

# CMU's IWSLT 2023 Simultaneous Speech Translation System

Brian Yan, Jiatong Shi, Soumi Maiti, William Chen, Xinjian Li, Yifan Peng, Siddhant Arora and Shinji Watanabe

## System Overview

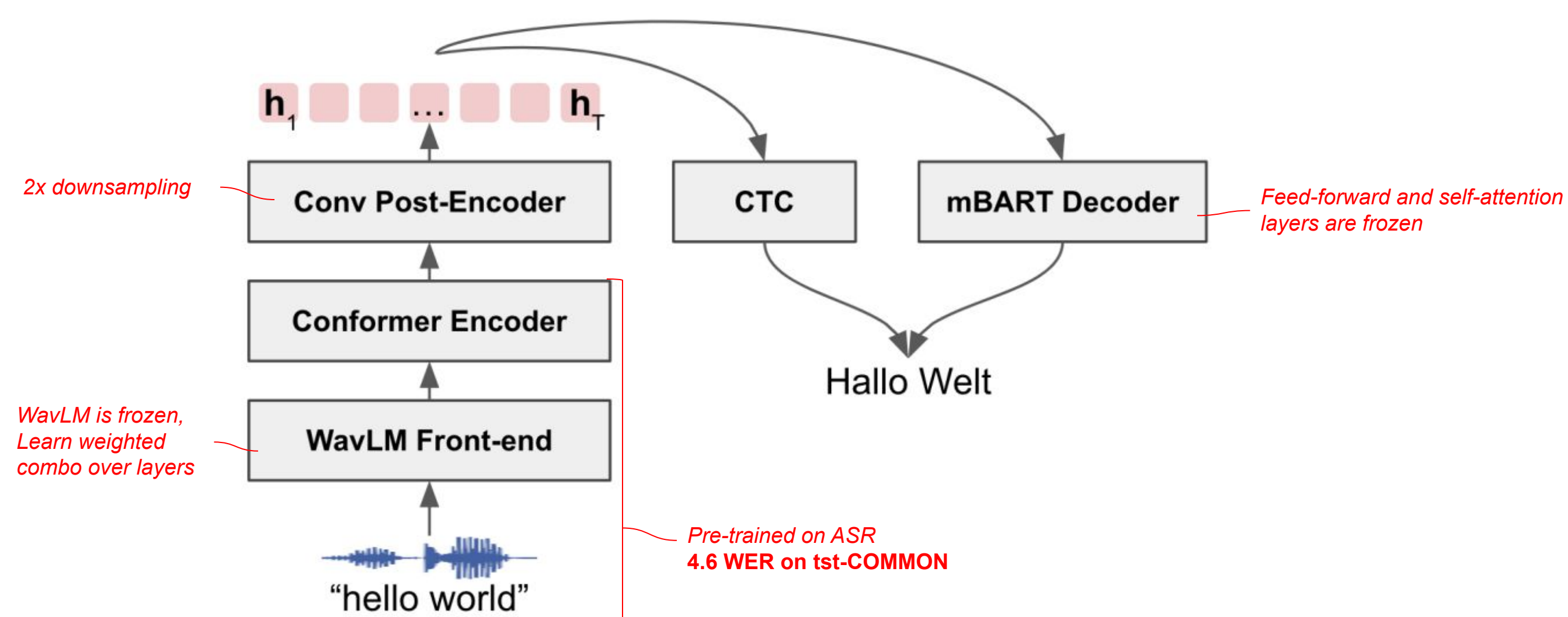
1. **Offline ST using joint CTC/attention** with self-supervised speech/text representations
2. **Offline-to-online adaptation** via chunk-based encoding and incremental beam search
3. **Cascaded simultaneous S2ST** by feeding incremental text outputs to a TTS model

### Data:

- MuST-C + Tedlium (w/ MT pseudo-labels) for ST
- CommonVoice subset for TTS

## Offline ST: Joint CTC/Attention + SSL/LLM

- Large scale model (~900M trainable params) with WavLM and mBART initializations



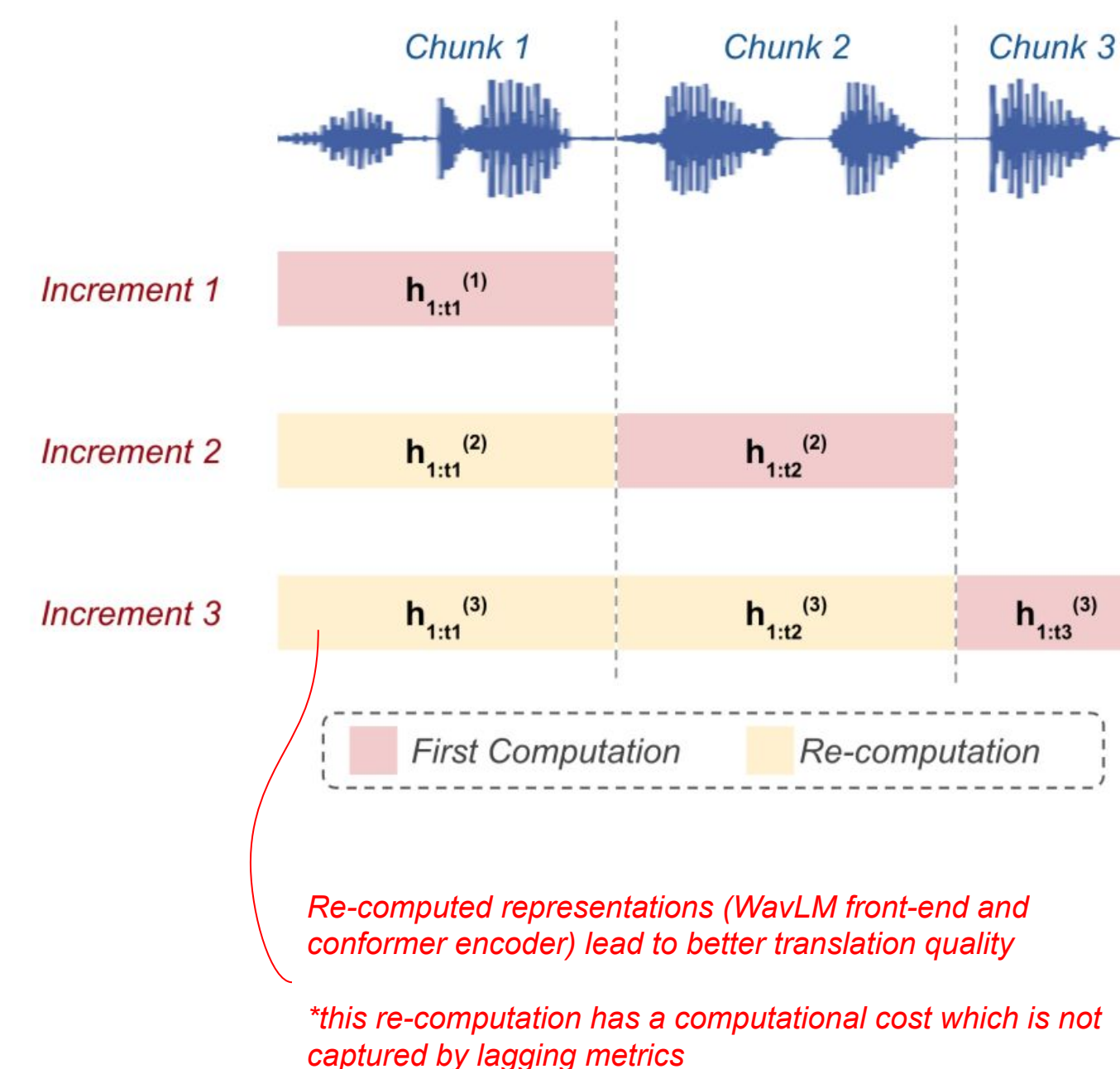
## Offline-to-Online: Incremental Encoding/Decoding

- Directly use offline ST model for online inference via incremental strategies

### Incremental Encoding

Re-compute encoder representations for each incremental chunk of speech

- 2s for speech-to-text
- 2.5s for speech-to-speech



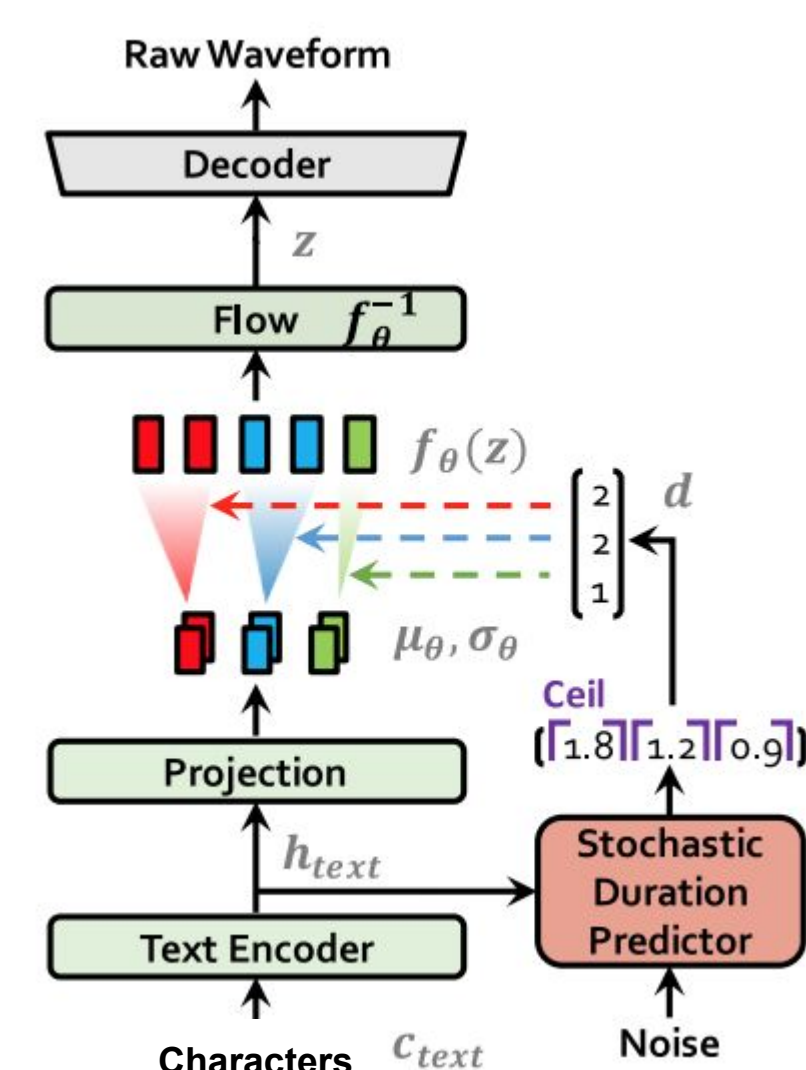
### Incremental Decoding

**Algorithm 1** Beam search step with rewinding of unreliable hypotheses on non-final chunks and incremental pruning upon end-detection.

```

1: procedure BEAMSTEP(hyps, prevHyps, isFinal)
2:   newHyps = {}; endDetected = False
3:   for  $y_{1:l-1} \in \text{prtHs}$  do
4:     attnCnds = top-k( $P_{\text{Attn}}(y_l|X, y_{1:l-1})$ ,  $k = p$ )
5:     for  $c \in \text{attnCnds}$  do
6:        $y_{1:l} = y_{1:l-1} \oplus c$ 
7:        $\alpha_{\text{CTC}} = \text{CTCScore}(y_{1:l}, X_{1:T})$ 
8:        $\alpha_{\text{Attn}} = \text{AttnScore}(y_{1:l}, X_{1:T})$ 
9:        $\beta = \text{LengthPen}(y_{1:l})$  Joint CTC/attn scoring
10:       $P_{\text{Beam}}(y_{1:l}|X) = \alpha_{\text{CTC}} + \alpha_{\text{Attn}} + \beta$ 
11:      newHyps[ $y_{1:l}$ ] =  $P_{\text{Beam}}(\cdot)$ 
12:      if (isFinal) and (c is <eos> or repeat) then
13:        endDetected = True
14:        newHyps = prevHyps  $\triangleright$  rewind
15:      else if l is maxL then
16:        endDetected = True End-of-chunk detection
17:      end if
18:    end for
19:  end for Pruning
20:  if endDetected then  $\triangleright$  incremental pruning
21:    newHyps = top-k( $P_{\text{Beam}}(\cdot)$ ,  $k = 1$ )
22:  else  $\triangleright$  standard pruning
23:    newHyps = top-k( $P_{\text{Beam}}(\cdot)$ ,  $k = b$ )
24:  end if
25:  return newHyps, endDetected
26: end procedure
    
```

## Cascaded Simultaneous S2ST



- Single-speaker VITS TTS model with character input
- CommonVoice data selection:
  - Evaluated the speech quality (DNSMOS) of top 5 most speakers by # of utterances
  - Set a threshold of 4.0 for selecting utterances
  - Choose the single speaker with most hours (12h)
- Cascaded inference on incremental text decodings (which correspond to ~2.5s of input speech)

Image Source: Conditional Variational Autoencoder with Adversarial Learning for End-to-End Text-to-Speech

## Results

- 6% quality degradation from ST to SST (~60% lagging reduction)
- 12% quality degradation from SST to SS2ST (includes ASR errors from ASR-BLEU)

MODEL	QUALITY	LATENCY	
OFFLINE SPEECH TRANSLATION (ST)	BLEU $\uparrow$	-	
Multi-Decoder CTC/Attn (Yan et al., 2023b)	30.1	-	-
WavLM-mBART CTC/Attn (Ours)	32.5	-	-
SIMUL SPEECH TRANSLATION (SST)	BLEU $\uparrow$	AL $\downarrow$	LAAL $\downarrow$
Time-Sync Blockwise CTC/Attn (Yan et al., 2023b)	26.6	1.93	1.98
WavLM-mBART CTC/Attn (Ours)	30.4	1.92	1.99
SIMUL SPEECH-TO-SPEECH TRANSLATION (SS2T)	ASR-BLEU $\uparrow$	SO $\downarrow$	EO $\downarrow$
WavLM-mBART CTC/Attn + VITS (Ours)	26.7	2.33	5.67

-2.1 BLEU

-3.7 BLEU