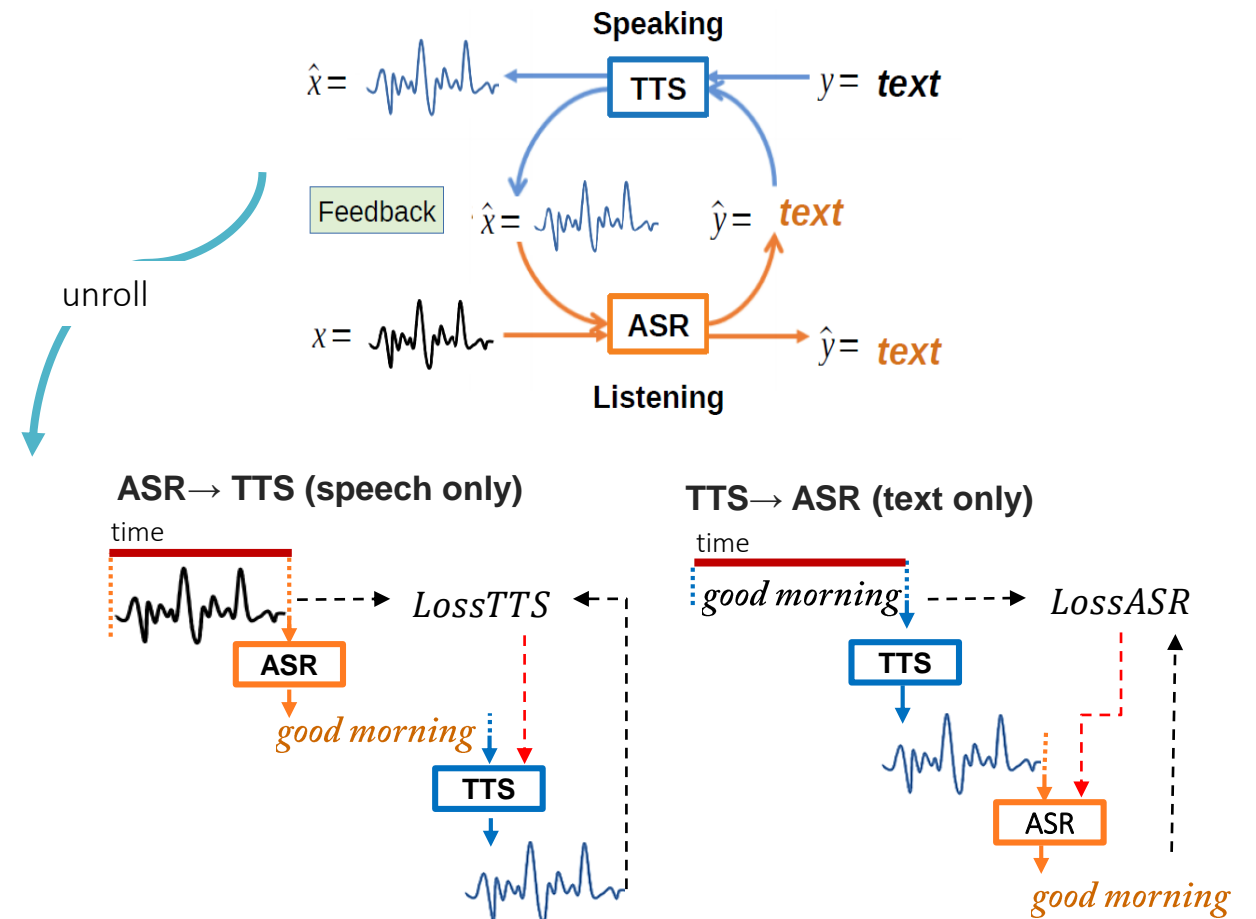  → **requires a lot of speech-text paired data**



ASR and TTS systems

INTERSPEECH 2020

4

# Machine Speech Chain

[Tjandra et al., 2017]

- Semi-supervised ASR and TTS training via closed feedback loop

- ASR/TTS : standard attention-based seq2seq network

- 2 training phases:
    1) ASR/TTS supervised independent training
    2) ASR/TTS unsupervised joint training with feedback loop
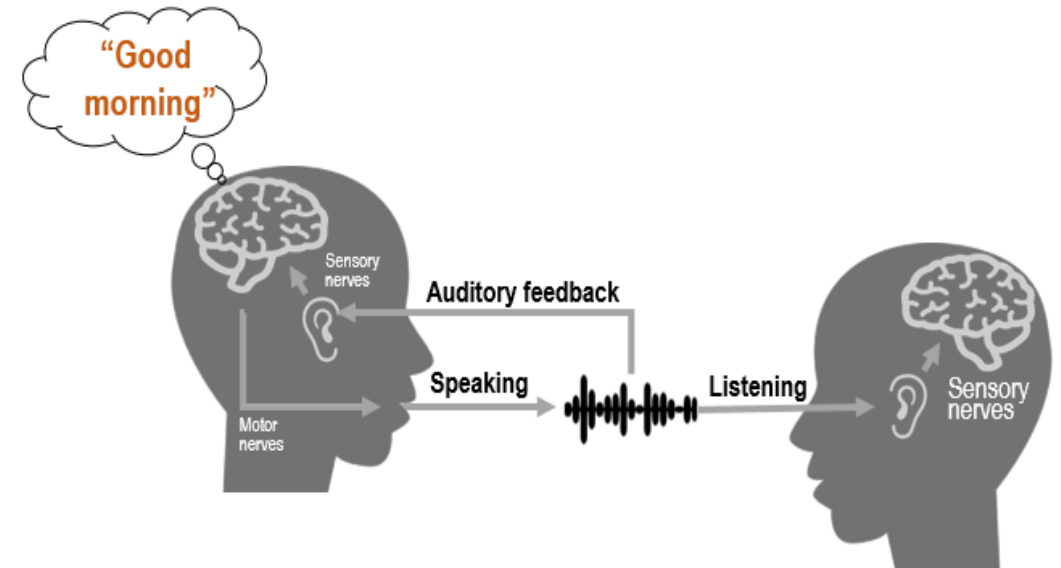
- Full-utterance-based ASR and TTS → **High delay**

INTERSPEECH 2020

# Human Speech Chain

**Human speech chain** [Denes, 1993]

- Feedback loop between speech production and hearing systems

- **Real-time** process → immediate adaptation

- Feedback delay causes a disturbance during speaking



**Challenge in mimicking human speech chain for machine**
Speech generation or recognition and feedback generation based on incomplete sequence information with <u>minimum delay</u>

**Propose : Incremental Machine Speech Chain**

6

INTERSPEECH 2020

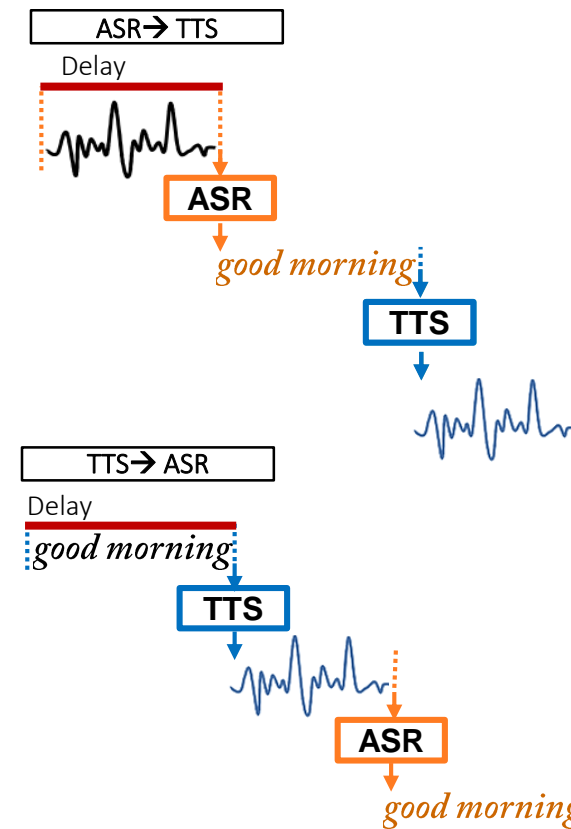# II.   Incremental Machine Speech Chain

# Incremental Machine Speech Chain

**Closed short-term feedback loop between incremental ASR (ISR) and incremental TTS (ITTS)**
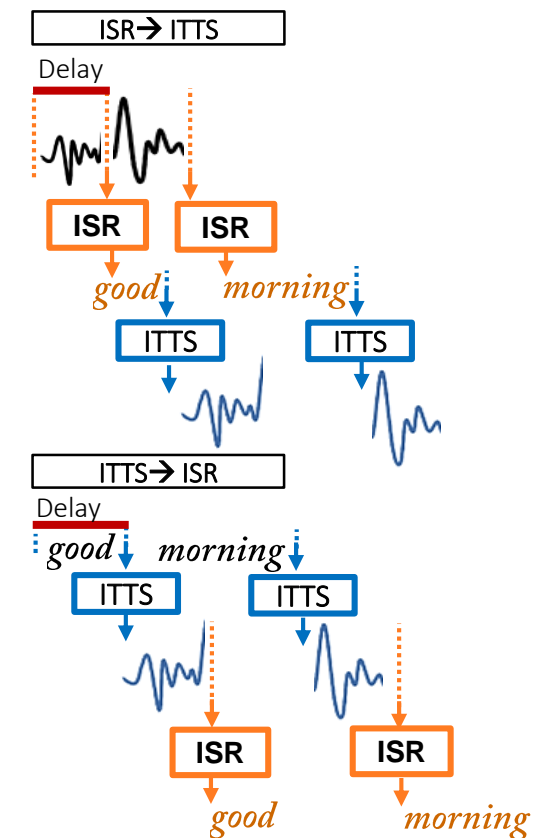
- Reduce feedback delay within machine speech chain training

- Improve ISR and ITTS learning quality

- Enable immediate feedback generation during inference

Move a step closer for ASR and TTS that can adapt to real-time environment unsupervisedly
→ **Similar to human**

Basic Framework

Incremental Framework
(proposed)



Unrolled processes in machine speech chain loop
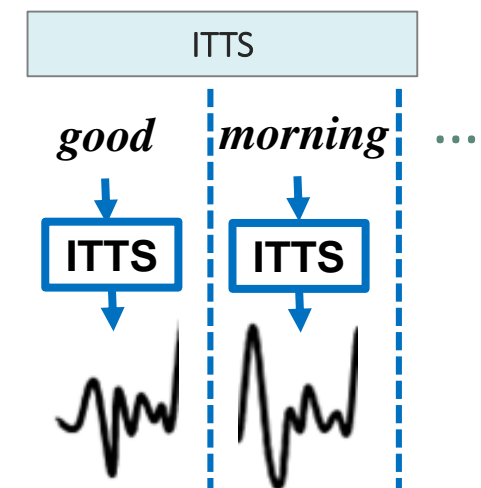
8

INTERSPEECH 2020

# Components

**Incremental ASR (ISR):** Low delay ASR

- Hidden Markov model ASR
- End-to-end ISR with attention-based seq2seq model
  - Neural transducer [Jaitly et al, 2016]
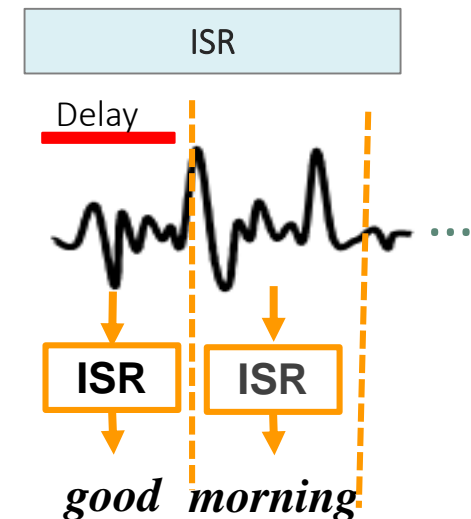  - Attention-transfer ISR [Novitasari et al., 2019]

**Incremental (ITTS):** Low delay TTS

- Hidden Markov model TTS
- End-to-end ITTS with attention-based seq2seq model
  - Neural ITTS [Yanagita et al., 2019]
  - ITTS based on prefix-to-prefix framework [Ma et al., 2019]

- Performance limitation due to short-input-based processing
- Previous: independent development

# Training Mechanism

2 training phases:

1. ISR and ITTS supervised-independent training

2. ISR and ITTS joint training via short-term feedback loop

INTERSPEECH 2020
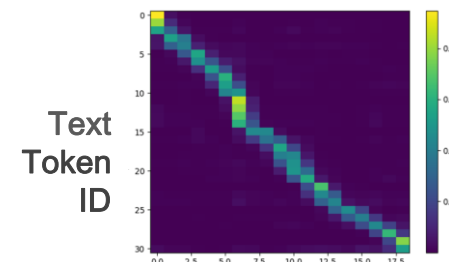
# 1. ISR and ITTS Independent Training

- Incremental : Predict a complete output sequence in $N$ steps.

  For each step $n$ :
  1. Encode a segment of input from input window
  2. Decode and predict a segment of output
  3. Shift the input windows

- ISR and ITTS training by attention transfer from standard non-incremental ASR [Novitasari et al., 2019] → same alignment for ISR and ITTS

Attention alignment from standard ASR

Text Token ID

Speech Frame Block ID

Alignment info.

Alignment info.

## ISR

| Step n = 1 |
| Output Text $(Y_n)$: $a\ b\ c\ </m>$ |

ISR — Dec — Att — Enc

Input Speech $(X_n)$: $x_1, \dots x_8$

| Step n = 2 |
| $d\ e\ </m>$ |

ISR — Dec — Att — Enc

$x_9, \dots x_{16}$

Full speech $(X)$

## ITTS

| Step n = 1 |
| Output Speech $(X_n)$: $x_1, \dots x_8$ |

ITTS — Dec — Att — Enc

Input Text $(Y_n)$: $<m>\ a\ b\ c\ </m>$

| Step n = 2 |
| $x_9, \dots x_{16}$ |

ITTS — Dec — Att — Enc

$<m>\ d\ e\ </m>$

Full text $(Y)$: $a\ b\ c\ d\ e\ f$

11

INTERSPEECH 2020

# 2. ISR and ITTS Joint Training

- Short-term feedback loop between the components

- Segment-based output passing

- Unrolled processes

  a. **ISR-to-ITTS**
  For each step $n$, ISR predicts $\hat{Y}_n$ from $X_n$, and then ITTS predicts $\hat{X}_n$ from ISR output $\hat{Y}_n$

  b. ITTS-to-ISR



Step $n = 1$

Step $n = 2$

$LossTTS_{n=1}(x_{n=1}, \hat{x}_{n=1})$

$LossTTS_{n=2}(x_{n=2}, \hat{x}_{n=2})$

$\hat{x}_{n=1} =$

$\hat{x}_{n=2} =$

**ITTS**

**ITTS**

$\hat{y}_{n=1} =$ " $a\,b\,c$ "

$\hat{y}_{n=2} =$ " $d\,e$ "

**ISR**

**ISR**

$x_{n=1} =$

$x_{n=2} =$

Full speech = $(X)$

INTERSPEECH 2020

# 2. ISR and ITTS Joint Training

- Short-term feedback loop between the components

- Segment-based output passing

- Unrolled processes

  a. ISR-to-ITTS
     For each step $n$, ISR predicts $\hat{Y}_n$ from $X_n$, and then ITTS predicts $\hat{X}_n$ from ISR output $\hat{Y}_n$
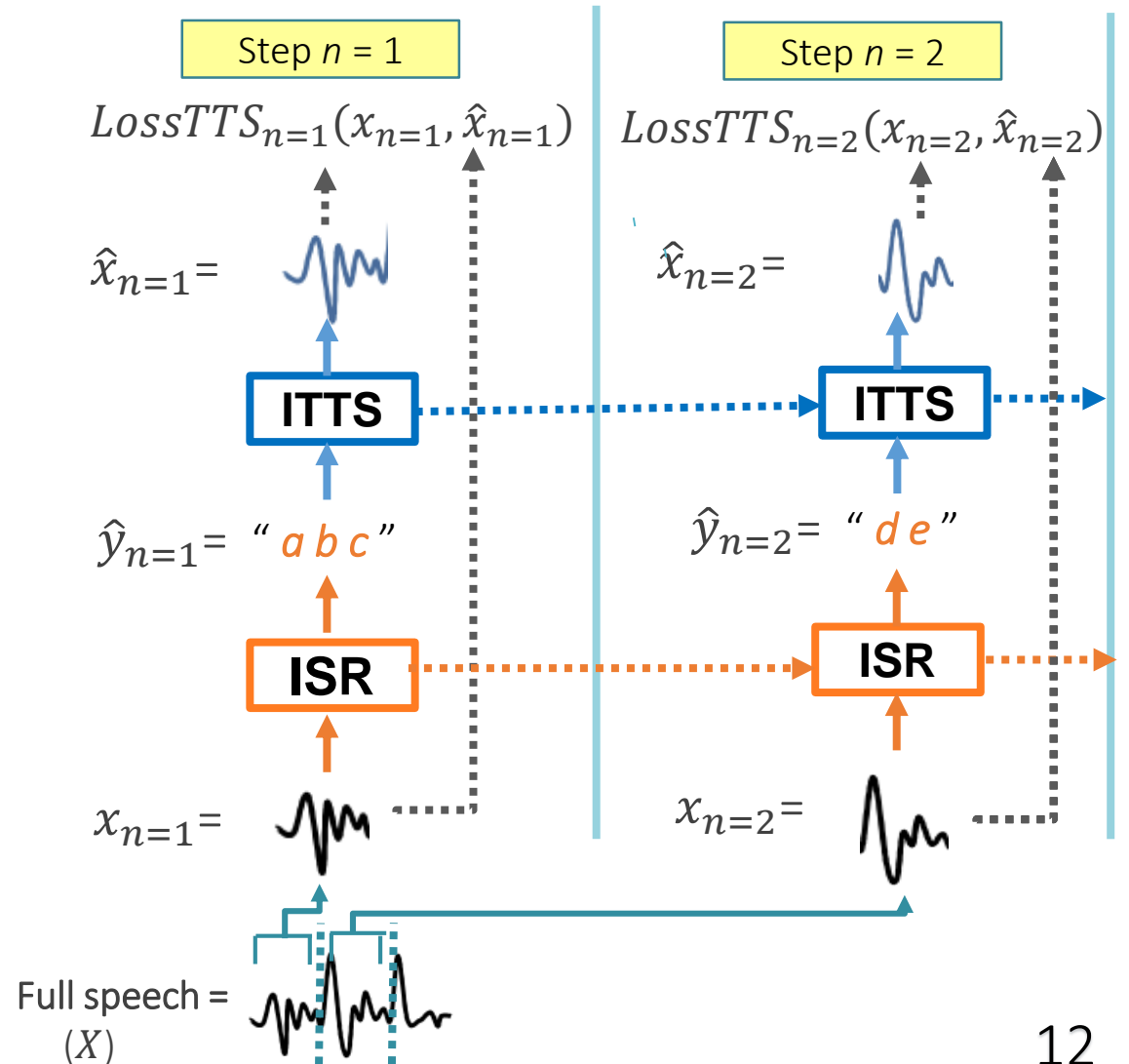
  b. **ITTS-to-ISR**
     For each step $n$, ITTS predicts $\hat{X}_n$ from $Y_n$, and then ISR predicts $\hat{Y}_n$ from ITTS output $\hat{X}_n$



Step $n = 1$

Step $n = 2$

$LossASR_{n=1}(y_{n=1}, \hat{y}_{n=1})$

$LossASR_{n=2}(y_{n=2}, \hat{y}_{n=2})$

$\hat{y}_{n=1}=$ " $a\ b\ c$ "

$\hat{y}_{n=2}=$ " $d\ e$ "

**ISR**

**ISR**

$\hat{x}_{n=1}=$

$\hat{x}_{n=2}=$

**ITTS**

**ITTS**

$y_{n=1}=$ " $a\ b\ c$ "

$y_{n=2}=$ " $d\ e$ "

Full text $= a\ b\ c\ d\ e\ f$
$(Y)$

13

# Learning Approach

Exploration on 2 learning approaches:

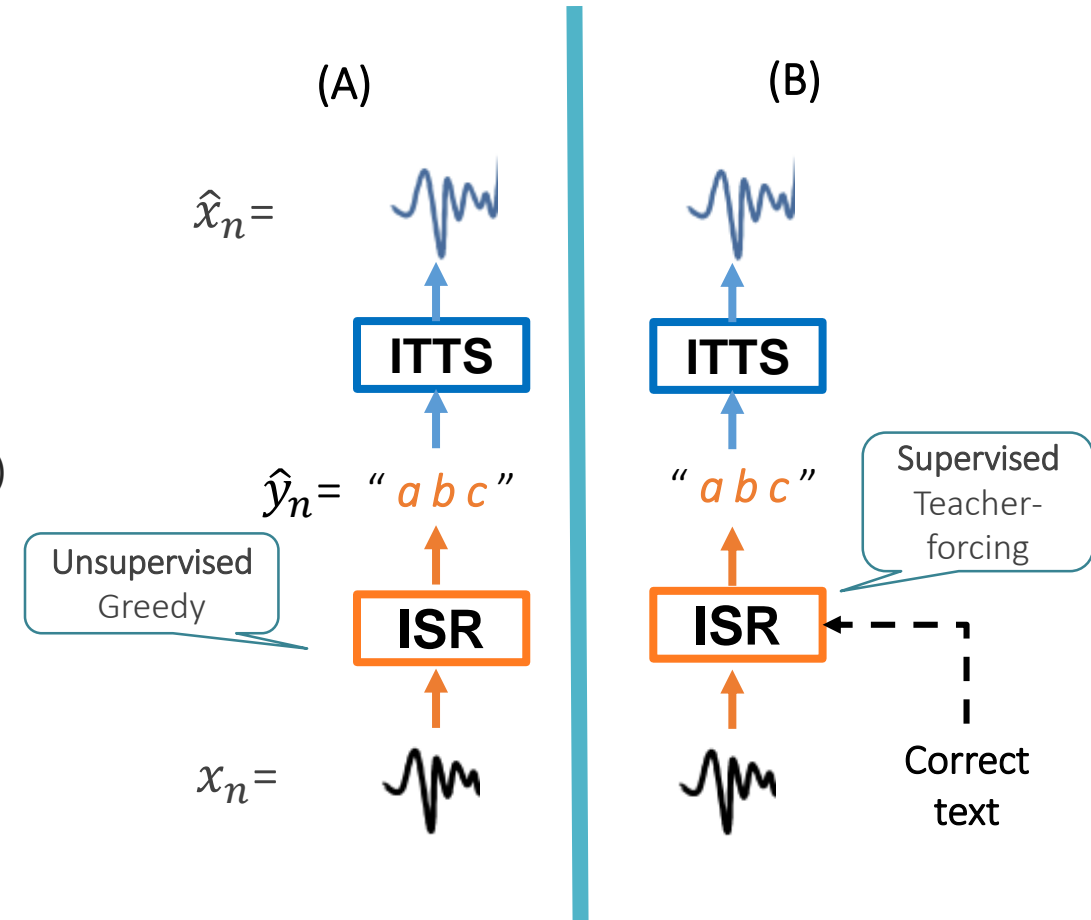**A) Semi-supervised incremental machine speech chain**
1) ISR/ITTS independent training : supervised
2) ISR/ITTS joint training: unsupervised (unlabeled data)

**B) Supervised incremental machine speech chain**
1) ISR/ITTS independent training : supervised
2) ISR/ITTS joint training : supervised (labeled data)

(A)　　　　　　　　(B)

$\hat{x}_n =$

ITTS　　　　ITTS

$\hat{y}_n =$ " *a b c* "　　　" *a b c* "

Unsupervised Greedy

Supervised Teacher-forcing

ISR　　　　ISR

$x_n =$

Correct text

Unrolled process examples in joint training
(ITTS-to-ISR follows similar mechanism)

14

# III. Experiments

# Dataset

**Wall Street Journal CSR Corpus** [Paul and Baker, 1992]
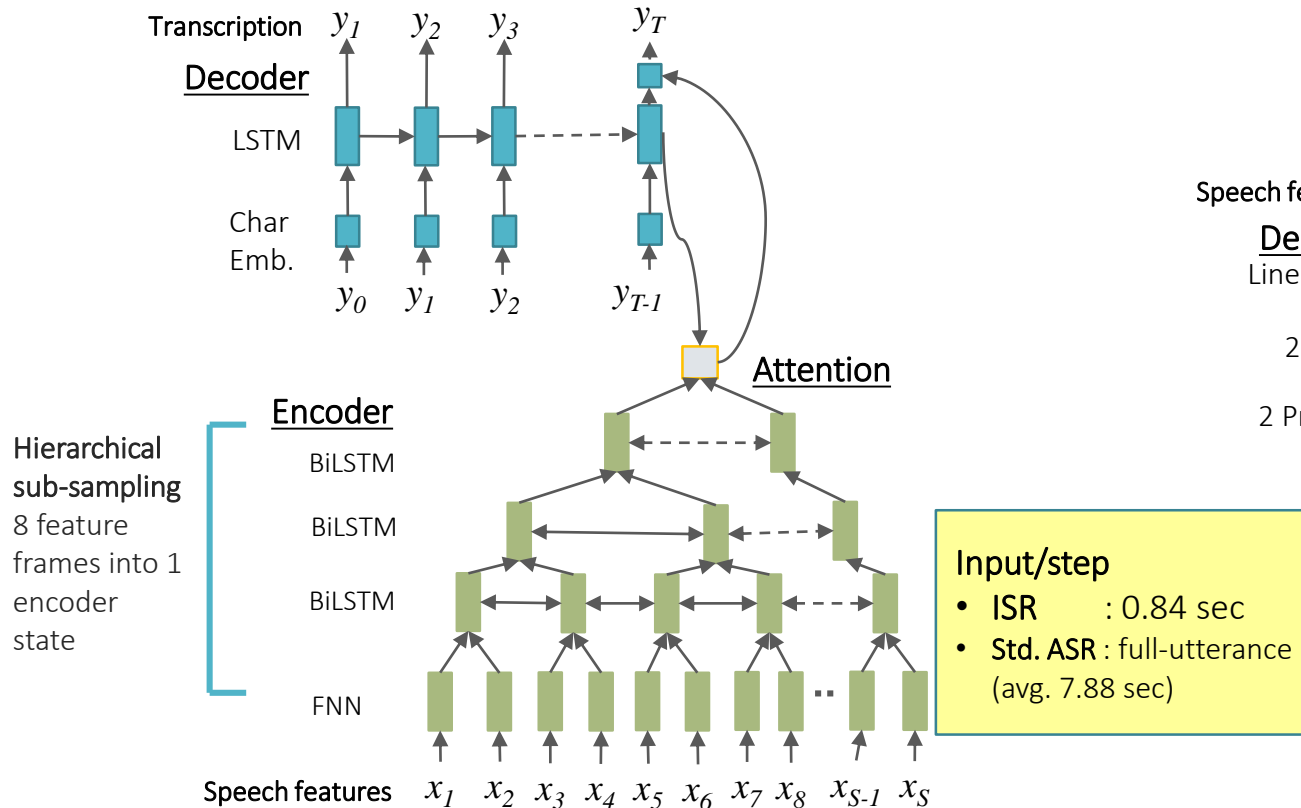
- Language : English
  - ❖ Training sets:
    - o *SI-84*        : 16 hours of speech, 83 speakers
    - o *SI-200*      : 66 hours of speech, 200 speakers
    - o *SI-284*      : *si84 + si200*
  - ❖ Dev. set    : *dev93*
  - ❖ Eval. set    : *eval92*
- Character-level
- Speech features: 80-dims log Mel spectrogram (window: 50 msec, shift: 12.5 msec)

         INTERSPEECH 2020

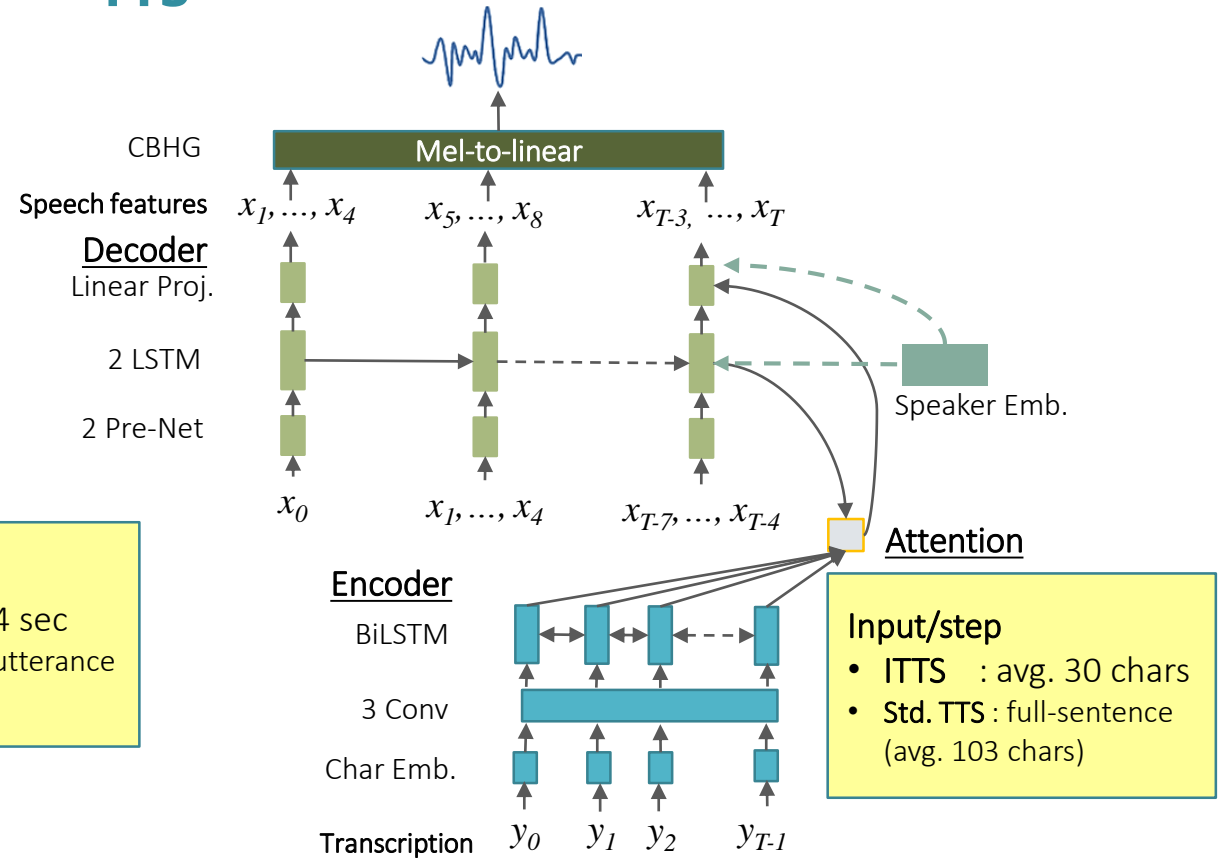# Experiments
## Model Configuration

* Same architecture for standard (non-incremental) and incremental models



**ASR**

**TTS** Tacotron 2 [Wang et al., 2017] structure with speaker embedding [Tjandra et al., 2018]
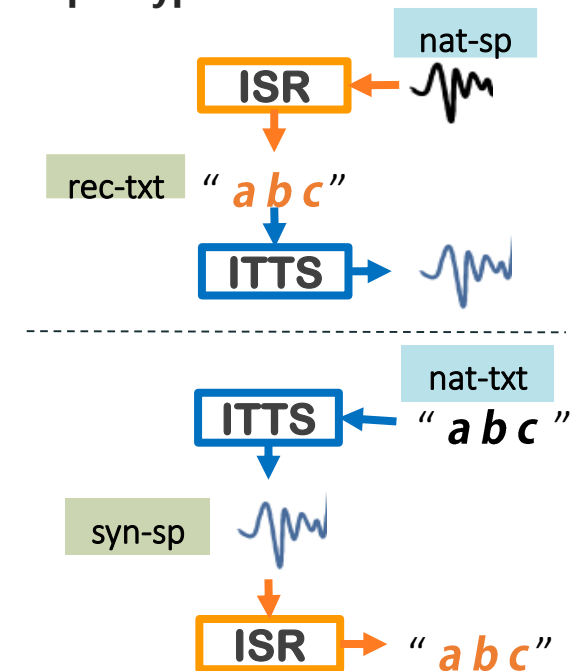
Input/step (ASR)
- ISR : 0.84 sec
- Std. ASR : full-utterance (avg. 7.88 sec)

Input/step (TTS)
- ITTS : avg. 30 chars
- Std. TTS : full-sentence (avg. 103 chars)

INTERSPEECH 2020

# Result

## ASR (CER%) and TTS (log Mel-spectrogram L2 loss) performances

| Data | ASR (CER%) | | | | TTS (L2-norm)$^2$ | | | |
|------|------------|--|--|--|-------------------|--|--|--|
| | Standard (delay: 7.88 sec) | | Incremental (delay: 0.84 sec) | | Standard (delay: 103 chars) | | Incremental (delay: 30 chars) | |
| | *nat-sp* | *syn-sp* | *nat-sp* | *syn-sp* | *nat-txt* | *rec-txt* | *nat-txt* | *rec-txt* |
| **Independent Training** | | | | | | | | |
| Indep-trn *SI-84* | 17.33 | 27.03 | 17.81 | 44.54 | 0.99 | 1.02 | 1.04 | 3.62 |
| Indep-trn *SI-284* | 7.16 | 9.60 | 7.97 | 19.99 | 0.75 | 0.77 | 0.84 | 1.31 |
| **Machine Speech Chain** | | | | | | | | |
| Indep-trn (*SI-84*) + chain-trn-greedy (*SI-200*) | 11.21 | 11.52 | 14.23 | 32.43 | 0.80 | 0.82 | 0.86 | 1.35 |
| Indep-trn (*SI-84*) + chain-trn-teachforce(*SI-200*) | 7.27 | 6.30 | 9.43 | 12.78 | 0.77 | 0.80 | 0.79 | 1.26 |

- **Baseline**
  - ☐ ISR and ITTS *indep-trn SI-84*
- **Topline**
  - ☐ Standard systems *indep-trn SI-284*
- **Proposed**
  - ☐ Incremental machine speech chain
- **Input type:**



- Incremental machine speech chain
  - o Improved ISR and ITTS
  - o Shorter delay with a close performance to the standard system

18

# IV. Conclusion

INTERSPEECH 2020

# Conclusion

**Incremental machine speech chain**

Short-term feedback loop for ISR/ITTS development by mimicking human speech chain

- o Reduced the delay with a close performance to the basic framework
- o Improve ISR and ITTS (natural/synthetic input)
- o Synthetic input processing: demonstration of real-time feedback generation

INTERSPEECH 2020

# Thank you

INTERSPEECH 2020