# 291K
# Deep Learning for Machine Translation
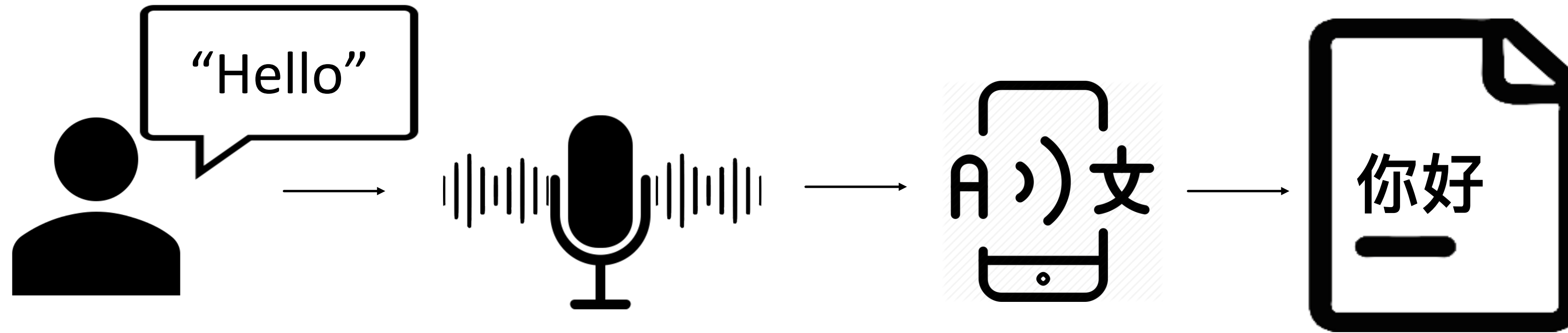# Speech Translation

Lei Li

UCSB

11/17/2021

# Outline

1. Overview: ST Problem and Challenge
2. Basic Model for Speech Translation
3. Break the Challenge of Data Scarcity
4. Better training strategy for ST
5. New ST-powered Products

# Speech-to-Text Translation(ST)

- source language *speech(audio)* ➞ target lang *text*



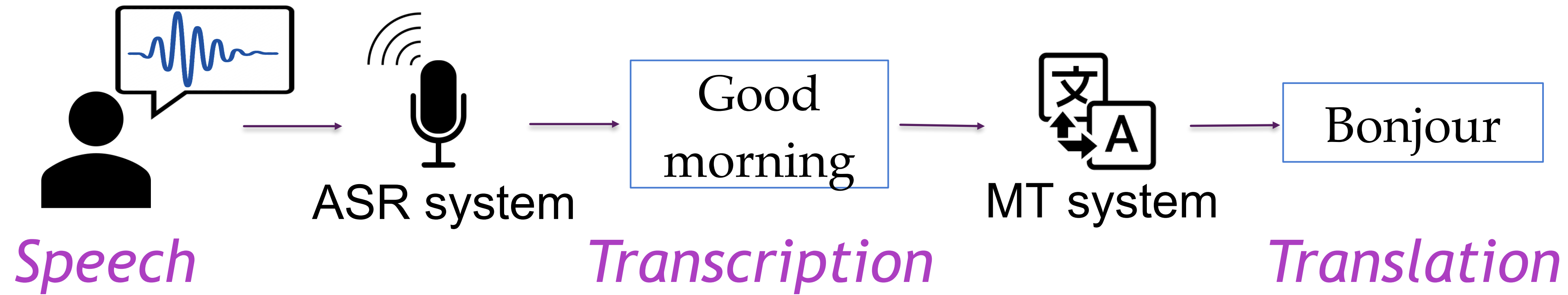| **Application Type** |
| --- |
| • (Non-streaming) ST e.g. video translation<br>• Streaming ST        e.g. realtime conference translation |

| **System** |
| --- |
| • Cascaded ST<br>• End-to-end ST |

# Cascaded ST System

- Challenges:

**1. Computationally inefficient**

**2. Error propagation**:  Wrong transcription ➡️ Wrong translation



*Speech*                    *Transcription*                    *Translation*

*do at this* *and see if it works for you* ➡️ 这样做，看看它是否对你有用
*duet this* *and see if it works for you* ➡️ 二重奏一下，看看它是否对你有用
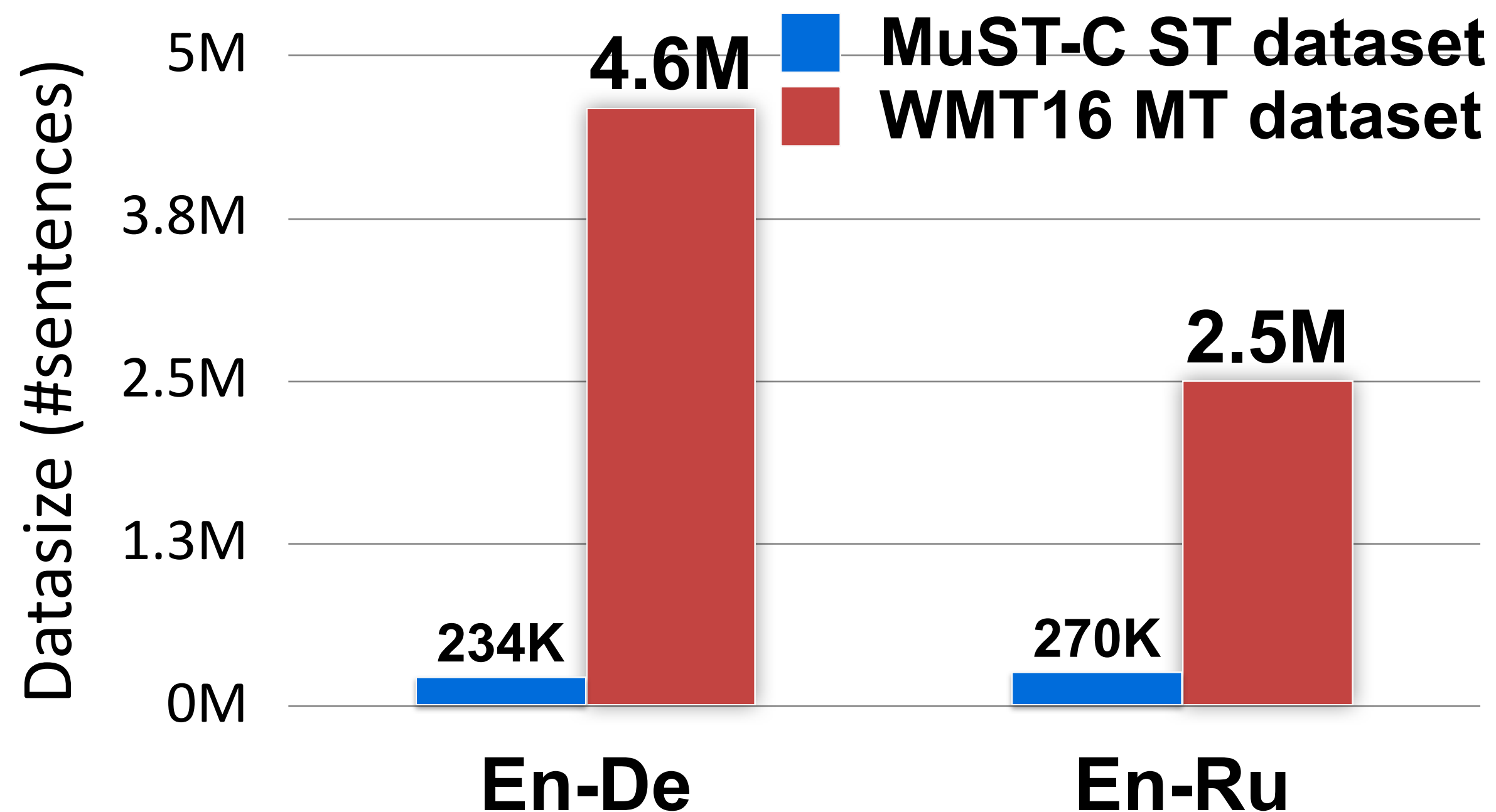
# End-to-end ST Model



- Single model to produce text translation from speech
- Basic model: Encoder-Decoder architecture (e.g. Transformer)
- Advantage:
  - Reduced latency, simpler deployment
  - Avoid error propagation

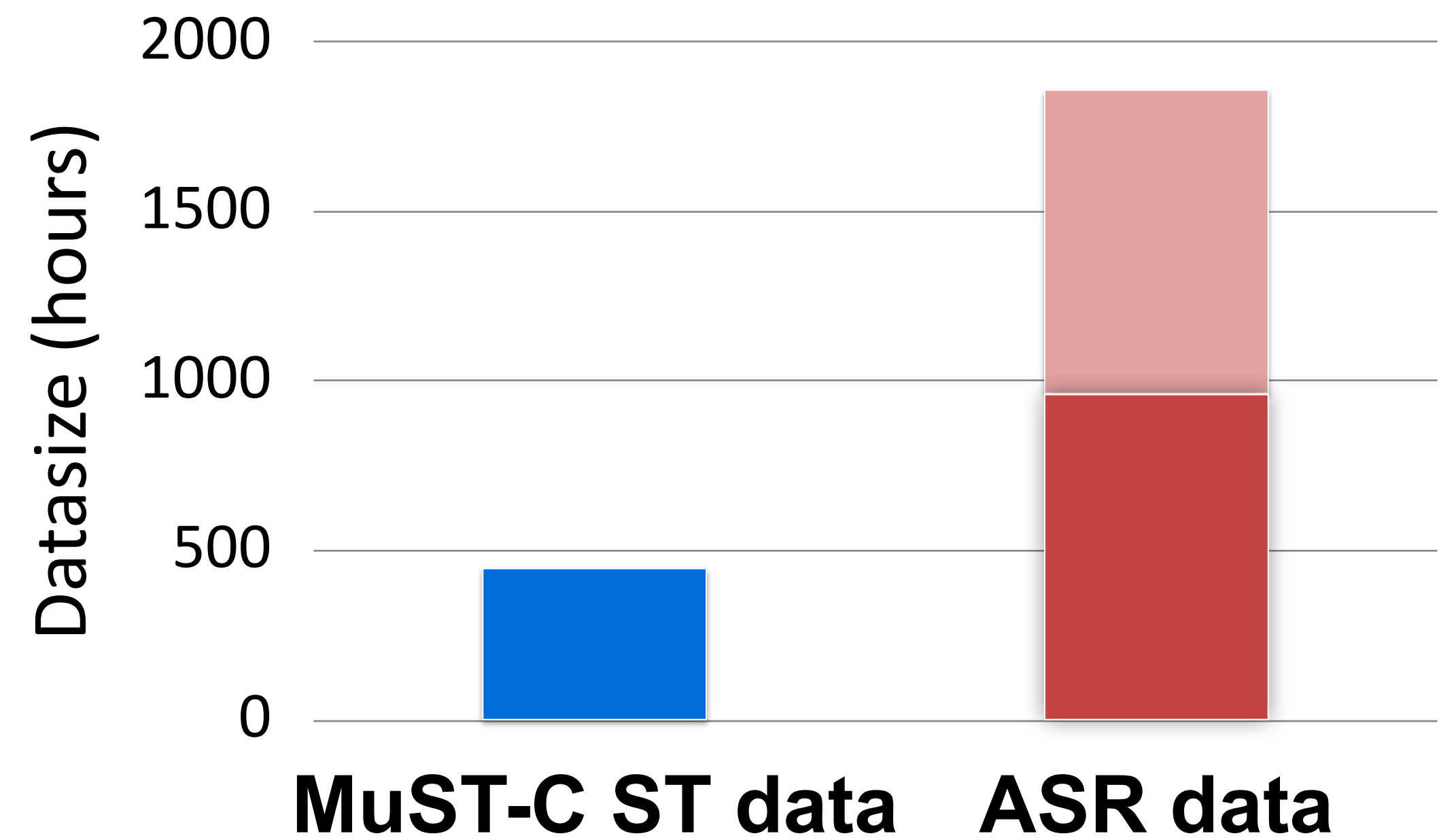[1] Bérard et al., Listen and translate: A proof of concept for end-to-end speech-to-text translation. 2016

# Challenge

- Data scarcity - lack of large parallel audio-translation corpus

- Modality Disparity between speech and text

**Dataset size (Text)**
**ST vs MT**

**Dataset size**
**ST vs ASR**

# Challenge

- Modality Disparity between speech and text
  - Disfluencies
    - Hesitations: "uh", "uhm", "hmm",
    - Discourse markers: "you know", "I mean",…
    - Repetitions: "It had, it had been a good day"
    - Corrections: "no, it cannot, I cannot go there"
  - Unlike (Text) MT, No punctuation
    - Let's eat Grandpa !
    - Let's eat, Grandpa !

# Basic End-to-end ST Model

# Basic ST model



Main differences to text machine translation

Input: Audio signal are continuous and much longer!

# Audio Signal

- Following best-practice from ASR

- Signal Sampling
  - Measure Amplitude of signal at time t
  - Typically 8kHz or 16 kHz

- Windowing — Frame
  - Split signal in different windows, called Frame
    ‣ Length: ~ 20-30 ms (typically 25ms)
    ‣ Stride: ~ 10 ms

# Audio Feature Extraction

- Speech feature extraction:
  - Most common:
    - Mel-Frequency Cepstral Coefficients (MFCC)
    - Log mel-filterbank features (FBANK)
  - Idea:
    - Analyse frequencies of the signal
  - Steps:
    - Discrete Fourier Transformation
    - Mel filter-banks
    - Log scale
    - (Inverse Discrete Fourier Transformation)
  - Size:
    - 20-100 features per frame
- Learned Feature: wav2vec

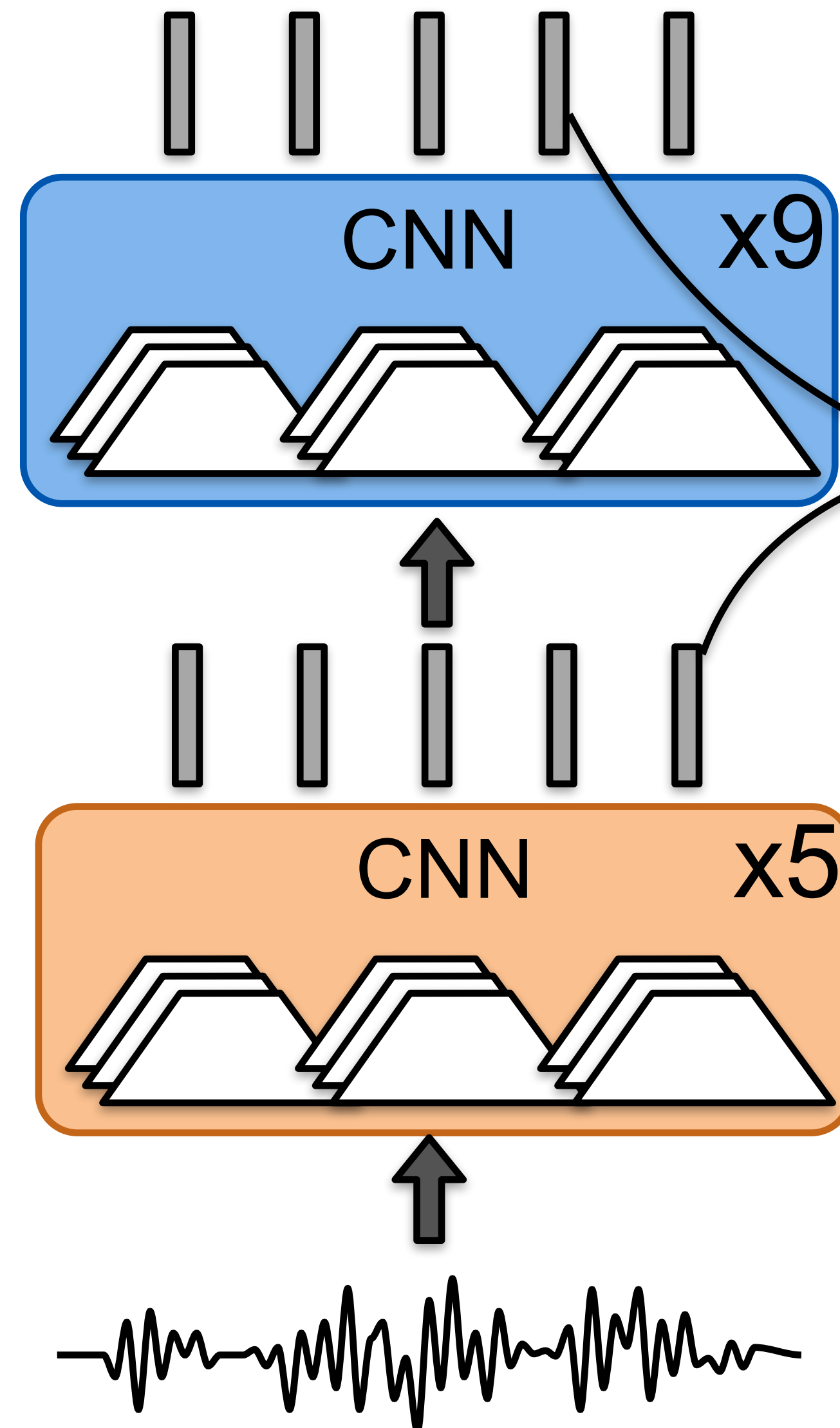# Wav2Vec: Self-supervised Speech Representation Learning

high-level context state c, each frame ~ 210ms, stride10ms

CNN x9

Training data: LibriSpeech 960 hrs audio only

Minimize contrastive loss

$$L = -\sum \left( \log \sigma(z_{t+1} \cdot h_t) + \sum \log \sigma(-z_- \cdot h_t) \right)$$

Low level acoustic state h, each frame ~ 30ms, stride10ms

CNN x5

Bring closer context and acoustic state

Bring further context and negative sampled acoustic state

wav2vec: Unsupervised Pre-training for Speech Recognition [Schneider et al, Interspeech 2010]

13

# Wav2Vec2.0: Contrastive on quantized acoustic state

Training data: (audio only)
LibriSpeech 960 hrs

LibriVox 53k hrs

Masked context during training

Quantized low-level acoustic state, each frame ~ 25ms, stride 20ms

Transformer Encoder x12

CNN x7

Minimize contrastive loss

$$L = - \sum \log \frac{\exp Sim(c_t, q_t)}{\sum \exp Sim(c_t, q_-)} + \text{penalty}$$

Bring closer masked context and quantized acoustic state

Wav2vec2.0: a Framework for Self-Supervised Learning of Speech Representations [Baevski et al, NeurIPS 2020] 14

# Basic Speech Translation Model (Similar to MT)

Transformer-based: N-layer convolution + attention encoder, M-layer decoder

Training data: <audio seq., translation text>

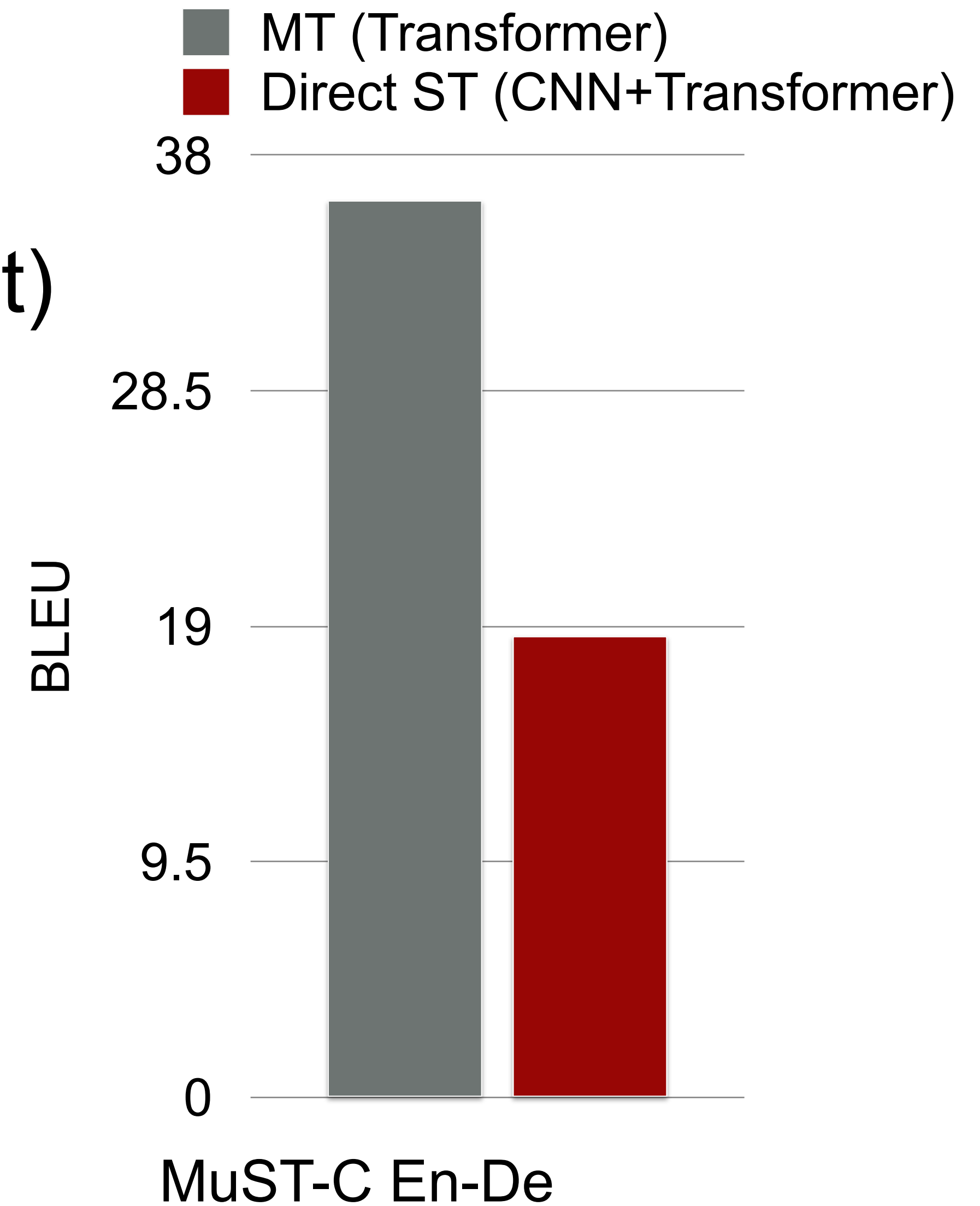# Speech Translation model lags behind MT

- Performance on MuST-C En-De:
  - ST 18.6
  - MT 36.2 (taking correct transcript as input)



MuST-C En-De

# Approaches for Speech Translation

- Utilizing additional parallel text from MT corpus - MT pretraining
  - Decoder initialization from separately trained MT model
  - Single-modal(audio) Encoder-Decoder: COSTT[Dong et al, AAAI 2021b]
- Using Additional ASR data - ASR Pre-training
  - Curriculum Pre-training [Wang et al, ACL 2020]
  - LUT [Dong et al, AAAI 2021a]
- Using additional raw audio data
  - Wav2vec & Wav2Vec2.0 [Schneider et al. Interspeech 2019, Baevski et al NeurIPS2020]
  - Apply to ST [Wang et al, 2021, Zhao et al, ACL 2021, Wang et al, Interspeech 2021]
- Distilling knowledge from Pre-trained Language Model (BERT)
  - LUT [Dong et al, AAAI 2021a]
- Learning Better Speech-text cross-modal representation for ST
  - TCEN-LSTM [Wang et al, AAAI 2020]
  - Chimera [Han et al, ACL 2021a]
  - Wav2vec2.0 + mBart + Self-training [Li et al, ACL 2021b]
  - FAT-ST [Zheng et al, ICML 2021]
- Better Fine-tuning Strategy
  - XSTNet [Ye et al, Interspeech 2021]

# Using external Parallel Text

**Dataset size
ST vs MT**



🤔 How to use <u>MT data</u> *with much larger scale* to improve ST performance?

# Separate Encoder-Decoder Pre-train

Speech Recognition
LibriSpeech corpus

Speech Translation
fine-tune on ST data

Machine Translation
WMT corpus

# Knowledge Distillation from MT model

MT pre-training    KL loss + ST Cross-entropy loss

Comment allez-vous ?

Comment allez-vous ?

**Decoder**

**Decoder**

**Encoder**

**Encoder**

How are you ?

How are you ?

End-to-End Speech Translation with Knowledge Distillation [Liu et al, Interspeech 2019]  20

# Motivation of Better Decoding

**Problem1:** How to give the decoder hints?
**Idea 1**: Introduce a consecutive decoder for trans-trans.

Compressed
Encoder

Consecutive
Decoder

**Problem2:** Long acoustic sequence is challenging for the encoder!
**Idea 2**: Introduce a compressed encoder to relief the model memory.

# Pre-train ST's decoder with full MT

How to make a single model's decoder to perform text translation?

Decoder   ==>  translation

Encoder -> Decoder  ==> transcribe and translation

**Trans**cription  **–**  **Trans**lation

(apples)

apples    pommes

| a | p | p | l | e | s |  | p | o | m | m | e | s |

Compressed
Encoder

Consecutive Decoder

Consecutive Decoding for Speech-to-text Translation [Q. Dong, M. Wang, H. Zhou, S. Xu, B. Xu, Lei Li, AAAI 2021]

# COSTT for ST

Step1: Pre-train using external MT corpus

Semantic represent:

CTC loss

Acoustic represent:

*Shrinking*

Transcript : *"Good morning"*   Translation: *"Bonjour"*

*Cross-Entropy loss*

Input : *Log-mel fbank feature of audio*

Acoustic-Semantic Encoder

Transcription-Translation Decoder

Step 2: Train encoder w/ shrinking module using CTC

Step 3: Train full model on ST data <audio, transcript, translation>

Consecutive Decoding for Speech-to-text Translation [Q. Dong, M. Wang, H. Zhou, S. Xu, B. Xu, Lei Li, AAAI 2021]

# Advantages of COSTT

- Unified training with both transcript and translation text

- Reduced encoder output size with CTC-guided shrinking

- Able to pre-train the decoder with external MT parallel data

Semantic
~10

Phoneme spikes

Acoustic
~1000

Consecutive Decoding for Speech-to-text Translation [Q. Dong, M. Wang, H. Zhou, S. Xu, B. Xu, Lei Li, AAAI 2021]

# Using external ASR data

**Dataset size**
**ST vs ASR**



How to use larger external <u>ASR data</u> to improve ST performance? 🤔

# Curriculum Pre-training with ASR data

① ASR Cross entropy + ASR CTC loss

I like to eat apple

Transformer Decoder

② Masked LM KL loss + Bilingual lexicon KL loss

eat

Ich esse gerne Apfel

I like to eat apple

Transformer Encoder

Transformer Encoder

| 0 | 1 | 2 | 3 |

+ + + +

2D Convolution

③ Translation cross entropy

Ich esse gerne Apfel

Transformer Decoder

**S** "I like to eat apple"

**S** "I like to eat apple"

Curriculum Pre-training for End-to-end Speech Translation [Wang et al, ACL 2020]

26

# ASR Pre-training helps ST

IWSLT & Librispeech

# Raw Text Pre-training



**Dataset size ST vs Raw text**

Datasize (million words)

3500
2800
2100
1400
700
0

3300M

400x

8.3M

English Wiki

BookCorpus

**MuST-C ST data**  **Raw text**

Using pre-trained LM in decoding weighting is easy!

**But**

🤔 How to use pre-trained BERT to improve ST performance?

# Drawbacks of the Encoder-Decoder Structure



**1.** A single encoder is hard to capture the representation of audio for the translation.
**2.** Limited in utilizing the information of "*transcription*" in the training.

# Motivation: Mimic human's behavior

**Question**: How human translate?



Listen — apples → Understand → Translate — pommes

"Listen-Understand-Translate"(LUT) model based motivated by human's behavior

# Motivation of Better Encoding

**Drawback 1:** A single encoder is not enough.

**Idea 1**: Introduce a semantic encoder

| Acoustic Encoder | → | Semantic Encoder (Understand) | → | Decoder (Translate) |

*supervise*

*"transcript"*

*supervise*

BERT of *"transcript"*

**Drawback 2:** Limit in using "transcript" info.

**Idea 2**: Utilizing the pre-trained representation (e.g. BERT) of the "transcript" to learn the semantic feature.

Listen, Understand and Translate [Q. Dong, R. Ye, M. Wang, H. Zhou, S. Xu, B. Xu, Lei Li, AAAI 2021]

# LUT: Utilizing Pre-trained Model on Raw Text

Training data: triples of

<speech, transcript_text, translate_text>

Transcript ($z$):
*"Good morning"*

BERT representation

Translation($y$):
*"Bonjour"*

*CTC loss*

*Distance loss*

*CE loss*

Input ($x$):
*Log-mel fbank feature*

Acoustic Encoder (Listen) → Semantic Encoder (Understand) → Translation Decoder (Translate)

Listen, Understand and Translate [Q. Dong, R. Ye, M. Wang, H. Zhou, S. Xu, B. Xu, Lei Li, AAAI 2021]

# ST Benefits from BERT, with Raw Text Pre-training



IWSLT & Librispeech

Legend:
- Transformer ST
- Transformer+ASR
- Transformer+Curriculum
- COSTT
- LUT

En-De: 12.5, 13.1, 18.2, 18.6, 18.6
En-Fr: 13.2, 16.9, 18, 18.2, 18.3

BLEU

# Audio Pre-training

## Dataset size
## ST vs raw Audio



🤔 How to use larger  raw audio data to improve ST performance?

# Speech Translation with Audio-Pretrain

## Wav2vec Pretrain + Fine-tune on ST



Comment allez-vous ?

Decoder

Encoder

Wav2vec 2.0

How are you ?

MuST-C ST results

LSTM [1]
Wav2vec-LSTM [1]
Transformer [2]
Wav2vec2.0-Transformer [3]

BLEU

En-De: 22.8, 23.6
En-Fr: 27.8, 29.8, 33.3, 34.6
En-Ru: 15.1, 17
En-Ro: 17.1, 18.2, 22.2, 22.4

[1] Self-supervised Representations improve end-to-end speech translation [Wu et al. InterSpeech 2020]
[2] NeurST toolkit [Zhao et al ACL2021 demo]    [3] End-to-end Speech Translation [Ye et al. InterSpeech 2021]

# Self-training with Audio data

Comment allez-vous ?

**Transformer Decoder**

Wav2vec 2.0
Transformer
CNN

How are you ?

Step 0. Audio-only pre-training for Wav2vec2.0

Step 1. Freeze Wav2vec2.0, train on ST

Step 2. Self-train on generated pseudo-translation with LibriVox audio

## CoVoST2 Results



Legend:
- Transformer [1]
- Transformer w/ ASR pre-train [1]
- Wav2vec2.0-Transformer [2]
- Wav2vec2.0-Transformer + Self-train [2]

BLEU axis: 0, 10, 20, 30, 40

En-De: 13.6, 16.3, 23.8, 26.5
En-Ca: 20.2, 21.8, 32.4, 34.1
En-Ar: 8.7, 12.1, 17.4, 20.2
En-Tr: 8.9, 10, 15.4, 17.5

[1] CoVoST 2 and Massively Multilingual Speech-to-Text Translation, [Wang et al InterSpeech 2021]
[2] Large-Scale Self- and Semi-Supervised Learning for Speech Translation [Wang et al. 2021]

# Fine-tuning Strategy for ST

# Cross Speech-Text Network (XSTNet)

# Supports to train MT data

☑ Transformer MT model

☑ We can add **more external MT data** to train Transformer encoder & decoder

# Supports inputs of two modalities

☑ Wav2vec2.0[1] as the acoustic encoder

☑ We add two convolution layers with 2-stride to shrink the length.



[1] wav2vec 2.0: A framework for self-supervised learning of speech representations, 2020

# Language indicator strategy

- We use language indicators to distinguish different tasks.

| Tasks | Source input | Target output |
|-------|--------------|---------------|
| MT | **<en>** This is a book. | **<fr>** c'est un livre. |
| ASR | **<audio>** 〜〜〜 | **<en>** This is a book. |
| ST | **<audio>** 〜〜〜 | **<fr>** c'est un livre. |

End-to-end Speech Translation via Cross-modal Progressive Training [Rong Ye, Mingxuan Wang, Lei Li, Interspeech 2021]

# Progressive Multi-task Training

**#** **L**arge-scale MT pre-training

Using **external MT** $D_{MT-ext}$

↓

**#** **M**ulti-task Finetune

Using **(1) external MT** $D_{MT-ext}$

(2) $D_{ST}$ with *<speech, translation>*

(3) $D_{ASR}$ with *<speech, transcript>*

**Progressive:**

*Don't stop*

*training $D_{MT-ext}$*

End-to-end Speech Translation via Cross-modal Progressive Training [Rong Ye, Mingxuan Wang, Lei Li, Interspeech 2021]

# XSTNet achieves State-of-the-art Performance

| Models | External Data | Pre-train Tasks | De | Es | Fr | It | Nl | Pt | Ro | Ru | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Transformer ST [13] | × | ASR | 22.8 | 27.4 | 33.3 | 22.9 | 27.2 | 28.7 | 22.2 | 15.1 | 24.9 |
| AFS [31] | × | × | 22.4 | 26.9 | 31.6 | 23.0 | 24.9 | 26.3 | 21.0 | 14.7 | 23.9 |
| Dual-Decoder Transf. [15] | × | × | 23.6 | 28.1 | 33.5 | 24.2 | 27.6 | 30.0 | 22.9 | 15.2 | 25.6 |
| Tang et al. [2] | MT | ASR, MT | 23.9 | 28.6 | 33.1 | - | - | - | - | - | - |
| FAT-ST (Big) [6] | ASR, MT, mono-data$^{\dagger}$ | FAT-MLM | 25.5 | 30.8 | - | - | 30.1 | - | - | - | - |
| W-Transf. | audio-only* | SSL* | 23.6 | 28.4 | 34.6 | 24.0 | 29.0 | 29.6 | 22.4 | 14.4 | 25.7 |
| **XSTNet (Base)** | audio-only* | SSL* | 25.5 | 29.6 | 36.0 | 25.5 | 30.0 | 31.3 | 25.1 | 16.9 | 27.5 |
| **XSTNet (Expand)** | MT, audio-only* | SSL*, MT | **27.8$^{\S}$** | **30.8** | **38.0** | **26.4** | **31.2** | **32.4** | **25.7** | **18.5** | **28.8** |

Table 1: *Performance (case-sensitive detokenized BLEU) on MuST-C test sets.* $\dagger$: *"Mono-data" means audio-only data from Librispeech, Libri-Light, and text-only data from Europarl/Wiki Text;* *: *"Audio-only" data from LibriSpeech is used in the pre-training of wav2vec2.0-base module, and "SSL" means the self-supervised learning from unlabeled audio data.* $\S$ *uses OpenSubtitles as external MT data.*

**XSTNet-Base**: Achieves the SOTA in the restricted setup

**XSTNet-Expand**: Goes better by using extra MT data

# XSTNet better than cascaded ST! a gain of 2.6 BLEU



What is "Cascaded-Strong" system?

Strong ASR model + Large-scale MT data

| Cascaded - Strong | Model | Training data | Performance (En-De) |
|---|---|---|---|
| ASR | W2V2+ Transformer | MuST-C $D_{ASR}$ | WER=13.0 |
| MT | Transformer-base | WMT + MuST-C $D_{MT}$ | BLEU=31.7 |

# Learning Better Speech-Text Bimodal Representation

- Chimera: Learning Fixed-size Shared Space for both audio and text, audio+MT pretraining [Han et al. 2021]

- Wav2vec2.0-mTransformer LNA: Use both audio pertaining + multilingual pertained language model, and selective efficient fine-tuning [Li et al. ACL 2021]

- FAT-ST: Masked pre-training for fused audio and text [Zheng et al. ICML 2021]

# Bi-modal Encoding Architecture for ST

Text Input

Word Embedding

Translation text

Common Encoder

Decoder

Bonjour

Speech Encoder

Audio input

Challenges: gap between text and audio
1. Length: ~20 (text) vs. ~ 1k-10k (audio)
2. Embedding space disparity

# Insights from Cognitive Neuroscience

Speech and text interfere with each other in brain[1]



**activation map**        **processing paths**
**Convergence sites** of *speech* (blue) and *text* (yellow)

[1] Van Atteveldt, Nienke, et al. "Integration of letters and speech sounds in the human brain." *Neuron* 43.2 (2004): 271-282.

[2] Spitsyna, Galina, et al. "Converging language streams in the human temporal lobe." *Journal of Neuroscience* 26.28 (2006): 7328-7336.

# Idea: Bridging the Speech-Text modality gap with Shared Semantic Representation

## ST triple data:

<speech, transcript_text, translate_text>

# Chimera Model for ST

Training with auxiliary objectives: ST + MT + Contrastive loss

Benefit: able to exploit large external MT data

# Chimera achieves the best (so far) BLEU on all languages in MuST-C

| Model | External Data | | | MuST-C EN-X | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Speech | ASR | MT | EN-DE | EN-FR | EN-RU | EN-ES | EN-IT | EN-RO | EN-PT | EN-NL |
| FairSeq ST [†] | × | × | × | 22.7 | 32.9 | 15.3 | 27.2 | 22.7 | 21.9 | 28.1 | 27.3 |
| Espnet ST [‡] | × | × | × | 22.9 | 32.8 | 15.8 | 28.0 | 23.8 | 21.9 | 28.0 | 27.4 |
| AFS [*] | × | × | × | 22.4 | 31.6 | 14.7 | 26.9 | 23.0 | 21.0 | 26.3 | 24.9 |
| Dual-Decoder [◇] | × | × | × | 23.6 | 33.5 | 15.2 | 28.1 | 24.2 | 22.9 | **30.0** | 27.6 |
| STATST [♯] | × | × | × | 23.1 | - | - | - | - | - | - | - |
| MAML [♭] | × | × | ✓ | 22.1 | 34.1 | - | - | - | - | - | - |
| Self-Training [○] | ✓ | ✓ | × | 25.2 | 34.5 | - | - | - | - | - | - |
| W2V2-Transformer [*] | ✓ | × | × | 22.3 | 34.3 | 15.8 | 28.7 | 24.2 | 22.4 | 29.3 | 28.2 |
| Chimera Mem-16 | ✓ | × | ✓ | 25.6 | 35.0 | 16.7 | 30.2 | 24.0 | 23.2 | 29.7 | 28.5 |
| Chimera | ✓ | × | ✓ | **27.1 •** | **35.6** | **17.4** | **30.6** | **25.0** | **24.0** | **30.2** | **29.2** |

# Audio and Multilingual Text Pretrain for Multilingual ST

Comment allez-vous ?

**Transformer Decoder**

**CNN**

Wav2vec 2.0

Transformer

CNN

How are you ?

- Encoder uses Wav2vec2.0 pre-trained on LibriVox-60k audio
- Decoder: mBart pre-trained on 50 monolingual text and 49 bitext
- ST finetune strategy (LNA):
  - Only fine-tune layer-norm and attention layers
- MT+ST multitask joint train with further parallel bitext data

Multilingual Speech Translation with Efficient Finetuning of Pretrained Models [Li et al, ACL 2021]

# Wav2vec2.0 retraining + Multilingual training effectively transfers to low resource source language

## CoVoST2 Results



Legend:
- Transformer
- m-Transformer
- Wav2vec2.0-mTransformer LNA

| | Fr-En | De-En | Es-En | Ca-En | It-En | Ru-En | Pt-En |
|---|---|---|---|---|---|---|---|
| Transformer | 24.3 | 8.4 | 12 | 14.4 | 0.2 | 1.2 | 0.5 |
| m-Transformer | 26.5 | 17.5 | 27 | 23.1 | 18.5 | 4.7 | 6.3 |
| Wav2vec2.0-mTransformer LNA | 35 | 28.2 | 35.2 | 31.1 | 27.6 | 22.8 | 24.1 |

## CoVoST2 Results



Legend:
- Transformer
- m-Transformer
- Wav2vec2.0-mTransformer joint train

| | En-De | En-Ca | En-Ar | En-Tr | En-Zh |
|---|---|---|---|---|---|
| Transformer | 13.6 | 20.2 | 8.7 | 8.9 | 20.6 |
| m-Transformer | 17.3 | 22.3 | 13 | 10.7 | 28.2 |
| Wav2vec2.0-mTransformer joint train | 25.8 | 30.9 | 18 | 17 | 33.3 |

Multilingual Speech Translation with Efficient Finetuning of Pretrained Models [Li et al, ACL 2021]

# Fused Acoustic and Text Masked Language Model (FAT-MLM)



L2 loss

2D Deconvolution

Cross-entropy

Good

Cross-entropy

Tag

Transformer Encoder

| En | En | En | En | En | En | En | En | De | De | De | De |
| + | + | + | + | + | + | + | + | + | + | + | + |

Acoustic embedding

| 0 | 1 | 2 | 3 | 0 | 1 | 2 | 3 |
| + | + | + | + | + | + | + | + |

Transformer Encoder

| <s> | [Mask] | Morning | </s> | <s> | Guten | [Mask] | </s> |

x                    y

| 0 | 1 | 2 | 3 |
| + | + | + | + |

2D Convolution

s   Mask   Mask

Pre-training data

1. Librispeech ASR 960h
2. Libri-light audio 3,748h
3. Europarl/wiki text 2.3M
4. MuST-C 408h
5. Europarl MT 1.9M

Fused Acoustic and Text Encoding for Multimodal Bilingual
Pretraining and Speech Translation, [Zheng et al ICML 2021]

53

$l_{ST}(s, y)$ ← | `<s>` | Guten | Tag | `</s>` | → $l_{MT}(x, y)$

Transformer Decoder

Transformer Encoder

Acoustic embedding

Transformer Encoder

| 0 | 1 | 2 | 3 |

+ + + +

2D Convolution

s

| 0 | 1 | 2 | 3 |

+ + + +

| `<s>` | Good | Morning | `</s>` |

x

Training:

- Pre-train FAT-MLM with all data
- Init FAT-ST with FAT-MLM, decoder copy encoder
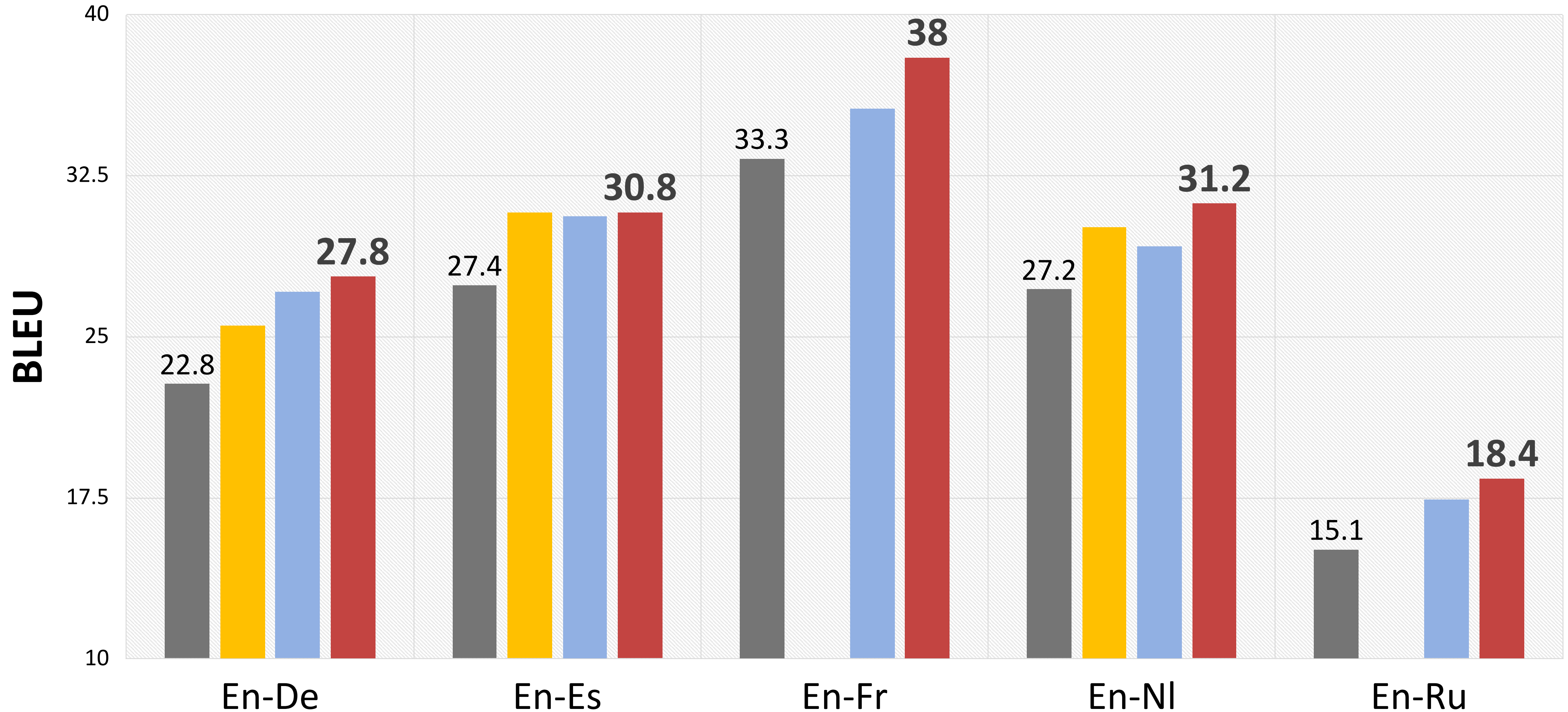- Further fine-tune on MuST-C ST data.

Fused Acoustic and Text Encoding for Multimodal Bilingual Pretraining and Speech Translation, [Zheng et al ICML 2021]

54

# Joint audio&text Pre-training task helps ST

| Pretrain Method | Models | En→De | En→Es | En→Nl | Avg. | Model Size |
|---|---|---|---|---|---|---|
| No Pretraining | ST | 19.64 | 23.68 | 23.01 | 22.11 | 31.25M |
| | ST + ASR | 21.70 | 26.83 | 25.44 | 24.66 (+2.55) | 44.82M |
| | ST + ASR & MT | 21.58 | 26.37 | 26.17 | 24.71 (+2.60) | 56.81M |
| | ST + MAM | 20.78 | 25.34 | 24.46 | 23.53 (+1.42) | 33.15M |
| | ST + MAM + ASR | 22.41 | 26.89 | 26.49 | 25.26 (+3.15) | 46.72M |
| | Liu et al. (2020b) | 22.55 | - | - | - | - |
| | Le et al. (2020) | 23.63 | 28.12 | 27.55 | 26.43 (+4.32) | 51.20M |
| | Cascade[§] | 23.65 | 28.68 | 27.91 | 26.75 (+4.64) | 83.79M |
| | FAT-ST (base). | 22.70 | 27.86 | 27.03 | 25.86 (+3.75) | 39.34M |
| ASR & MT | ST | 21.95 | 26.83 | 26.03 | 24.94 (+2.83) | 31.25M |
| | ST + ASR & MT | 22.05 | 26.95 | 26.15 | 25.05 (+2.94) | 56.81M |
| MAM | FAT-ST (base) | 22.29 | 27.21 | 26.26 | 25.25 (+3.14) | 39.34M |
| FAT-MLM | FAT-ST (base) | **23.68** | 28.61 | **27.84** | 26.71 (+4.60) | 39.34M |
| | FAT-ST (big) | 23.64 | **29.00** | 27.64 | **26.76** (+4.65) | 58.25M |

Fused Acoustic and Text Encoding for Multimodal Bilingual Pretraining and Speech Translation, [Zheng et al ICML 2021]
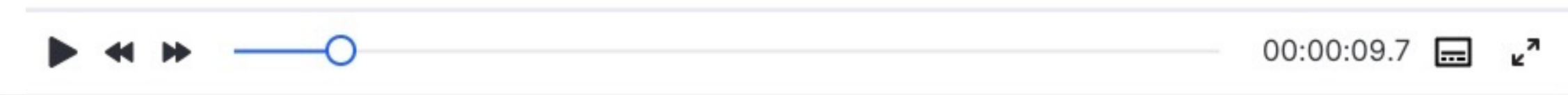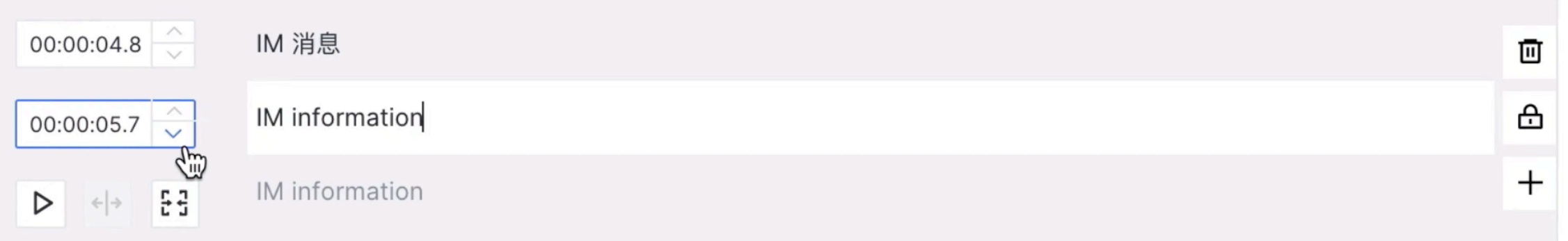
# Pre-training Improves ST Performance

- MuST-C Results

# Summary

| | Direct Supervision | Contrastive | Masked LM | Knowledge distillation | Progressive train | Selective Fine-tune | Self-training |
|---|---|---|---|---|---|---|---|
| **MT Parallel Text** | COSTT | | | [Liu et al. 2019] | XSTNet | | |
| **ASR Speech-Transcript** | LUT | | | | | | |
| **Audio-only** | | Wav2vec Wav2vec 2.0 | | | | | [Wang et al. 2021] |
| **Raw text** | | | | LUT | | | |
| **Speech+Text** | | Chimera | FAT-ST | | XSTNet | LNA | |

# Speech Translation Product Demo

# VolcTransStudio: Video Translation Platform



实时翻译，自动提示 & 交互式修改

Correct-and-Memorize: Learning to translation from interactive revisions [Rongxiang Weng, Hao Zhou, Shujian Huang, Yifan Xia, Lei Li, Jiajun Chen. IJCAI 19]

# Summary for Speech Translation Pre-training

- Parallel speech translation data is scarce
- Pre-training to utilize external large data
  - MT data (Parallel text)
  - ASR data (Speech-transcript)
  - Raw text (Monolingual and Multilingual)
  - Audio-only
- Network architecture to solve modality disparity
  - CNN-Transformer
  - Fixed-size shared memory module
  - Bimodal input with length shrinking for audio
- Techniques to better pre-train and better fine-tune
  - Contrastive prediction
  - Masked LM
  - Quantization of audio representation
  - Knowledge distillation
  - Progressive pre-training

# Language Presentation

# Reference

- Speech Translation
  - wav2vec: Unsupervised Pre-training for Speech Recognition
  - wav2vec 2.0: A framework for self-supervised learning of speech representations
  - Investigating self-supervised pre-training for end-to-end speech translation
  - Self-supervised representations improve end-to-end speech translation (wav2vec + LSTM seq2seq)
  - Large-Scale Self-and Semi-Supervised Learning for Speech Translation
  - Consecutive Decoding for Speech-to-text Translation
  - "Listen, Understand and Translate": Triple Supervision Decouples End-to-end Speech-to-text Translation
  - Learning Shared Semantic Space for Speech-to-Text Translation [ACL 21]
  - Multilingual Speech Translation with Efficient Finetuning of Pretrained Models [ACL 21]
  - Fused Acoustic and Text Encoding for Multimodal Bilingual Pretraining and Speech Translation [ICML 21]
  - End-to-end Speech Translation via Cross-modal Progressive Training [Interspeech 21]
  - Curriculum Pre-training for End-to-end Speech Translation [ACL 20]
  - End-to-End Speech Translation with Knowledge Distillation [Interspeech 19]