# ESPnet-ST-v2: Multipurpose Spoken Language Translation Toolkit

Brian Yan\*1 Jiatong Shi\*1 Yun Tang<sup>2</sup> Hirofumi Inaguma<sup>2</sup> Peter Polák<sup>3</sup> Yifan Peng<sup>1</sup> Siddharth Dalmia<sup>1</sup> Patrick Fernandes<sup>1</sup> Dan Berrebbi<sup>1</sup> Tomoki Hayashi<sup>4</sup> Xiaohui Zhang<sup>2</sup> Zhaoheng Ni<sup>2</sup> Moto Hira<sup>2</sup> Juan Pino<sup>2</sup> Shinji Watanabe<sup>1,5</sup> Soumi Maiti<sup>1</sup> <sup>1</sup>Carnegie Mellon University <sup>2</sup>Meta AI <sup>3</sup>Charles University <sup>4</sup>Nagoya University <sup>5</sup>Johns Hopkins University {byan, jiatongs}@cs.cmu.edu

# Abstract

ESPnet-ST-v2 is a revamp of the open-source ESPnet-ST toolkit necessitated by the broadening interests of the spoken language translation community. ESPnet-ST-v2 supports 1) offline speech-to-text translation (ST), 2) simultaneous speech-to-text translation (SST), and 3) offline speech-to-speech translation (S2ST) each task is supported with a wide variety of approaches, differentiating ESPnet-ST-v2 from other open source spoken language translation toolkits. This toolkit offers state-ofthe-art architectures such as transducers, hybrid CTC/attention, multi-decoders with searchable intermediates, time-synchronous blockwise CTC/attention, Translatotron models, and direct discrete unit models. In this paper, we describe the overall design, example models for each task, and performance benchmarking behind ESPnet-ST-v2, which is publicly available at https://github.com/espnet/espnet.1

#### 1 Introduction

The objective of this project is to contribute to the diversity of the open-source spoken language translation ecosystem. Toward this, we launched this ESPnet-ST-v2 update in collaboration with researchers working on Fairseq (Ott et al., 2019) and TorchAudio (Yang et al., 2021b). This project focuses on: offline speech-to-text (ST), simultaneous speech-to-text (SST), and offline speech-to-speech (S2ST). These three spoken language translation tasks have drawn significant interest, as evidenced by rising IWSLT<sup>2</sup> shared task participation.

The ST task can be considered a base form of spoken language translation. Early approaches to ST stemmed from coupling statistical automatic speech recognition (ASR) (Huang et al., 2014) and text-to-text translation (MT) (Al-Onaizan et al., 1999), and this type of cascaded approach is still

common in the neural network era (Bentivogli et al., 2021; Zhang et al., 2022). End-to-end differentiable (E2E) approaches have recently emerged as an alternative offering greater simplicity and superior performance in some cases (Inaguma et al., 2021b); however, E2E approaches still benefit from techniques originating from ASR and MT (Gaido et al., 2021; Inaguma et al., 2021a).

SST modifies ST by imposing an additional streaming requirement, where systems are expected to produce textual outputs while incrementally ingesting speech input. Both the aforementioned cascaded and end-to-end approaches to ST have been adapted for SST (Ma et al., 2020b; Iranzo-Sánchez et al., 2021; Chen et al., 2021), although the more direct nature of the latter may be advantageous for latency-sensitive applications. On the other hand, S2ST extends ST by producing target speech rather than target text. Again, cascaded approaches of ST followed by text-to-speech (TTS) came first (Waibel et al., 1991; Black et al., 2002) and E2E approaches followed (Jia et al., 2019; Lee et al., 2022a; Jia et al., 2022a; Inaguma et al., 2022), with the latter offering smaller footprints and greater potential to retain source speech characteristics.

Given the recent swell in E2E ST, SST, and S2ST research, we have revamped ESPnet-ST (Inaguma et al., 2020) which previously only supported E2E ST. In particular, this work:

- Implements ST, SST, and S2ST using common Pytorch-based modules, including encoders, decoders, loss functions, search algorithms, and self-supervised representations.
- Builds a variety of example E2E models: attentional encoder-decoders, CTC/attention, multi-decoders with searchable intermediates, and transducers for ST. Blockwise attentional encoder-decoders, time-synchronous blockwise CTC/attention and blockwise transducers for SST. Spectral models (i.e. Translatotron) and

<sup>&</sup>lt;sup>1</sup>Please see our documentation for ST/SST and S2ST to get started. Example models and tutorials are provided.

<sup>&</sup>lt;sup>2</sup>International Workshop on Spoken Language Translation

discrete unit based models for S2ST.

 Benchmarks the ST, SST, and S2ST performance of ESPnet-ST-v2 against top IWSLT shared task systems and other prior works.

With this major update, ESPnet-ST-v2 keeps pace with the interests of the community and offers a variety of unique features, making it a valuable complement to Fairseq (Wang et al., 2020), NeurST (Zhao et al., 2021), and other spoken language translation toolkits.

#### 2 Related Works

ESPnet-ST-v2 follows a long line of open-source speech processing toolkits which can support spoken language translation (Zenkel et al., 2018; Shen et al., 2019; Kuchaiev et al., 2019; Hayashi et al., 2020; Wang et al., 2020; Zhao et al., 2021).

In Table 1 we compare ESPnet-ST-v2 to Fairseq (Wang et al., 2020) and NeurST (Zhao et al., 2021), two toolkits which also cover multiple types of spoken language translation. Fairseq and NeurST offer cascaded and E2E approaches to ST and SST (some of which are not offered by ESPnet-ST-v2). Meanwhile, ESPnet-ST-v2 focuses on E2E approaches and offers multiple unique core architectures not covered by the other toolkits. For S2ST, Fairseq and ESPnet-ST-v2 both offer a range of approaches. All told, ESPnet-ST-v2 offers the greatest variety across ST, SST, and S2ST – however, we view these toolkits as complementary. The following section elaborates on the unique features of ESPnet-ST-v2.

# 3 ESPnet-ST-v2

In this section, we first describe the overall design and then introduce a few key features.

# 3.1 Modular Design

Figure 1 illustrates the software architecture of ESPnet-ST-v2. This modular design is an improvement over the ESPnet-ST-v1 where monolithic model and task definitions made it more difficult to extend and modify the toolkit. We also designed ESPnet-ST-v2 such that modules developed for adjacent tasks (e.g. ASR, TTS, MT) can also be readily used for spoken language translation.

In ESPnet-ST-v2 major neural network modules, such as frontends, encoders, decoders, search, and loss functions, inherit from common abstract classes making them easy to interchange. These modules, which are detailed further in the next

		cs'	ESPARIST		
	R	AET.	AET.	estort	
FEATURES	E.S.	\$3	· EV	AEC	
Offline ST	1	1	1	1	
End-to-End Architecture(s)	1	1	1	/	
Attentional Enc-Dec	1	1	1	1	
CTC/Attention	1	-	-	-	
Transducer	1	-	-	-	
Hierarchical Encoders	1	-	-	-	
Multi-Decoder	1	-	-	-	
Cascaded Architectures	1	1	1	✓	
Speech SSL Representations	$\checkmark^1$	-	1	-	
Speech & Text Pre-training	1	1	1	✓	
Joint Speech/Text Pre-training	-	-	1	-	
Simultaneous ST	1	-	1	$\checkmark^3$	
End-to-End Architecture(s)	1	-	1	-	
Contextual Block Encoders	1	-	-	-	
Blockwise Attn Enc-Dec	1	-	-	-	
Blockwise CTC/Attention	1	-	-	-	
Blockwise Transducer	1	-	-	-	
Wait-K Attn Enc-Dec	-	-	1	-	
Monotonic Attn Enc-Dec	-	-	1	-	
Cascaded Architectures	-	-	1	$\checkmark^3$	
Offline S2ST	1	-	1	-	
End-to-End Architecture(s)	1	-	1	-	
Spec Enc-Dec (Translatotron)		-	1	-	
Spec Multi-Dec (Translatotron 2)	1	-	1	-	
Discrete Enc-Dec (Speech-to-Unit)	1	-	1	-	
Discrete Multi-Decoder (UnitY)	1	-	1	-	
Speech SSL Representations	$\checkmark^1$	-	1	-	
Neural Vocoder Support	$\checkmark^2$	1	1	-	

Table 1: Key features of ESPnet-ST-v2 compared to ESPnet-ST-v1 (Inaguma et al., 2020), Fairseq (Wang et al., 2020), and NeurST (Zhao et al., 2021). Comparison intends to highlight unique features of ESPnet-ST-v2 and not to comprehensively review all toolkits. <sup>1</sup>Supports S3PRL (Yang et al., 2021a). <sup>2</sup>Supports both spectral and discrete. <sup>3</sup>Only supports text-to-text.

subsection, are used as building blocks in wrapper classes which are used to construct model architectures. Then the fully constructed models are fed to task wrappers which prepare data loaders, initialize models, and handle training/validation. For inference, pythonic APIs invoke search algorithms over the trained models and direct outputs to scoring scripts. For instance, the third-party SimulEval tool for evaluating SST latency (Ma et al., 2020a) is integrated via this API layer. We are also integrating with TorchAudio (Yang et al., 2021b) in the same manner. Finally, recipe scripts define experimental pipelines from data preparation to evaluation.

# 3.2 Key Features

Each of the following modeling components feature a variety of interchangeable approaches.

#### Modules (Python/Pytorch) Wrappers (Python) # Frontend modules # Model wrappers for AbsFrontend connecting modules DefaultFrontend ESPnetSTModel # ST/SST S3prlFrontend ESPnetS2STModel FusedFrontends # Task wrappers for # Encoder modules model init & training AbsEncoder AbsTask ConformerEncoder STTask TransformerEncoder S2STTask EBranchfromerEncoder ContextBlkConformerEncoder ContextBlkTransformerEncode APIs (Python) # Inference APIs # Decoder modules Speech2Text AbsDecoder Speech2TextStreaming TransformerDecoder Speech2Speech RNNDecoder TransducerDecoder HuggingFaceTransformerDecoder Scripts (Bash/Python) # Search modules simuleval agent.py nbest\_mbr.py BeamSearch BeamSearchOnline # Misc. modules Recipes (Bash) JointNetwork st.sh # ST/SST s2st.sh

Figure 1: Software architecture of ESPnet-ST-v2.

Frontends & Targets Spectral features (e.g. FBANK) and features extracted from speech self-supervised learning (SSL) representations are supported, as well as fusions over multiple features (Berrebbi et al., 2022). For speech SSL features, ESPnet-ST-v2 integrates with the S3PRL toolkit (Yang et al., 2021a). These speech SSL representations are also used to generate discrete targets for S2ST (Lee et al., 2022a).

Encoder Architectures Conformer (Gulati et al., 2020; Guo et al., 2021), Branchformer (Peng et al., 2022), EBranchformer (Kim et al., 2023), and Transformer (Vaswani et al., 2017; Karita et al., 2019) encoder architectures are supported for ST and S2ST. For SST, a blockwise scheme is adopted following (Tsunoo et al., 2021; Deng et al., 2022) to form contextual block Conformer and Transformer encoders. Intermediate CTC (Lee and Watanabe, 2021) and Hierachical CTC (Sanabria and Metze, 2018) encoding are also supported; these techniques have been shown to stabilize deep encoder optimization (Lee and Watanabe, 2021) and improve representations for sequence tasks involving source-to-target re-ordering (Yan et al., 2023).

**Decoder Architectures** Attentional Transformer and recurrent neural network decoders are supported (Karita et al., 2019). Multi-decoder schemes which allow for E2E differentiable decoder cascades via searchable hidden intermediates (Dalmia et al., 2021), are also supported; this technique

has been shown to improve sequence modeling for tasks which naturally decompose into sub-tasks. Finally, large language model decoders (e.g. mBART (Liu et al., 2020b)) can be adopted through an integration with HuggingFace (Wolf et al., 2020).

**Loss Functions** Cross-entropy (for attentional decoders), CTC, and Transducer are supported for ST and SST. Multi-objective training with CTC/attention and CTC/transducer as well as multi-tasked training (e.g. ASR/MT/ST) is also supported. For S2ST, L1 and mean square error losses are also supported for spectral models.

Search Algorithms For offline attentional decoder models, label-synchronous beam search is supported with optional CTC joint decoding for multi-objective models (Watanabe et al., 2017). For offline Transducer models, the original Graves beam search (Graves, 2012) as well as timesynchronous and alignment-synchronous beam search (Saon et al., 2020) beam searches are supported. For SST, both incremental decoding and non-incremental (allowing re-translation) decoding (Liu et al., 2020a) are supported, along with stable hypothesis detection methods (Polák et al., 2022). Blockwise attentional decoder models use a label-synchronous beam search or timesynchronous beam search if a CTC branch is available. Blockwise transducer models use timesynchronous beam search.

Synthesis & Post-processing For ST, Minimum Bayes Risk (MBR) ensembling (Fernandes et al., 2022) is supported for leveraging quality-metrics (e.g. BLEU) to compare and rank n-best outputs from one or more models. For S2ST, neural vocoders are supported for both spectral and discrete inputs (Hayashi et al., 2020, 2021).

#### 4 Example Models

In this section, we introduce example models which are pre-built in ESPnet-ST-v2 using the neural network components described in the previous section. These examples include state-of-the-art core architectures, as evidenced by prior studies and our performance benchmarking (presented in §5).

### 4.1 ST Models

**CTC/Attention** (**CA**) Following Yan et al. (2023), we use Conformer encoders with hierarchical CTC encoding and Transformer decoders. The hierarchical CTC encoding, which aligns the first

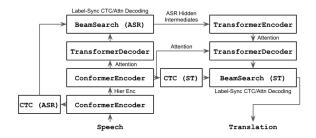


Figure 2: Multi-Decoder CTC/Attention for ST.

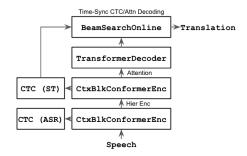


Figure 3: Time-Sync Blockwise CTC/Attn for SST.

N layers of the encoder towards ASR targets and the last M layers towards ST targets, regularizes the final encoder representations to be monotonic with respect to the target. CTC/attention models are jointly decoded using either label-synchronous (wherein the attention branch is primary) or time-synchronous (wherein the CTC branch is primary) beam search. For offline tasks, label-synchrony has shown greater performance (Yan et al., 2023).

**Multi-Decoder CTC/Attention (MCA)** As shown in Figure 2, the Multi-decoder decomposes ST into two sub-tasks, logically corresponding to ASR and MT encoder-decoder models, while maintaining E2E differentiability (Dalmia et al., 2021). This Multi-decoder scheme is also combined with the CTC/attention scheme described in the blurb above, following Yan et al. (2022). We use Conformer encoders with hierarchical CTC for encoding speech and Transformer encoders for encoding intermediate ASR text. We use Transformer decoders for both ASR and ST. During inference, the ASR stage is decoded first and then the final MT/ST stage is decoded; both stages use label-synchronous joint CTC/attention beam search.

#### 4.2 SST Models

**Time-Synchronous Blockwise CTC/Attention** (**TBCA**) As shown in Figure 3, we adapt the aforementioned CTC/attention model for ST (§4.1)

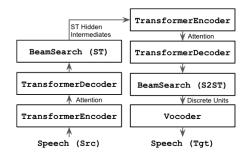


Figure 4: Discrete Multi-Decoder (UnitY) for S2ST.

to SST by replacing the Conformer encoder with a contextual block Conformer (Tsunoo et al., 2021). During inference, we initially followed Deng et al. (2022) and used the label-synchronous CTC/attention beam search originally proposed for ASR by Tsunoo et al. (2021). However, we found that label-synchrony results in overly conservative boundary block detection for SST. Therefore we opt instead for the time-synchronous variant which relies on CTC's more robust end-detection (Yan et al., 2023) to control boundary block detection; this change reduces latency without sacrificing quality. To perform incremental decoding without re-translation (as expected by SimulEval), hypotheses are pruned after processing all of the time steps for each encoder block.

Blockwise Transducer (BT) As demonstrated by Xue et al. (2022), Transducers can be effectively applied to SST despite the monotonic nature of their underlying alignment model. We build Transducers for SST using contextual block Conformer encoders and unidirectional LSTM decoders. We found that the aforementioned hierarchical CTC encoding (§4.1) improves training stability and convergence rate. During inference, we found that the time-synchronous algorithm described by Saon et al. (2020) outperformed the original Graves decoding (Graves, 2012) and the later proposed alignment-synchronous algorithms (Saon et al., 2020). We also found that length normalization is required to avoid overly short outputs. Incremental decoding is applied in the same manner as for TBCA.

#### 4.3 S2ST Models

**Spectral Multi-Decoder (Translatotron 2)** Similar to the MCA model for ST (§4.1), the spectral Multi-decoder (Jia et al., 2022a) decomposes S2ST into ST and TTS sub-tasks. The ST sub-

Toolkit	MODEL TYPE	DE ES FR			avg	
OFFLINE SPEECE	H TRANSLATION (ST)	BLEU↑				
NeurST (Zhao et al., 2021)	Attentional Enc-Dec (AED)	22.8 27.4 33.3			27.8	
Fairseq (Wang et al., 2020)	Attentional Enc-Dec (AED)	22.7	27.2	32.9	27.6	
ESPnet-ST-v1 (Inaguma et al., 2020)	Attentional Enc-Dec (AED)	22.9	28.0	32.8	27.9	
ESPnet-ST-v2 (this work)	Multi-Decoder CTC/Attn (MCA)	27.9 32.1 38.5				
SIMULTANEOUS SPE	ECH TRANSLATION (SST)	BLEU↑/ AL↓				
Fairseq (Wang et al., 2020) Wait-K Attentional Enc-Dec (WAED)		18.6 / 6.8	22.9 / 6.9	28.5 / 6.7	23.3 / 6.8	
ESPnet-ST-v2 (this work)	Time-Sync Blockwise CTC/Attn (TBCA)	23.5 / 2.3	29.2 / 2.4	32.7 / 2.3	28.5 / 2.3	
OFFLINE SPEECH-TO-SPEECH TRANSLATION (S2ST)			ASR-E	BLEU ↑		
Fairseq (Inaguma et al., 2022)	a et al., 2022) Discrete Multi-Decoder (UnitY)		32.3	30.9	29.6	
ESPnet-ST-v2 (this work) Discrete Multi-Decoder (UnitY)		23.7	32.0	33.1	29.6	

Table 2: Overview of ESPnet-ST-v2's ST, SST, and S2ST performances compared to other open-source toolkits. Results are presented on MuST-C-v1 (English-to-X) for ST/SST and on CVSS-C (X-to-English) for S2ST.

task is modeled with an encoder-decoder network while the TTS sub-task is modeled with an autoregressive synthesizer. The synthesizer attends over both the ST-encoder and ST-decoder hidden states. We use Transformers for the ST encoder-decoder and a Tacotron-style (Wang et al., 2017) decoder as the synthesizer. During inference, we first use beam search for the ST sub-task and then autoregressively generate Mel-spectrograms. The final waveform speech is generated with a HiFi-GAN vocoder (Kong et al., 2020).

**Discrete Multi-Decoder (UnitY)** The UnitY model (Inaguma et al., 2022) is similar to Translatotron 2, but critically predicts discrete units of speech SSL representations rather than spectral information in the final stage. In other words, UnitY is Multi-decoder consisting of a ST subtask followed by a text-to-unit (T2U) sub-task (see Figure 4). We use Transformer-based encoderdecoders for both sub-tasks. During inference, the ST stage is first decoded and then followed by the T2U stage. Both stages use label synchronous beam search. The final speech is generated with a unit HiFi-GAN vocoder with Fastspeech-like duration prediction (Polyak et al., 2021; Lee et al., 2022a), which is separately trained in the Parallel-WaveGAN toolkit (Hayashi et al., 2020, 2021).

# 5 Performance Benchmarking

In this section, we 1) compare open-source toolkits 2) compare our different example models and 3) compare our models with top IWSLT shared task systems and state-of-the-art prior works.

MODEL	HIERENC	BLEU↑
Attn Enc-Dec (AED)	-	25.7
Multi-Decoder Attn Enc-Dec (MAED)	-	27.6
CTC/Attention (CA)	✓	28.6
Multi-Decoder CTC/Attn (MCA)	✓	28.8
Transducer (T)	✓	27.6

Table 3: Example ST models – results on MuST-C-v2 En-De tst-COMMON.

#### **5.1** Experimental Setup

Please refer to §A.1 for reproducibility details. The following is only a summary of our setup.

**Data** We use MuST-C-v1 or MuST-C-v2 (Di Gangi et al., 2019) for ST/SST and CVSS-C for S2ST (Jia et al., 2022b). For IWSLT comparisons, we combine MuST-C-v1, MuST-C-v2, and ST-TED (Niehues et al., 2018) for ST/SST.

**Models** Unless otherwise indicated, we use a "base" setting for our models. Our base models have 40-80M trainable parameters across all tasks and are trained on a  $\sim$ 400h of single language pair data from a single corpus. For ST/SST, we also use a "large" setting for benchmarking against IWSLT submissions. Our large models have 150-200M trainable parameters and are trained on  $\sim$ 1000h of single language pair data from multiple corpora.

**Scoring** For ST/SST, we evaluate detokenized case-sensitive BLEU (Post, 2018). For SST, we additionally evaluate Average Lagging (AL) (Ma et al., 2020a). For S2ST, we evaluate ASR-BLEU by transcribing the generated speech and then evaluating the BLEU of this transcription.

Model	KD	BT	Ens	BLEU↑
IWSLT'21 (Top 3 of 6)				
1 Volctrans E2E <sup>†</sup>	/	-	1	24.3
2 OPPO Cascade <sup>†</sup>	✓	1	1	22.6
3 Volctrans Cascade <sup>†</sup>	✓	✓	✓	22.2
ESPnet-ST-v2				
A Base CA	-	-	-	23.2
B Base MCA	-	-	-	23.6
C Large CA	-	-	-	24.3
D Large MCA	-	-	-	25.1
E MBR (A+B+C+D)	-	-	✓	25.4

Table 4: Base and large CTC/attention (CA) and Multidecoder CTC/attention (MCA) models compared to top IWSLT 2021 systems for the given segmentation tst2020 En-De test set. KD=Knowledge Distillation, BT=Back-Translation, Ens=Ensemble. †Uses WMT MT data.

Model	BSz	BLEU†/AL↓
Blockwise Attn Enc-Dec (BAED)	40	22.8 / 3.23
Label-Sync Blockwise CTC/Attn (LBCA)	40	24.4 / 3.23
Time-Sync Blockwise CTC/Attn (TBCA)	40	24.6 / 2.34
Blockwise Transducer (BT)	40	22.9 / 2.37
Blockwise Attn Enc-Dec (BAED)	20	21.0 / 2.77
Label-Sync Blockwise CTC/Attn (LBCA)	20	<b>22.9</b> / 2.77
Time-Sync Blockwise CTC/Attn (TBCA)	20	22.8 / <b>1.63</b>
Blockwise Transducer (BT)	20	20.9 / 1.71

Table 5: Example SST models – results on MuST-C-v2 En-De tst-COMMON. BSz=Block Size.

### 5.2 Results

**Toolkit Comparison** Table 2 summarizes ESPnet-ST-v2 performance, showing one best example model (§4) for each task. ESPnet-ST-v1, Fairseq, and NeurST models are also referenced for comparison. On ST/SST, ESPnet-ST-v2 is 4-7 BLEU higher with 4.5 sec lower AL.<sup>3</sup> On S2ST ESPnet-ST-v2 is on par with Fairseq.

**ST** Table 3 shows a variety of approaches, of which the CTC/attention and Multi-decoder CTC/attention (MCA) models show the strongest performances. In Table 4, we scale these two approaches by training on larger corpora and increasing model capacity – our large MCA model outperforms the best IWSLT 2021 offline track submission on the 2020 test set with given segmentation.

**SST** Table 5 shows a variety of approaches, of which the blockwise Transducer (BT) and time-synchronous blockwise CTC/attention (TBCA)

MODEL	SSL	LLM	KD	BLEU↑ / AL↓
IWSLT'22 (Top 3 of 5)				
1 CUNI-KIT E2E	1	✓	-	31.5 / 1.93
2 UPV Cascade <sup>†</sup>	-	-	-	27.8 / 1.93
3 FBK E2E <sup>†</sup>	-	-	✓	25.0 / 1.99
ESPnet-ST-v2				
A Base TBCA	-	-	-	24.7 / 1.93
B Large TBCA	-	-	-	26.6 / 1.93

Table 6: Base and large time-sync CTC/attention (TBCA) models compared to top IWSLT 2022 systems for the medium latency regime. Evaluated on EnDe tst-COMMON-v2. SSL=Speech Self-Supervised Learning, LLM=Large Pre-trained Language Model, KD=Knowledge Distillation. †Uses WMT MT data.

MODEL	TYPE	ASR-BLEU↑
Prior Works		
1 Translatotron (Jia et al., 2019)	Spectral	14.4
2 Translatotron2 (Jia et al., 2022a)	Spectral	30.3
4 Speech-to-Unit (Lee et al., 2022a)	Discrete	30.8
5 UnitY (Inaguma et al., 2022)	Discrete	32.3
ESPnet-ST-v2		
A Attn Enc-Dec (Translatotron)	Spectral	16.6
B Multi-Decoder (Translatotron2)	Spectral	24.3
C Attn Enc-Dec (Speech-to-Unit)	Discrete	31.3
D Multi-Decoder (UnitY)	Discrete	32.0

Table 7: Example S2ST models – results on CVSS-C Es-En test set. Prior works shown for comparison.

models have the lowest AL. We choose to scale the TBCA to compare with IWSLT submissions due to its superior translation quality, but note that the BT has lower computational overhead due primarily to the lack of source-target computation; AL is non-computation aware. In Table 6, we fit the TBCA to the 2 second AL latency regime by selecting a blocksize of 32 and scale it with more data and model capacity – our large TBCA model would have ranked 3rd out of 6 amongst IWSLT 2022 submissions without using any SSL / LLM representations or knowledge distillation.

**S2ST** Table 7 shows a variety of approaches compared to prior works with comparable architectures – our S2ST models are generally on par with prior works which are considered state-of-the-art. In fact, all of our models slightly outperform their respective prior works except for Translatotron 2. Further, in Table 8 we ablate a range of SSL types for both the frontend and discrete units demonstrating the flexibility of our toolkit.

<sup>&</sup>lt;sup>3</sup>This comparison refers to the originally published results from the toolkit description papers. Note that subsequent works using these toolkits have improved the performance.

FRONTEND	DISCRETE UNIT	ASR-BLEU↑
FBANK	HuBERT	14.8
wav2vec2†	HuBERT	21.2
HuBERT†	HuBERT	21.4
mHuBERT	HuBERT	21.5
WavLM†	HuBERT	22.8
FBANK	WavLM	15.0
wav2vec2†	WavLM	21.6
HuBERT†	WavLM	22.1
mHuBERT	WavLM	22.0
WavLM†	WavLM	23.1

Table 8: Ablation on different types of SSL for the frontend and discrete unit portions of S2ST models. †Trained with large settings, others with base settings.

#### 6 Conclusion

We presented ESPnet-ST-v2 which now supports offline speech translation, simultaneous speech translation, and offline speech-to-speech translation. ESPnet-ST-v2 will continue to grow to support the community's interests. Future updates may include more new tasks, such as simultaneous speech-to-speech translation, and cross-toolkit integrations via TorchAudio.

## Limitations

The first set of limitations to be aware of are datarelated. Although prior works have shown the feasibility of building E2E systems without source language transcriptions (Lee et al., 2022b; Chen et al., 2022; Zhang et al., 2021), in this work we only investigate cases where triplet data (source speech, source transcript, target translation) is available for ST/SST and where quadruplet data (source speech, source transcript, target translation, target speech) is available for S2ST.

The second set of limitations to be aware of are evaluation-related. For SST, we follow prior works (Ma et al., 2020a; Wang et al., 2020; Anastasopoulos et al., 2022) and evaluate AL which is a measure of how much the system outputs lags behind the amount of input read. Notably, this does not consider the actual computation time and only the input-to-output ratio. For S2ST, we follow prior works (Jia et al., 2022a; Inaguma et al., 2022) and evaluate ASR-BLEU. This evaluation is dependent on an ASR system, which is not standardized across prior works. And further, our evaluation of S2ST outputs does not include naturalness. Finally, in this work we have not conducted any human evaluation of translation outputs.

# Acknowledgements

Brian Yan and Shinji Watanabe are supported by the Human Language Technology Center of Excellence. This work also used the Extreme Science and Engineering Discovery Environment (XSEDE) (Towns et al., 2014), which is supported by National Science Foundation grant number ACI-1548562; specifically, the Bridges system (Nystrom et al., 2015), as part of project cis210027p, which is supported by NSF award number ACI-1445606, at the Pittsburgh Supercomputing Center. This work also used GPUs donated by the NVIDIA Corporation.

#### References

Yaser Al-Onaizan, Jan Curin, Michael Jahr, Kevin Knight, John Lafferty, Dan Melamed, Franz-Josef Och, David Purdy, Noah A Smith, and David Yarowsky. 1999. Statistical machine translation. In *Final Report, JHU Summer Workshop*, volume 30.

Antonios Anastasopoulos, Loïc Barrault, Luisa Bentivogli, Marcely Zanon Boito, Ondřej Bojar, Roldano Cattoni, Anna Currey, Georgiana Dinu, Kevin Duh, Maha Elbayad, Clara Emmanuel, Yannick Estève, Marcello Federico, Christian Federmann, Souhir Gahbiche, Hongyu Gong, Roman Grundkiewicz, Barry Haddow, Benjamin Hsu, Dávid Javorský, Věra Kloudová, Surafel Lakew, Xutai Ma, Prashant Mathur, Paul McNamee, Kenton Murray, Maria Nădejde, Satoshi Nakamura, Matteo Negri, Jan Niehues, Xing Niu, John Ortega, Juan Pino, Elizabeth Salesky, Jiatong Shi, Matthias Sperber, Sebastian Stüker, Katsuhito Sudoh, Marco Turchi, Yogesh Virkar, Alexander Waibel, Changhan Wang, and Shinji Watanabe. 2022. Findings of the IWSLT 2022 evaluation campaign. In Proceedings of the 19th International Conference on Spoken Language Translation (IWSLT 2022), pages 98-157, Dublin, Ireland (in-person and online). Association for Computational Linguistics.

Luisa Bentivogli, Mauro Cettolo, Marco Gaido, Alina Karakanta, Alberto Martinelli, Matteo Negri, and Marco Turchi. 2021. Cascade versus direct speech translation: Do the differences still make a difference? In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 2873–2887.

Dan Berrebbi, Jiatong Shi, Brian Yan, Osbel López-Francisco, Jonathan Amith, and Shinji Watanabe. 2022. Combining Spectral and Self-Supervised Features for Low Resource Speech Recognition and Translation. In *Proc. Interspeech* 2022, pages 3533–3537.

- Alan W Black, Ralf D Brown, Robert Frederking, Rita Singh, John Moody, and Eric Steinbrecher. 2002. Tongues: Rapid development of a speech-to-speech translation system. In *Proceedings of Second International Conference on Human Language Technology Research HLT*, pages 183–186.
- Junkun Chen, Mingbo Ma, Renjie Zheng, and Liang Huang. 2021. Direct simultaneous speech-to-text translation assisted by synchronized streaming asr. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4618–4624.
- Peng-Jen Chen, Kevin Tran, Yilin Yang, Jingfei Du, Justine Kao, Yu-An Chung, Paden Tomasello, Paul-Ambroise Duquenne, Holger Schwenk, Hongyu Gong, et al. 2022. Speech-to-speech translation for a real-world unwritten language. *arXiv* preprint *arXiv*:2211.06474.
- Siddharth Dalmia, Brian Yan, Vikas Raunak, Florian Metze, and Shinji Watanabe. 2021. Searchable hidden intermediates for end-to-end models of decomposable sequence tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1882–1896, Online. Association for Computational Linguistics.
- Keqi Deng, Shinji Watanabe, Jiatong Shi, and Siddhant Arora. 2022. Blockwise Streaming Transformer for Spoken Language Understanding and Simultaneous Speech Translation. In *Proc. Interspeech 2022*, pages 1746–1750.
- Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2019. MuST-C: a Multilingual Speech Translation Corpus. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 2012–2017, Minneapolis, Minnesota. Association for Computational Linguistics.
- Patrick Fernandes, António Farinhas, Ricardo Rei, José De Souza, Perez Ogayo, Graham Neubig, and André FT Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412.
- Marco Gaido, Mauro Cettolo, Matteo Negri, and Marco Turchi. 2021. CTC-based compression for direct speech translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 690–696, Online. Association for Computational Linguistics.
- Alex Graves. 2012. Sequence transduction with recurrent neural networks. *ArXiv*, abs/1211.3711.

- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. 2020. Conformer: Convolution-augmented transformer for speech recognition. *Proc. Interspeech 2020*, pages 5036–5040.
- Pengcheng Guo, Florian Boyer, Xuankai Chang, Tomoki Hayashi, Yosuke Higuchi, Hirofumi Inaguma, Naoyuki Kamo, Chenda Li, Daniel Garcia-Romero, Jiatong Shi, et al. 2021. Recent developments on espnet toolkit boosted by conformer. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5874–5878. IEEE.
- Tomoki Hayashi, Ryuichi Yamamoto, Katsuki Inoue, Takenori Yoshimura, Shinji Watanabe, Tomoki Toda, Kazuya Takeda, Yu Zhang, and Xu Tan. 2020. ESPnet-TTS: Unified, reproducible, and integratable open source end-to-end text-to-speech toolkit. In ICASSP 2020-2020 IEEE international conference on acoustics, speech and signal processing (ICASSP), pages 7654–7658. IEEE.
- Tomoki Hayashi, Ryuichi Yamamoto, Takenori Yoshimura, Peter Wu, Jiatong Shi, Takaaki Saeki, Yooncheol Ju, Yusuke Yasuda, Shinnosuke Takamichi, and Shinji Watanabe. 2021. ESPnet2-TTS: Extending the edge of TTS research. *arXiv* preprint arXiv:2110.07840.
- Xuedong Huang, James Baker, and Raj Reddy. 2014. A historical perspective of speech recognition. *Communications of the ACM*, 57(1):94–103.
- Hirofumi Inaguma, Tatsuya Kawahara, and Shinji Watanabe. 2021a. Source and target bidirectional knowledge distillation for end-to-end speech translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1872–1881.
- Hirofumi Inaguma, Shun Kiyono, Kevin Duh, Shigeki Karita, Nelson Enrique Yalta Soplin, Tomoki Hayashi, and Shinji Watanabe. 2020. ESPnet-ST: All-in-one speech translation toolkit. *arXiv preprint arXiv:2004.10234*.
- Hirofumi Inaguma, Sravya Popuri, Ilia Kulikov, Peng-Jen Chen, Changhan Wang, Yu-An Chung, Yun Tang, Ann Lee, Shinji Watanabe, and Juan Pino. 2022. UnitY: Two-pass direct speech-to-speech translation with discrete units. arXiv preprint arXiv:2212.08055.
- Hirofumi Inaguma, Brian Yan, Siddharth Dalmia, Pengcheng Guo, Jiatong Shi, Kevin Duh, and Shinji Watanabe. 2021b. ESPnet-ST IWSLT 2021 offline speech translation system. In *Proceedings of the 18th International Conference on Spoken Language Translation (IWSLT 2021)*, pages 100–109, Bangkok, Thailand (online). Association for Computational Linguistics.

- Javier Iranzo-Sánchez, Javier Jorge, Pau Baquero-Arnal, Joan Albert Silvestre-Cerdà, Adrià Giménez, Jorge Civera, Albert Sanchis, and Alfons Juan. 2021. Streaming cascade-based speech translation leveraged by a direct segmentation model. *Neural Networks*, 142:303–315.
- Ye Jia, Michelle Tadmor Ramanovich, Tal Remez, and Roi Pomerantz. 2022a. Translatotron 2: High-quality direct speech-to-speech translation with voice preservation. In *International Conference on Machine Learning*, pages 10120–10134. PMLR.
- Ye Jia, Michelle Tadmor Ramanovich, Quan Wang, and Heiga Zen. 2022b. CVSS corpus and massively multilingual speech-to-speech translation. In *Proceedings of Language Resources and Evaluation Conference (LREC)*, pages 6691–6703.
- Ye Jia, Ron J Weiss, Fadi Biadsy, Wolfgang Macherey, Melvin Johnson, Zhifeng Chen, and Yonghui Wu. 2019. Direct speech-to-speech translation with a sequence-to-sequence model. *Proc. Interspeech* 2019, pages 1123–1127.
- Shigeki Karita, Nanxin Chen, Tomoki Hayashi, Takaaki Hori, Hirofumi Inaguma, Ziyan Jiang, Masao Someki, Nelson Enrique Yalta Soplin, Ryuichi Yamamoto, Xiaofei Wang, et al. 2019. A comparative study on transformer vs RNN in speech applications. In 2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), pages 449–456. IEEE.
- Kwangyoun Kim, Felix Wu, Yifan Peng, Jing Pan, Prashant Sridhar, Kyu J Han, and Shinji Watanabe. 2023. E-branchformer: Branchformer with enhanced merging for speech recognition. In 2022 IEEE Spoken Language Technology Workshop (SLT), pages 84–91. IEEE.
- Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. Hifi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis. *Advances in Neural Information Processing Systems*, 33:17022–17033.
- Oleksii Kuchaiev, Jason Li, Huyen Nguyen, Oleksii Hrinchuk, Ryan Leary, Boris Ginsburg, Samuel Kriman, Stanislav Beliaev, Vitaly Lavrukhin, Jack Cook, et al. 2019. Nemo: a toolkit for building ai applications using neural modules. *arXiv preprint arXiv:1909.09577*.
- Ann Lee, Peng-Jen Chen, Changhan Wang, Jiatao Gu, Sravya Popuri, Xutai Ma, Adam Polyak, Yossi Adi, Qing He, Yun Tang, et al. 2022a. Direct speech-to-speech translation with discrete units. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3327–3339.
- Ann Lee, Hongyu Gong, Paul-Ambroise Duquenne, Holger Schwenk, Peng-Jen Chen, Changhan Wang, Sravya Popuri, Yossi Adi, Juan Pino, Jiatao Gu, and

- Wei-Ning Hsu. 2022b. Textless speech-to-speech translation on real data. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 860–872, Seattle, United States. Association for Computational Linguistics.
- Jaesong Lee and Shinji Watanabe. 2021. Intermediate loss regularization for CTC-based speech recognition. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6224–6228. IEEE.
- Danni Liu, Gerasimos Spanakis, and Jan Niehues. 2020a. Low-Latency Sequence-to-Sequence Speech Recognition and Translation by Partial Hypothesis Selection. In *Proc. Interspeech 2020*, pages 3620–3624.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020b. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguis*tics. 8:726–742.
- Xutai Ma, Mohammad Javad Dousti, Changhan Wang, Jiatao Gu, and Juan Pino. 2020a. Simuleval: An evaluation toolkit for simultaneous translation. In *Proceedings of the EMNLP*.
- Xutai Ma, Juan Pino, and Philipp Koehn. 2020b. SimulMT to simulST: Adapting simultaneous text translation to end-to-end simultaneous speech translation. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 582–587.
- Jan Niehues, Rolando Cattoni, Sebastian Stüker, Mauro Cettolo, Marco Turchi, and Marcello Federico. 2018. The IWSLT 2018 evaluation campaign. In *Proceedings of the 15th International Conference on Spoken Language Translation*, pages 2–6, Brussels. International Conference on Spoken Language Translation.
- Nicholas A Nystrom, Michael J Levine, Ralph Z Roskies, and J Ray Scott. 2015. Bridges: a uniquely flexible hpc resource for new communities and data analytics. In *Proceedings of the 2015 XSEDE Conference: Scientific Advancements Enabled by Enhanced Cyberinfrastructure*, pages 1–8.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*.
- Yifan Peng, Siddharth Dalmia, Ian Lane, and Shinji Watanabe. 2022. Branchformer: Parallel mlp-attention architectures to capture local and global context for speech recognition and understanding.

- In *International Conference on Machine Learning*, pages 17627–17643. PMLR.
- Peter Polák, Ngoc-Quan Ngoc, Tuan-Nam Nguyen, Danni Liu, Carlos Mullov, Jan Niehues, Ondrej Bojar, and Alexander Waibel. 2022. Cuni-kit system for simultaneous speech translation task at iwslt 2022. *IWSLT* 2022, page 277.
- Adam Polyak, Yossi Adi, Jade Copet, Eugene Kharitonov, Kushal Lakhotia, Wei-Ning Hsu, Abdelrahman Mohamed, and Emmanuel Dupoux. 2021. Speech Resynthesis from Discrete Disentangled Self-Supervised Representations. In *Proc. Interspeech* 2021.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Ramon Sanabria and Florian Metze. 2018. Hierarchical multitask learning with CTC. In 2018 IEEE Spoken Language Technology Workshop (SLT), pages 485–490. IEEE.
- George Saon, Zoltán Tüske, and Kartik Audhkhasi. 2020. Alignment-length synchronous decoding for RNN transducer. In *ICASSP 2020 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7804–7808.
- Jonathan Shen, Patrick Nguyen, Yonghui Wu, Zhifeng Chen, Mia X Chen, Ye Jia, Anjuli Kannan, Tara Sainath, Yuan Cao, Chung-Cheng Chiu, et al. 2019. Lingvo: a modular and scalable framework for sequence-to-sequence modeling. arXiv preprint arXiv:1902.08295.
- J. Towns, T. Cockerill, M. Dahan, I. Foster, K. Gaither, A. Grimshaw, V. Hazlewood, S. Lathrop, D. Lifka, G. D. Peterson, R. Roskies, J. R. Scott, and N. Wilkins-Diehr. 2014. Xsede: Accelerating scientific discovery. *Computing in Science & Engineering*, 16(5):62–74.
- Emiru Tsunoo, Yosuke Kashiwagi, and Shinji Watanabe. 2021. Streaming transformer ASR with blockwise synchronous beam search. In 2021 IEEE Spoken Language Technology Workshop (SLT), pages 22–29. IEEE.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Alex Waibel, Ajay N Jain, Arthur E McNair, Hiroaki Saito, Alexander G Hauptmann, and Joe Tebelskis. 1991. JANUS: a speech-to-speech translation system using connectionist and symbolic processing strategies. In [Proceedings] ICASSP 91: 1991 International Conference on Acoustics, Speech, and Signal Processing, pages 793–796 vol.2.

- Changhan Wang, Yun Tang, Xutai Ma, Anne Wu, Dmytro Okhonko, and Juan Pino. 2020. Fairseq S2T: Fast speech-to-text modeling with fairseq. In Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing: System Demonstrations, pages 33–39.
- Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly, Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. 2017. Tacotron: Towards end-to-end speech synthesis. *Proc. Interspeech* 2017, pages 4006–4010.
- Shinji Watanabe, Takaaki Hori, Suyoun Kim, John R. Hershey, and Tomoki Hayashi. 2017. Hybrid CTC/attention architecture for end-to-end speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 11(8):1240–1253.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.
- Jian Xue, Peidong Wang, Jinyu Li, Matt Post, and Yashesh Gaur. 2022. Large-Scale Streaming Endto-End Speech Translation with Neural Transducers. In *Proc. Interspeech* 2022, pages 3263–3267.
- Brian Yan, Siddharth Dalmia, Yosuke Higuchi, Graham Neubig, Florian Metze, Alan W Black, and Shinji Watanabe. 2023. CTC alignments improve autoregressive translation. *The 17th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*.
- Brian Yan, Patrick Fernandes, Siddharth Dalmia, Jiatong Shi, Yifan Peng, Dan Berrebbi, Xinyi Wang, Graham Neubig, and Shinji Watanabe. 2022. CMU's iwslt 2022 dialect speech translation system. In *Proceedings of the 19th International Conference on Spoken Language Translation (IWSLT 2022)*, pages 298–307.
- Shu-wen Yang, Po-Han Chi, Yung-Sung Chuang, Cheng-I Jeff Lai, Kushal Lakhotia, Yist Y. Lin, Andy T. Liu, Jiatong Shi, Xuankai Chang, Guan-Ting Lin, Tzu-Hsien Huang, Wei-Cheng Tseng, Ko tik Lee, Da-Rong Liu, Zili Huang, Shuyan Dong, Shang-Wen Li, Shinji Watanabe, Abdelrahman Mohamed, and Hung yi Lee. 2021a. SUPERB: Speech processing universal performance benchmark. *Proc. Interspeech* 2021.
- Yao-Yuan Yang, Moto Hira, Zhaoheng Ni, Anjali Chourdia, Artyom Astafurov, Caroline Chen, Ching-Feng

Yeh, Christian Puhrsch, David Pollack, Dmitriy Genzel, Donny Greenberg, Edward Z. Yang, Jason Lian, Jay Mahadeokar, Jeff Hwang, Ji Chen, Peter Goldsborough, Prabhat Roy, Sean Narenthiran, Shinji Watanabe, Soumith Chintala, Vincent Quenneville-Bélair, and Yangyang Shi. 2021b. Torchaudio: Building blocks for audio and speech processing. *arXiv* preprint arXiv:2110.15018.

Thomas Zenkel, Matthias Sperber, Jan Niehues, Markus Müller, Ngoc-Quan Pham, Sebastian Stüker, and Alex Waibel. 2018. Open source toolkit for speech to text translation. *Prague Bull. Math. Linguistics*, 111:125–135.

Chen Zhang, Xu Tan, Yi Ren, Tao Qin, Kejun Zhang, and Tie-Yan Liu. 2021. UWSpeech: Speech to speech translation for unwritten languages. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14319–14327.

Weitai Zhang, Zhongyi Ye, Haitao Tang, Xiaoxi Li, Xinyuan Zhou, Jing Yang, Jianwei Cui, Pan Deng, Mohan Shi, Yifan Song, et al. 2022. The USTC-NELSLIP offline speech translation systems for IWSLT 2022. In *Proceedings of the 19th International Conference on Spoken Language Translation (IWSLT 2022)*, pages 198–207.

Chengqi Zhao, Mingxuan Wang, Qianqian Dong, Rong Ye, and Lei Li. 2021. NeurST: Neural speech translation toolkit. In the 59th Annual Meeting of the Association for Computational Linguistics (ACL): System Demonstrations.

# A Appendix

# A.1 Reproducibility

Table 9 shows the hyperparameters for the models presented in §5. All of our data preparation scripts are available in ESPnet: https://github.com/espnet/espnet/tree/master/egs2.

Model	Task	Encoder(s)	Decoder(s)	Frontend	Pre-Train Init	Multi-Obj	Src BPE	Tgt BPE	# Params
AED (Table 3)	ST	12 lyr, 4 head, 256 adim	(ASR) 6 lyr, 4 head (ST) 6 lyr, 4 head	FBANK	ASR Enc/Dec	ASR	4k	4k	60M
MAED (Table 3)	ST	(ASR) 12 lyr, 4 head, 256 adim (MT) 2 lyr, 4 head, 256 adim	(ASR) 6 lyr, 4 head (MT) 6 lyr, 4 head	FBANK	ASR Enc/Dec	ASR	4k	4k	60M
CA (Table 3)	ST	18 lyr, 4 head, 256 adim	(ASR) 6 lyr, 4 head (ST) 6 lyr, 4 head	FBANK	ASR Enc/Dec	ASR	4k	4k	70M
MCA (Table 3)	ST	(ASR) 18 lyr, 4 head, 256 adim (MT) 4 lyr, 4 head, 256 adim	(ASR) 6 lyr, 4 head (MT) 6 lyr, 4 head	FBANK	ASR Enc/Dec/CTC	ASR	4k	4k	70M
T (Table 3)	ST	18 lyr, 4 head, 256 adim	1 lyr, 512 dim, 640 joint (ASR) 6 lyr, 4 head	FBANK	ASR Enc/Dec/CTC	ASR	4k	4k	70M
Large CA (Table 4)	ST	18 lyr, 8 head, 512 adim	(ASR) 6 lyr, 4 head (ST) 6 lyr, 4 head	HuBERT	ASR Enc/Dec/CTC	ASR	8k	16k	210M
Large MCA (Table 4)	ST	(ASR) 18 lyr, 8 head, 512 adim (MT) 4 lyr, 8 head, 512 adim	(ASR) 6 lyr, 8 head (MT) 6 lyr, 8 head	HuBERT	ASR Enc/Dec/CTC	ASR	8k	8k	210M
BAED (Table 5)	SST	18 lyr, 4 head, 256 adim	6 lyr, 4 head	FBANK	ASR Enc lyr 1-12	-	4k	4k	70M
LBCA (Table 5)	SST	18 lyr, 4 head, 256 adim	6 lyr, 4 head	FBANK	ASR Enc lyr 1-12	-	4k	4k	70M
TBCA (Table 5)	SST	18 lyr, 4 head, 256 adim	6 lyr, 4 head	FBANK	ASR Enc lyr 1-12	-	4k	4k	70M
BT (Table 5)	SST	18 lyr, 4 head, 256 adim	1 lyr, 4 head, 640 joint	FBANK	ASR Enc lyr 1-12	-	4k	4k	40M
Large TBCA (Table 6)	SST	18 lyr, 8 head, 512 adim	6 lyr, 8 head	FBANK	ASR Enc lyr 1-12	-	8k	8k	150M
Translatotron (Table 7)	S2ST	12 lyr, 4 head, 256 adim	6 lyr, 1024 dim	FBANK	-	ASR, ST	7k	500	80M
Translatotron2 (Table 7)	S2ST	16 lyr, 4 head, 256 adim	(ST) 6 lyr, 4 head (TTS) 2 lyr, 1024 dim	FBANK	-	ASR, ST	7k	500	50M
Speech-to-Unit (Table 7)	S2ST	12 lyr, 4 head, 512 adim	6 lyr, 8 head	FBANK	-	ASR, ST	7k	500	40M
UnitY (Table 7)	S2ST	(ST) 16 lyr, 4 head, 256 adim (T2U) 2 lyr, 4 head, 256 adim	(ST) 4 lyr, 4 head (T2U) 2 lyr, 8 head	FBANK	-	ASR, ST	7k	500	40M

Table 9: ST, SST, and S2ST model hyperparameters. Parameter counts are rounded to the nearest 10 million.