

# Unified Speech-Text Pre-training for Speech Translation and Recognition

Anonymous ACL submission

## Abstract

In this work, we describe a method to jointly pre-train speech and text in an encoder-decoder modeling framework for speech translation and recognition. The proposed method utilizes multi-task learning to integrate four self-supervised and supervised subtasks for cross modality learning. A self-supervised speech subtask, which leverages unlabelled speech data, and a (self-)supervised text to text subtask, which makes use of abundant text training data, take up the majority of the pre-training time. Two auxiliary supervised speech tasks are included to unify speech and text modeling space. Detailed analysis reveals learning interference among subtasks. In order to alleviate the subtask interference, two pre-training configurations are proposed for speech translation and speech recognition respectively. Our experiments show the proposed method can effectively fuse speech and text information into one model. It achieves between 1.7 and 2.3 BLEU improvement above the state of the art on the MUST-C speech translation dataset and comparable WERs to wav2vec 2.0 on the LIBRISPEECH speech recognition task.

## 1 Introduction

Pre-training can learn universal feature representations from a large training corpus and is beneficial for downstream tasks with limited amounts of training data (Peters et al., 2018; van den Oord et al., 2018; Chung et al., 2018; Zoph et al., 2020). With the advancement of computational power and self-supervised pre-training approaches, large volumes of unlabeled data may now be used in pre-training. Methods, such as BERT (Devlin et al., 2019), BART (Lewis et al., 2020b) and wav2vec2.0 (Baevski et al., 2020b), have emerged as the backbone of many speech and natural language processing tasks.

The aforementioned pre-training methods focus on learning feature representation either from text

or speech. Many speech applications combine information learnt from both speech and text corpora to achieve state of the art results. In speech processing, transcribed speech training data is generally very scarce for many languages. It is difficult to build robust linguistic knowledge representation solely based on labeled speech training data. Jia et al. (2019); Chen et al. (2021) propose to generate synthetic data from text to augment speech training corpora corpus. Li et al. (2021) demonstrate that models initialized with pre-trained wav2vec2.0 and mBART (Liu et al., 2020) modules are competitive for the multilingual speech to text translation task. Chuang et al. (2020) propose to concatenate the acoustic model and BERT model for speech Q&A. Chung et al. (2021b) align speech utterance representation to the corresponding text sentence representation, in which both representations are generated from unsupervised pre-trained models, for speech understanding.

In this study, we are interested in pre-training for speech to text tasks based on the Encoder-Attention-Decoder (EAD) framework. In particular, we seek to answer the question whether the integration of data from different modalities is beneficial for representation learning. To answer this question, we propose Speech and Text joint Pre-Training (STPT), a multi-task learning framework to combine different modalities, i.e., speech and text, in the pre-training stage. A self-supervised speech subtask and a (self-)supervised text to text subtask dominate the pre-training computation to leverage large amounts of unlabelled speech data and abundant text training corpus. Two auxiliary supervised speech subtasks are used to unify different modalities in the same modeling space. The proposed method fuses information from the text and speech training corpus into a single model, and it effectively improves the performance of downstream tasks, such as speech to text translation (ST) and automatic speech recognition (ASR). Our con-

084 tributions are summarized as follows:

- 085 1. We propose a multi-task learning framework  
086 to jointly pre-train speech and text in one  
087 model.
- 088 2. We conduct detailed analyses on the proposed  
089 pre-training method, which reveal the interfer-  
090 ence among different subtasks.
- 091 3. Two joint pre-training configurations are pro-  
092 posed to alleviate learning interference for  
093 ASR and ST respectively.
- 094 4. State-of-the-art results are achieved on the  
095 downstream tasks. We obtain at least 1.7  
096 BLEU improvement compared with the best  
097 MUST-C ST system reported and comparable  
098 WERs as wav2vec 2.0 in the LIBRISPEECH  
099 ASR task.

## 100 2 Related work

101 **Pre-training:** Self-supervised pre-training is usu-  
102 ally optimized with two different criteria: con-  
103 trastive loss (van den Oord et al., 2018; Chung and  
104 Glass, 2020; Baevski et al., 2020b) and masked  
105 prediction loss (Devlin et al., 2019). Contrastive  
106 loss focuses on distinguishing the positive samples  
107 from the negative ones given the reference sam-  
108 ple and it has achieved great success for speech  
109 recognition (Baevski et al., 2020b). Masked predic-  
110 tion loss has been first studied for natural language  
111 processing tasks (Devlin et al., 2019; Lewis et al.,  
112 2020b) with subsequent application to speech pro-  
113 cessing (Baevski et al., 2020a; Hsu et al., 2021).  
114 Chung et al. (2021a) combine contrastive loss and  
115 masked prediction loss, which shows good perfor-  
116 mance for the downstream ASR task. The opti-  
117 mization of our self-supervised speech task is more  
118 related to the masked prediction loss. Instead of  
119 predicting the hard discretized label for the masked  
120 frames, which is error prone, we use KL divergence  
121 to minimize the distribution difference between the  
122 same feature frames with and without masking.  
123 Please refer to subsection 3.2 for more details.

124 **Self-training (or iterative pseudo labelling):**  
125 self-training is another widely used approach to  
126 take advantage of unlabelled speech data to im-  
127 prove the ASR performance (Kahn et al., 2020; Xu  
128 et al., 2020; Pino et al., 2020; Zhang et al., 2020;  
129 Wang et al., 2021a; Xiao et al., 2021; Wang et al.,  
130 2021b). A seed model, which usually is trained

with a small amount of supervised speech train-  
ing data, is employed to generate pseudo labels  
for the unlabelled speech data. The speech data  
with pseudo labels is augmented into the training  
dataset to build another model, which is expected  
to outperform the seed model due to more train-  
ing data exposure. Similar to self-training, we also  
use small amounts of supervised data to unify the  
speech and text modeling space. However, the  
self-supervised speech training in this work avoids  
making hard predictions and uses KL divergence to  
maximize the mutual information between masked  
span and observed feature frames.

**Multi-task learning:** Due to data scarcity, multi-  
task learning is widely adopted to leverage parallel  
text training data for ST (Weiss et al., 2017; Anasta-  
sopoulos and Chiang, 2018; Tang et al., 2021b; Ye  
et al., 2021). Those methods primarily use super-  
vised speech data sets during multi-task learning,  
whereas our method can leverage large amounts of  
unlabeled speech data during the pre-training stage,  
which has the potential to improve performance  
even further.

A concurrent work from Ao et al. (2021) also  
proposes to jointly pre-train speech and text for  
ASR and text to speech application, which is fully  
unsupervised. Our method focuses on taking ad-  
vantage of the supervised speech data, which is the  
same data used for fine-tuning, to improve the joint  
speech text pre-training, and the results demon-  
strate the efficacy of supervised speech data in pre-  
training. Another concurrent work is from Bapna  
et al. (2021), which focuses on speech encoder  
pre-training using both speech and text data. Our  
method emphasizes the encoder-decoder frame-  
work and training both encoder and decoder in the  
pre-training stage.

## 168 3 Method

ASR and ST are the two main downstream tasks for  
the proposed pre-training method. Figure 1 depicts  
our joint pre-training framework, which consists of  
four subtasks:

- 173 1. (Self-)supervised text to text subtask (T2T)
- 174 2. Self-supervised speech learning subtask  
175 (SSL)
- 176 3. Supervised speech phoneme classification  
177 subtask (S2P)

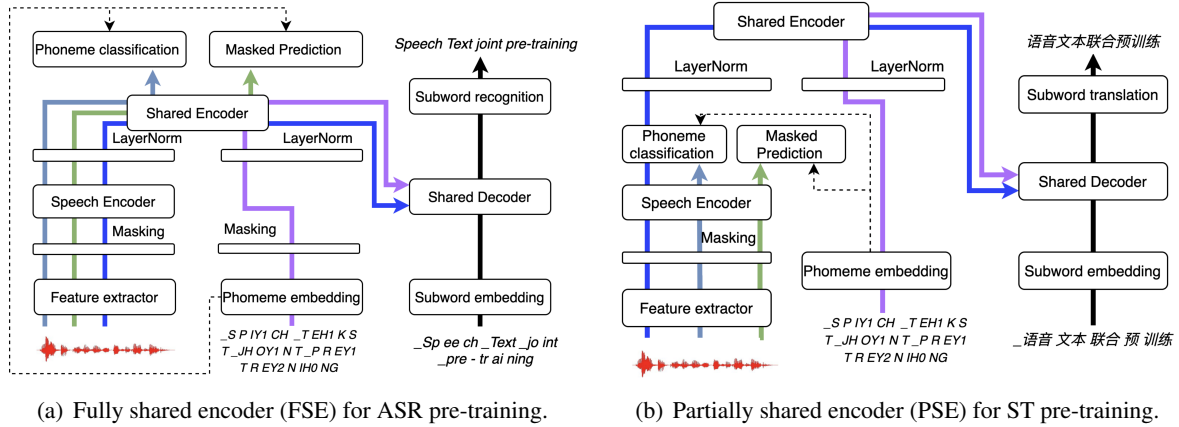


Figure 1: Speech text joint pre-training framework. The purple, green, steelblue and blue lines depict the encoders for the T2T, SSL, S2P and S2T subtasks respectively. The black lines show the decoder model for the T2T and S2T subtasks. The dotted lines indicate the phoneme embedding is applied in the SSL and S2P subtasks.

#### 4. Supervised EAD based speech to text subtask, which is the same as the downstream task, i.e., ST or ASR (S2T)

The choice of the T2T subtask depends on the downstream task. For ASR, the T2T subtask is a denoising autoencoder task (BART) (Lewis et al., 2020a) while ST utilizes a text based neural machine translation task. The SSL subtask is a self-supervised speech learning task to leverage large amounts of unlabelled speech data optimized by the masked prediction loss. The last two supervised speech tasks (S2P and S2T) are used to unify two modalities, i.e., speech and text, into one modeling space.

In this study, we find that the subtasks for the ASR pre-training are complementary, while interference is observed in subtasks of the ST pre-training at some encoder layers. We propose two different configurations: fully shared encoder (FSE) (Figure 1(a)) for the ASR pre-training, and partially shared encoder (PSE) (Figure 1(b)) for the ST pre-training. The FSE configuration aims to encourage information sharing between different subtasks while the PSE configuration tries to minimize the information sharing between encoder only subtasks, i.e., SSL and S2P, and sequence to sequence EAD tasks, i.e., subtask T2T and S2T. More subtask interference analysis is presented in subsection 5.2. We describe the details of each subtask in the following subsections.

### 3.1 (Self-)supervised text to text subtask

In the sequence to sequence ASR and ST tasks, the decoder is a text generator conditioned on the en-

coder outputs. Large amounts of training samples are required to cover different linguistic aspects of the target language. The abundant text corpus is an ideal supplement to the limited supervised speech data corpus. Assume the target text sequence is  $Y = (y_1, y_2, \dots, y_N)$ , its corresponding corrupted version,  $X = \text{NOISE}(Y) = (x_1, x_2, \dots, x_M)$ , can be created by masking or replacing token spans in  $Y$  (Lewis et al., 2020a) for the ASR pre-training. If the downstream task is ST,  $X$  is the corresponding source token sequence. The task is optimized by maximizing cross entropy  $\mathcal{L}_{T2T}$

$$\mathcal{L}_{T2T} = - \sum_i^N \log p(y_i | y_{1:i-1}, X) \quad (1)$$

In this subtask, we also convert the input text into the corresponding pronunciation form, i.e., phoneme sequence, as it would be easier to align the encoder outputs from speech and text. The purple and black lines in Figure 1 describe the encoder and decoder of the T2T subtask.

### 3.2 Self-supervised speech subtask

The SSL subtask aims to leverage vast amounts of unlabelled speech data and learn general speech representations. The model configuration follows wav2vec2.0 (Baevski et al., 2020b) where the speech model includes a feature extractor and a context encoder. The context encoder corresponds to the speech encoder in Figure 1(b) in the ST pre-training. If ASR is the downstream task, the context encoder includes one extra shared encoder as shown in Figure 1(a).

The SSL subtask is optimized via masked prediction loss and it consists of two-pass computation. Given the speech input  $S = (s_1, s_2, \dots, s_T)$ , the feature extractor and context encoder outputs are  $Z = (z_1, z_2, \dots, z_{T'})$  and  $O = (o_1, o_2, \dots, o_{T'})$  respectively, where the speech input is down-sampled by the feature extractor and  $T > T'$ . In the first pass, the output  $O$  is compared with the phoneme embedding  $E = (e_1, e_2, \dots, e_I)$ , which is from the T2T subtask described in subsection 3.1.  $I$  is the phoneme vocabulary size. The predicted phoneme distribution  $p(o_j|e_i)$  is defined as

$$p(o_j|e_i) = \frac{\exp(o_j^\top \cdot e_i)}{\sum_{i'} \exp(o_j^\top \cdot e_{i'})} \quad (2)$$

In the second pass, speech feature spans  $\hat{Z} \subset Z$  are selected and corrupted as wav2vec2.0 (Baevski et al., 2020b).  $\hat{O}$  is the corresponding context encoder output from  $\hat{Z}$ . We train the model to infer the corrupted  $p(\hat{o}_j|e_i)$  to be similar as  $p(o_j|e_i)$  by minimizing KL divergence.

$$\mathcal{L}_{\text{SSL}} = - \sum_{\hat{o}_j \in \hat{O}} \sum_i p(o_j|e_i) \log \frac{p(\hat{o}_j|e_i)}{p(o_j|e_i)} \quad (3)$$

### 3.3 Supervised speech phoneme classification

The S2P subtask is employed to unify the self-supervised trained speech and text models. It shares the same model as in the SSL subtask. In this subtask, a transcribed ASR data set is used and the goal of this task is to predict the frame level phoneme labels. A HMM-GMM model is trained with the same transcribed dataset using Kaldi (Povey et al., 2011) to generate the frame-level labels with forced-alignment. The phoneme classification task is optimized with the cross entropy loss

$$\mathcal{L}_{\text{S2P}} = - \sum_{o_j \in O} \log p(o_j|e_{a(j)}) \quad (4)$$

where  $a(j)$  is the phoneme label associated with the context encoder output  $o_j$ . The S2P subtask is depicted with steelblue lines in Figure 1.

### 3.4 Supervised EAD based speech to text subtask

Besides the S2P subtask mentioned in the previous subsection, we include the potential downstream EAD based task, i.e. ASR or ST, as another auxiliary subtask during the pre-training stage. In

many speech translation datasets, such as MuST-C (Gangi et al., 2019) or CoVoST (Wang et al., 2020), we have both speech transcription and translation labels. The speech transcription is used in the S2P subtask while the S2T subtask can make use of the corresponding translation labels. We hope this auxiliary task would make the transition from pre-training to fine-tuning smooth and result in better performance in the downstream task. The components involved during optimization include feature extractor, speech encoder, shared encoder (connected with blue lines in Figure 1), and decoder (black lines in Figure 1). They are trained with cross entropy criterion

$$\mathcal{L}_{\text{S2T}} = - \sum_t \log p(y_t|y_{t-1}, O) \quad (5)$$

where  $O$  is the input speech and  $Y = (y_1, \dots, y_N)$  is the target labels.

The overall pre-training loss is defined as the combination of four losses discussed above

$$\mathcal{L} = \mathcal{L}_{\text{T2T}} + \alpha \mathcal{L}_{\text{SSL}} + \beta \mathcal{L}_{\text{S2P}} + \gamma \mathcal{L}_{\text{S2T}} \quad (6)$$

where  $\alpha$ ,  $\beta$  and  $\gamma$  are task weights for the SSL, S2P and S2T subtasks respectively.

### 3.5 Implementation details

During the pre-training, the shared encoder inputs come from two sources, either from speech encoder outputs in the S2T subtask or phoneme embeddings in the T2T subtask. The shared encoder inputs might be in different numerical scales. In order to stabilize the multi-task training, a LayerNorm (Ba et al., 2016) is applied to the shared encoder inputs and places those inputs in the same numerical scale as shown in Figure 1.

The S2P and SSL subtasks rely on the phoneme embeddings from the T2T subtask. The optimization for the SSL and S2P subtasks are not stable at the early pre-training stage if the phoneme embeddings are randomly initialized. In our implementation, we pre-train the shared encoder and decoder with the T2T subtask. It helps the stabilization of the training and achieve a better feature representation. Also, the joint pre-training is time consuming since it takes care of four different subtasks. Pre-training encoder and decoder via the T2T subtask can make use of more text training data given the same training time.

## 4 Experimental setting

In the pre-training, we first train modules with the T2T subtask until they are converged. Then the entire model is jointly optimized with all subtasks mentioned in section 3. Finally, the pre-trained model is fine-tuned on the downstream tasks. Two downstream tasks, ASR and ST, are examined in this study. The training data for each subtask for pre-training and fine-tuning is described below.

### 4.1 Pre-training

**T2T:** For ASR pre-training, the language model (LM) training dataset<sup>1</sup> for LIBRISPEECH is used to build the monolingual BART model. It has about 800 million words. For ST pre-training, we take the parallel training corpus from WMT. We examine our methods on two translation directions in MUST-C: English-Spanish (EN-ES), which uses WMT13 training corpus, and English-French (EN-FR), which takes the WMT14 training data. There are 370 million and 1 billion English words in the EN-ES and EN-FR parallel training datasets respectively.

We use “g2p\_en” Python package (Lee and Kim, 2018) to convert the training text into the corresponding phoneme representation, which is based on the CMU English dictionary. We further extend the phoneme set by distinguishing the first phoneme in the word with an additional “\_” mark appended, which is similar to the notation in the SentencePiece process. The input phoneme vocabulary size is 134.

**SSL:** For both ASR and ST pre-training, 60k hours of unlabelled English speech data from Libri-light (Kahn et al., 2020) is used to build the self-supervised speech task. We set the maximum utterance duration to 37.5 seconds and minimum duration to 4 seconds. We randomly sample audio segments with maximum duration if utterances are longer than the maximum duration. No voice activity detection is applied.

**S2P:** We use the transcribed LIBRISPEECH dataset for ASR pre-training. In ST pre-training, the MUST-C training dataset is used, where the corresponding English transcription is used as the training target labels after it is converted into phoneme representation. The phoneme level segmentation is obtained via force-alignment, which is conducted using HMM/GMM trained from the same speech data with the Kaldi toolkit (Povey et al., 2011).

<sup>1</sup><https://www.openslr.org/11/>

**S2T:** We use the same labelled data in the S2P subtask for the S2T subtask, i.e., LIBRISPEECH training data for the ASR pre-training and MUST-C data for the ST pre-training. Instead of using phoneme representation, the target labels are encoded with SentencePiece (Kudo and Richardson, 2018). For both ASR and ST tasks, the vocabulary is an Unigram model with size 10k and full character coverage on the training text data.

### 4.2 Fine-tuning

In the fine-tuning stage, we keep optimizing the model with the T2T and S2T subtasks. Two encoder-only subtasks (SSL and S2P) are dropped, since the model has learnt good speech representation from the unlabeled speech data in the pre-training. The ASR system is evaluated on four LIBRISPEECH testsets: dev-clean, dev-other, test-clean and test-other. WER is reported in the experiments. ST models are evaluated on the tst-COMMON testset from MUST-C. Case-sensitive detokenized SACREBLEU (Post, 2018) is used to measure the ST performance.

### 4.3 Model configuration

The model takes raw speech audio as input. The feature encoder contains seven blocks and the temporal convolutions in each block have 512 channels with strides (5,2,2,2,2,2,2) and kernel widths (10,3,3,3,3,2,2). The speech encoder, shared encoder and shared decoder are all with 6 transformer layers, model dimension 768, inner dimension (FFN) 3,072 and 8 attention heads. We adopt Pre-LN in the transformer block as Xiong et al. (2020). The total number of parameters is 169 millions.

The task weight for each subtask is set by the number of mini-batches used during training. In the pre-training, the ratio of mini-batch numbers for each subtasks are 1.0, 7.0, 0.5 and 0.5 for the T2T, SSL, S2P and S2T subtasks respectively.

We mask 30% tokens in the T2T BART subtask in ASR pre-training, and no masking is applied for the T2T NMT subtask in ST pre-training. 7% of the feature frames in the SSL subtask and 3% of the feature frames in the two supervised speech subtasks are randomly selected as the mask span starting time-step. The mask span length is 10.

The models are optimized with Adam (Kingma and Ba, 2014) for both pre-training and fine-tuning. The final results are evaluated using an averaged model from checkpoints of the last 10 epochs. Additional experimental details such as learning rate

and mini-batch sizes are included in [Appendix A<sup>2</sup>](#).

## 5 Experimental results

### 5.1 Main results

We present the 960 hours LIBRISPEECH recognition results in [Table 1](#). We include results from the literature from row one to four and list both decoding results without and with an external LM. The WERs obtained with LM are displayed within “()”.

The first part of the table shows results from the wav2vec 2.0 base model, which is a CTC based ASR system. Second part of the table presents results from three ASR systems reported using the EAD modeling framework (row two to four). We mainly compare the proposed method with systems based on the EAD modeling framework. LAS is a LSTM based system trained with the LIBRISPEECH data only. Transformer ([Tang et al., 2021b](#)) and SpeechT5 ([Ao et al., 2021](#)) are based on multi-task learning and jointly trained with text tasks. Besides joint training in the fine-tuning stage, SpeechT5 also utilizes unsupervised joint pre-training to incorporate text data in the early training stage.

The results from the proposed STPT is presented in the third part of the table (row five). STPT without an external LM outperforms all previous reported EAD-based systems. On average, there is a 17.5% relative WER reduction compared to SpeechT5 with LM. When LM is applied, model with STPT only reduces WER slightly, with an average WER reduction of less than 0.1. The decoding LM is trained with the text training corpus from LIBRISPEECH, the same as the T2T subtask in the pre-training and fine-tuning. Other systems, on the other hand, show a considerable WER reduction when the LM is applied during decoding. It indicates that our multi-task learning in the pre-training and fine-tuning stages can effectively fuse linguistic information in the text data corpus into the ASR model. The external LM might not be required if it is trained on the same text corpus.

In [Table 2](#), we present the speech translation results on the MuST-C datasets. Row one to four are the latest results from literature. Row one shows the results by training a speech to text translation task alone. Row two and three present results from two multi-task systems with speech and text jointly trained together. Row four is the best system reported, which is initialized with the pre-trained

<sup>2</sup>We will open-source the code after the ACL review.

wav2vec 2.0 and machine translation model, then fine-tuned with joint speech and text training. Our method achieves 2.3 and 1.7 more BLEU scores for EN-ES and EN-FR translation directions compared with the results in row four ([Ye et al., 2021](#)).

### 5.2 Impact of model structure

Interference among subtasks may impede the progress of multi-task learning and result in inferior results. In this work, we examine the task interference via comparing the gradient similarity between pair subtasks. We choose the pre-trained models using the FSE configuration and accumulate gradients from one of four jointly trained subtasks discussed in [section 3](#). We prepare 20 batches of training samples for each subtask, and retrieve the accumulated gradients by sending these batches to the models. Then we calculate the pairwise cosine similarity between gradients from any two subtasks.

The pairwise subtask gradient similarity from the shared encoder are presented in [Figure 2](#). The [Figure 2\(a\)](#) demonstrates the gradient similarity in ASR pre-training. In most layers, the gradient similarities are small. No serious gradient interference is observed. The [Figure 2\(b\)](#) depicts the gradient similarity from the ST pre-training. Compared with the ASR pre-training, the S2T and T2T subtasks are replaced by sequence to sequence speech translation and text based neural machine translation subtasks in the ST pre-training. The interference between different subtasks is significant as large positive and negative gradient similarities are observed in the third and fifth layers, as shown in [Figure 2](#).

Similarly, we compare task gradients in the speech encoder and no obvious task interference is observed within the speech encoder for both ASR and ST pre-training. Detailed analysis on the speech encoder is included in the [Appendix B](#).

In order to alleviate the task interference, we propose the PSE configuration for the ST pre-training instead of the FSE configuration. [Table 3](#) presents the performance comparison between two configurations on both ASR and ST pre-training. On the left part of the table, we list the ASR results using 100 hours labelled speech data (train-clean-100) in the pre-training and fine-tuning. While the right part of the table shows the BLEU from the speech translation STPT evaluated on the MUST-C dataset. As we expected, the FSE configuration encourages information sharing among tasks and it achieves

Data set	Dev		Test		
	clean	other	clean	other	ave.
wav2vec 2.0 (Baevski et al., 2020b)	3.2 (1.8)	8.9 (4.7)	3.4 (2.1)	8.5 (4.8)	6.0 (3.4)
LAS (Park et al., 2019)	-	-	2.8 (2.5)	6.8 (5.8)	-
Transformer (Tang et al., 2021b)	2.8	7.0	3.1	7.2	5.0
SpeechT5 (Ao et al., 2021)	2.7 (2.2)	6.9 (5.6)	2.9 (2.3)	7.1 (5.7)	4.9 (4.0)
STPT	2.0 (2.1)	4.4 (4.2)	2.1 (2.1)	4.6 (4.5)	3.3 (3.2)

Table 1: WER results on Librispeech. “()” indicates the WER is measured with an external LM.

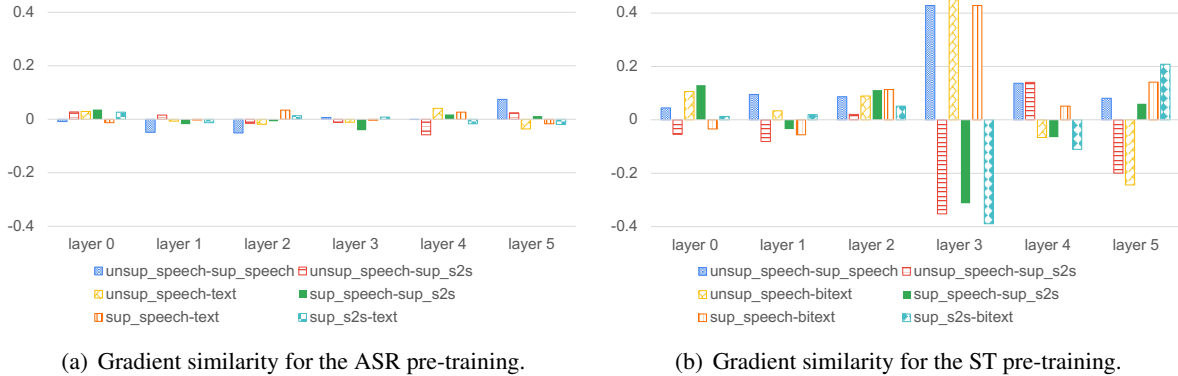


Figure 2: Gradient similarity for different subtasks on the shared text encoder.

Data corpus	EN-ES	EN-FR
Inaguma et al. (2020)	28.0	32.7
Tang et al. (2021a)	31.0	37.4
Zheng et al. (2021)	30.8	-
Ye et al. (2021)	30.8	38.0
STPT	33.1	39.7

Table 2: BLEU results of three language pairs on the MuST-C tst-COMMON.

Config.	Librispeech (WER ↓)		MuST-C (BLEU ↑)	
	dev clean	dev other	EN-ES	EN-FR
FSE	3.2	6.8	31.4	38.3
PSE	3.1	8.3	33.1	39.7

Table 3: Comparison of two pre-training configurations for ASR and ST.

lower WER for the ASR task, which indicates subtasks in the ASR pre-training are complementary to each other. On the other hand, the PSE configuration minimizes the information sharing between EAD subtasks and encoder only subtasks, and it leads to higher BLEU for the ST task.

### 5.3 Impact of supervised data

The supervised speech data connects the text and speech modeling and unifies the representation from different modalities. An interesting question we want to investigate is how much super-

vised data is enough to learn a good cross modality representation. In this experiment, we choose different amounts of labelled speech data for ASR pre-training and fine-tuning, varied from 960 hours (the full dataset), 100 hours (train-clean-100) and 10 hours as (Kahn et al., 2020), to answer this question.

In Table 4, the first column shows the amounts of supervised speech data available during the pre-training and the second column presents the amount of labelled data used in the fine-tuning stage. In pre-training, the same supervised speech data is used in the S2P and S2T subtasks.

The first observation is that more supervised speech data in the pre-training stage is always helpful to get smaller WER. For example, if the models are fine-tuned with the full LIBRISPEECH training dataset, the average WER are 3.3 (row one), 3.6 (row two) and 4.0 (row four) for experiments with 960, 100 and 10 hours labelled data in the pre-training stage. The second observation is that we are still able to obtain good speech presentations even with small amounts of labelled data. In row four, the model is pre-trained with 10 hours labelled data, then fine-tuned with 960 hours supervised speech data. It can achieve an average 4.0 WER, which is as good as the previously reported EAD systems in Table 1 if not better. However, we also

538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565

PT (h)	FT (h)	Dev		Test	
		clean	other	clean	other
960	960	2.0	4.4	2.1	4.6
100	960	2.3	4.9	2.2	5.1
	100	3.2	6.8	3.5	7.2
10	960	2.7	5.3	2.8	5.3
	100	3.8	7.8	4.0	7.7
	10	19.9	27.5	22.0	28.8

Table 4: Impact of the amounts of supervised data. “PT” and “FT” stand for pre-training and fine-tuning respectively.

notice the performance degrades quickly if only small amounts of labelled speech data are available. The average WER is increased to 24.6 (row six) when only 10 hours of supervised speech data is employed in both pre-training and fine-tuning.

#### 5.4 Ablation study

In Table 5, we present an ablation study by removing different steps/tasks in the pre-training stage.

In order to make the pre-training more stable, the model training adopts a three-stage optimization strategy as discussed in subsection 3.5: 1) pre-training the T2T subtask to have a good initialization on the phoneme embeddings 2) joint pre-training with four subtasks to leverage large amounts of unlabelled speech data and abundant text data and 3) fine-tuning the model on the downstream task for best performance. In the second row, we skip the T2T pre-training step and initialize the model randomly for the joint pre-training. 0.5 WER increase is observed in average on two LIBRISPEECH dev sets. It also has more impact on the EN-ES speech translation direction where 1.2 BLEU score is lost without proper initialization.

In the third row, we present the results without the S2T subtask. For both ASR and ST, significant performance degradation is observed, with an average 1.1 WER increase for two ASR tests and 1.8 BLEU decrease for two ST directions. We also try removing the S2P subtask while still keeping the S2T subtask. The training doesn’t converge. The SSL subtask is with very small or zero cost since all predictions collapse into one or two target phonemes. Also, little progress has been made for the S2T subtask even though it is co-trained with the SSL and T2T subtasks.

In the last row, the model is trained without pre-training, i.e., only the T2T and S2T subtasks are optimized. Compared with the STPT results, there is about 1.4 WER increase for two LIBRISPEECH

Config.	Librispeech (WER ↓)		MuST-C (BLEU ↑)	
	dev clean	dev other	EN-ES	EN-FR
STPT	2.0	4.4	33.1	39.7
- T2T pre-training	2.4	5.0	31.9	39.2
- EAD task	2.9	5.6	31.3	38.0
- pre-training	2.8	6.4	30.6	35.4

Table 5: Ablation study for STPT.

test sets and 3.4 BLEU decrease for the two ST directions on average.

#### 5.5 Discussion

In the ST pre-training, we use bitext data in the T2T subtask while the monolingual text data is only employed in the ASR pre-training. It is possible to include the BART task as another subtask in the ST pre-training stage, since the monolingual text data would be useful for the low resource speech translation directions, where only limited speech and bitext training data is available. In the S2T subtask, we take the downstream task as the auxiliary subtask, i.e., EAD based ASR for the ASR pre-training and speech to text translation for the ST pre-training. It will be interesting to extend this work for the multilingual scenario, where ASR could be treated as a special translation direction and we could have a uniformed pre-training framework for both ASR and ST tasks. We will leave these two extensions as our future work.

#### 6 Conclusion

In this work, we present a method to jointly pre-train speech and text in one model for speech translation and recognition under the EAD framework. It includes four self-supervised and supervised subtasks from two different input modalities, hence the proposed method can leverage large amounts of unlabelled speech data and abundant text data in the pre-training stage. We conduct detailed analysis on the interference among different subtasks and propose two model configurations for the ASR and ST pre-training respectively to alleviate the subtask interference. Our experimental results show the proposed method can effectively fuse information within text and speech training data into one model. We achieves between 1.7 and 2.3 BLEU improvement over the state of the art on the MUST-C EN-FR and EN-ES speech translation tasks, and comparable WERs as wav2vec 2.0 in the LIBRISPEECH ASR task.



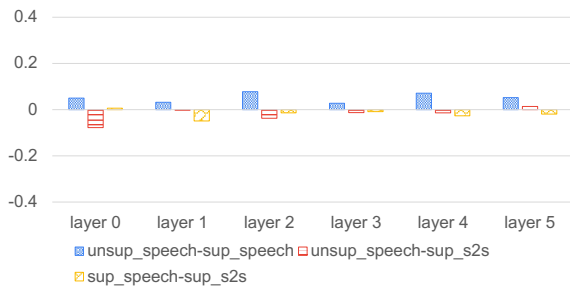
645	<b>References</b>		
646	Antonios Anastasopoulos and David Chiang. 2018.		
647	Tied multitask learning for neural speech translation.		
648	In <i>NAACL-HLT</i> .		
649	Junyi Ao, Rui Wang, Long Zhou, Shujie Liu, Shuo Ren,		
650	Yu Wu, Tom Ko, Qing Li, Yu Zhang, Zhihua Wei,		
651	Yao Qian, Jinyu Li, and Furu Wei. 2021. Speecht5:		
652	Unified-modal encoder-decoder pre-training for spo-		
653	ken language processing.		
654	Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton.		
655	2016. Layer normalization. <i>ArXiv</i> , abs/1607.06450.		
656	Alexei Baevski, Steffen Schneider, and Michael Auli.		
657	2020a. vq-wav2vec: Self-supervised learning of dis-		
658	crete speech representations. In <i>ICLR</i> .		
659	Alexei Baevski, Henry Zhou, Abdelrahman Mohamed,		
660	and Michael Auli. 2020b. wav2vec 2.0: A frame-		
661	work for self-supervised learning of speech represen-		
662	tations. In <i>NeurIPS</i> .		
663	Ankur Bapna, Yu an Chung, Nan Wu, Anmol Gu-		
664	lati, Ye Jia, Jonathan H. Clark, Melvin Johnson, Ja-		
665	son Riesa, Alexis Conneau, and Yu Zhang. 2021.		
666	Slam: A unified encoder for speech and language		
667	modeling via speech-text joint pre-training. <i>ArXiv</i> ,		
668	abs/2110.10329.		
669	Zhehuai Chen, Yu Zhang, Andrew Rosenberg, Bhu-		
670	vana Ramabhadran, Gary Wang, and Pedro J.		
671	Moreno. 2021. Injecting text in self-supervised		
672	speech pretraining. <i>ArXiv</i> , abs/2108.12226.		
673	Yung-Sung Chuang, Chi-Liang Liu, Hung yi Lee, and		
674	Lin-Shan Lee. 2020. Speechbert: An audio-and-text		
675	jointly learned language model for end-to-end spo-		
676	ken question answering. In <i>INTERSPEECH</i> .		
677	Yu-An Chung and James Glass. 2020. Improved		
678	speech representations with multi-target autoregres-		
679	sive predictive coding. In <i>ACL</i> .		
680	Yu-An Chung, Wei-Hung Weng, Schrasing Tong, and		
681	James R. Glass. 2018. Unsupervised cross-modal		
682	alignment of speech and text embedding spaces. In		
683	<i>NeurIPS</i> .		
684	Yu-An Chung, Yu Zhang, Wei Han, Chung-Cheng		
685	Chiu, James Qin, Ruoming Pang, and Yonghui Wu.		
686	2021a. W2v-bert: Combining contrastive learning		
687	and masked language modeling for self-supervised		
688	speech pre-training. <i>ArXiv</i> , abs/2108.06209.		
689	Yu-An Chung, Chenguang Zhu, and Michael Zeng.		
690	2021b. Splat: Speech-language joint pre-training		
691	for spoken language understanding. In <i>NAACL</i> .		
692	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and		
693	Kristina Toutanova. 2019. Bert: Pre-training of deep		
694	bidirectional transformers for language understand-		
695	ing. In <i>NAACL-HLT</i> .		
	Mattia Antonino Di Gangi, Roldano Cattoni, Luisa	696	
	Bentivogli, Matteo Negri, and Marco Turchi. 2019.	697	
	MuST-C: a multilingual speech translation corpus.	698	
	In <i>NAACL-HLT</i> .	699	
	Wei-Ning Hsu, Yao-Hung Hubert Tsai, Benjamin	700	
	Bolte, Ruslan Salakhutdinov, and Abdelrahman Mo-	701	
	hamed1. 2021. Hubert: How much can a bad teacher	702	
	benefit asr pre-training. In <i>ICASSP</i> .	703	
	H. Inaguma, S. Kiyono, K. Duh, S. Karita, N. Soplín,	704	
	T. Hayashi, and S. Watanabe. 2020. Espnet-st: All-	705	
	in-one speech translation toolkit. In <i>ACL</i> .	706	
	Ye Jia, Melvin Johnson, Wolfgang Macherey, Ron J.	707	
	Weiss, Yuan Cao, Chung-Cheng Chiu, Naveen Ari,	708	
	Stella Laurenzo, and Yonghui Wu. 2019. Leverag-	709	
	ing weakly supervised data to improve end-to-end	710	
	speech-to-text translation. <i>ICASSP</i> , pages 7180–	711	
	7184.	712	
	J. Kahn, A. Lee, and A. Hannun. 2020. Self-training	713	
	for end-to-end speech recognition. In <i>Proc. of</i>	714	
	<i>ICASSP</i> .	715	
	J. Kahn, M. Rivière, W. Zheng, E. Kharitonov,	716	
	Q. Xu, P. E. Mazaré, J. Karadayi, V. Liptchinsk-	717	
	sky, R. Collobert, C. Fuegen, T. Likhomanenko,	718	
	G. Synnaeve, A. Joulin, A. Mohamed, and	719	
	E. Dupoux. 2020. Libri-light: A benchmark	720	
	for asr with limited or no supervision. In	721	
	<i>ICASSP</i> , pages 7669–7673. <a href="https://github.com/facebookresearch/libri-light">https://github.com/facebookresearch/libri-light</a> .	722	
		723	
	Diederik P Kingma and Jimmy Ba. 2014. Adam: A	724	
	method for stochastic optimization. In <i>ICLR</i> .	725	
	T. Kudo and J. Richardson. 2018. Sentencepiece:	726	
	A simple and language independent subword tok-	727	
	enizer and detokenizer for neural text processing. In	728	
	<i>EMNLP</i> .	729	
	Y. Lee and T. Kim. 2018. Learning pronunciation from	730	
	a foreign language in speech synthesis networks.	731	
	<i>ArXiv</i> .	732	
	Mike Lewis, Yinhan Liu, Naman Goyal, Mar-	733	
	jan Ghazvininejad, Abdelrahman Mohamed, Omer	734	
	Levy, Ves Stoyanov, and Luke Zettlemoyer. 2020a.	735	
	Bart: Denoising sequence-to-sequence pre-training	736	
	for natural language generation, translation, and	737	
	comprehension. In <i>ACL</i> .	738	
	Mike Lewis, Yinhan Liu, Naman Goyal, Mar-	739	
	jan Ghazvininejad, Abdelrahman Mohamed, Omer	740	
	Levy, Veselin Stoyanov, and Luke Zettlemoyer.	741	
	2020b. BART: Denoising sequence-to-sequence	742	
	pre-training for natural language generation, trans-	743	
	lation, and comprehension. In <i>Proceedings of the</i>	744	
	<i>58th Annual Meeting of the Association for Compu-</i>	745	
	<i>tational Linguistics</i> , pages 7871–7880, Online. As-	746	
	sociation for Computational Linguistics.	747	
	Xian Li, Changhan Wang, Yun Tang, C. Tran, Yuqing	748	
	Tang, Juan Miguel Pino, Alexei Baevski, Alexis	749	
	Conneau, and Michael Auli. 2021. Multilingual	750	

751	speech translation from efficient finetuning of pre-trained models. In <i>ACL/IJCNLP</i> .	Ron J. Weiss, Jan Chorowski, Navdeep Jaitly, Yonghui Wu, and Zhifeng Chen. 2017. <a href="#">Sequence-to-sequence models can directly translate foreign speech</a> . In <i>INTERSPEECH</i> .	803
752			804
753	Yinhan Liu, Jiatao Gu, Naman Goyal, X. Li, Sergey Edunov, Marjan Ghazvininejad, M. Lewis, and L. Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. <i>ArXiv</i> , abs/2001.08210.	Alex Xiao, C. Fuegen, and Abdel rahman Mohamed. 2021. <a href="#">Contrastive semi-supervised learning for asr</a> . In <i>ICASSP</i> .	805
754			806
755			807
756			808
757			809
758	D. Park, W. Chan, Y. Zhang, C. Chiu, B. Zoph, E. Cubuk, and Q. Le. 2019. <a href="#">SpecAugment: A simple data augmentation method for automatic speech recognition</a> . <i>Interspeech</i> .	Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tie-Yan Liu. 2020. On layer normalization in the transformer architecture. In <i>ICML</i> , volume abs/2002.04745.	810
759			811
760			812
761			813
762			814
763	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. <a href="#">Deep contextualized word representations</a> . In <i>NAACL-HLT</i> .	Qiantong Xu, Tatiana Likhomanenko, Jacob Kahn, Awni Y. Hannun, Gabriel Synnaeve, and Roman Collobert. 2020. Iterative pseudo-labeling for speech recognition. In <i>Interspeech</i> , volume abs/2005.09267.	815
764			816
765			817
766	Juan Miguel Pino, Qiantong Xu, Xutai Ma, Mohammad Javad Dousti, and Yun Tang. 2020. Self-training for end-to-end speech translation. In <i>INTERSPEECH</i> .	Rong Ye, Mingxuan Wang, and Lei Li. 2021. <a href="#">End-to-end speech translation via cross-modal progressive training</a> . <i>ArXiv</i> .	818
767			819
768			820
769			821
770	Matt Post. 2018. <a href="#">A call for clarity in reporting BLEU scores</a> . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	Yu Zhang, James Qin, Daniel S. Park, Wei Han, Chung-Cheng Chiu, Ruoming Pang, Quoc V. Le, and Yonghui Wu. 2020. Pushing the limits of semi-supervised learning for automatic speech recognition. <i>Proc. of NeurIPS SAS Workshop</i> .	822
771			823
772			824
773			825
774			826
775	Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely. 2011. The kaldi speech recognition toolkit. In <i>ASRU</i> .	Renjie Zheng, Junkun Chen, M. Ma, and Liang Huang. 2021. <a href="#">Fused acoustic and text encoding for multi-modal bilingual pretraining and speech translation</a> . <i>ArXiv</i> .	827
776			828
777			829
778			830
779			831
780			832
781	Yun Tang, Juan Miguel Pino, Xian Li, Changan Wang, and Dmitriy Genzel. 2021a. Improving speech translation by understanding and learning from the auxiliary text translation task. In <i>ACL</i> .	Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and Quoc V. Le. 2020. Rethinking pre-training and self-training. <i>ArXiv</i> , abs/2006.06882.	833
782			834
783			835
784			836
785	Yun Tang, Juan Miguel Pino, Changan Wang, Xutai Ma, and Dmitriy Genzel. 2021b. A general multi-task learning framework to leverage text data for speech to text tasks. <i>ICASSP</i> .	<b>A Optimization setting</b>	837
786			838
787			839
788			840
789	Aäron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. <a href="#">Representation learning with contrastive predictive coding</a> . <i>ArXiv</i> .	<b>T2T subtask pre-training</b> The T2T model is pre-trained with learning rate 0.01 using Adam optimization. The maximum tokens per mini-batch is 2048 with 8 V100 GPU cards. The model is updated 400,000 until fully converged.	841
790			842
791			843
792	Changan Wang, Anne Wu, and Juan Pino. 2020. <a href="#">Covost 2 and massively multilingual speech-to-text translation</a> .	<b>Pre-training with all subtasks</b> The model then keeps optimizing with all four subtasks: T2T, SSL, S2P and S2T, with learning rate 0.001. The model is trained using 16 A100 GPU cards with update frequency 12. The maximum token number per batch for the T2T subtask is 2048 while the maximum sample number is 750,000 (46s) for the speech input in three speech subtasks. The maximum update number is 800,000 and 200,000 for the ASR pre-training and the ST pre-training respectively.	844
793			845
794			846
795	Changan Wang, Anne Wu, Juan Miguel Pino, Alexei Baevski, Michael Auli, and Alexis Conneau. 2021a. <a href="#">Large-scale self- and semi-supervised learning for speech translation</a> . In <i>Interspeech</i> .	<b>Fine-tuning</b> The model is fine-tuned on the downstream task with learning rate 0.0003 and 8 V100 GPU cards. The update frequency set to 3. The	847
796			848
797			849
798			850
799	Chengyi Wang, Yu Wu, Shujie Liu, Jinyu Li, Yao Qian, K. Kumatani, and Furu Wei. 2021b. <a href="#">Unispeech at scale: An empirical study of pre-training method on large-scale speech recognition dataset</a> . In <i>ICML</i> .		851
800			852
801			853
802			854

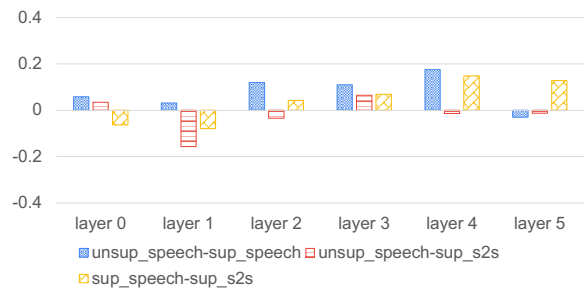
855 maximum update numbers are dependent on the  
856 amounts of supervised speech data available. We  
857 choose 100,000 for the ASR task with 960 hours  
858 training data and 20,000 for 100 or 10 hours train-  
859 ing data. For the ST task, the maximum update  
860 number is set to 50,000.

## 861 **B Gradient similarity of the speech** 862 **encoder**

863 Three subtasks: SSL, S2P, and S2T, share the  
864 speech encoder during the joint pre-training. Sim-  
865 ilar pairwise gradient similarity analysis is con-  
866 ducted on these three subtasks at the speech en-  
867 coder, as shown in [Figure 3](#). The gradient similarity  
868 analysis for the ASR pre-training is presented in  
869 the left subfigure while the ST-pretraining is listed  
870 in the right. In both cases, the gradient similarities  
871 for different subtask pairs are small, i.e., absolute  
872 values of the gradient similarities are all below 0.2.  
873 It indicates the task interference between different  
874 subtasks are not significant.



(a) Gradient similarity for the ASR pre-training.



(b) Gradient similarity for the ST pre-training.

Figure 3: Gradient similarity for different subtasks on the speech encoder.