

# Speech-to-Speech Translation with Discrete-Unit-Based Style Transfer

Anonymous ACL submission

## Abstract

Direct speech-to-speech translation (S2ST) with discrete self-supervised representations has achieved remarkable accuracy, but is unable to preserve the speaker timbre of the source speech. Meanwhile, the scarcity of high-quality speaker-parallel data poses a challenge for learning style transfer during translation. We propose an S2ST framework with style-transfer capability on the basis of discrete self-supervised speech representations and codec units. The acoustic language model we introduce for style transfer leverages self-supervised in-context learning, acquiring style transfer ability without relying on any speaker-parallel data, thereby overcoming the issue of data scarcity. By using extensive training data, our model achieves zero-shot cross-lingual style transfer on previously unseen source languages. Experiments show that our model generates translated speeches with high fidelity and style similarity.<sup>1</sup>

## 1 Introduction

Speech-to-speech translation (S2ST) aims to translate spoken utterances from one language to another, which can bring immense convenience to international communication. Compared to conventional cascaded systems comprising ASR, text translation and TTS models (Lavie et al., 1997; Nakamura et al., 2006; Wahlster, 2013), direct S2ST models without intermediate text generation have a more concise pipeline with less computation cost and error propagation while facilitating application to unwritten languages, and thus spark widespread interest in the community.

Mainstream approaches of direct S2ST (Lee et al., 2021a,b; Huang et al., 2022; Popuri et al., 2022) utilize discrete speech representation from self-supervised models (such as HuBERT (Hsu

et al., 2021)) as prediction target, and then use them to reconstruct the waveform. Such representation eliminates speaker identity and prosody of the speeches and retains only semantic contents, which simplifies the target distribution and makes the translation less challenging. However, it also has the drawback of losing the style information of the source speech. Extra voice conversion systems are needed if users want to keep the source speaker timbre, which may cause degradation in audio quality and content accuracy.

Some works propose direct S2ST with style transfer (Jia et al., 2021; Song et al., 2023). These methods depend on paired data that source and target speech share the same speaker. However, such data from the real world is extremely scarce as it requires a large number of multilingual speakers, while simulated data from multilingual TTS systems suffers from less diversity and extra data collection costs. Recent large-scale S2ST models (Rubenstein et al., 2023; Barrault et al., 2023) have also incorporated the capability of style transfer, yet their sub-modules are highly coupled and are difficult to apply to other S2ST models.

Inspired by recent progress in spoken language models (Borsos et al., 2023; Wang et al., 2023), we propose a novel approach for direct S2ST with the ability of cross-lingual style transfer, and does not rely on any speaker-parallel data. We utilize two types of discrete representations, namely semantic and acoustic units, from a self-supervised speech model and a neural codec, separately. Our method encompasses three stages: 1) speech-to-semantic-unit translation, which translates source speech to target semantic units; 2) acoustic unit modeling, which generates target acoustic units from translated semantic units using style information in the source speech; and 3) unit-to-wave generation, which reconstructs high-fidelity translated speech from the acoustic units. The modules of the three stages are trained independently and decou-

<sup>1</sup>Audio samples are available at <http://stylelm.github.io/>

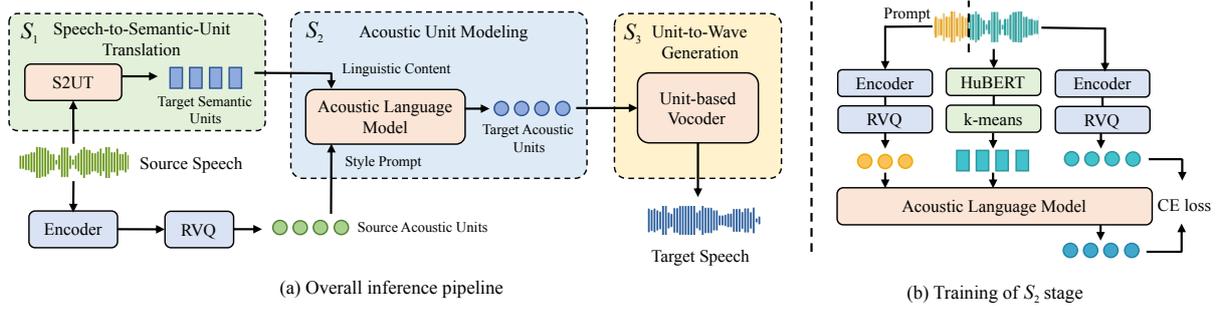


Figure 1: We propose an S2ST approach with style transfer based on discrete representations from a self-supervised speech model and a neural codec. Figure (a) shows the inference pipeline of our method; figure (b) illustrates the self-supervised training process of the acoustic language model of  $S_2$ .

080 pled from each other, allowing our framework to  
 081 be applied to various existing speech-to-unit  
 082 translation models.

083 For the acoustic unit modeling stage, we intro-  
 084 duce an acoustic language model. It employs a self-  
 085 supervised training approach and learns style trans-  
 086 fer through in-context learning, which relies on no  
 087 speaker-parallel data, and thus addresses the issue  
 088 of data scarcity. By utilizing extensive training  
 089 data, our model achieves zero-shot cross-lingual  
 090 style transfer with source languages not included  
 091 in the training. Experiments show that our model  
 092 generates results with superior audio quality and  
 093 style similarity while maintaining accurate content.

094 Our contributions can be summarized as follows:

- 095 • We propose an S2ST approach with cross-  
 096 lingual style transfer capability, even on previ-  
 097 ously unseen source languages.
- 098 • By employing self-supervised training, our  
 099 model does not rely on any speaker-parallel  
 100 data, thus addressing the issue of data scarcity.
- 101 • The decoupling nature of the sub-modules en-  
 102 ables our framework to be adopted by various  
 103 existing speech-to-unit translation models.
- 104 • Experiments show that our method generates  
 105 translated speeches with high quality and style  
 106 similarity.

## 107 2 Method

108 The overall inference pipeline of our method is il-  
 109 lustrated in Fig.1 (a). Our method comprises three  
 110 consecutive stages, utilizing two distinct types of  
 111 discrete units: 1) speech-to-semantic-unit transla-  
 112 tion stage  $S_1$ , which converts source audio into  
 113 semantic units of the translated speech; 2) acoustic

114 unit modeling stage  $S_2$ , generating target acoustic  
 115 units conditioned on the semantic output from the  
 116 preceding stage and the acoustic units of the source  
 117 speech as style prompt; 3) unit-to-wave genera-  
 118 tion stage  $S_3$ , producing translated speech that main-  
 119 tains consistent style with the source. We provide  
 120 details about these two types of units and the three  
 121 stages in the following subsections.

### 122 2.1 Semantic and Acoustic Units

123 Discrete HuBERT (Hsu et al., 2021) units obtained  
 124 from the clustering of self-supervised speech rep-  
 125 resentations are shown (Lee et al., 2021b; Huang  
 126 et al., 2022) to be effective in providing seman-  
 127 tic content information and are widely adopted  
 128 in S2ST as prediction target (Lee et al., 2021a,b;  
 129 Huang et al., 2022; Popuri et al., 2022). HuBERT  
 130 encodes the target speech into continuous represen-  
 131 tations with a frame length of 20 ms, and these rep-  
 132 resentations are then discretized with the k-means  
 133 algorithm to get the semantic units.

134 On the other hand, audio codec models with  
 135 encoder-decoder architecture such as SoundStream  
 136 (Zeghidour et al., 2021) have recently shown out-  
 137 standing performance in learning acoustic infor-  
 138 mation. Such a codec model can produce discrete  
 139 representations (i.e. the acoustic units) of audio  
 140 by employing a convolutional encoder followed by  
 141 a residual vector quantizer. These representations  
 142 contain detailed acoustic information and can be  
 143 used to reconstruct waveforms with the correspond-  
 144 ing decoder or an additional vocoder.

### 145 2.2 Speech-to-Semantic-Unit Translation

146 The speech-to-semantic-unit translation stage gen-  
 147 erates translated semantic units conditioned on  
 148 source speech input, achieving translation of lin-  
 149 guistic content. Various models (Lee et al., 2021a;

Huang et al., 2022; Popuri et al., 2022) have been proposed for this procedure. These models share a common basic architecture of a convolutional speech encoder followed by an encoder-decoder architecture based on a transformer or conformer. Due to the decoupling nature of the sub-modules of the three stages, we have the flexibility to adopt different S2UT models in this stage, and we attempted two of them in our experiments (See Section 3.1).

### 2.3 Acoustic Unit Modeling

The acoustic unit modeling stage  $S_2$  generates translated acoustic units from semantic tokens and style prompts. The core component of  $S_2$  is an acoustic language model, which is basically a decoder-only transformer. The model takes a prefix sequence formed by concatenating acoustic unit sequence  $\mathbf{a}_p$ , which serves as a style prompt, and the target semantic sequence  $\mathbf{s}$ , and generates the target acoustic sequence  $\mathbf{a}$  autoregressively. This procedure can be formulated as

$$p(\mathbf{a} | \mathbf{a}_p, \mathbf{s}; \theta_{AR}) = \prod_{t=1}^T \prod_{c=1}^C p(\mathbf{a}_t^c | \mathbf{a}_{<t}, \mathbf{a}_t^{<c}, \mathbf{a}_p, \mathbf{s}; \theta_{AR}) \quad (1)$$

The entire sequence is in the format of  $[\mathbf{a}_p | \mathbf{s} | \mathbf{a}]$ , with a separator token between each pair of adjacent parts. 3 codebooks are used for  $\mathbf{a}_p$  and  $\mathbf{a}$ .

The training procedure of  $S_2$  is illustrated in Figure 1(b). It adopts a self-supervised training paradigm, where the first three seconds of each audio sample is truncated as prompt, and the acoustic language model is trained to predict the acoustic units of the remaining part conditioned on its semantic units and the prompt acoustic units with cross-entropy loss. This in-context learning approach enables the model to grasp the correspondence in acoustic characteristics between the two parts and acquire style transfer ability. During inference, we use semantic tokens from the previous stage and acoustic units of source speech as the style prompt to realize cross-lingual style transfer.

### 2.4 Unit-to-Wave Generation

In the waveform generation stage  $S_3$ , we adopt a GAN-based unit vocoder to map the target acoustic units to high-fidelity waveforms. Our vocoder is derived from BigVGAN (Lee et al., 2022), with a generator built from a set of look-up tables (LUT) that embed the discrete units, and a series of blocks composed of transposed convolution and a residual block with dilated layers. Multi-period discrim-

inator (MPD) and multi-resolution discriminator (MRD) are used for adversarial training.

## 3 Experiments

### 3.1 Setup

**Datasets** We use two language pairs in the CVSS dataset (Jia et al., 2022) as the translation benchmark, which are French-English (Fr-En) and Spanish-English (Es-En). For  $S_2$  and  $S_3$  stages, we use the *unlab-60k* subset of Libri-Light (Kahn et al., 2020) to train the acoustic language model, and use LibriTTS (Zen et al., 2019) to train the SoundStream model and the vocoder. All audio is processed at a 16 kHz sampling rate. We provide more details about the datasets in Appendix A.

**Model Configurations** We apply the publicly available multilingual HuBERT (mHuBERT) model with the k-means model of 1000 clusters for the 11th-layer features and train a SoundStream model with a size of 1024 for each codebook and an overall downsampling rate of 320. For stage  $S_1$ , we train an S2UT-conformer for Fr-En following (Lee et al., 2021a), and an xm-transformer for Es-En following (Popuri et al., 2022) but without mbart-decoder initialization. The decoder-only transformer of  $S_2$  has about 760M parameters, with details of its architecture provided in Appendix B. **Baselines** Considering that previous S2ST models with style transfer (Jia et al., 2021; Song et al., 2023; Rubenstein et al., 2023; Barrault et al., 2023) either differ from ours in settings or are not open-sourced, we mainly compare our model with S2UT models used in  $S_1$  followed by a single-speaker vocoder, and cascaded pipelines formed by appending various voice conversion models after the vocoder, which are PPG-VC(Liu et al., 2021), NANSY(Choi et al., 2021) and YourTTS(Casanova et al., 2022).

**Evaluation Metrics** We employ both objective and subjective metrics to measure the model performance in terms of translation accuracy, speech quality, and style similarity with the source speech. For objective evaluation, we calculate the BLEU score between the ASR-transcripts of the translated speech and reference text as well as speaker cosine similarity (Cos). For subjective metrics, we use crowd-sourced human evaluation with 1-5 Likert scales and report mean opinion scores on speech quality (MOS) and style similarity (SMOS) with 95% confidence intervals (CI). More details are provided in Appendix C.

Table 1: Translation Quality and Audio Similarity on CVSS Dataset.

ID	Model	BLEU (Fr-En) ( $\uparrow$ )	BLEU (Es-En) ( $\uparrow$ )	MOS ( $\uparrow$ )	SMOS( $\uparrow$ )	Cos ( $\uparrow$ )
1	S2UT	18.08	23.78	$3.73 \pm 0.05$	/	/
2	S2UT + PPG-VC	17.05	23.03	$3.37 \pm 0.07$	$3.30 \pm 0.06$	0.65
3	S2UT + NANSY	17.21	23.36	$3.56 \pm 0.06$	$3.47 \pm 0.05$	0.68
4	S2UT + YourTTS	16.73	22.09	$3.74 \pm 0.05$	$3.60 \pm 0.06$	0.69
5	Ours	17.64	23.41	$3.86 \pm 0.06$	<b><math>3.69 \pm 0.05</math></b>	<b>0.74</b>
6	GT (CVSS-C)	84.52	88.54	$3.92 \pm 0.05$	/	/
7	GT (CVSS-T)	81.48	84.81	$3.95 \pm 0.05$	$3.56 \pm 0.06$	0.68

Table 2: Ablation on Training Data Volume and Sizes of  $S_2$  Model.

ID	Model	BLEU (Fr-En) ( $\uparrow$ )	BLEU (Es-En) ( $\uparrow$ )	MOS ( $\uparrow$ )	SMOS ( $\uparrow$ )	Cos ( $\uparrow$ )
Ablation on Training Data Volume						
1	LibriTTS	17.62	23.37	$3.84 \pm 0.05$	$3.55 \pm 0.05$	0.67
2	Libri-Light unlab-60k	17.64	23.41	$3.86 \pm 0.05$	$3.69 \pm 0.05$	0.74
3	+ CVSS source	17.25	23.49	$3.85 \pm 0.05$	$3.71 \pm 0.05$	0.76
Ablation on Model Size						
4	Small (160M)	16.55	21.78	$3.73 \pm 0.06$	$3.58 \pm 0.05$	0.70
5	Base (430M)	16.87	22.36	$3.81 \pm 0.05$	$3.64 \pm 0.05$	0.73
6	Large (760M)	17.64	23.41	$3.86 \pm 0.05$	$3.69 \pm 0.05$	0.74

### 3.2 Results and Analysis

Table 1 summarizes the main experiment results. We observe a comprehensive decrease in BLEU scores for 2-5 compared to 1, indicating that additional style transfer processes lead to a loss in semantic content. Nevertheless, our model achieves the slightest decrease of 0.44 and 0.37 in BLEU, together with the highest MOS of 3.86. This indicates that in comparison to cascaded voice conversion, our style transfer mechanism based on discrete intermediate representations can mitigate quality and content losses during the transfer and produce high-quality audio.

On the other hand, our model achieves the highest speaker similarity, with SMOS being 3.69 and Cos being 0.74, which surpasses all three cascaded systems and even the CVSS-T target, demonstrating the outstanding performance in zero-shot cross-lingual style transfer of our model. This can be attributed to the large model size and extensive training data, through which our model acquires strong zero-shot style transfer capability and can generalize effectively to unseen source languages.

### 3.3 Ablation Studies

We further conduct ablations on the training data volume and model size of  $S_2$ , and the results are summarized in Table 2. We observe that when using LibriTTS with shorter duration and fewer

speakers, there is a significant decrease in SMOS and Cos of 0.14 and 0.07, with only a minor decrease in BLEU and MOS of 0.02, 0.04, and 0.02. This suggests that the model’s style transfer performance relies on a large amount of speech data from multiple speakers, while achieving high-quality speech generation does not require as much data. We also add part of the speech from the CVSS source to the training data, obtaining a marginal improvement of 0.02 on both Cos and SMOS. This indicates that with extensive training data, the performance of  $S_2$  on unseen source languages is close to that on seen languages. Furthermore, we observe a comprehensive improvement in all metrics as the model size increases in 4-6, proving that the superior performance of our acoustic language model is closely linked to its large parameter size.

## 4 Conclusions

We propose an S2ST approach with style transfer capability by adopting an acoustic language model that learns style transfer through in-context learning. By adopting self-supervised training and large-scale training data, our method addresses the scarcity of speaker-parallel data and achieves cross-lingual style transfer with unseen source languages. Experiments indicate that our approach achieves outstanding results in terms of translation accuracy, speech quality and style similarity.

## 5 Limitations and Potential Risks

Despite that our model excels in style transfer and generating high-quality translated speech, it still suffers from several limitations: 1) Our evaluation (especially the objective evaluation) of style transfer capability mainly focuses on the global speaker timbre, and we have not yet delved deeply into other stylistic characteristics such as prosody and emotion. We leave the exploration of these aspects for future work. 2) The large model size and the autoregressive generation paradigm may lead to efficiency issues, such as long inference latency. 3) The BLEU scores heavily depend on the ASR quality, which may not accurately reflect the speech translation performance. Future directions could be improving ASR quality or exploring other evaluation metrics without reliance on ASR models. Besides, due to the speaker timbre transfer capability of our model, it may be misused to disinform, defame, or commit fraud. We will add some constraints to guarantee people who use our code or pre-trained model will not use the model in illegal cases.

## References

Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenhaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, et al. 2023. Seamless: Multilingual expressive and streaming speech translation. *arXiv preprint arXiv:2312.05187*.

Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, et al. 2023. Audiolm: a language modeling approach to audio generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.

Edresson Casanova, Julian Weber, Christopher D Shulby, Arnaldo Candido Junior, Eren Gölge, and Moacir A Ponti. 2022. Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone. In *International Conference on Machine Learning*, pages 2709–2720. PMLR.

Hyeong-Seok Choi, Juheon Lee, Wansoo Kim, Jie Lee, Hoon Heo, and Kyogu Lee. 2021. Neural analysis and synthesis: Reconstructing speech from self-supervised representations. *Advances in Neural Information Processing Systems*, 34:16251–16265.

Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction

of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3451–3460.

Rongjie Huang, Jinglin Liu, Huadai Liu, Yi Ren, Lichao Zhang, Jinzheng He, and Zhou Zhao. 2022. Transpeech: Speech-to-speech translation with bilateral perturbation. *arXiv preprint arXiv:2205.12523*.

Ye Jia, Michelle Tadmor Ramanovich, Tal Remez, and Roi Pomerantz. 2021. Translatotron 2: Robust direct speech-to-speech translation.

Ye Jia, Michelle Tadmor Ramanovich, Quan Wang, and Heiga Zen. 2022. Cvss corpus and massively multilingual speech-to-speech translation. *arXiv preprint arXiv:2201.03713*.

Jacob Kahn, Morgane Rivièrè, Weiyi Zheng, Evgeny Kharitonov, Qiantong Xu, Pierre-Emmanuel Mazaré, Julien Karadayi, Vitaliy Liptchinsky, Ronan Collobert, Christian Fuegen, et al. 2020. Libri-light: A benchmark for asr with limited or no supervision. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7669–7673. IEEE.

Alon Lavie, Alex Waibel, Lori Levin, Michael Finke, Donna Gates, Marsal Gavaldà, Torsten Zeppenfeld, and Puming Zhan. 1997. Janus-iii: Speech-to-speech translation in multiple languages. In *1997 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 1, pages 99–102. IEEE.

Ann Lee, Peng-Jen Chen, Changhan Wang, Jiatao Gu, Sravya Popuri, Xutai Ma, Adam Polyak, Yossi Adi, Qing He, Yun Tang, et al. 2021a. Direct speech-to-speech translation with discrete units. *arXiv preprint arXiv:2107.05604*.

Ann Lee, Hongyu Gong, Paul-Ambroise Duquenne, Holger Schwenk, Peng-Jen Chen, Changhan Wang, Sravya Popuri, Yossi Adi, Juan Pino, Jiatao Gu, et al. 2021b. Textless speech-to-speech translation on real data. *arXiv preprint arXiv:2112.08352*.

Sang-gil Lee, Wei Ping, Boris Ginsburg, Bryan Catanzaro, and Sungroh Yoon. 2022. Bigvgan: A universal neural vocoder with large-scale training. *arXiv preprint arXiv:2206.04658*.

Songxiang Liu, Yewen Cao, Disong Wang, Xixin Wu, Xunying Liu, and Helen Meng. 2021. Any-to-many voice conversion with location-relative sequence-to-sequence modeling. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:1717–1728.

Satoshi Nakamura, Konstantin Markov, Hiromi Nakaiwa, Gen-ichiro Kikui, Hisashi Kawai, Takatoshi Jitsuhiro, J-S Zhang, Hirofumi Yamamoto, Eiichiro Sumita, and Seiichi Yamamoto. 2006. The atr multilingual speech-to-speech translation system. *IEEE Transactions on Audio, Speech, and Language Processing*, 14(2):365–376.

409 Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan,  
 410 Sam Gross, Nathan Ng, David Grangier, and Michael  
 411 Auli. 2019. fairseq: A fast, extensible toolkit for  
 412 sequence modeling. In *Proceedings of NAACL-HLT*  
 413 *2019: Demonstrations*.

414 Sravya Popuri, Peng-Jen Chen, Changhan Wang, Juan  
 415 Pino, Yossi Adi, Jiatao Gu, Wei-Ning Hsu, and Ann  
 416 Lee. 2022. Enhanced direct speech-to-speech trans-  
 417 lation using self-supervised pre-training and data aug-  
 418 mentation. *arXiv preprint arXiv:2204.02967*.

419 Paul K Rubenstein, Chulayuth Asawaroengchai,  
 420 Duc Dung Nguyen, Ankur Bapna, Zalán Borsos,  
 421 Félix de Chaumont Quitry, Peter Chen, Dalia El  
 422 Badawy, Wei Han, Eugene Kharitonov, et al. 2023.  
 423 Audiopalm: A large language model that can speak  
 424 and listen. *arXiv preprint arXiv:2306.12925*.

425 Kun Song, Yi Ren, Yi Lei, Chunfeng Wang, Kun Wei,  
 426 Lei Xie, Xiang Yin, and Zejun Ma. 2023. Styles2st:  
 427 Zero-shot style transfer for direct speech-to-speech  
 428 translation. *arXiv preprint arXiv:2305.17732*.

429 Wolfgang Wahlster. 2013. *Verbmobil: foundations of*  
 430 *speech-to-speech translation*. Springer Science &  
 431 Business Media.

432 Changhan Wang, Anne Wu, and Juan Pino. 2020. Cov-  
 433 ost 2 and massively multilingual speech-to-text trans-  
 434 lation. *arXiv preprint arXiv:2007.10310*.

435 Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang,  
 436 Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu,  
 437 Huaming Wang, Jinyu Li, et al. 2023. Neural codec  
 438 language models are zero-shot text to speech synthe-  
 439 sizers. *arXiv preprint arXiv:2301.02111*.

440 Neil Zeghidour, Alejandro Luebs, Ahmed Omran,  
 441 Jan Skoglund, and Marco Tagliasacchi. 2021.  
 442 Soundstream: An end-to-end neural audio codec.  
 443 *IEEE/ACM Transactions on Audio, Speech, and Lan-  
 444 guage Processing*, 30:495–507.

445 Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J  
 446 Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. 2019.  
 447 LibriTTS: A corpus derived from librispeech for text-  
 448 to-speech. *arXiv preprint arXiv:1904.02882*.

## 449 A Datasets

450 In this section, we provide details of the translation  
 451 benchmark dataset and the corpora for training  $S_2$   
 452 and  $S_3$  models.

453 **CVSS** CVSS (Jia et al., 2022) is an S2ST bench-  
 454 mark dataset derived from the CoVoST 2 (Wang  
 455 et al., 2020) speech-to-text translation corpus by  
 456 synthesizing the translation text into speech us-  
 457 ing TTS systems. It comprises two sub-versions  
 458 of CVSS-C and CVSS-T, where the target speech  
 459 in CVSS-C is generated by a single-speaker TTS  
 460 system while that of CVSS-T is generated by a  
 461 multi-speaker TTS system with speaker timbre

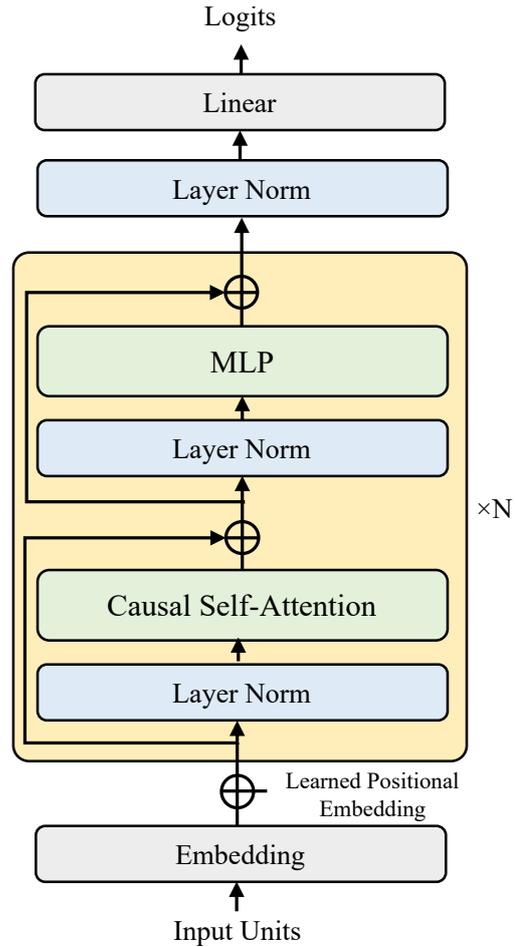


Figure 2: Structure of the Acoustic Language Model.

462 transferred from the source speech. We use CVSS-  
 463 C for training and evaluating the translation models,  
 464 and provide results of ground truth target audios  
 465 in CVSS-T as a reference for style transfer perfor-  
 466 mance.

467 **Libri-Light** Libri-Light is a large-scale corpus con-  
 468 taining unlabelled speech from audiobooks in En-  
 469 glish. The *unlab-60k* subset we use consists of  
 470 57.7k hours of audio with 7,439 speakers.

471 **LibriTTS** LibriTTS is a multi-speaker English  
 472 TTS dataset. It comprises 585.5 hours of audio  
 473 with 2,456 speakers.

## 474 B Model Settings

475 We illustrate the structure of the acoustic language  
 476 model in Figure 2, and provide hyperparameters  
 477 of our  $S_2$  and  $S_3$  stage models in Table 3. We  
 478 also refer the readers to the original papers (Lee  
 479 et al., 2021a; Popuri et al., 2022) for details of  $S_1$   
 480 models used. Each sub-module is trained with 4  
 481 NVIDIA-V100 GPUs for about a week.

Hyperparameter		Prompt-Singer
Acoustic Language Model (Small)	Layers	22
	Hidden Dim	768
	Attention Headers	12
	FFN Dim	3,072
	Number of Parameters	160.5M
Acoustic Language Model (Base)	Layers	26
	Hidden Dim	1,152
	Attention Headers	16
	FFN Dim	4,608
	Number of Parameters	420.2M
Acoustic Language Model (Large)	Layers	26
	Hidden Dim	1,536
	Attention Headers	16
	FFN Dim	6,144
	Number of Parameters	763.1M
Unit Vocoder	Upsample Rates	[5,4,2,2,2,2]
	Hop Size	320
Unit Vocoder	Upsample Kernel Sizes	[9,8,4,4,4,4]
	Number of Parameters	121.6M

Table 3: Hyperparameters of  $S_2$  and  $S_3$  Stage Models.

## C Evaluation Metrics

For translation accuracy, we use an open-sourced ASR model in *fairseq*<sup>2</sup> (Ott et al., 2019) framework to transcribe the audios and then calculate the BLEU score between the transcripts and the reference text. For speaker similarity, we use Rysemblyzer<sup>3</sup>, which is a public-available speaker encoder to extract speaker embeddings of the synthesized and source speech and calculate their cosine similarity.

Our subjective evaluation tests are crowd-sourced and conducted via Amazon Mechanical Turk. For audio quality evaluation, we ask the testers to examine the audio quality and naturalness. For style similarity, we instruct the testers to evaluate the style similarity between the synthesized and source speech while ignoring the content. The testers rate scores on 1-5 Likert scales. We provide screenshots of the testing interfaces in Figure 3 and 4. Each data item is rated by 2 testers, and the testers are paid \$8 hourly.

We calculate BLEU scores over the entire test split and randomly sample 500 items from each language pair for other metrics, which represents approximately 3% of the test set.

<sup>2</sup>[https://github.com/facebookresearch/fairseq/tree/main/examples/speech\\_to\\_speech/asr\\_bleu](https://github.com/facebookresearch/fairseq/tree/main/examples/speech_to_speech/asr_bleu)

<sup>3</sup><https://github.com/resemble-ai/Resemblyzer>

**Previewing Answers Submitted by Workers**  
 This message is only visible to you and will not be shown to Workers.  
 You can test completing the task below and click "Submit" in order to preview the data and format of the submitted results.

**Instructions** | **Shortcuts** | How natural (i.e. human-sounding) is this recording? Please focus on examining the audio quality and naturalness, and ignore the differences of style (timbre, emotion and prosody).

▶ 0:00 / 0:01

**Select an option**

Excellent - Completely natural speech - 5	1
4.5	2
Good - Mostly natural speech - 4	3
3.5	4
Fair - Equally natural and unnatural speech - 3	5
2.5	6
Poor - Mostly unnatural speech - 2	7
1.5	8
Bad - Completely unnatural speech - 1	9

Figure 3: Screenshot of MOS testing.

**Previewing Answers Submitted by Workers**  
 This message is only visible to you and will not be shown to Workers.  
 You can test completing the task below and click "Submit" in order to preview the data and format of the submitted results.

**Instructions** | **Shortcuts** | How similar is this recording to the reference audio? Please focus on the similarity of the style (speaker identity, emotion and prosody) to the reference, and ignore the differences of content, grammar, or audio quality.

**Reference audio:**  
 ▶ 0:00 / 0:06

**Testing audio:**  
 ▶ 0:00 / 0:03

**Select an option**

Excellent - Completely similar speech - 5	1
4.5	2
Good - Mostly similar speech - 4	3
3.5	4
Fair - Equally similar and dissimilar speech - 3	5
2.5	6
Poor - Mostly dissimilar speech - 2	7
1.5	8
Bad - Completely dissimilar speech - 1	9

Figure 4: Screenshot of SMOS testing.