# Improving Speech Translation by Cross-modal Multi-grained Contrastive Learning

Hao Zhang, Nianwen Si, Yaqi Chen, Wenlin Zhang, Xukui Yang, Dan Qu, Wei-Qiang Zhang, *Senior Member, IEEE*

*Abstract*—The end-to-end speech translation (E2E-ST) model has gradually become a mainstream paradigm due to its low latency and less error propagation. However, it is non-trivial to train such a model well due to the task complexity and data scarcity. The speech-and-text modality differences result in the E2E-ST model performance usually inferior to the corresponding machine translation (MT) model. Based on the above observation, existing methods often use sharing mechanisms to carry out *implicit knowledge transfer* by imposing various constraints. However, the final model often performs worse on the MT task than the MT model trained alone, which means that the knowledge transfer ability of this method is also limited. To deal with these problems, we propose the FCCL (<u>F</u>ine- and <u>C</u>oarse- Granularity <u>C</u>ontrastive <u>L</u>earning) approach for E2E-ST, which makes *explicit knowledge transfer* through cross-modal multi-grained contrastive learning. A key ingredient of our approach is applying contrastive learning at both sentence- and frame-level to give the comprehensive guide for extracting speech representations containing rich semantic information. In addition, we adopt a simple whitening method to alleviate the representation degeneration in the MT model, which adversely affects contrast learning. Experiments on the MuST-C benchmark show that our proposed approach significantly outperforms the state-of-the-art E2E-ST baselines on all eight language pairs. Further analysis indicates that FCCL can free up its capacity from learning grammatical structure information and force more layers to learn semantic information.

*Index Terms*—Speech Translation, Contrastive Learning, End-to-End.

## I. INTRODUCTION

SPEECH Translation (ST) takes the speech in one language (source) as input and outputs the translated text in another language (target). Traditional ST systems [1], [2] cascade automatic speech recognition (ASR) and machine translation (MT), which might suffer from error propagation and high latency. With the rapid progress of deep learning, end-to-end speech translation (E2E-ST) has attracted more attention [3]–[12]. Different from the traditional cascading method,

Hao Zhang, Nianwen Si, Yaqi Chen, Wenlin Zhang, Xukui Yang, and Dan Qu are with the School of the Information Engineering University, Zhengzhou 450000, China. Email: (haozhang012, snw1608, chyaqi163, wenlinzzz, gzyangxk, qudan_xd)@163.com.

Wei-Qiang Zhang is with the Department of Electronic Engineering, Tsinghua University, Beijing 100084, China. Email: (wqzhang@tsinghua.edu.cn).
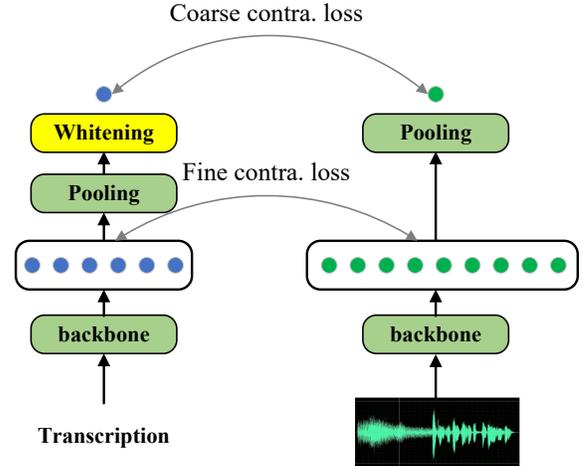
Fig. 1. Schematic diagram of multi-grained contrastive learning. Text and speech are inputted into the backbone to get the token-level (frame-level) representations. The number of blue and green points before pooling represents the sequence length. We then pool the token-level representations over the time dimension to obtain the sentence-level representations. We conduct fine and coarse granularity contrastive learning at the token- and sentence-level representations, respectively.

which decomposes ST into two sub-tasks, E2E-ST jointly handles them in a single neural network, which endows it with unique advantages, such as less error propagation and fewer parameters.

However, compared with MT and ASR tasks, E2E-ST is a cross-modal translation task, and its training data is more challenging to collect. The task complexity and data scarcity mean that it is non-trivial to train such a model well. Although both ST[1] and MT are translation tasks, the performance of the ST model is usually much inferior to the corresponding text-based MT model [7] with the same source and target languages. This can be attributed to the modality gap between speech and text. Compared to discrete text, speech is fine-grained and contains more noise, making it more challenging to extract speech representations containing rich semantic information. Based on the above observation, existing methods focus on transferring knowledge from the MT to the ST model with sophisticated techniques, such as progressive training [10], triangular decomposition agreement [8], and manifold mixup [9]. The common idea behind these methods is to implicitly constrain the parameter space of the ST model by treating the MT task as a constraint term with the sharing mechanism, which can be called ***implicit knowledge***

---

[1]If not special specified, ST mentioned below refers to E2E-ST by default.

*transfer*. Although impressive performance improvement has been achieved, the performance of the final model in the MT task is often inferior to that of the MT model trained alone [9], [13], which means the auxiliary knowledge from the MT task will be less and less as the training progresses. Consequently, implicitly transferring knowledge may not be optimal for the ST task.

Recalling the original problem mentioned above, the modality gap makes it hard to extract semantically rich representation from speech. On the contrary, it is easy to obtain semantic information from text due to its structural advantage. Why not use text representation to provide direct guidance for extracting speech representation? In this paper, we propose FCCL (Fine- and Coarse- Granularity Contrastive Learning) to transfer explicit knowledge from the MT to the ST model. The fundamental motivation behind our method is that since the semantic space is shared between different modalities, the speech and text with similar semantics should be close in the semantic space. In contrast, data pairs with different meanings should be pushed far away. Furthermore, the continuous semantic space of speech should be linked with the discrete symbolic space of text, which is more flexible in transferring knowledge.

Specifically, as shown in Figure 1, we provide the comprehensive guidance for extracting speech representations from both frame- and sentence-level to bridge the modality gap. The encoder output of the ST model is regarded as the frame-level speech representation. We average the encoder output along the time dimension to get the sentence-level speech representation. Similar to the ST model, we can get the token- and sentence-level text representation for the MT model. Coarse granularity contrastive learning is conducted at sentence-level. ST is a generative task, which means each speech frame in the encoder should have precise semantics. The representation learned from sentence-level can be suboptimal for frame-level. Thus, we conduct fine granularity contrastive at frame-level (token-level) to learn more refined representations.

It is straight to implement coarse granularity contrast. The fine granularity contrast, however, is not the same. To implement fine granularity contrast, we should obtain the alignment between the speech frame and text token. A common way to find this correspondence is to carry out force alignment through an additional alignment model [9]. However, its computational complexity is high, and the alignment model is not always available. To deal with this problem, we propose a maximum similarity method to get the alignment in an unsupervised manner. Based on the obtained alignment, we conduct fine granularity contrastive to find the optimal the token-level representation for the decoder.

Our main contributions are summarized as follows:

1) We propose FCCL, a cross-modal multi-grained contrastive learning method, to conduct *explicit knowledge transfer* from the MT to the ST model.
2) We propose a maximum similarity method to effectively and efficiently find the correspondence between speech frames and text tokens in an unsupervised manner, which needs negligible computation overhead.

3) We use a whitening method to alleviate the representation degeneration of the MT model by transforming the sentence representations into a standard normal distribution, which satisfies isotropy.
4) We show through CCA analysis that FCCL can free up its capacity from learning grammatical structure information and force more layers to learn semantic information.
5) We conduct experiments on the MuST-C benchmark on all eight language pairs. The experiment results and detailed analysis verify the effectiveness of our proposed method.

## II. RELATED WORKS

**End-to-end ST** Early ST system cascaded the ASR and MT system. Benefiting from the development of deep learning in recent years, E2E-ST [14], [15] has become a mainstream paradigm due to its advantages, such as lower latency, alleviation of error propagation, and fewer parameters. However, due to the scarcity of triplet training data and the complexity of cross-modal translation task, it is non-trivial to train such a model well. To overcome the scarcity of training data, various methods have been explored, such as data augmentation [16]–[18], multi-task learning [19], [20], sub-module pre-training [21]–[25], self-training [26], meta-learning [27], and interactive decoding [20]. Meanwhile, some works focus on alleviating the task complexity. One branch notices that the encoder of the ST model is overburdened. They decouple the ST encoder into an acoustic encoder and a semantic encoder to improve the ability to extract information from the speech feature [6], [7], [11]. Another branch aims to bridge the modality gap between speech and text. They focus on transferring knowledge from the MT to the ST model with sophisticated techniques, such as progressive training [10], triangular decomposition agreement [8] and manifold mixup [9]. In this work, we explore how to bridge the modality gap via *explicit knowledge transfer* based on contrastive learning.

Contrastive learning aims to learn general representations on many unlabeled data. It has extensively promoted progress in computer vision [28], [29], natural language processing (NLP) [30]–[32], and speech [33]–[35]. However, details of view generation are crucial and require careful design [36]. In contrast, using multiple modalities as different views is simpler and more natural. Some studies extend contrastive learning to the multimodal domain [36]–[38]. [12], [39] introduce a similar idea to design a speech-text cross-modal contrastive learning module in speech translation. However, they suffer from two unique problems.

First, they trained ST and MT tasks together. Unless unique methods are used [13], the joint-trained model often performs worse on the MT task than the MT model trained alone. Additionally, joint training will also limit the use of advanced techniques to improve the MT performance, which is directly related to the text representation quality.

Second, they ignore the characteristics of the ST task. ST is a generative task, which means each speech frame in the encoder should have precise semantics. However, the
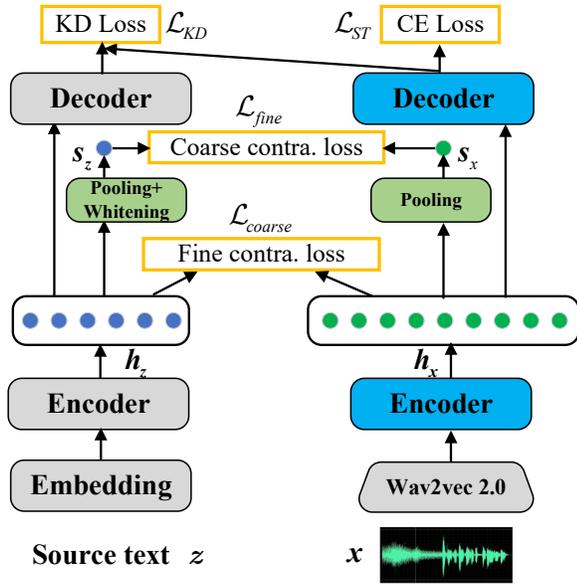
Fig. 2. Overview of our proposed method. Pooling means averaging the encoder output over the time dimension. FCCL contains three modules, a pretrained MT model, an ST model, and a contrastive module. The grey modules in the figure indicate that its parameters are no longer updated during training. The contra. in the picture is the abbreviation of contrastive. During inference, only the ST model is preserved, and all other modules are discarded.

representation learned from sentence-level can be sub-optimal for frame-level [40]. In this work, we propose cross-modal multi-grained contrastive learning to give the comprehensive guide for extracting speech representations containing rich semantic information. We obtain the text representation from a pretrained MT model. This can ensure high-quality text representation during the training.

**Fine Granularity Contrastive learning** Contrastive learning is usually conducted on the overall representation of the input. Nevertheless, better overall representation does not guarantee more accurate fine granularity representation [40]. It is suboptimal for generative tasks or dense prediction tasks. In computer vision, some studies perform contrastive learning at pixel-level to learn finer input representations [40]–[42]. In the ST task, we need to conduct fine granularity contrastive at frame-level so that each speech frame in the encoder has precise semantics. Different from [9], which gets the correspondence between speech frame and text token by force alignment, we propose a maximum method to obtain the correspondence in an unsupervised manner with negligible latency overhead.

## III. METHOD

In this section, we begin with the fundamental problem formulation of E2E-ST. Then we introduce the coarse and fine granularity contrastive learning in Sections III-C and III-D, respectively. The overall structure is shown in Figure 2.

### A. Problem Formulation

The speech translation corpus usually contains *speech-transcription-translation* triples, denoted as

$\mathcal{D} = \{(x^{(n)}, y^{(n)}, z^{(n)})\}_{n=1}^N$, where $x, y, z$ are the audio, the translation in the target language, and the corresponding transcription in the source language, respectively. During training, we optimize the maximum likelihood estimation loss on the training set:

$$\mathcal{L}_{ST}(\theta) = - \sum_{(x,y) \in \mathcal{D}} \log p(y|x) \qquad (1)$$

### B. Model Architecture

FCCL contains an MT model, a contrastive learning module, and an ST model. We freeze the parameters of the MT model to avoid performance drop during training. And this manner can also facilitate using sophisticated methods to improve the MT model performance, which can provide better text representations.

**Wav2vec 2.0 as feature extractor** Traditional acoustic features cause model performance degradation when training data is insufficient, since the manually customized process will inevitably lead to information loss [43]. Following the previous works [9], [10], [12], we adopt Wav2vec 2.0 [44] without finetune to extract speech representation for raw speech $x$. Its parameters are frozen to facilitate quick experiments.

### C. Coarse Granularity Contrastive Learning

Contrastive learning aims to learn general representations from two views of the same input [28], [29]. However, details of view generation are crucial and require careful design [36]. In contrast, using multiple modalities as different views is simpler and more natural. In this paper, we treat speech and the corresponding transcription as expressions of the same semantics in different modalities.

Given the input speech-transcription pairs $\{(x^{(n)}, z^{(n)})\}_{n=1}^N$, the encoded representations are $\{(h_x^{(n)}, h_z^{(n)})\}_{n=1}^N$. We average them in terms of the time dimension to get the sentence-level representation. The natural idea is to compute contrastive loss on the obtained sentence-level representation. However, the representation degeneration [30], [45] of the MT model makes it problematic. As shown in Figure 3a, affected by word frequency, the representations finally learned by the MT model are squeezed into a cone and are not uniformly distributed with respect to direction. The sentence-level representation - as average of the encoder output - suffers from the same issues [46]. Thus, the sentence-level representation space is semantically non-smoothing and poorly defined in some areas. This will lead to the phenomenon that some negative samples are not similar to the speech samples, but the calculated cosine similarity is relatively significant, which will further affect the calculation of contrastive loss [47]. Inspired by research in NLP [48], [49], we adopt a simple whitening strategy to make the sentence-level representations anisotropic (Figure 3b).

The representations for contrastive loss calculation are $\{(s_x^{(n)}, s_z^{(n)})\}_{n=1}^N$, where $s_x^{(i)} = \text{AveragePooling}(h_x^{(i)})$, $s_z^{(i)} = \text{Whitening}(\text{AveragePooling}(h_z^{(i)}))$, $s_x^{(i)} \in 1 \times d$, $s_z^{(i)} \in 1 \times d$.
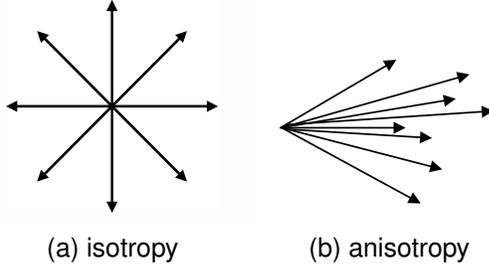
Fig. 3. Schematic diagram of isotropy and anisotropy. The arrow represents the directions of the word representations.

The text representation is fixed, while the speech representation is dynamic. So, we only treat the text representations $s_z^{(k),k \neq i}$ as negative samples to ensure the consistency of negative sample representation. Moreover, we set up a First-In-First-Out (FIFO) queue [28] to store the text representations of the previous batch to decouple the relationship between the number of negative samples and the batch size. The contrastive loss for coarse granularity is as follows:

$$\mathcal{L}_{coarse} = -\frac{1}{N}\sum_{i=1}^{N} \log \frac{e^{sim(s_x^i \cdot s_z^i/\tau)}}{\sum_{j=1}^{N} e^{sim(s_x^i \cdot s_z^j/\tau)} + \sum_{k=1}^{K} e^{sim(s_x^i \cdot s_z^k/\tau)}} \quad (2)$$

where $sim(\cdot)$ is the cosine similarity function, $N$ is the batch size, $\tau$ is the temperature and $K$ is the number of negative samples in the queue.

**Whitening** Given a set of sentence representations $\{q_i\}_{i=1}^{N}, q_i = \text{AveragePooling}(h^{(i)})$, a linear transformation is performed to ensure the mean value is zero and the covariance is the identity matrix.

$$\tilde{q}_i = (q_i - \mu)W \quad (3)$$

where $W$ is the transform matrix, and $\mu$ is defined as the mean vector of the entire embedding set $\{q_i\}_{i=1}^{N}$, i.e., $\mu = \frac{1}{N}\sum_{i=1}^{N} q_i$. We denote the original covariance matrix of $\{q_i\}_{i=1}^{N}$ as $\Sigma$.

The transformed covariance matrix $\widetilde{\Sigma}$ of $\{\tilde{q}_i\}_{i=1}^{N}$ need to be the identity matrix:

$$\widetilde{\Sigma} = W^T \Sigma W = I \quad (4)$$

Therefore,

$$\Sigma = (W^T)^{-1}W$$
$$= (W^{-1})^T W^{-1} \quad (5)$$

where $\Sigma$ is a positive definite symmetric matrix. SVD decomposition is calculated to get the solution:

$$\Sigma = U\Lambda U^T$$
$$W = U\sqrt{\Lambda^{-1}} \quad (6)$$

### D. Fine Granularity Contrastive Learning

ST is a generative task, which means each speech frame in the encoder should have precise semantics. Nevertheless, better overall representation does not guarantee more accurate fine granularity representation. We need more advanced processing. The length inconsistency of speech and text makes fine granularity contrast less straightforward than coarse granularity. The alignment between speech frames and text tokens must be obtained. We propose a maximum similarity method to find this correspondence in an unsupervised manner to avoid force alignment, as the alignment model is not always available.

We get the cosine similarity matrix $\Delta \in {}^{T_x \times T_z}$ based on the encoded representations $\{(h_x^{(n)}, h_z^{(n)})\}_{n=1}^{N}$. $T_x, T_z$ are the length of speech frames and text tokens, respectively. Each speech frame is matched to the text token with the maximum similarity. We denote the $j$-th frame representation of the $i$-th speech as $h_{x,j}^{(i)}$. This process can be described as follows:

$$pos_{x,j}^{(i)} = \{h_{z,m}^{(i)} | \arg\max_m sim(h_{x,j}^{(i)} \cdot h_{z,m}^{(i)}), m = 1, 2, \cdots T_z\} \quad (7)$$

where $h_{x,j}^{(i)}$ is the $j$-th frame representation of the $i$-th speech, $h_{z,m}^{(i)}$ is the $m$-th token representation of the corresponding transcription. The entire matching process can be quickly calculated through the matrix with a small amount of calculation. Notes that we do not use whitening here to ensure computational instability because the length of the token sequence is limited. Given the correspondence between the speech frame and text token, we can quickly get the fine granularity contrastive loss:

$$\mathcal{L}_{fine} = -\frac{1}{N}\sum_{i=1}^{N}\sum_{j=1}^{T_x} \log \frac{e^{sim(h_{x,j}^{(i)} \cdot pos_{x,j}^{(i)}/\tau)}}{e^{sim(h_{x,j}^{(i)} \cdot pos_{x,j}^{(i)}/\tau)} + \sum_{k=1}^{K} e^{sim(h_{x,j}^{(i)} \cdot s_z^k/\tau)}} \quad (8)$$

where $K$ is the number of negative samples, $T_x$ is the length of the frame-level speech and $N$ is the batch size. We use the same negative sample strategy as we use for coarse contrast, because it's conceptually and computationally simpler.

### E. Dropdim

Currently commonly used data augmentation strategies in speech include Speed Perturb [50] and SpecAugment [51]. However, when using Wav2vec 2.0 as the acoustic feature extractor, the effect of these data augmentation strategies is usually limited. In this paper, we adopt dropdim [52], a recently proposed structured dropout method, as a data augmentation strategy. The key idea is to broke the excessive co-adapting between different embedding dimensions and force the self-attention to encode meaningful features with a certain number of embedding dimensions erased.

### F. Training

The contrastive learning is conducted on the encoder. We combine the word-level knowledge distillation [5] in the

decoder to give multi-level guidance for the training of the ST model. The overall loss is the weighted sum of all previous losses:

$$\mathcal{L} = (1 - \gamma)\mathcal{L}_{ST} + \alpha\mathcal{L}_{coarse} + \beta\mathcal{L}_{fine} + \gamma\mathcal{L}_{KD} \quad (9)$$

where $\alpha, \beta, \gamma$, are hyper-parameter to adjust the weight of each loss. $\mathcal{L}_{KD}$ represents the word-level knowledge distillation loss.

## IV. EXPERIMENT

TABLE I
STATISTICS OF ALL DATASETS

| Language (EN-) | MuST-C | | External MT | |
|---|---|---|---|---|
| | hours | #sent | Source | #sent |
| Germany (DE) | 408h | 234K | WMT16 | 4.6M |
| French (FR) | 492h | 280K | WMT14 | 40.8M |
| Russian (RU) | 489h | 270K | WMT16 | 2.5M |
| Spanish (ES) | 504h | 270K | WMT13 | 15.2M |
| Italian (IT) | 465h | 258K | OPUS100 | 1.0M |
| Romanian (RO) | 432h | 240K | WMT16 | 0.6M |
| Portuguese (PT) | 385h | 244K | OPUS100 | 1.0M |
| Dutch (NL) | 442h | 253K | OPUS100 | 1.0M |

### A. Dataset and Processing

**MuST-C** MuST-C [53] is a multilingual dataset extracted from TED talks, including source audio, transcriptions, and text translations. Its source language is English, and the target language cover eight language direction: German (De), French (Fr), Russian (Ru), Spanish (Es), Italian (It), Romanian (Ro), Portuguese (Pt), and Dutch (NL). It is the most extensive training data for speech translation. We select the model according to its performance on the validation set and test it on the tst-COMMON set.

**External MT Datasets** The MT model is trained separately and has the same structure as the ST model. It allows us to use parallel sentence pairs in the external MT datasets in addition to the transcription-translation pairs in the ST corpus. Table I lists the statistics of all the datasets included.

**Processing** We use the original 16-bit 16kHz mono-channel audio waveform as speech input. We tokenize and true case all texts via Moses[2]. In each language direction, we apply BPE [54] on the combination of source and target text to obtain shared sub-word units with a vocabulary size of 8K.

### B. Experimental setups

**Model Configuration** We use Wav2vec 2.0 following the large configuration in [44], which is self-supervised pretrained on Librispeech [55] audio data only[3]. We use Transformer [56] as the backbone of the model, including 6 encoder layers and 6 decoder layers. We train both small and medium size models. For the small size model, each layer comprises 256 hidden units, 4 attention heads, and 2048 feed-forward hidden units. For the medium size model, the above parameters are set to

512, 8, 2048. The queue size $K$ and temperature $\tau$ are set to 1000 and 0.12, respectively, according to our pilot study.

**MT model Pretrain** The previous work [11] shows a domain mismatch between extra datasets and MuST-C. Therefore, the MT model is first trained on the extra MT datasets and then finetuned on transcription-translation pairs of MuST-C.

**E2E-ST Training and Inference** Following the previous work, we initialize the ST model encoder and decoder with the pretrained MT model. During training, we use the Adam [57] optimizer with $\beta_1 = 0.9, \beta_2 = 0.98$ and adopt the default learning schedule in ESPnet [58]. The dropout rate and the value of label smoothing are all set to 0.1. Regarding dropdim described in Section III-E, we adopt the random mask strategy described in their paper with a mask rate of 0.05. An early stop strategy is adopted during training with three epochs patiences. We set $\alpha$, $\beta$ and $\gamma$ to 1.0, 1.0, 0.6 respectively. All models are trained on 2 Nvidia Tesla-V100 GPUs. It takes about one day to converge.

During inference, we average the best 5 checkpoints for evaluation. We use beam search with a beam size of 10, and the length penalty is 0.6. We report the case-sensitive SacreBLEU[4] [59] for fair comparison with previous work.

**Baseline Systems** To verify the effectiveness of our method, we compare with the following E2E-ST systems: STAST [7], AFS [60], SATE [11], Dual-Decoder [61], XSTNet [10], TDA [8], STEMM [9], JT-S-MT [13], Chimera [12], STPT [25], SpeechUT [24] and ConST [31]. We implement FCCL_base which has the same model architecture as our proposed FCCL. The only difference is that it is trained without fine and coarse granularity contrastive learning.

### C. Main Results

**Comparison with E2E Baselines** Table II shows the results on the MuST-C dataset. The previous works have shown that the additional datasets can significantly improve the ST model performance. Therefore, to be fair, we mainly compare prior results both with and without additional MT datasets. (a) Without external MT datasets. FCCL$^s$ obtain an average improvement of 0.7 BLEU compared with FCCL$^s$_base. This show that contrastive learning can effectively guide the learning of the ST model, leading to better translation performance. When including the results from previous work, FCCL$^m$ outperforms the previous models and achieves new state-of-the-art. Different from STEMM [9], which adopts the manifold mixup to alleviate the representation discrepancy in an implicit manner, we bridge the modality gap through explicit knowledge transfer with the help of contrastive learning. Thus, FCCL achieves better performance, and the improvement over STEMM is also remarkable. (b) With external MT datasets. Compared to itself, FCCL$^m$ can achieve 1.2/1.5 BLEU improvement in different model sizes when additional MT datasets are available. This demonstrates FCCL ability to leverage additional datasets. Chimera [12] designed a shared semantic memory through contrastive learning to learn the semantic information shared between modalities. However, it

---

[2]https://www.statmt.org/moses/
[3]https://huggingface.co/facebook/wav2vec2-large-960h

[4]sacreBLEU signature: nrefs:1 — bs:1000 — seed:12345 — case:mixed — eff:no — tok:13a — smooth:exp — version:2.0.0

TABLE II
BLEU SCORES ON MUST-C TST-COMMON SET. "EXTERNAL DATA" INDICATES WHETHER THE METHOD USES ADDITIONAL DATA. THE SUPERSCRIPTS $s$ AND $m$ REPRESENT THE SMALL MODEL AND MEDIUM MODEL, RESPECTIVELY.

| Model | External Data | | BLEU | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Speech | MT | EN-DE | EN-FR | EN-RU | EN-ES | EN-IT | EN-RO | EN-PT | EN-NL | Avg. |
| w/o external MT data | | | | | | | | | | | |
| STAST [7] | × | × | 23.1 | - | - | - | - | - | - | - | - |
| AFS [60] | × | × | 22.4 | 31.6 | 14.7 | 26.9 | 23.0 | 21.0 | 30.0 | 24.9 | 23.9 |
| SATE [11] | × | × | 25.2 | - | - | - | - | - | - | - | - |
| Dual-Decoder [61] | × | × | 23.6 | 33.5 | 15.2 | 28.1 | 24.2 | 22.9 | 30.0 | 27.6 | 25.7 |
| STPT [25] | ✓ | × | - | **39.7** | - | **33.1** | - | - | - | - | - |
| XSTNet [10] | ✓ | × | 25.5 | 36.0 | 16.9 | 29.6 | 25.5 | **25.1** | 31.3 | 30.0 | 27.5 |
| TDA [8] | × | × | 25.4 | 36.1 | 16.4 | 29.6 | 25.1 | 23.9 | 31.1 | 29.6 | 27.2 |
| STEMM [9] | ✓ | × | 25.6 | 36.1 | 17.1 | 30.3 | 25.6 | 24.3 | 31.0 | 30.1 | 27.5 |
| ConST [31] | ✓ | × | 25.7 | 36.8 | 17.3 | 30.4 | 26.3 | 24.8 | **32.0** | **30.6** | 28.0 |
| FCCL$^s$_base | ✓ | × | 24.9 | 35.7 | 17.1 | 29.9 | 25.2 | 23.7 | 30.5 | 29.6 | 27.1 |
| FCCL$^s$ | ✓ | × | 25.7 | 36.5 | 17.5 | 30.4 | 26.0 | 24.6 | 31.4 | 30.3 | 27.8 |
| FCCL$^m$ | ✓ | × | **25.9** | 36.8 | **17.6** | 30.7 | **26.4** | 25.0 | 31.8 | 30.5 | **28.1** |
| w/ external MT data | | | | | | | | | | | |
| SATE [11] | ✓ | ✓ | 28.1 | - | - | - | - | - | - | - | - |
| JT-S-MT [13] | × | ✓ | 26.8 | 37.4 | - | 31.0 | - | - | - | - | - |
| SpeechUT [24] | × | ✓ | **30.1** | **41.4** | - | **33.6** | - | - | - | - | - |
| XSTNet [10] | ✓ | ✓ | 27.8 | 38.0 | 18.5 | 30.8 | 26.4 | 25.7 | **32.4** | **31.2** | 28.8 |
| Chimera [12] | ✓ | ✓ | 26.3 | 35.6 | 17.4 | 30.6 | 25.0 | 24.0 | 30.2 | 29.2 | 27.3 |
| TDA [8] | ✓ | ✓ | 27.1 | 37.4 | - | - | - | - | - | - | - |
| STEMM [9] | ✓ | ✓ | 28.7 | 37.4 | 17.8 | 31.0 | 25.8 | 24.5 | 31.7 | 30.5 | 28.4 |
| ConST [31] | ✓ | ✓ | 28.3 | 38.3 | 18.9 | 32.0 | 27.2 | 25.6 | **33.1** | **31.7** | 29.4 |
| FCCL$^s$_base | ✓ | ✓ | 27.6 | 36.8 | 18.0 | 30.3 | 25.7 | 25.1 | 31.0 | 30.2 | 28.0 |
| FCCL$^s$ | ✓ | ✓ | 28.7 | 37.5 | 19.1 | 31.2 | 26.5 | 26.0 | 32.1 | 31.0 | 29.0 |
| FCCL$^m$ | ✓ | ✓ | 29.0 | 38.3 | **19.7** | 31.9 | **27.3** | **26.8** | 32.7 | 31.6 | **29.6** |

TABLE III
COMPARISON WITH CASCADED MODELS ON MUST-C EN-DE AND EN-FR TST-COMMON SET.

| Model | | BLEU | |
|---|---|---|---|
| | | En-De | En-Er |
| | XSTNet [10] | 25.2 | 34.9 |
| Cascaded | STEMM [9] | 27.5 | - |
| | SATE [11] | 28.2 | - |
| End-to-end | FCCL$^m$ | **29.0** | **38.3** |

TABLE IV
MODEL PERFORMANCE UNDER DIFFERENT TRAINING METHODS.

| Models | Task | |
|---|---|---|
| | ST | MT |
| Baseline | 23.79 | 28.10 |
| joint | 24.43 | 27.47 |
| pretrain | 25.71 | 28.10 |

limits the feature output lengths of the two modalities to be consistent, which will sacrifice the MT model performance [31]. Our proposed method does not have this limitation. ConST [31] introduces a similar idea to bridge the modality gap. Nevertheless, they ignore the nature of ST task and only conduct contrastive learning at sentence-level. In contrast, we propose to give more clear guidance at both sentence- and frame-level, which achieves better performance. Our model is worse than SpeechUT [24] and STPT [25]. However, they mainly focus on the pre-training procedure. Our proposed method is orthogonal with theirs, and we will investigate how to combine them together in the future.

**Comparison with Cascaded Baselines** To further validate the effectiveness of our proposed method, we compare with several strong cascaded baseline systems, all of which are trained with additional datasets. As described in Table III, our proposed method outperform the cascade model and achieve better performance.

## V. ANALYSIS

### A. How to get the MT model?

When transferring knowledge from the MT to the ST model, an important question is how to get the MT model. In this section, we conduct some analysis of this. Pretrain represents we pretrain the MT model and freeze its parameters, which is used in this paper. Joint represents that we jointly train the MT and ST task in a shared model. Baseline means the ST model is trained only with the cross-entropy loss.

As shown in Table IV, the performance of the joint-trained model on the ST task is better than baseline, but worse than pretrain. Although joint training can use the MT task as an additional task to constrain the parameter solution space of the ST task, the opposite is not always true. In the absence of sophisticated techniques, the performance of the joint-trained model on the MT task is often inferior to the MT model trained alone. Although a tuned semantic space benefits the early learning of ST, the performance decreasing of the joint-trained model on the MT task in later training will have a more significant negative impact. In the future, we will study how to combine the advantages of pretrain and joint training.

### B. Is it necessary to initialize the ST model with a pretrained MT model?

In this paper, we propose to improve the ST model performance by explicit knowledge transfer. However, the ST encoder and decoder are initialized from a pretrained MT model.
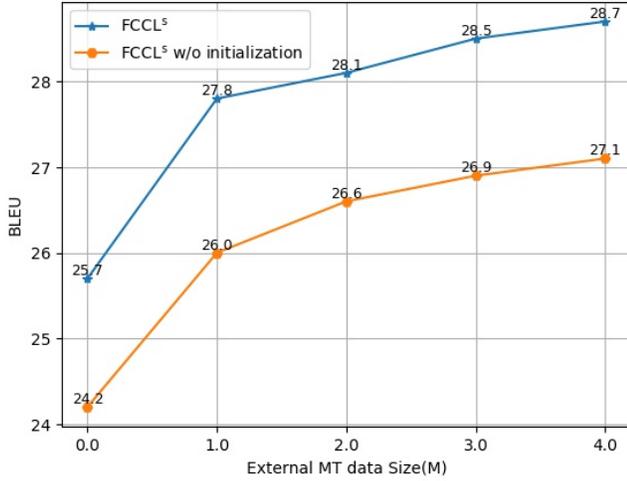
Fig. 4. BLEU scores on MuST-C En-De tst-COMMON against the size of external MT data.



Fig. 5. Histogram and probability density function (pdf) of the cosine similarity.

Initialization corresponds to implicit knowledge transfer. In this section, we study the effect of initialization on FCCL.

As shown in Figure 4, initialization can yield improvements over 1.0 BLEU. This can be attributed to the fact that initialization from a pretrained MT model can eliminate the randomness caused by different starting points of optimization [62], making the parameter space of the ST model to be simpler than that of the model trained from scratch. Although some studies in NLP show that initialization from a pretrained model (such as BART [63] and mBART [64]) will undermine the downstream task performance under high-resource settings (data pairs $>= 10M$) [65], its condition is too rigorous to observe a similar phenomenon in the ST task. On the one hand, as shown in Table I, the amount of data in the MuST-C dataset is limited. On the other hand, the external supervision signal used to guide the ST model learning may still not be strong enough. That results in FCCL not being completely free of the effects of initialization. We will conduct further research in the future, making FCCL fully "explicit".

### C. How does the whitening method work?

In this paper, we use the whitening method to alleviate the representation degeneration of the MT model. A natural question is how it works and affects the calculation of contrastive loss. We analyzed the cosine similarity distribution to make some explanations. Specifically, we randomly select a speech sample from the MuST-C En-De tst-COMMON set as the anchor sample, and 1000 transcription samples. For each sample, we average its encoder output over the time dimension to get the overall representation. The cosine similarity between 1000 samples and anchor sample (speech sample) is $[s_x^1 \cdot s_z^1, s_x^1 \cdot s_z^2, \cdots, s_x^1 \cdot s_z^j], j = 1, \cdots 1000$, where $s_x^1$ denotes the speech representation, $s_z^j$ is $j$-th text representation. $s_z^1$ is the corresponding positive sample of $s_x^1$, and the rest are treated as negative. Then we normalize the cosine similarity $[s_x^1 \cdot s_z^1, s_x^1 \cdot s_z^2, \cdots, s_x^1 \cdot s_z^j]/s_x^1 \cdot s_z^1, j = 1, \cdots 1000$.

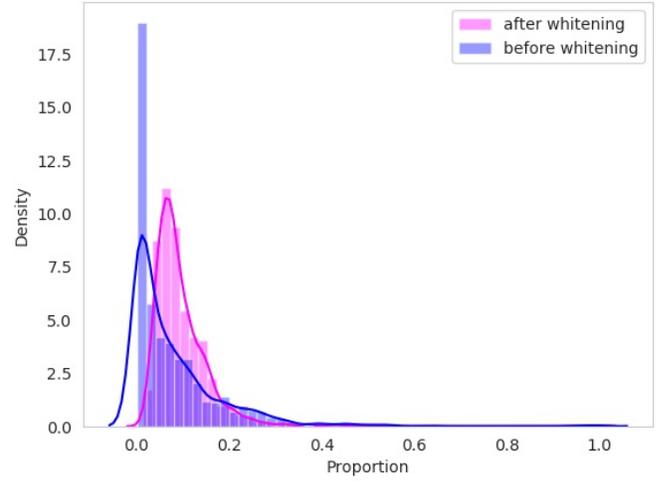The histogram and probability density function (pdf) of the cosine similarity is shown in Figure 5. Before whitening,

the pdf is a long-tailed distribution, which means that there exist many samples with high resemblance to the anchor sample except the positive sample. Unfortunately, many of these samples are not really similar to the anchor sample. The anisotropy of text representation leads to the appearance of spurious correlation. These long-tailed samples can be eliminated after whitening, indicating that the whitening operation can alleviate the long-tailed problem of cosine similarity distribution. The contrastive loss is based on cosine similarity. Thus, the whitening method can make contrastive learning more precise.

### D. Effectiveness of Each Learning Objective

In order to analyze the effectiveness of different modules in our method, we conduct an ablation study on the MuST-C En-De dataset. As shown in Table V, each part of FCCL is necessary and has a positive effect. (1) Removing the fine and coarse granularity contrastive loss brings about a decrease of 0.45 and 0.68 BLEU, respectively, indicating that the contrastive loss can guide the learning of the ST model. (2) Moreover, removing both of them will cause further performance degradation (0.8 BLEU), indicating that fine and coarse contrastive learning are complementary. Without coarse granularity contrast, fine granularity contrast produces slight effect on model performance, because good features will not be learned if incorrect correspondence is extracted. (3) When the contrast loss is calculated directly without whitening, 0.69 BLEU reduction was observed. This valid the analysis in Section V-C, showing that whitening can alleviate the long-tailed problem of cosine similarity distribution and make contrastive learning more precise. (4) When knowledge distillation is removed, the model performance drops by 0.67 BLEU. In FCCL, knowledge distillation and contrastive learning guide the ST model from the decoder and encoder outputs, respectively. They are complementary, and removing either one will adversely affect the model. (5) It is worth noting that when dropdim is removed, the model performance has

TABLE V
BLEU SCORES ON MUST-C EN-DE TST-COMMON SET WHEN DIFFERENT PARTS ARE REMOVED.

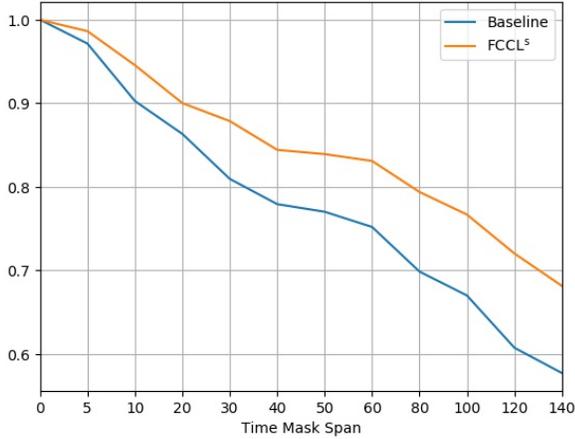| Model | BLEU |
|---|---|
| $FCCL^s$ | 25.71 |
| -fine contra. | 25.26(-0.45) |
| -coarse contra. | 25.03(-0.68) |
| -fine contra. | 24.91(-0.80) |
| -whitening | 25.02(-0.69) |
| -knowledge distillation | 25.04(-0.67) |
| -dropdim | 24.76(-0.95) |



Fig. 6. Model performance on MuST-C En-De tst-COMMON under different strength perturbations. The horizontal coordinate represents the time mask span. The vertical coordinate represents the ratio between the BLEU value under the perturbed input and the BLEU value under the unperturbed input.

a significant drop, about 0.95 BLEU. This can be attributed to the fact that dropdim in FCCL can not only enhance the generalization ability of the model, but also generate more challenging representations to enhance the effect of contrastive learning.

### E. Model Hallucinations

The attention mechanism of a model might not reflect a model actual inner reasoning. In MT, Lee [66] proposed the concept of hallucinations. A model hallucinates if a small perturbation in the input causes a sharp change in the output, indicating that the model does not really pay attention to the input. He selects the most common words from the corpus as perturbations. However, giving the exact definition in speech requires force alignment, which is costly. Instead, a simple perturbation method is considered in this paper. We adopt the time mask in SpecAugment [51] as a perturbation method to test the model performance under perturbation.

As shown in Figure 6, the blue and yellow curves represent the performance of the Baseline and $FCCL^s$ under different strength perturbations compared to the unperturbed input. The baseline represents an ST model trained only with the end-to-end cross-entropy loss $\mathcal{L}_{ST}$. The larger the time mask span, the more noise in the input. It can be seen that the performance of $FCCL^s$ decreases more slowly when the perturbation strength gradually increases, suggesting it is more resilient against perturbations and more attentive to the content of its input.

### F. Canonical Correlation Analysis

To further analyze FCCL, we turn to canonical correlation analysis (CCA), which finds a linear transformation that maximizes the correlation between two high-dimensional representations. Raghu [67] defined each neuron activation vector as its response over a finite set of inputs, and the amount of data determines the dimension of the activation vector. For a given dataset $X = \{x_1, \cdots, x_m\}$, the activation vector of the neuron $i$ in the layer $l$ is defined as $z_i^l$ i.e., $z_i^l = (z_i^l(x_1), \cdots, z_i^l(x_m))$.

In CCA, enough data is essential. On the one hand, it can reduce the occurrence of spurious correlation, and on the other hand, it can also ensure the stability of the calculation. We mix the MuST-C En-De tst-COMMON and dev sets to increase the amount of data used to compute CCA. Additionally, we normalize the activation vector values to [-1, 1] as suggested by [68]. For each input, we average the representation over the time dimension as the neuron response to the corresponding input. We save the activations for each encoder layer for three types ST models, namely Baseline, $FCCL^s$, and Random models. The parameters of random model are entirely random. We use it as a comparison to exclude the influence of spurious correlations caused by insufficient data. The baseline model is the same as defined in Section V-E. We also save activations for the MT model. Then we compute the projection weighted CCA (PWCCA) [69] between activations from those layers. A high value indicates that the representations between layers are linearly related, implying that the layers capture similar information.

Figure 7 plots the PWCCA between different layers of different networks, and we can observe three exciting phenomena.

(1) For the Random model, it already has certain similarities with the MT model. On the one hand, the features extracted by Wav2vec 2.0 already contain rich acoustic information. On the other hand, the limited amount of data used to calculate PWCCA leads to spurious correlations.

(2) FCCL forces more layers to learn semantic information. The output of the sixth layer of the MT model encoder contains the high-level semantic information of the text, and the lower layer output contains grammatical structure information. In Baseline, only the last two layers are used to learn semantic information (only the fifth and sixth layers in the Baseline model have the max correlations with the sixth layer of the MT model). In contrast, the last four layers are used to learn semantic information in the $FCCL^s$ (the last four layers in $FCCL^s$ model have the max correlations with the sixth layer of the MT model).

(3) Compared with the Baseline, FCCL can effectively improve the correlation between each layer and the semantic representation of the text (the first row of $FCCL^s$ has significantly higher values than the Baseline). Since a network capacity is limited, $FCCL^s$ can free up its capacity from learning grammatical structure information and force more layers to learn semantic information. In this way, the model can extract representations containing more semantic information, and its performance also gets improvement. This proves that
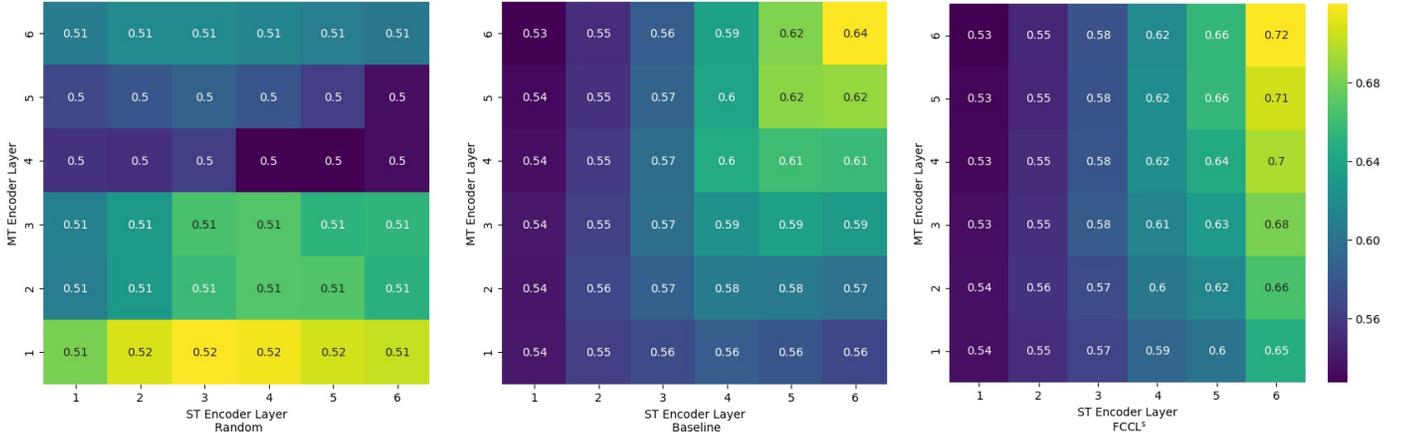
Fig. 7.  PWCCA between encoder states of different layers of the ST and the MT model.

our method can effectively bridge the gap between the two modalities, validating our model design.

### G. Is the maximum similarity method reasonable?

When conducting fine granularity contrastive learning, we need to get the correspondence between the speech frames and text tokens. To avoid the use of an extra alignment model, we propose a maximum similarity method. Two natural questions are whether this approach makes sense and whether it incurs additional computing costs because it is online.

**We answer the first question by performing fine granularity visualization.** We randomly select paired speech-transcription from MuST-C En-De tst-COMMON set, then calculate the similarity matrix $\Delta \in {}^{T_x \times T_z}$ between speech and transcription representation. As shown in Figure 8, the overall correspondence is monotonic. Some non-monotonic alignment can be attributed to the presence of silence and noise frames in speech. The model learns to correspond these frames to unique text token.

**We answer the second question by computing the FLOPs.** The FLOPs with and without fine-grained contrastive learning are 21.67G and 21.66G, respectively. With an increase of only 0.01G FLOPs, fine granularity contrastive learning improves the performance by 0.45 BLEU (as shown in Table V). In general, by adopting this method, we can effectively and efficiently find the correspondence between speech frame and text token in an unsupervised manner with negligible latency overhead.

### H. Visualization of Coarse Granularity Alignment

We randomly select 30 speech-transcription pairs from MuST-C En-De tst-COMMON set, and then apply T-SNE [70] to the vector representations of these samples to reduce the dimension to two. Note that these vector representations are obtained by averaging the encoder outputs over the time dimension.

The results are visualized in Figure 9. Each speech-transcription pair is connected by a solid line. It can be
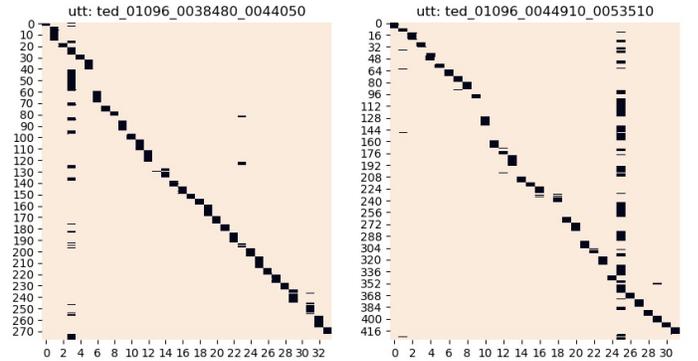


Fig. 8.  Visualization of similarity matrix after masking.

intuitively seen from the figure that most paired speech-transcription are projected together, and some even overlap with each other. This proves that FCCL is capable of bridging the representation divergence of the two modalities. In addition, some speech representations in the figure are still clustered together, mainly because we compute the contrastive loss across modalities and not within the modal. Thus, the speech representations do not show good uniformity.

### I. Visualization of Model Learning Dynamic

Figure 10 shows the learning dynamics. The blue and yellow curves represent the accuracy of the Baseline and FCCL$^s$ on the validation set, respectively. Although the FCCL$^s$ achieves better performance, its performance in the early stage is worse than the Baseline. One reason is that we use dropdim to increase the learning difficulty of the model. In addition, the quality of the speech representation in the early stage is poor, so the maximum similarity method cannot find the correct correspondence. A possible solution is to discard the fine granularity contrastive loss in the early stage of training and add it until the ST model has a specific representation ability.
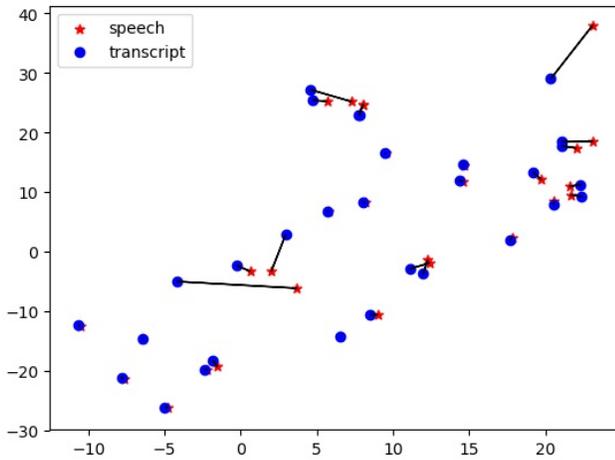
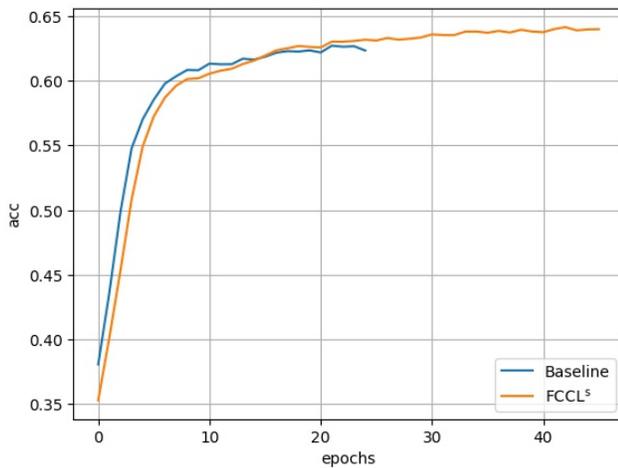Fig. 9.  Visualization the sentence-level representation.



Fig. 10.  Visualization of training process.

## VI. CONCLUSION AND FUTURE WORK

In this paper, we propose a cross-modal multi-grained contrast learning method, FCCL, for *explicit knowledge transfer* from the MT to the ST model. In addition, we propose whitening to solve the representation degeneration of text representation in the MT model. Experiments on the MuST-C dataset in all 8 languages demonstrate the effectiveness of our method. Additional experiment analysis and visualization show that FCCL is capable of bridging the speech-text representation gap and exhibits stronger robustness.

Although our method exhibits the desired effect, it relied on transcription to guide the extract of speech representation during training. For the more than 7000 languages and dialects worldwide, most of them do not have corresponding translations or even transcriptions, and our method does not work in such scenarios. More effective strategies to improve the quality of E2E-ST need to be explored in the future.

## REFERENCES

[1] H. Ney, "Speech translation: Coupling of recognition and translation," in *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 1.  IEEE, 1999, pp. 517–520.

[2] L. Mathias and W. Byrne, "Statistical phrase-based speech translation," in *Proceedings of IEEE International Conference on Acoustics Speech and Signal Processing*, vol. 1.  IEEE, 2006, pp. I–I.

[3] A. Bérard, L. Besacier, A. C. Kocabiyikoglu, and O. Pietquin, "End-to-end automatic speech translation of audiobooks," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2018, pp. 6224–6228.

[4] M. Sperber, G. Neubig, J. Niehues, and A. Waibel, "Attention-passing models for robust and data-efficient end-to-end speech translation," *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 313–325, 2019.

[5] Y. Liu, H. Xiong, Z. He, J. Zhang, H. Wu, H. Wang, and C. Zong, "End-to-end speech translation with knowledge distillation," *arXiv preprint arXiv:1904.08075*, 2019.

[6] Q. Dong, R. Ye, M. Wang, H. Zhou, S. Xu, B. Xu, and L. Li, "Listen, understand and translate: Triple supervision decouples end-to-end speech-to-text translation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 14, 2021, pp. 12 749–12 759.

[7] Y. Liu, J. Zhu, J. Zhang, and C. Zong, "Bridging the modality gap for speech-to-text translation," *arXiv preprint arXiv:2010.14920*, 2020.

[8] Y. Du, Z. Zhang, W. Wang, B. Chen, J. Xie, and T. Xu, "Regularizing end-to-end speech translation with triangular decomposition agreement," in *arXiv preprint arXiv:2112.10991*, 2021.

[9] Q. Fang, R. Ye, L. Li, Y. Feng, and M. Wang, "Stemm: Self-learning with speech-text manifold mixup for speech translation," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022, pp. 7050–7062.

[10] R. Ye, M. Wang, and L. Li, "End-to-end speech translation via cross-modal progressive training," *arXiv preprint arXiv:2104.10380*, 2021.

[11] C. Xu, B. Hu, Y. Li, Y. Zhang, S. Huang, Q. Ju, T. Xiao, and J. Zhu, "Stacked acoustic-and-textual encoding: Integrating the pre-trained models into speech translation encoders," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, 2021, pp. 2619–2630.

[12] C. Han, M. Wang, H. Ji, and L. Li, "Learning shared semantic space for speech-to-text translation," in *Findings of the Association for Computational Linguistics: ACL-IJCNLP*, 2021, pp. 2214–2225.

[13] Y. Tang, J. Pino, X. Li, C. Wang, and D. Genzel, "Improving speech translation by understanding and learning from the auxiliary text translation task," *arXiv preprint arXiv:2107.05782*, 2021.

[14] A. Bérard, O. Pietquin, C. Servan, and L. Besacier, "Listen and translate: A proof of concept for end-to-end speech-to-text translation," *arXiv preprint arXiv:1612.01744*, 2016.

[15] L. Duong, A. Anastasopoulos, D. Chiang, S. Bird, and T. Cohn, "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, pp. 949–959.

[16] A. C. Kocabiyikoglu, L. Besacier, and O. Kraif, "Augmenting librispeech with french translations: A multimodal corpus for direct speech translation evaluation," *arXiv preprint arXiv:1802.03142*, 2018.

[17] Y. Jia, M. Johnson, W. Macherey, R. J. Weiss, Y. Cao, C.-C. Chiu, N. Ari, S. Laurenzo, and Y. Wu, "Leveraging weakly supervised data to improve end-to-end speech-to-text translation," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.  IEEE, 2019, pp. 7180–7184.

[18] J. Pino, L. Puzon, J. Gu, X. Ma, A. D. McCarthy, and D. Gopinath, "Harnessing indirect training data for end-to-end automatic speech translation: Tricks of the trade," in *Proceedings of the 16th International Conference on Spoken Language Translation*, 2019.

[19] R. J. Weiss, J. Chorowski, N. Jaitly, Y. Wu, and Z. Chen, "Sequence-to-sequence models can directly translate foreign speech," *arXiv preprint arXiv:1703.08581*, 2017.

[20] Y. Liu, J. Zhang, H. Xiong, L. Zhou, Z. He, H. Wu, H. Wang, and C. Zong, "Synchronous speech recognition and speech-to-text translation with interactive decoding," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 05, 2020, pp. 8417–8424.

[21] S. Bansal, H. Kamper, K. Livescu, A. Lopez, and S. Goldwater, "Pre-training on high-resource speech recognition improves low-resource speech-to-text translation," in *2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, 2019, pp. 58–68.

[22] M. C. Stoian, S. Bansal, and S. Goldwater, "Analyzing asr pretraining for low-resource speech-to-text translation," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.  IEEE, 2020, pp. 7909–7913.

[23] C. Wang, Y. Wu, S. Liu, M. Zhou, and Z. Yang, "Curriculum pre-training for end-to-end speech translation," *arXiv preprint arXiv:2004.10093*, 2020.

[24] Z. Zhang, L. Zhou, J. Ao, S. Liu, L. Dai, J. Li, and F. Wei, "Speechut: Bridging speech and text with hidden-unit for encoder-decoder based speech-text pre-training," *arXiv preprint arXiv:2210.03730*, 2022.

[25] Y. Tang, H. Gong, N. Dong, C. Wang, W.-N. Hsu, J. Gu, A. Baevski, X. Li, A. Mohamed, M. Auli *et al.*, "Unified speech-text pre-training for speech translation and recognition," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 1488–1499.

[26] J. Pino, Q. Xu, X. Ma, M. J. Dousti, and Y. Tang, "Self-training for end-to-end speech translation," *Proc. Interspeech 2020*, pp. 1476–1480, 2020.

[27] S. Indurthi, H. Han, N. K. Lakumarapu, B. Lee, I. Chung, S. Kim, and C. Kim, "Data efficient direct speech-to-text translation with modality agnostic meta-learning," *arXiv preprint arXiv:1911.04283*, 2019.

[28] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9729–9738.

[29] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proceedings of International conference on machine learning*. PMLR, 2020, pp. 1597–1607.

[30] T. Gao, X. Yao, and D. Chen, "Simcse: Simple contrastive learning of sentence embeddings," *arXiv preprint arXiv:2104.08821*, 2021.

[31] Y. Yan, R. Li, S. Wang, F. Zhang, W. Wu, and W. Xu, "Consert: A contrastive framework for self-supervised sentence representation transfer," *arXiv preprint arXiv:2105.11741*, 2021.

[32] T. Z. W. Ye, B. Yang, L. Zhang, X. Ren, D. Liu, J. Sun, S. Zhang, H. Zhang, and W. Zhao, "Frequency-aware contrastive learning for neural machine translation," *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.

[33] L. Wang and A. v. d. Oord, "Multi-format contrastive learning of audio representations," *arXiv preprint arXiv:2103.06508*, 2021.

[34] Z.-Q. Zhang, Y. Song, M.-H. Wu, X. Fang, and L.-R. Dai, "Xlst: Cross-lingual self-training to learn multilingual representation for low resource speech recognition," *arXiv preprint arXiv:2103.08207*, 2021.

[35] A. Xiao, C. Fuegen, and A. Mohamed, "Contrastive semi-supervised learning for asr," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2021, pp. 3870–3874.

[36] J.-B. Alayrac, A. Recasens, R. Schneider, R. Arandjelović, J. Rama-puram, J. De Fauw, L. Smaira, S. Dieleman, and A. Zisserman, "Self-supervised multimodal versatile networks," *Proceedings of Advances in Neural Information Processing Systems*, vol. 33, pp. 25–37, 2020.

[37] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *Proceedings of International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.

[38] H.-H. Wu, P. Seetharaman, K. Kumar, and J. P. Bello, "Wav2clip: Learning robust audio representations from clip," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2022, pp. 4563–4567.

[39] R. Ye, M. Wang, and L. Li, "Cross-modal contrastive learning for speech translation," in *Proc. of NAACL*, 2022.

[40] X. Wang, R. Zhang, C. Shen, T. Kong, and L. Li, "Dense contrastive learning for self-supervised visual pre-training," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3024–3033.

[41] Y. Zeng, X. Zhang, and H. Li, "Multi-grained vision language pre-training: Aligning texts with visual concepts," *arXiv preprint arXiv:2111.08276*, 2021.

[42] D. Wang and S. Karout, "Fine-grained multi-modal self-supervised learning," *arXiv preprint arXiv:2112.12182*, 2021.

[43] Q. Dong, Y. Zhu, M. Wang, and L. Li, "Unist: Unified end-to-end model for streaming and non-streaming speech translation," *arXiv preprint arXiv:2109.07368*, 2021.

[44] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Proceedings of Advances in Neural Information Processing Systems*, vol. 33, pp. 12 449–12 460, 2020.

[45] B. Li, H. Zhou, J. He, M. Wang, Y. Yang, and L. Li, "On the sentence embeddings from pre-trained language models," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 9119–9130.

[46] J. Gao, D. He, X. Tan, T. Qin, L. Wang, and T.-Y. Liu, "Representation degeneration problem in training natural language generation models," *arXiv preprint arXiv:1907.12009*, 2019.

[47] R. Cao, Y. Wang, Y. Liang, L. Gao, J. Zheng, J. Ren, and Z. Wang, "Exploring the impact of negative samples of contrastive learning: A case study of sentence embedding," in *Findings of the Association for Computational Linguistics*, 2022, pp. 3138–3152.

[48] J. Huang, D. Tang, W. Zhong, S. Lu, L. Shou, M. Gong, D. Jiang, and N. Duan, "Whiteningbert: An easy unsupervised sentence embedding approach," in *Findings of the Association for Computational Linguistics: EMNLP 2021*, 2021, pp. 238–244.

[49] J. Su, J. Cao, W. Liu, and Y. Ou, "Whitening sentence representations for better semantics and faster retrieval," *arXiv preprint arXiv:2103.15316*, 2021.

[50] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur, "Audio augmentation for speech recognition." in *Proceedings of INTERSPEECH*, 2015, pp. 3586–3589.

[51] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "Specaugment: A simple data augmentation method for automatic speech recognition," *arXiv preprint arXiv:1904.08779*, 2019.

[52] H. Zhang, D. Qu, K. Shao, and X. Yang, "Dropdim: A regularization method for transformer networks," *IEEE Signal Processing Letters*, vol. 29, pp. 474–478, 2022.

[53] M. A. Di Gangi, R. Cattoni, L. Bentivogli, M. Negri, and M. Turchi, "Must-c: a multilingual speech translation corpus," in *Proceedings of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2019, pp. 2012–2017.

[54] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," *arXiv preprint arXiv:1508.07909*, 2015.

[55] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An asr corpus based on public domain audio books," in *Proceedings of International Conference on Acoustics, Speech and Signal Processing*, 2015, pp. 5206–5210.

[56] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Proceedings of Advances in neural information processing systems*, vol. 30, 2017.

[57] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[58] H. Inaguma, S. Kiyono, K. Duh, S. Karita, N. Yalta, T. Hayashi, and S. Watanabe, "ESPnet-ST: All-in-one speech translation toolkit," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Online: Association for Computational Linguistics, Jul. 2020, pp. 302–311. [Online]. Available: https://www.aclweb.org/anthology/2020.acl-demos.34

[59] M. Post, "A call for clarity in reporting BLEU scores," in *Proceedings of the Third Conference on Machine Translation: Research Papers*. Belgium, Brussels: Association for Computational Linguistics, Oct. 2018, pp. 186–191. [Online]. Available: https://www.aclweb.org/anthology/W18-6319

[60] B. Zhang, I. Titov, B. Haddow, and R. Sennrich, "Adaptive feature selection for end-to-end speech translation," *arXiv preprint arXiv:2010.08518*, 2020.

[61] H. Le, J. Pino, C. Wang, J. Gu, D. Schwab, and L. Besacier, "Dual-decoder transformer for joint automatic speech recognition and multi-lingual speech translation," *arXiv preprint arXiv:2011.00747*, 2020.

[62] M. Rofin, N. Balagansky, and D. Gavrilov, "Linear interpolation in parameter space is good enough for fine-tuned language models," *arXiv preprint arXiv:2211.12092*, 2022.

[63] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, "Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7871–7880.

[64] Y. Liu, J. Gu, N. Goyal, X. Li, S. Edunov, M. Ghazvininejad, M. Lewis, and L. Zettlemoyer, "Multilingual denoising pre-training for neural machine translation," *Transactions of the Association for Computational Linguistics*, vol. 8, pp. 726–742, 2020.

[65] L. Ding, K. Peng, and D. Tao, "Improving neural machine translation by denoising training," *arXiv preprint arXiv:2201.07365*, 2022.

[66] K. Lee, "Hallucinations in neural machine translation," in *Proceedings of International Conference on Learning Representations*, 2019.

[67] M. Raghu, J. Gilmer, J. Yosinski, and J. Sohl-Dickstein, "Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability," *Proceedings of Advances in neural information processing systems*, vol. 30, 2017.

[68] N. Kambhatla, L. Born, and A. Sarkar, "Cipherdaug: Ciphertext based data augmentation for neural machine translation," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022, pp. 201–218.

[69] A. Morcos, M. Raghu, and S. Bengio, "Insights on representational similarity in neural networks with canonical correlation," *Proceedings of Advances in Neural Information Processing Systems*, vol. 31, 2018.

[70] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of machine learning research*, vol. 9, no. 11, 2008.