# MAM: Masked Acoustic Modeling for End-to-End Speech-to-Text Translation

**Junkun Chen** [* 1 2]  **Mingbo Ma** [* 1]  **Renjie Zheng** [1]  **Kaibo Liu** [1]  **Liang Huang** [1 2]

## Abstract

End-to-end Speech-to-text Translation (E2E-ST), which directly translates source language speech to target language text, is widely useful in practice, but traditional cascaded approaches (ASR+MT) often suffer from error propagation in the pipeline. On the other hand, existing end-to-end solutions heavily depend on the source language transcriptions for pre-training or multi-task training with Automatic Speech Recognition (ASR). We instead propose a simple technique to learn a robust speech encoder in a self-supervised fashion only on the speech side, which can utilize speech data without transcription. This technique termed Masked Acoustic Modeling (MAM), not only provides an alternative solution to improving E2E-ST, but also can perform pre-training on any acoustic signals (including non-speech ones) without annotation. We conduct our experiments over 8 different translation directions. In the setting without using any transcriptions, our technique achieves an average improvement of +1.1 BLEU, and +2.3 BLEU with MAM pre-training. Pre-training of MAM with arbitrary acoustic signals also has an average improvement with +1.6 BLEU for those languages. Compared with ASR multi-task learning solution, which replies on transcription during training, our pre-trained MAM model, which does not use transcription, achieves similar accuracy.

## 1 Introduction

Speech-to-text translation (ST), which translates the source language speech to target language text, is useful in many scenarios such as international conferences, travels, foreign-language video subtitling, etc. Conventional cascaded approaches to ST (Ney, 1999; Matusov et al., 2005; Mathias & Byrne, 2006; Berard et al., 2016) first transcribe the speech audio into source language text (ASR) and then perform text-to-text machine translation (MT), which inevitably suffers from error propagation in the pipeline. To alleviate
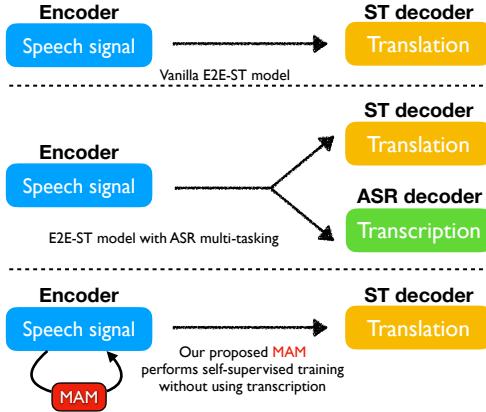
*Figure 1.* Comparisons with different existing solutions and our proposed Masked Acoustic Modeling (MAM).

this problem, recent efforts explore end-to-end approaches (E2E-ST) (Weiss et al., 2017; Berard et al., 2018; Vila et al., 2018; Gangi et al., 2019), which are computationally more efficient at inference time and mitigate the risk of error propagation from imperfect ASR.

To improve the translation accuracy of E2E-ST models, researchers either initialize the encoder of ST with a pre-trained ASR encoder (Berard et al., 2018; Bansal et al., 2019; Wang et al., 2020b) to get better representations of the speech signal, or perform Multi-Task Learning (MTL) with ASR to bring more training and supervision signals to the shared encoder (Anastasopoulos et al., 2016; Anastasopoulos & Chiang, 2018; Sperber et al., 2019; Liu et al., 2019b) (see Fig. 1). These methods improve the translation quality by providing more training signals to the encoder to learn better phonetic information and hidden representation correspondence (Stoian et al., 2020).

However, both above solutions assume the existence of substantial speech transcriptions of the source language. Unfortunately, this assumption is problematic. On the one hand, for certain low-resource languages, especially endangered ones (Bird, 2010; Bird et al., 2014), the source speech transcriptions are expensive to collect. Moreover, according to the report from Ethnologue, there are more than 3000 languages that have no written form or no standard orthography, making phonetic transcription impossible (Duong et al., 2016). On the other hand, the amount of speech audios with

transcriptions are limited (as they are expensive to collect), and there exist far more audios without any annotations. It will be much more straightforward and cheaper to leverage these raw audios to train a robust encoder directly.

To relieve from the dependency on source language transcriptions, we present a straightforward yet effective solution, Masked Acoustic Modeling (MAM), to utilize the speech data in a self-supervised fashion without using any source language transcription, unlike other speech pre-training models (Chuang et al., 2019; Wang et al., 2020c). Aside from the regular training of E2E-ST (without ASR as MTL or pre-training), MAM masks certain portions of the speech input randomly and aims to recover the masked speech signals with their context on the encoder side. MAM not merely provides an alternative solution to improving E2E-ST, but also is a general technique that can be used as a pre-training module on arbitrary acoustic signals, e.g., multilingual speech, music, animal sounds. The contributions of our paper are as follows:

- We demonstrate the importance of a self-supervising module for E2E-ST. Unlike all previous attempts, which heavily depend on transcription, MAM improves the capacity of the encoder by recovering masked speech signals merely based on their context. MAM also can be used together with transcriptions in ASR pre-training and MTL settings to further boost the translation accuracy.

- MAM also can be used as a pre-training module solely by itself. During pre-training, MAM is capable to utilize arbitrary acoustic signal (e.g., music, animal sound) other than regular speech audio. Considering there are much more acoustic data than human speech, MAM has better potential to be used for pre-training. To the best of our knowledge, MAM is the first technique that is able to perform pre-training with any form of the audio signal.

- For 8 different translation directions, when we do not use any transcription, MAM demonstrates an average BLEU improvements of 1.09 in the basic setting and 2.26 with pre-training.

- We show that the success of MAM does not rely on intensive or expensive computation. MAM only has 6.5% more parameters than the baseline model.

## 2   Preliminaries: ASR and ST

We first briefly review the standard E2E-ST and E2E-ST with ASR MTL to set up the notations.

### 2.1   Vanilla E2E-ST Training with Seq2Seq

Regardless of particular design of Seq2Seq models for different tasks, the encoder always takes the source input sequence $\boldsymbol{x} = (x_1, ..., x_n)$ of $n$ elements where each

$x_i \in \mathbb{R}^{d_x}$ is a $d_x$-dimension vector and produces a sequence of hidden representations $\boldsymbol{h} = f(\boldsymbol{x}) = (h_1, ..., h_n)$ where $h_i = f(\boldsymbol{x})$. The encoding function $f$ can be implemented by a mixture between Convolution, RNN and Transformer. More specifically, $\boldsymbol{x}$ can be the spectrogram or mel-spectrogram of the source speech, and each $x_i$ represents the frame-level speech feature with certain duration.

On the other hand, the decoder greedily predicts a new output word $y_t$ given both the source sequence $\boldsymbol{x}$ and the prefix of decoded tokens, denoted $\boldsymbol{y}_{<t} = (y_1, ..., y_{t-1})$. The decoder continues the generation until it emits <eos> and finishes the entire decoding process. Finally, we obtain the hypothesis $\boldsymbol{y} = (y_1, ..., \text{<eos>})$ with the model score which defined as following:

$$p(\boldsymbol{y} \mid \boldsymbol{x}) = \prod_{t=1}^{|\boldsymbol{y}|} p(y_t \mid \boldsymbol{x}, \boldsymbol{y}_{<t}) \tag{1}$$

During the training time, the entire model aims to maximize the conditional probability of each ground-truth target sentence $\boldsymbol{y}^\star$ given input $\boldsymbol{x}$ over the entire training corpus $D_{\boldsymbol{x},\boldsymbol{y}^\star}$, or equivalently minimizing the following loss:

$$\ell_{\text{ST}}(D_{\boldsymbol{x},\boldsymbol{y}^\star}) = -\sum_{(\boldsymbol{x},\boldsymbol{y}^\star) \in D_{\boldsymbol{x},\boldsymbol{y}^\star}} \log p(\boldsymbol{y}^\star \mid \boldsymbol{x}) \tag{2}$$

### 2.2   Multi-task Learning with ASR

To further boost the performance of E2E-ST, researchers proposed to either use pre-trained ASR encoder to initialize ST encoder, or to perform ASR MTL together with ST training. We only discuss the MTL since pre-training does not require significant change to Seq2Seq model.

During multi-task training, there are two decoders sharing one encoder. Besides the MT decoder, there is also another decoder for generating transcriptions. With the help of ASR training, the encoder is able to learn more accurate speech segmentations (similar to forced alignment) making the global reordering of those segments for MT relatively easier. We defined the following training loss for ASR:

$$\ell_{\text{ASR}}(D_{\boldsymbol{x},\mathbf{z}^\star}) = -\sum_{(\boldsymbol{x},\mathbf{z}^\star) \in D_{\boldsymbol{x},\mathbf{z}^\star}} \log p(\mathbf{z}^\star \mid \boldsymbol{x}) \tag{3}$$

where $\mathbf{z}^\star$ represents the annotated, ground-truth transcription for speech audio $\boldsymbol{x}$. In our baseline setting, we also hybrid CTC/Attention framework (Watanabe et al., 2017) on the encoder side. In the case of multi-task training with ASR for ST, the total loss is defined as

$$\ell_{\text{MTL}}(D_{\boldsymbol{x},\boldsymbol{y}^\star,\mathbf{z}^\star}) = \ell_{\text{ST}}(D_{\boldsymbol{x},\boldsymbol{y}^\star}) + \ell_{\text{ASR}}(D_{\boldsymbol{x},\mathbf{z}^\star}) \tag{4}$$

where $D_{\boldsymbol{x},\boldsymbol{y}^\star,\mathbf{z}^\star}$ is the training dataset which contains speech, translation and transcription triplets.

## 3   Masked Acoustic Modeling

All the existing solutions to boost the current E2E-ST performance heavily depend on the availability of the transcription
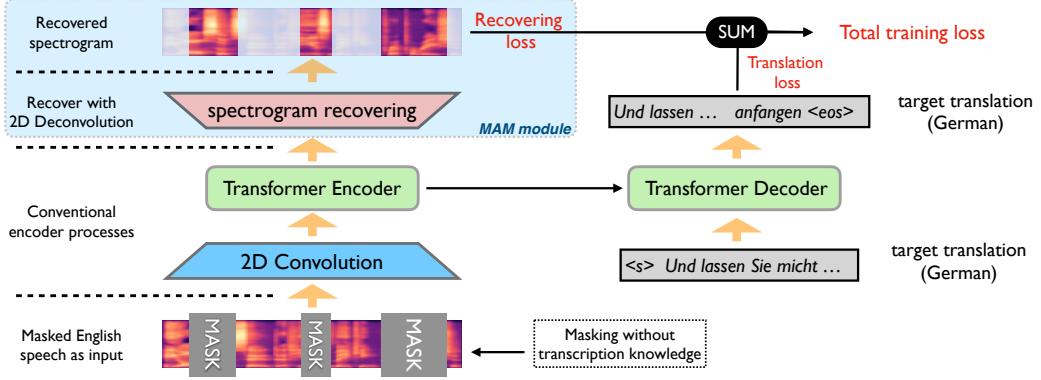
*Figure 2.* MAM (in blue box) can be treated as one extra module besides standard Transformer encoder-decoder and convolution layers for processing speech signals.

of the source language. Those solutions are not able to take advantage of large amount of speeches without any annotations. They also become inapplicable when the source language is low-resource or even does not have a standard orthography system. Therefore, the ideal solution should not be constrained by source language transcription and still achieves similar translation quality. Thus, we introduce MAM in this section.

## 3.1 MAM as Part of Training Objective

We propose to perform self-supervised training on the encoder side by reconstructing sabotaged speech signals from the input. Note that MAM is totally different from another self-supervised training (Chuang et al., 2019; Wang et al., 2020a;c) which rely on transcription to segment the speech audio with forced alignment tools(Povey et al., 2011; McAuliffe et al., 2017). We directly apply random masks with different widths over speech audio, eliminating the dependency of transcription. Therefore, MAM can be easily applied to speech audio without transcription and even to any non-human speech audio, e.g., music and animal sound.

Formally, we define a random replacement function over the original speech input $\boldsymbol{x}$:

$$\hat{\boldsymbol{x}} \sim \text{Mask}_{\text{frame}}(\boldsymbol{x}, \lambda), \qquad (5)$$

where $\text{Mask}(\cdot)_{\text{frame}}$ randomly replaces some certain frames in $\boldsymbol{x}$ with the same random initialized vector, $\epsilon \in \mathbb{R}^{d_x}$, with a probability of $\lambda$ (30% in our experiments). Note that we use the same vector $\epsilon$ to represent all the corrupted frames (see one example in Fig.4b). Then we obtain a corrupted input $\hat{\boldsymbol{x}}$ and its corresponding latent representation $\hat{\boldsymbol{h}}$.

For MAM module, we have the following training objective to reconstruct the original speech signal with the surrounding context information with self-supervised fashion:

$$\ell_{\text{Rec}}(D_{\boldsymbol{x}}) = \sum_{\boldsymbol{x} \in D_{\boldsymbol{x}}} ||\boldsymbol{x} - \phi(f(\hat{\boldsymbol{x}}))||_2^2 \qquad (6)$$

where $\phi$ is a reconstruction function which tries to recover the original signal from the hidden representation $f(\hat{\boldsymbol{x}})$ with corrupted inputs. For simplicity, we use regular 2D deconvolution as $\phi$, and mean squared error for measuring the difference between original input and recovered signal. Finally, we have the following total loss of our model

$$\ell_{\text{MAM}}(D_{\boldsymbol{x},\boldsymbol{y}^\star}) = \ell_{\text{ST}}(D_{\boldsymbol{x},\boldsymbol{y}^\star}) + \ell_{\text{Rec}}(D_{\boldsymbol{x}})$$

To further boost the performance of E2E-ST, we can train MAM with ASR MTL when transcription is available:

$$\ell_{\text{MAM + MTL}}(D_{\boldsymbol{x},\boldsymbol{y}^\star,\boldsymbol{z}^\star}) = \ell_{\text{MTL}}(D_{\boldsymbol{x},\boldsymbol{y}^\star,\boldsymbol{z}^\star}) + \ell_{\text{Rec}}(D_{\boldsymbol{x}})$$

## 3.2 Different Masking Strategies

MAM aims at much harder tasks than pure textual pretraining models, e.g., BERT or ERINE, which only perform semantic learning over missing tokens. In our case, we not only try to recover semantic meaning, but also acoustic characteristics of given audio. MAM simultaneously predicts the missing words and generates spectrograms like speech synthesis tasks.

To ensure the masked segments contain different levels of granularity of speech semantic, we propose the following masking methods.

**Single Frame Masking**  Uniformly mask $\lambda\%$ frames out of $\boldsymbol{x}$ to construct $\hat{\boldsymbol{x}}$. Note that we might have continuous frames that were masked.

**Span Masking**  Similar with SpanBERT (Joshi et al., 2020b), we first sample a serial of span widths and then apply those spans randomly to different positions of the input signal. Note that we do not allow overlap in this case. Our span masking is defined as $\hat{\boldsymbol{x}} \sim \text{Mask}_{\text{span}}(\boldsymbol{x}, \lambda)$.

## 3.3 Pre-training MAM

MAM is a powerful technique that is not only beneficial to the conventional training procedure, but also can be used as a pre-training framework that does not need any annotation.

The bottleneck of current speech-related tasks, e.g., ASR, ST, is lacking of the annotated training corpus. For some languages that do not even have a standard orthography system, these annotations are even impossible to obtain. Although current speech-related, pre-training frameworks (Chuang et al., 2019; Wang et al., 2020c) indeed relieve certain needs of large scale parallel training corpus for E2E-ST, all of these pre-training methods still need intense transcription annotation for the source speech.

During pre-training time, we only use the encoder part of MAM. Thanks to our flexible masking techniques, MAM is able to perform pre-training with any kind of audio signal. This allows us to perform pre-training with MAM with three different settings, pre-training with source language speech, with multilingual speech, and arbitrary audios. To the best of our knowledge, MAM is the first pre-training technique that can be applied to arbitrary audios. Considering about the vast arbitrary acoustic signals existing on the Internet (e.g., youtube), MAM has great potential to further boost the downstream tasks. MAM that pre-trained with arbitrary acoustic signals does not differentiate languages and provides the unified pre-trained model for any downstream, fine-tuning task. This is different from the multilingual pre-training setting since the downstream task's source speech language is not necessary to be included in the pre-training stage, which is essential to the low-resource and zero-resource languages.

## 4 Experiments

In this section, we conducted MAM pre-training experiment on three corpora, Libri-Light (only English speech, medium version) (Kahn et al., 2020), Common Voice [1] (We select 14 languages, which contains ca, de, en, es, fr, it, kab, nl, pl, pt, ro, ru, zh-CN, zh-TW) (Ardila et al., 2020), and Audioset (arbitrary acoustic data) (Gemmeke et al., 2017). The statistical results of the dataset are shown in Table. 2. Note that Audioset includes a wide range of arbitrary sounds, from human and animal sounds, to natural and environmental sounds, to musical and miscellaneous sounds.

Then, we analyze the performance of MAM in E2E-ST with 8 different language translation directions using English as the source speech on MuST-C dataset (Di Gangi et al., 2019). All raw audio files are processed by Kaldi (Povey et al., 2011) to extract 80-dimensional log-Mel filterbanks stacked with 3-dimensional pitch features using a window size of 25 ms and step size of 10 ms. We train sentencepiece (Kudo & Richardson, 2018) models with a joint vocabulary size of 8K for each dataset. We remove samples that have more than 3000 frames for GPU efficiency. Our basic Transformer based E2E-ST framework has similar

[1]https://commonvoice.mozilla.org/en/datasets

|  | ST | ST+ASR | ST+MAM |
|---|---|---|---|
| # of parameters | 31M | 47M | 33M |

Table 1. MAM only has 6.5% more parameters than the baseline model while ASR multi-tasking needs to use 51.6% more parameters.

|  | MuST-C | Libri-Speech | Libri-Light | Common Voice | Audioset |
|---|---|---|---|---|---|
| Type | ♦★ | ♦ |  | ♦ |  |
| Hours | 408h | 960h | 3748h | 4421h | 6873h |

Table 2. The statistical results of corpora. ♦ and ★ denote the corpus has transcripts and translations, respectively. Note that although Common Voice has transcripts, we do not use them.

settings with ESPnet-ST(Inaguma et al., 2020). We first downsample the speech input with 2 layers of 2D convolution of size 3 with stride size of 2. Then there is a standard 12-layers Transformer with 2048 hidden size to bridge the source and target side. We only use 4 attention heads on each side of the transformer and each of them has a dimensionality of 256. For MAM module, we simply linearly project the outputs of the Transformer encoder to another latent space, then upsample the latent representation with 2-layers deconvolution to match the size of the original input signal. For the random masking ratio $\lambda$, we choose 30% across all the experiments including pre-training. During inference, we do not perform any masking over the speech input. We average the last 5 checkpoints for testing. For decoding, we use a beam search with setting beam size and length penalty to 5 and 0.6, respectively.

Our MAM is very easy to replicate as we do not perform any parameters and architecture search upon the baseline system. Due to the simple, but effective design of MAM, MAM does not rely on intensive computation. It converges within 2 days of training with 8 1080Ti GPUs for the basic model. We showcase the comparison of parameters between different solutions to E2E-ST in Table. 1 This makes a big difference with current popular intensive computations frameworks such as BERT(Devlin et al., 2019) (340M parameters) and GPT3(Brown et al., 2020) (175B parameters), making this technique is accessible to regular users.

### 4.1 Analyzing ASR and MAM

Aside from the extra training signal that is introduced by transcriptions, there is a deeper reason why ASR and MAM are beneficial to E2E-ST. In this section, we first discuss the difficulties and challenges in E2E-ST. Then we analyze the reasons why ASR MTL and MAM are helpful for E2E-ST by visualizing the self-attention over source side encoder.

Compared with other tasks, e.g., MT or ASR, which also employ Seq2Seq framework for E2E training, E2E-ST is a more difficult and challenging task in many ways. Firstly,
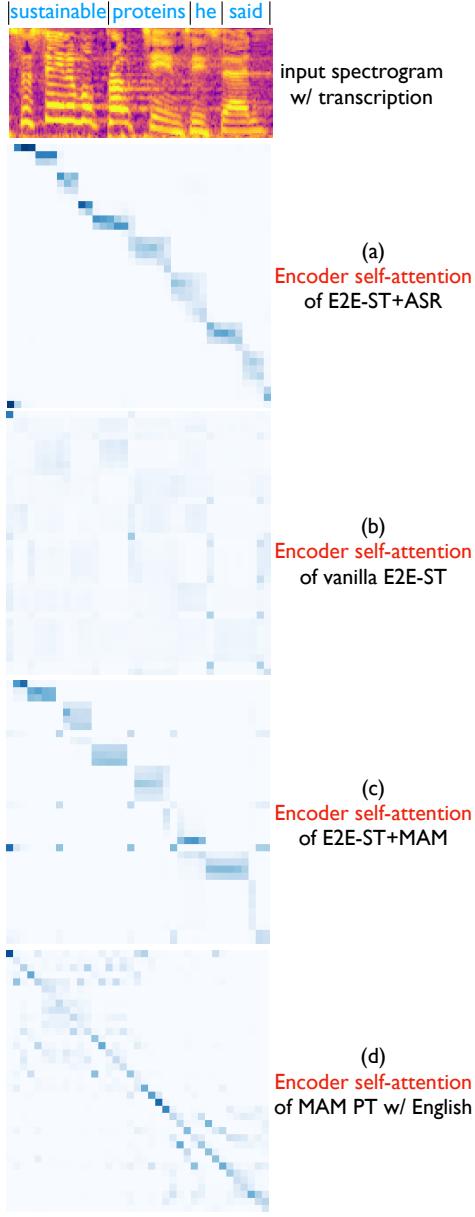
*Figure 3.* One head of the last layer self-attention comparison between different models. ASR MTL and MAM help the encoder learns similar self-attentions. See detailed discussion in 4.1.

data modalities are different on the source and target sides. For ST, the encoder deals with speech signals and tries to learn word presentations on the decoder side, while MT has text format on both sides. Secondly, due to the nature of the high sampling rate of speech signals, speech inputs are generally multiple (e.g. 4 to 7) times longer than the target sequence, which increases the difficulties of learning the correspondence between source and target. Thirdly, compared with the monotonicity natural of the alignment of ASR, ST usually needs to learn the global reordering

between speech signal and translation, and this raises the difficulties to another level. Especially in ST, since source and target are in different languages, it is very challenging to obtain the corresponding phoneme or syllable segments given the training signal from a different language.
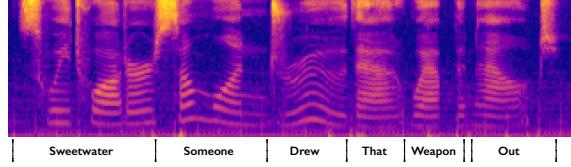
Fig. 3 tries to explain and analyze the difference between E2E-ST (a) and E2E-ST with ASR MTL (b). We extract the most top layer from the encoder for comparison. We notice that E2E-ST (a) tends to get more meaningful self-attention on the encoder with the training signal from ASR. With help from ASR, the source input spectrogram is chunked into segments that contain phoneme-level information. During training, the monotonicity natural of the ASR alignment functions as a forced alignment to group a set of adjacent frames to represent certain phonemes or syllables from source speech. With a larger scale of segmented spectrograms, the target side decoder only needs to perform reordering on those segments instead of frames. Our observations also align with the analysis from Stoian et al. (2020).

We also visualize the self-attention on encoder for E2E-ST with MAM (without pre-training) in (c) of Fig. 3. We find that MAM has the similar ability with ASR to segment the source speech into chunks. As it is shown in (d) of Fig. 3, when we only perform pre-training on the English speech (Libri-Light dataset), without E2E-ST training, self-attentions that are generated by pre-trained MAM are mostly monotonic on source side. Recovering local frames usually needs the information from surrounding context, especially for the speaker and environment-related characteristic. But we still observe that self-attention sometimes focuses on longer distance frames as well. This type of attention is very similar with low to mid layer self-attention of ASR. When there is a down streaming task (e.g., ASR or ST) is used for fine tuning, the top layer's self-attention will get chunked attention which is similar to (a) and (c).
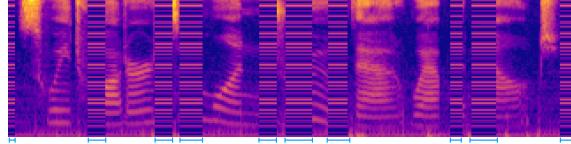
To conclude, we observe that MAM functions very similar to ASR on the encoder side. Hence, MAM is a reliable framework that can be used as an alternative solution when there is no transcription available. Especially, with the help of a large scale acoustic dataset, which does not have transcription annotation, MAM provides the E2E-ST a much better encoder initialization.

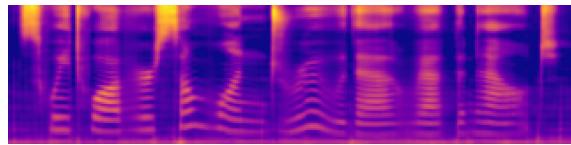### 4.2 Visualizing Reconstruction

To demonstrate what MAM has learned from pre-training step, we first showcase the reconstruction ability of MAM by visualizing the differences of spectrograms between the original and recovered inputs. This experiment was conducted on two corpora, Libri-Light and the Free Music Archive (FMA) (Defferrard et al., 2016) dataset. We use
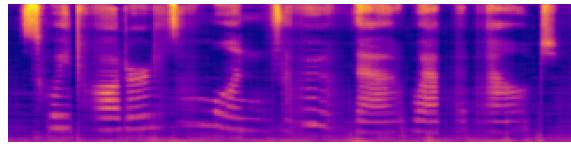
(a) The original speech spectrogram. Note that though we annotate the transcription underneath, we do not use transcription information at all during pre-training.



(b) We mask the selected frames (underlined with blue lines) with the same random initialized vector.



(c) Recovered spectrogram with MAM, pre-trained with Libri-Light corpus.



(d) MAM that pre-trains with FMA music corpus still have the ability to reconstruct corrupted speech signal.

*Figure 4.* One speech example to showcase the reconstruction ability of pre-trained MAM. We notice that MAM reconstructs the corrupted audio signal in both pre-training with ordinary speech and music dataset.

the "fma-medium" setting [2] which contains about 25,000 tracks of 30 seconds music within 16 unbalanced genres. The total music length is about 208 hours. We use FMA dataset for reconstruction visualization since FMA only contains music data and the characteristic of the music signal is very different from pure human speech. Note that our reconstructed spectrograms are a little blur compared with the original input since there are some downsampling steps in the E2E-ST baseline framework.

To verify the pre-trained results of MAM, we demonstrate the reconstruction ability of MAM by visualizing the results in Fig. 4. We show the original spectrogram of input speech in Fig. 4a. Then we corrupted the original spectrogram by replacing the selected mask frames with $\epsilon$, which is a random initialized vector, to form $\hat{x}$ (see Fig. 4b). In Fig. 4c, we show that our proposed MAM is able to recover the missing segments of input speech by pre-training over the

Libri-Light dataset. More interestingly, since MAM does not need any transcription to perform pre-training, we also pre-train MAM with FMA dataset. Surprisingly, as shown in Fig. 4d, MAM performs very similar reconstruction ability compared with the one that is pre-trained with speech dataset considering the corrupted audio is only about pure speech. This might be because some music tracks include human singing voices and MAM learns human speech characteristics from those samples though human singing voice can be quite different from speech. We also conduct reconstruction with speech pre-trained MAM for corrupted FMA data (see Fig. A1 in Appendix).

### 4.3 Translation Accuracy Comparisons

We showcase the translation accuracy of MAM comparing against to 6 baselines from Table 3 to Table 5:

- **Cascade**: cascade framework first transcribes the speech into transcription then passes the results to later machines translation system.
- **MT with ASR annotation**: an MT system which directly generates the target translation from the human-annotated transcription.
- **E2E-ST**: this is the vanilla translation system which does not use transcriptions in MuST-C.
- **E2E-ST + ASR MTL**: ST trained with ASR MTL using the transcription in MuST-C.
- **ST + SpecAugment**: a data augmentation method (Park et al., 2019; Bahar et al., 2019) by performing random masking over input speech.
- **E2E-ST + ASR PT**: the encoder of ST is initialized by a pre-trained ASR encoder which is trained from the speech and transcription pairs in Libri-Speech (Panayotov et al., 2015).

To better make a conclusion of our results from Table 3 to Table 5, we organize the comparisons as follows.

### 4.3.1 Comparison in the settings without transcription and pre-training

In Table 3, we first compare MAM against E2E-ST where there is no transcription and pre-training. Both MAM with single and span masking methods achieve averagely +0.44 (single) and +1.09 (span) improvements in BLEU score correspondingly against to E2E-ST in 8 different translation directions. Span masking consistently outperforms single frame masking as it is a more difficult self-supervised task.

### 4.3.2 Comparison in the pre-training settings without transcription

In Table 4, we have three different pre-training settings for MAM, which are pre-training with English speech (same

|  | Models | De | Es | Fr | It | Nl | Pt | Ro | Ru | Avg. Δ |
|---|---|---|---|---|---|---|---|---|---|---|
| Baselines | MT with ASR annotation (Di Gangi et al., 2019) | 28.09 | 34.16 | 42.23 | 30.40 | 33.43 | 32.44 | 28.16 | 18.30 | |
| | Cascaded methods (Inaguma et al., 2020) | 23.65 | 28.68 | 33.84 | 24.04 | 27.91 | 29.04 | 22.68 | 16.39 | |
| | E2E-ST | 19.64 | 23.68 | 28.91 | 19.95 | 23.01 | 24.00 | 21.06 | 12.05 | - |
| | ST + SpecAug | 20.06 | 24.51 | 29.26 | 20.27 | 23.73 | 24.40 | 21.21 | 12.84 | +0.49 |
| Ours | MAM (single) | 20.34 | 24.46 | 29.18 | 19.52 | 23.81 | 24.56 | 21.37 | 12.57 | +0.44 |
| | MAM (span) | 20.78 | 25.34 | 30.26 | 20.51 | 24.46 | 24.90 | 21.62 | 13.14 | **+1.09** |

*Table 3.* Comparisons between MAM and other baselines over 8 languages on MuST-C. In this setting, we use MAM as an extra training module for E2E-ST, and there is no pre-training involved. We notice that MAM with span masking achieves better performance and there is 1.09 BLEU score improvements upon E2E-ST. The column starts with "Avg. Δ" summarizes the average improvements upon baseline method, E2E-ST. See more discussions in Sec. 4.3.1.

|  | Models | De | Es | Fr | It | Nl | Pt | Ro | Ru | Avg. Δ |
|---|---|---|---|---|---|---|---|---|---|---|
| Baselines | E2E-ST | 19.64 | 23.68 | 28.91 | 19.95 | 23.01 | 24.00 | 21.06 | 12.05 | - |
| | E2E-ST+ASR PT* | 20.75 | 25.57 | 30.75 | 20.62 | 24.31 | 25.33 | 22.50 | 14.24 | †+1.47 |
| | E2E-ST+ASR MTL | 21.70 | 26.83 | 31.36 | 21.45 | 25.44 | 26.52 | 23.71 | 14.54 | †+2.41 |
| Ours | MAM w/ English PT | 21.44 | 26.48 | 31.21 | 21.28 | 25.22 | 26.41 | 23.83 | 14.53 | **+2.26** |
| | MAM w/ multilingual PT | 21.02 | 25.93 | 30.62 | 21.05 | 24.87 | 25.64 | 22.94 | 13.90 | +1.71 |
| | MAM w/ acoustic PT | 20.81 | 25.85 | 30.48 | 20.52 | 24.81 | 25.46 | 22.90 | 13.83 | +1.55 |

*Table 4.* Comparisons between span-based MAM pre-training with different pre-training corpus and other baselines over 8 languages on MuST-C. PT is short for pre-training. * denotes pretrained with Librispeech corpus. We use Libri-Light for English pre-training, Common Voice for multi-lingual pre-training and Audioset for arbitrary acoustic data pre-training. The methods denote with † use transcription in pre-training or MTL, but all our MAM methods do not use transcription. MAM pre-training with English corpus achieves very similar performance with E2E-ST+ASR MTL. The column starts with "Avg. Δ" summarizes the average improvements upon baseline method. See more discussions in Sec. 4.3.2.

with the source language) from Libri-Light, multilingual speech data from Common Voice, and arbitrary acoustic data from Audioset corpus. Among those methods, MAM pre-trained with Libri-Light achieves the best results as it consistently outperforms the baseline. Averagely speaking, there is +2.26 improvements compared with E2E-ST. When we compare to "E2E-ST+ASR PT", there are about +0.79 improvements in BLEU score across 8 target languages.

Surprisingly, MAM trained with acoustic data still achieves about +1.55 improvements upon E2E-ST. Considering acoustic data does not need any annotation and this kind of dataset is much easier to collect, the results are very encouraging. With the help of vast acoustic data on the website (e.g., youtube), MAM trained with arbitrary acoustic data has great potential to further boost the performance. To the best of our knowledge, MAM is the first technique that performs pre-training with any form of the audio signal.

MAM trained with Common Voice does not have significant improvements with two following reasons: firstly, speech audios in Common Voice sometimes are very short (about 2 to 3 seconds) while MuST-C usually contains much longer speech (above 10 seconds) leading to very limited options for random masking; secondly there are much fewer English speech in this corpus.

### 4.3.3 Comparison in the settings using transcription

In this setting, we use "E2E-ST + ASR MTL" as the baseline. MAM MTL with pre-training over Libri-Light achieves +0.9 average improvements over 8 languages.

### 4.3.4 Comparison to Wav2vec

We also compare MAM against to other wav2vec-based methods (Wu et al., 2020) in Table 6. Due to the differences in baseline methods, to make a fair comparison, we only compare the relative improvements upon our own baseline on the same test data. MAM still achieves much larger improvements upon a much stronger baseline. Especially our baseline is already about 4 BLEU points better than the baseline in wav2vec, MAM still achieves +1.6 more BLEU points improvements compared with wav2vec-based pre-training methods making our performance on En-Ro about 5.6 BLEU points better than wav2vec-based methods.

### 4.4 Comparisons in Low and Mid-resource Settings

In Table 7, we reduce the size of MuST-C from 408 hours to 50 hours and 200 hours to mimic the low and mid-resource language speech translation.

In the scenario when the source language is extremely low-resource (no transcribed pre-training and fine-tuning data), we have "E2E-ST" as the baseline. MAM in both multi-

| Models | De | Es | Fr | It | Nl | Pt | Ro | Ru | Avg. Δ |
|---|---|---|---|---|---|---|---|---|---|
| E2E-ST+ASR MTL | 21.70 | 26.83 | 31.36 | 21.45 | 25.44 | 26.52 | 23.71 | 14.54 | - |
| MAM+ASR MTL | 22.41 | 26.89 | 32.55 | 22.12 | 26.49 | 27.22 | 24.45 | 14.90 | +0.69 |
| MAM w/ English PT + ASR MTL | 22.87 | 26.86 | 32.80 | 22.12 | 26.81 | 27.43 | 24.65 | 15.21 | **+0.90** |

*Table 5.* Comparisons between MAM with ASR MTL and E2E-ST with ASR MTL. MAM still achieves an improvement about +0.9 BLEU. The column starts with "Avg. Δ" summarizes the average improvements upon baseline method.

| Wav2vec-based Method (Wu et al., 2020) | | | | |
|---|---|---|---|---|
| | En-Fr | Δ | En-Ro | Δ |
| their baseline[†] | 27.8 | - | 17.1 | - |
| + wav2vec PT[†] | 29.8 | +2.0 | 18.2 | +1.1 |
| + vq-wav2vec PT[†] | 28.6 | +0.8 | 17.4 | +0.3 |
| MAM-based Method | | | | |
| our baseline | 28.9 | - | 21.1 | - |
| + MAM w/ English PT | 31.2 | **+2.3** | 23.8 | **+2.7** |

*Table 6.* Comparisons between wav2vec-based pre-train method for E2E-ST. Results that are decorated with [†] are from Wu et al. (2020). Our relative improvements over baseline methods are much larger than wav2vec-based pre-training methods. See more discussions in Sec. 4.3.4.

lingual and acoustic pre-training boosts the performance significantly.

When the transcription is only available at the fine-tuning stage (compare with "+ASR MTL"), MAM pre-trained with Libri-Light achieves similar performance without using transcription in fine-tuning.

In the cases when there is no transcription in the fine-tuning stage, but there exist a large scale annotated pre-training corpus, MAM still achieves similar performance in 200 hours training setting without using any transcription.

| Models | Fr | | Es | |
|---|---|---|---|---|
| | 50h | 200h | 50h | 200h |
| E2E-ST | 0.52 | 19.83 | 0.4 | 16.13 |
| E2E-ST+ASR MTL | 8.84 | 25.64 | 7.67 | 20.21 |
| E2E-ST+ASR PT* | 12.50 | 24.59 | 11.80 | 19.35 |
| MAM | 0.6 | 20.54 | 0.4 | 16.85 |
| MAM w/ English PT | 6.84 | 24.86 | 6.53 | 19.17 |
| MAM w/ acoustic PT | 3.29 | 22.22 | 2.46 | 17.98 |

*Table 7.* Experimental comparisons with difference training resource. * denotes pretrained with Librispeech corpus. See Section 4.4 for detailed discussion.

## 5 Related Work

Text-based BERT-style (Devlin et al., 2019; Liu et al., 2019a; Joshi et al., 2020a; Zhang et al., 2019) frameworks are extremely popular in recent years due to the remarkable improvement that they bring to the downstream tasks at fine-tuning stages. Inspired by techniques mentioned above, MAM also performs self-supervised training that masks cer-

tain portions randomly over the input signals. But different from BERT-style pre-training, MAM tries to recover the missing semantic information (e.g., words, subword units) and learns the capabilities to restore the missing speech characteristics and generate the original speech.

SpecAugment (Park et al., 2019) was originally proposed for ASR as a data augmentation method by applying a mask over input speech, then it is adapted to ST by Bahar et al. (2019). However, there is no recovering step in SpecAugment, and it can not be used as a pre-training framework.

For the self-supervised training in speech domain, Chuang et al. (2019); Wang et al. (2020c;a) proposed to use forced-alignment to segment speech audio into pieces at word level and masked some certain words during fine-tuning. Obviously the forced-alignment based approaches rely on the transcriptions of source speech, and can not be applied to low or zero resource source speech while MAM will relief the needs of large-scale, annotated speech and translation pairs during pre-training.

Baevski et al. (2020) proposed wav2vec 2.0 pre-training model for ASR task which masks the speech input in the latent space and pre-trains the model via contrastive task defined over a quantization of the latent representations. In contrast, as the objective of MAM is much simpler and straightforward, we don't need much extra fine-tuning efforts given an E2E-ST baseline and massive computational resource. Furthermore, wav2vec 2.0 is build upon discretized, fix-size, quantized codebooks, and it is not easy to be adapted to arbitrary acoustic signal pre-training. Lastly, wav2vec 2.0 is designed to have ASR as the downstream task, and the fine-tuning stage relies on CTC loss (Graves et al., 2006) which is not straightforward to be adapted in translation task since translation usually involves with many reordering between target and source side while CTC depends on monotonic transition function on source side.

## 6 Conclusions

We have presented a novel acoustic modeling framework MAM in this paper. MAM not only can be used as an extra component during training time, but also can be used as a separate pre-training framework with arbitrary acoustic signal. We demonstrate the effectiveness of MAM with multiple different experiment settings in 8 languages. Es-

pecially, for the first time, we show that pre-training with arbitrary acoustic data with MAM boosts the performance of speech translation.

# References

Anastasopoulos, A. and Chiang, D. Tied multitask learning for neural speech translation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2018.

Anastasopoulos, A., Chiang, D., and Duong, L. An unsupervised probability model for speech-to-translation alignment of low-resource languages. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016.

Ardila, R., Branson, M., Davis, K., Henretty, M., Kohler, M., Meyer, J., Morais, R., Saunders, L., Tyers, F. M., and Weber, G. Common voice: A massively-multilingual speech corpus. 2020.

Baevski, A., Zhou, H., Mohamed, A., and Auli, M. wav2vec 2.0: A framework for self-supervised learning of speech representations. *NeurIPS 2020*, 2020.

Bahar, P., Zeyer, A., Schlüter, R., and Ney, H. On using specaugment for end-to-end speech translation. 2019.

Bansal, S., Kamper, H., Livescu, K., Lopez, A., and Goldwater, S. Pre-training on high-resource speech recognition improves low-resource speech-to-text translation. 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)*, 2019.

Berard, A., Pietquin, O., Servan, C., and Besacier, L. Listen and translate: A proof of concept for end-to-end speech-to-text translation. *CoRR*, abs/1612.01744, 2016. URL http://arxiv.org/abs/1612.01744.

Berard, A., Besacier, L., Kocabiyikoglu, A., and Pietquin, O. End-to-end automatic speech translation of audiobooks. *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6224–6228, 2018.

Bird, S. A scalable method for preserving oral literature from small languages. In Chowdhury, G., Koo, C., and Hunter, J. (eds.), *The Role of Digital Libraries in a Time of Global Change*. Springer Berlin Heidelberg, 2010.

Bird, S., Gawne, L., Gelbart, K., and McAlister, I. Collecting bilingual audio in remote indigenous communities. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, 2014.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. 2020.

Chuang, Y.-S., Liu, C.-L., Lee, H.-Y., and Lee, L.-s. Speech-BERT: An Audio-and-text Jointly Learned Language Model for End-to-end Spoken Question Answering. *arXiv e-prints*, 2019.

Defferrard, M., Benzi, K., Vandergheynst, P., and Bresson, X. Fma: A dataset for music analysis. *arXiv preprint arXiv:1612.01840*, 2016.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.

Di Gangi, M. A., Cattoni, R., Bentivogli, L., Negri, M., and Turchi, M. MuST-C: a Multilingual Speech Translation Corpus. In *NAACL*, 2019.

Duong, L., Anastasopoulos, A., Chiang, D., Bird, S., and Cohn, T. An attentional model for speech translation without transcription. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2016.

Ethnologue. Ethnologue (21st edition). URL https://www.ethnologue.com/enterprise-faq/how-many-languages-world-are-unwritten-0.

Gangi, M. A. D., Negri, M., and Turchi, M. Adapting Transformer to End-to-End Spoken Language Translation. In *Proc. Interspeech 2019*, 2019.

Gemmeke, J. F., Ellis, D. P. W., Freedman, D., Jansen, A., Lawrence, W., Moore, R. C., Plakal, M., and Ritter, M. Audio set: An ontology and human-labeled dataset for audio events. In *Proc. IEEE ICASSP 2017*, 2017.

Graves, A., Fernández, S., Gomez, F., and Schmidhuber, J. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd International Conference on Machine Learning*, 2006.

Inaguma, H., Kiyono, S., Duh, K., Karita, S., Soplin, N. E. Y., Hayashi, T., and Watanabe, S. Espnet-st: All-in-one speech translation toolkit. *arXiv preprint arXiv:2004.10234*, 2020.
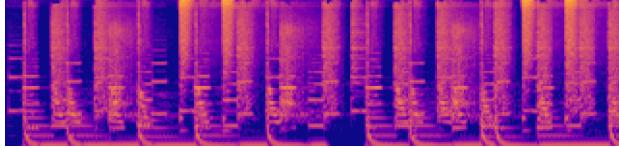
Joshi, M., Chen, D., Liu, Y., Weld, D. S., Zettlemoyer, L., and Levy, O. SpanBERT: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8, 2020a.

Joshi, M., Chen, D., Liu, Y., Weld, D. S., Zettlemoyer, L., and Levy, O. Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8:64–77, 2020b.

Kahn, J., Rivière, M., Zheng, W., Kharitonov, E., Xu, Q., Mazaré, P. E., Karadayi, J., Liptchinsky, V., Collobert, R., Fuegen, C., Likhomanenko, T., Synnaeve, G., Joulin, A., Mohamed, A., and Dupoux, E. Librilight: A benchmark for asr with limited or no supervision. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7669–7673, 2020. https://github.com/facebookresearch/libri-light.

Kudo, T. and Richardson, J. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Brussels, Belgium, November 2018. Association for Computational Linguistics.

Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692, 2019a. URL http://arxiv.org/abs/1907.11692.

Liu, Y., Xiong, H., Zhang, J., He, Z., Wu, H., Wang, H., and Zong, C. End-to-End Speech Translation with Knowledge Distillation. In *Proc. Interspeech 2019*, 2019b.

Mathias, L. and Byrne, W. Statistical phrase-based speech translation. In *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings*, 2006.

Matusov, E., Kanthak, S., and Ney, H. On the integration of speech recognition and statistical machine translation. In *INTERSPEECH*, 2005.

McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., and Sonderegger, M. Montreal forced aligner: Trainable text-speech alignment using kaldi. In *Proc. Interspeech 2017*, pp. 498–502, 2017. doi: 10.21437/Interspeech.2017-1386. URL http://dx.doi.org/10.21437/Interspeech.2017-1386.

Ney, H. Speech translation: coupling of recognition and translation. In *1999 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings.*

*ICASSP99 (Cat. No.99CH36258)*, volume 1, pp. 517–520 vol.1, 1999.

Panayotov, V., Chen, G., Povey, D., and Khudanpur, S. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5206–5210. IEEE, 2015.

Park, D. S., Chan, W., Zhang, Y., Chiu, C.-C., Zoph, B., Cubuk, E. D., and Le, Q. V. Specaugment: A simple data augmentation method for automatic speech recognition. In *Interspeech 2019*, 2019.

Povey, D., Ghoshal, A., Boulianne, G., Goel, N., Hannemann, M., Qian, Y., Schwarz, P., and Stemmer, G. The kaldi speech recognition toolkit. In *In IEEE 2011 workshop*, 2011.

Sperber, M., Neubig, G., Niehues, J., and Waibel, A. Attention-Passing Models for Robust and Data-Efficient End-to-End Speech Translation. *Transactions of the Association for Computational Linguistics (TACL)*, 2019. URL https://arxiv.org/abs/1904.07209.

Stoian, M., Bansal, S., and Goldwater, S. Analyzing asr pretraining for low-resource speech-to-text translation. In *ICASSP*, 2020.

Vila, L. C., Escolano, C., Fonollosa, J. A. R., and Costa-jussà, M. R. End-to-end speech translation with the transformer. In *IberSPEECH*, 2018.

Wang, C., Wu, Y., Du, Y., Li, J., Liu, S., Lu, L., Ren, S., Ye, G., Zhao, S., and Zhou, M. Semantic mask for transformer based end-to-end speech recognition. In *Interspeech*, 2020a.

Wang, C., Wu, Y., Liu, S., Yang, Z., and Zhou, M. Bridging the gap between pre-training and fine-tuning for end-to-end speech translation. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, 2020b.

Wang, C., Wu, Y., Liu, S., Zhou, M., and Yang, Z. Curriculum pre-training for end-to-end speech translation. In *ACL*, 2020c.

Watanabe, S., Hori, T., Kim, S., Hershey, J. R., and Hayashi, T. Hybrid ctc/attention architecture for end-to-end speech recognition. 2017.

Weiss, R. J., Chorowski, J., Jaitly, N., Wu, Y., and Chen, Z. Sequence-to-sequence models can directly translate foreign speech. In *Proc. Interspeech 2017*, 2017.

Wu, A., Wang, C., Pino, J., and Gu, J. Self-supervised representations improve end-to-end speech translation. In *Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020*, 2020.
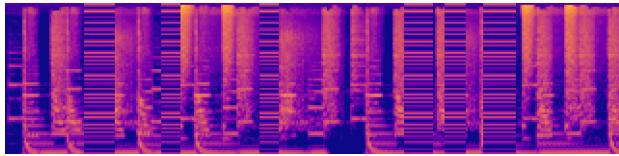
Zhang, Z., Han, X., Liu, Z., Jiang, X., Sun, M., and Liu, Q. ERNIE: enhanced language representation with informative entities. 2019.
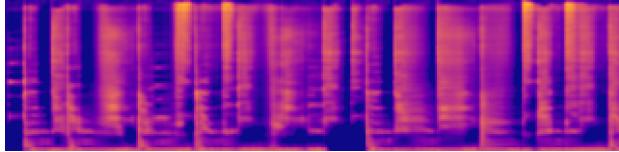
# Appendix

In the other way around, we also try to use Libri-Light pretrained MAM to recover the corrupted music in Fig. A1. MAM that pre-trained with human speech data does not show good reconstruction in Fig. A1c since there are many different musical instruments' sounds that are unseen in speech data.



(a) The original musical spectrogram that is mixed with different instruments' sound.



(b) We mask the selected frames (underlined with blue lines) with the same random initialized vector.



(c) Recovered spectrogram with MAM, pre-trained with Libri-Light corpus.

*Figure A1.* One speech example to showcase the reconstruction ability of pre-trained MAM. Pre-trained MAM with Libri-Light corpus (only human speech data) can not reconstruct the original music spectrogram accurately since there are many different musical instruments' sound that is unseen in speech data.