

COLLD: CONTRASTIVE LAYER-TO-LAYER DISTILLATION FOR COMPRESSING MULTILINGUAL PRE-TRAINED SPEECH ENCODERS

Heng-Jui Chang^{1,2,*}, Ning Dong², Ruslan Mavlyutov², Sravya Popuri², Yu-An Chung²

¹MIT CSAIL ²Meta AI

hengjui@mit.edu, andyyuan@meta.com

ABSTRACT

Large-scale self-supervised pre-trained speech encoders outperform conventional approaches in speech recognition and translation tasks. Due to the high cost of developing these large models, building new encoders for new tasks and deploying them to on-device applications are infeasible. Prior studies propose model compression methods to address this issue, but those works focus on smaller models and less realistic tasks. Thus, we propose Contrastive Layer-to-layer Distillation (CoLLD), a novel knowledge distillation method to compress pre-trained speech encoders by leveraging masked prediction and contrastive learning to train student models to copy the behavior of a large teacher model. CoLLD outperforms prior methods and closes the gap between small and large models on multilingual speech-to-text translation and recognition benchmarks.

Index Terms— Self-supervised learning, knowledge distillation, model compression, multilingual speech translation

1. INTRODUCTION

Self-supervised learning (SSL) for speech encoder pre-training benefits various speech processing tasks and outperforms conventional approaches [1]. SSL methods leverage large unlabeled speech corpus to train deep neural networks to encode useful representations and succeed in applications like speech translation [2] and automatic speech recognition (ASR) [3]. However, powerful speech encoders usually have many parameters, making real-time or on-device speech processing less feasible.

Researchers propose model compression techniques to address the issues of large speech encoders. The compressed SSL pre-trained encoders can be applied to various downstream tasks. These approaches can be categorized into knowledge distillation (KD) and parameter pruning. In KD, a lightweight student model learns to predict hidden representations to mimic the large teacher model’s behavior [4–10]. DistilHuBERT [4] predicts multiple hidden layers in a HuBERT teacher [11] using the student’s output with separate prediction heads. FitHuBERT [5] and Ashihara et al. [6] propose layer-to-layer (L2L) KD that uses narrow and deep students to layer-wise distill the teacher’s hidden representations. In unstructured pruning, parameters with small values are set to zero [12], while structured pruning removes submodules from a model [13–15] to reduce the parameters but requires complicated implementation. Other studies combine the above methods [16] or techniques like layer-skipping [17] and low-bit quantization [18].

Although existing methods succeed in many tasks, most works focus on compressing small SSL models and evaluating with unrealistic problem setups. Those works compress a HuBERT Base [11]

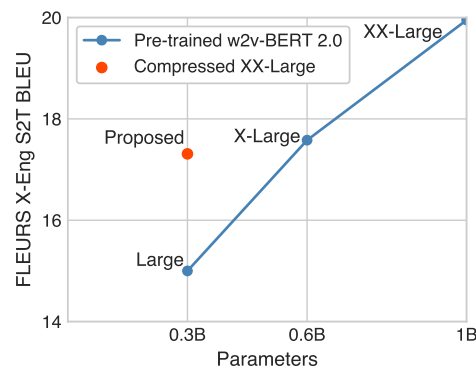


Fig. 1. Encoder sizes vs. X-Eng speech-to-text translation BLEU scores. The proposed model is a compressed XX-Large model.

model (95M parameters) to models around 20M to 30M parameters and evaluate with the Speech processing Universal PERFORMANCE Benchmark (SUPERB) [19, 20]. These compressed models are unsuitable for complex tasks that require fine-tuning because of the small model capacities, limiting application scenarios. Under this setting, the effectiveness of these methods for large-scale models and problems remains to be discovered.

To bridge the gap between academic research and real-world problems, we extend the speech encoder compression task to a large-scale pre-trained speech encoder (w2v-BERT 2.0 [2]) and apply the compressed model to multilingual speech-to-text translation (S2T). This problem is challenging because the original model is significantly larger (1B parameters), and the compressed model is fine-tuned with a more complicated yet realistic task. Following previous studies, we use unlabeled data to compress an SSL pre-trained teacher model because this setup allows flexible utilization and avoids fine-tuning huge encoders. Moreover, the compressed encoder has 300M parameters, which is currently the largest encoder size widely used in both production and academia [19].

Under this new problem setting, we propose Contrastive Layer-to-layer Distillation (CoLLD) by combining L2L KD [6] and a contrastive masked prediction learning objective [21]. First, some student model input frames are masked while the teacher remains unmasked. Then, each masked student’s hidden layer frame classifies the corresponding teacher’s hidden layer frame from a set of distractors, where the distractors are randomly sampled from other frames of the teacher’s representations. After distillation, we evaluate the student model with internal and public benchmarks, covering S2T and multilingual ASR. As shown in Fig. 1 and Sec. 3, CoLLD surpasses prior distillation methods, narrows the performance gap between large models (0.6B and 1.0B parameters) and outperforms strong baselines like XLS-R [22] and MMS [23].

* Work done during an internship at Meta AI.

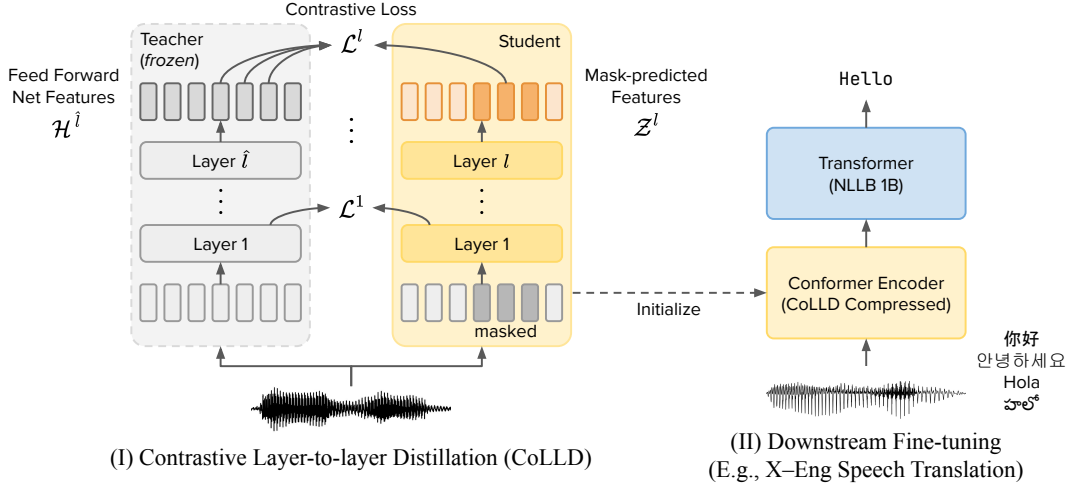


Fig. 2. An illustration of the proposed Contrastive Layer-to-layer Distillation (CoLLD) framework. (I) CoLLD feeds the same input to a frozen teacher and a learnable student model, where the student’s input frames are partially masked. For each student layer l , the masked representations learn to classify the corresponding teacher frame in layer \hat{l} from K distractor frames. (II) After distillation, the student model weights initialize downstream models and are fine-tuned with labeled data to perform tasks like multilingual speech translation.

2. METHOD

2.1. Overview

We propose the Contrastive Layer-to-layer Distillation (CoLLD) framework as shown in Fig. 2. First, the student’s layers are trained to predict teacher hidden layer representations (Sec. 2.2). Next, we incorporate masked prediction to encourage the student model to learn better representations (Sec. 2.3). Finally, a contrastive learning objective prevents the model from collapsing. (Sec. 2.4).

2.2. Layer-to-layer Distillation

Moreover, as Ashihara et al. [6] pointed out, deep and narrow student models better capture the teacher’s behavior. We follow [5] and [6] by assigning each student layer to predict a teacher’s hidden layer. The student-to-teacher layer mapping is obtained as follows. Let L^T and L^S as the numbers of teacher and student layers, with $L^T \geq L^S$. The l^{th} student layer learns to predict the \hat{l}^{th} teacher layer, where

$$\hat{l} = \text{round} \left((l - 1) \frac{L^T - 1}{L^S - 1} \right) + 1, \quad (1)$$

for $l = 1, 2, \dots, L^S$. Each student layer is assigned to predict a unique teacher layer, and the selected layers are uniformly distributed across the teacher model. This mapping rule allows flexible student architectures for different applications.

Previous works distill the final output of each teacher layer [4,5]. Inspired by data2vec [24], we let the student model predict each teacher layer’s feed-forward net (FFN) features for better learning targets. Specifically, the student learns from the outputs of the second FFN of each Conformer block in the teacher [25].

2.3. Masked Prediction

Prior KD methods usually keep the student’s inputs unmasked [4,5], but many SSL methods rely on masked language modeling [11, 21, 24], and studies have shown this technique useful for knowledge distillation [7,9]. Therefore, we only mask the student’s input frames and apply L2L distillation to the masked frames.

2.4. Contrastive Distillation Objective

We found that utilizing L1 or L2 losses for KD sometimes leads to collapsed representations when incorporating masked prediction if the hyperparameters are not carefully tuned. Hence, we propose a contrastive learning objective to mitigate this issue [21, 26]. For each masked timestep $t \in \mathcal{T}$ in an utterance, the student’s l^{th} layer output z_t^l predicts the \hat{l}^{th} teacher layer representation $h_t^{\hat{l}}$. The student minimizes the distance between z_t^l and $h_t^{\hat{l}}$. The conventional L2 regression loss is written as

$$\mathcal{L}^l = \sum_{t \in \mathcal{T}} \left\| z_t^l - h_t^{\hat{l}} \right\|_2^2, \quad (2)$$

while the proposed contrastive distillation objective is

$$\mathcal{L}^l = - \sum_{t \in \mathcal{T}} \log \frac{\exp \left(\cos \left(z_t^l, h_t^{\hat{l}} \right) / \tau \right)}{\sum_{h' \in \mathcal{H}_t^{\hat{l}}} \exp \left(\cos \left(z_t^l, h' \right) / \tau \right)}, \quad (3)$$

where $\mathcal{H}_t^{\hat{l}}$ is a set composed of $h_t^{\hat{l}}$ and K distractors [26] sampled from the \hat{l}^{th} teacher layer with indices also in \mathcal{T} . $\tau > 0$ is a hyperparameter and $\cos(\cdot, \cdot)$ denotes cosine similarity. With this objective, the model is expected to avoid collapsing.

3. EXPERIMENTS

3.1. Setup

3.1.1. Model

All experiments are based on w2v-BERT 2.0 [2], a series of SSL speech encoders trained with contrastive learning [21] and masked language modeling [11]. The Conformer [25] architectures and forward computing costs are listed in Table 2. A depth-wise convolution kernel size of 31 is used. Each model takes 80-dimensional filter bank features as input and downsamples each utterance by concatenating consecutive frames to reduce the frame rate from 100Hz to 50Hz. Excluding Large₄₀, all w2v-BERT 2.0 models are pre-trained from scratch with an internal corpus containing 4M hours of unlabeled speech, covering 143+ languages. Unless stated otherwise, students are randomly initialized Large₄₀ or Large₁₂ models that distill knowledge from the XX-Large teacher.

Table 1. BLEU scores of multilingual speech-to-text translation (X-Eng S2T) evaluated on CoVoST 2 [27] and FLEURS 101 languages test set [28]. Excluding pre-trained from scratch topline, each model has 0.3B parameters. Avg indicates an averaged score across all languages.

Method	CoVoST 2				FLEURS-101									
	High	Mid	Low	Avg	WE	EE	CMN	SSA	SA	SEA	CJK	Avg		
Pre-trained w2v-BERT 2.0														
XX-Large (1.0B Teacher)	37.6	35.8	29.7	32.6	27.5	25.8	20.2	10.0	19.1	16.2	14.3	20.0		
X-Large (0.6B)	36.4	34.0	28.0	31.0	25.7	23.7	17.0	8.0	16.4	13.9	9.9	17.6		
Large ₁₂ (0.3B)	33.5	31.3	23.7	27.4	22.4	20.7	14.3	7.1	12.9	11.4	8.1	15.0		
Layer Removal from Pre-trained XX-Large														
Layer Skipping Large ₁₂	29.6	26.4	15.9	21.0	15.6	14.0	9.1	3.5	8.2	7.1	4.2	9.7		
Bottom Layers Large ₁₂	31.4	28.9	21.8	25.3	20.0	18.1	12.4	5.9	10.8	10.2	6.9	13.1		
CoLLD from Pre-trained XX-Large														
Large ₁₂	34.6	32.8	25.5	29.0	24.3	23.2	17.3	7.5	16.2	14.0	12.3	17.2		
Large ₄₀	35.4	33.6	26.9	30.1	24.0	23.3	17.9	7.6	16.6	14.0	11.8	17.3		
Loss	Contrastive (Eq. 3) → L2 (Eq. 2)		33.7	32.2	25.2	28.5	22.6	22.0	16.8	6.9	15.1	13.3	11.2	16.2
Target	FFN → Layer Output		34.1	32.2	25.5	28.7	22.7	22.0	16.8	7.4	15.9	13.3	10.9	16.4
Data	Multilingual → Monolingual		33.7	31.7	24.8	28.2	22.4	21.4	15.1	5.8	13.8	11.8	9.5	15.2
Initialization	Random → Layer Skipping Large ₁₂		30.1	27.3	18.0	22.5	19.1	17.6	11.7	4.8	10.7	9.1	6.9	12.4
Initialization	Random → Bottom Layers Large ₁₂		29.1	25.8	17.8	21.8	18.6	17.0	11.3	4.9	10.3	8.7	6.3	12.0
CoLLD from Fine-tuned XX-Large														
Large ₄₀ + S2T Fine-tuned XX-Large Teacher	36.2	34.5	27.9	31.0	25.2	24.6	19.1	8.2	18.2	15.4	13.7	18.5		

Table 2. w2v-BERT 2.0 [2] architectures with different dimensions, feed-forward net sizes (FFN), and attention heads. The number of parameters (Param) and multiply-accumulate operation (MACs) during forward pass indicate required spaces and computation costs. MACs are calculated with an input utterance of 20 seconds long.

Model	Param	Dim	FFN	Layer	Head	GMACs \downarrow
XX-Large	1.0B	1024	4096	40	16	1214.1
X-Large	0.6B	1024	4096	24	16	728.5
Large ₁₂	0.3B	1024	4096	12	16	364.3
Large ₄₀	0.3B	768	1024	40	8	462.3

MACs computation: <https://github.com/zhijian-liu/torchprofile>

3.1.2. Knowledge Distillation

We implement experiments with fairseq [32]. Only 92k hours of audio data in the 4M hours corpus are used for distillation because KD requires fewer updates than pre-training, where the amount of used training data is calculated according to [33]. We set $\tau = 0.1$ and $K = 100$ in Eq. 3. Downsampled features are randomly masked with a span of 10 frames and a probability of 0.065, resulting in approximately 49% of masked frames. Each model is trained with 200k updates using an Adam optimizer [34] with a peak learning rate of 10^{-4} , $\beta_1 = 0.9$, $\beta_2 = 0.98$, $\epsilon = 10^{-6}$, and a weight decay of 10^{-2} . The learning rate ramps up linearly in the first 4k updates and linearly decays to 0 for the rest. Each model is compressed on 32 NVIDIA A100 80GB GPUs, with an effective batch size of 27.7 minutes of audio data in each update. Large₁₂ and Large₄₀ students take 2 and 4 days to distill from the XX-Large teacher. Although the parameters of 0.3B models are similar, the distillation time of the 40-layer student is higher because the forward operation of each hidden layer cannot be parallelized. Some prior KD methods are not included for comparison because they require complex implementation and hyperparameter search.

3.1.3. Multilingual Speech Translation

The speech-to-English-text translation (X-Eng S2T) model comprises a Conformer encoder, a length adaptor [35], and a 1.3B-parameter NLLB-200 machine translation model [36]. The fine-tuning data include approximately 60k hours of paired speech and translation text that cover 88 X-English directions. The Conformer encoder is fine-tuned entirely, but only the layer norm and self-attention for NLLB. The learning rate linearly increases to 10^{-4} in the first 5k updates (2×10^{-4} for XX-Large), and then follows the inverse square root schedule [37]. All models are trained with an effective batch size of 64 minutes of audio and 150k updates. We use 16 to 64 NVIDIA V100 32GB GPUs, depending on the model size. We evaluate fine-tuned S2T models on CoVoST 2 [27] and FLEURS [28] with a decoding beam size of 5.

3.2. Fine-tuning Results

This section reveals the effectiveness of CoLLD through fine-tuning w2v-BERT 2.0 models on X-Eng S2T. As shown in Table 1, we offer three pre-trained from scratch w2v-BERT 2.0 models, where the 0.3B model is served as a baseline. Layer removal baselines preserve 30% of the layers of the XX-Large model by either preserving the bottom layers or uniformly skipping layers following Eq. 1.

CoLLD Large₄₀ surpasses 0.3B baselines by at least two BLEUs in most subsets, indicating that the student successfully acquires knowledge from the XX-Large teacher. Although CoLLD Large₄₀ is incapable of reaching the same performance as the 1.0B teacher because of the model capacity, the gap between the 0.3B and 0.6B models is significantly reduced. Especially in FLEURS, CoLLD offers slightly superior BLEU scores in most subsets compared with the 0.6B topline. Hence, CoLLD Large₄₀ is comparable with the X-Large w2v-BERT 2.0 but requires only half of the parameters.

We offer ablation studies in the same table. The overall S2T performance is degraded by replacing each of the proposed components in CoLLD with prior methods, indicating the necessity of

Table 3. SSL pre-trained models with 0.3B parameters on the 10-minute set of the ML-SUPERB benchmark [29]. The metrics include accuracy (Acc%), character error rate (CER%), phone error rate (PER%), and SUPERB score (SUPERB_s) [30].

SSL Model	Pre-training / Distillation Data		Mono-ASR CER/PER↓	Multi-ASR		LID Normal Acc↑	Multi-ASR + LID		SUPERB _s ↑	
	#Hours	#Langs		Normal	Few-shot		Normal	Few-shot		
			CER↓	CER↓	Acc↑	CER↓	CER↓			
No Compression Baseline										
XLSR 53 [31]	56k	53	49.5	33.9	43.6	6.6	45.6	33.4	43.2	403.4
XLS-R 128 [22]	400k	128	39.7	29.2	40.9	66.9	55.6	28.4	42.1	734.1
MMS [23]	491k	1406	33.8	28.7	36.5	62.3	71.9	31.5	30.9	829.1
w2v-BERT 2.0 Large ₁₂	4M	143+	46.6	27.2	32.2	37.0	78.5	27.2	31.7	698.8
Proposed										
CoLLD Large ₄₀	92k	143+	35.5	22.2	29.6	82.8	85.7	21.9	28.7	988.7

the design of CoLLD. First, a shallow and wide student architecture (Large₁₂) drops one BLEU score in most test sets compared with the deeper model (Large₄₀), corroborating with prior studies [5, 6]. Still, Large₁₂ outperforms all baselines, and the fine-tuning and inference costs of the shallow model are lower than those of the deep model. Therefore, the choice between shallow and deep models depends on the application scenario. Second, optimizing with L2 loss or learning from each teacher layer’s output leads to 1 to 2 BLEU score degradation, showing that the proposed techniques distill better representations from the teacher. Third, replacing distillation data with a 1k hours English speech corpus decreases BLEU scores but performs better than the baselines, implying that CoLLD still works even when the training data diversity is reduced. Furthermore, initializing student models with some teacher layers results in significantly worse scores, so model initialization is unnecessary. Note that we do not compare with DistilHuBERT because prior works have shown L2L KD has superior performance [5, 6]. The ablation studies clearly show the importance of the proposed CoLLD.

To push the limit of CoLLD, we consider distilling from an S2T fine-tuned teacher for comparison. In the last part of Table 1, the results of a CoLLD Large₄₀ model distilled from an S2T fine-tuned XX-Large teacher are reported. This compressed model offers superior performance compared with the 0.6B topline in many evaluation subsets, showing that CoLLD is applicable to both pre-trained and fine-tuned w2v-BERT 2.0 models. Thus, if a teacher model fine-tuned with labeled data is available, CoLLD produces better-compressed models. Overall, CoLLD successfully compresses a pre-trained XX-Large w2v-BERT 2.0 by 70% while retaining good X-Eng S2T performance.

3.3. Multilingual SUPERB

This section evaluates CoLLD with Multilingual SUPERB (ML-SUPERB) [29], a standard multilingual speech processing benchmark, to offer a more comprehensive comparison with other SSL models. ML-SUPERB covers 143 languages and four tasks: monolingual ASR (Mono-ASR), multilingual ASR (Multi-ASR), language identification (LID), and Multi-ASR + LID. We use the 10-minute set of ML-SUPERB to show the performance of pre-trained models in a low-resource setting. For a fair comparison, the pre-trained and distilled models are frozen and serve as feature extractors during downstream model training. We follow the implementation as in ESPnet [38].

As shown in Table 3, w2v-BERT 2.0 offers a solid baseline compared to prior works because this model is trained with significantly

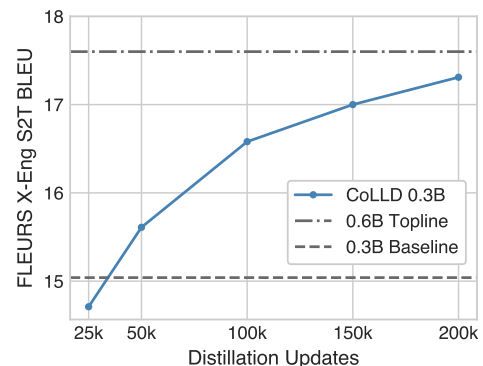


Fig. 3. Distillation updates vs. FLEURS-101 X-Eng BLEU scores.

more data. Next, CoLLD surpasses w2v-BERT 2.0 and other prior methods in most ML-SUPERB tasks and achieves the best overall SUPERB score by using only 92k hours of distillation data. The results again corroborate that CoLLD successfully distills knowledge from the XX-Large teacher.

3.4. Impact of Distillation Updates

This section investigates the impact of the data required for CoLLD by varying the total number of distillation updates. As shown in Fig. 3, CoLLD surpasses the 0.3B pre-trained from scratch baseline with only 50k of distillation updates. Meanwhile, when trained with 200k updates, CoLLD reaches a similar performance as the 0.6B topline model. Therefore, the amount of distillation data is highly correlated to downstream performance, and the distilled models offer better representations when more data and computation resources are available.

4. CONCLUSION

This paper proposes CoLLD, a novel model compression method by combining layer-to-layer knowledge distillation and contrastive learning for large-scale multilingual speech encoders. We show that CoLLD is superior over prior compression methods on multilingual speech recognition and speech-to-text translation by evaluating the proposed methods on internal and public benchmarks. This approach reduces model sizes of powerful pre-trained speech encoders while retaining good performance after fine-tuning, enabling on-device and streaming applications.

5. REFERENCES

- [1] A. Mohamed *et al.*, “Self-supervised speech representation learning: A review,” *IEEE JSTSP*, 2022.
- [2] Seamless Communication *et al.*, “Seamlessm4t—massively multilingual & multimodal machine translation,” *arXiv*, 2023.
- [3] Y. Zhang *et al.*, “Google usm: Scaling automatic speech recognition beyond 100 languages,” *arXiv*, 2023.
- [4] H.-J. Chang, S.-w. Yang, and H.-y. Lee, “DistilHuBERT: Speech representation learning by layer-wise distillation of hidden-unit bert,” in *ICASSP*, 2022.
- [5] Y. Lee, K. Jang, J. Goo, Y. Jung, and H. Kim, “Fithubert: Going thinner and deeper for knowledge distillation of speech self-supervised learning,” *Interspeech*, 2022.
- [6] T. Ashihara, T. Moriya, K. Matsuura, and T. Tanaka, “Deep versus wide: An analysis of student architectures for task-agnostic knowledge distillation of self-supervised speech models,” *Interspeech*, 2022.
- [7] R. Wang, Q. Bai, J. Ao, L. Zhou, Z. Xiong, Z. Wei, Y. Zhang, T. Ko, and H. Li, “Lighthubert: Lightweight and configurable speech representation learning with once-for-all hidden-unit bert,” *Interspeech*, 2022.
- [8] K.-P. Huang, T.-h. Feng, Y.-K. Fu, T.-Y. Hsu, P.-C. Yen, W.-C. Tseng, K.-W. Chang, and H.-y. Lee, “Ensemble knowledge distillation of self-supervised speech models,” in *ICASSP*, 2023.
- [9] K. Jang, S. Kim, S.-Y. Yun, and H. Kim, “Recycle-and-distill: Universal compression strategy for transformer-based speech ssl models with attention map reusing and masking distillation,” *Interspeech*, 2023.
- [10] H. Wang, S. Wang, W.-Q. Zhang, and J. Bai, “Distilxlr: A light weight cross-lingual speech representation model,” *Interspeech*, 2023.
- [11] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *TASLP*, vol. 29, 2021.
- [12] C.-I. J. Lai, Y. Zhang, A. H. Liu, S. Chang, Y.-L. Liao, Y.-S. Chuang, K. Qian, S. Khurana, D. Cox, and J. Glass, “PARP: Prune, adjust and re-prune for self-supervised speech recognition,” *NeurIPS*, 2021.
- [13] Y. Peng, K. Kim, F. Wu, P. Sridhar, and S. Watanabe, “Structured pruning of self-supervised pre-trained models for speech recognition and understanding,” in *ICASSP*, 2023.
- [14] H. Jiang, L. L. Zhang, Y. Li, Y. Wu, S. Cao, T. Cao, Y. Yang, J. Li, M. Yang, and L. Qiu, “Accurate and structured pruning for efficient automatic speech recognition,” *Interspeech*, 2023.
- [15] H. Wang, S. Wang, W.-Q. Zhang, H. Suo, and Y. Wan, “Task-agnostic structured pruning of speech representation models,” *Interspeech*, 2023.
- [16] Y. Peng, Y. Sudo, S. Muhammad, and S. Watanabe, “Dphubert: Joint distillation and pruning of self-supervised speech models,” *Interspeech*, 2023.
- [17] Y. Peng, J. Lee, and S. Watanabe, “I3d: Transformer architectures with input-dependent dynamic depth for speech recognition,” in *ICASSP*, 2023.
- [18] C.-F. Yeh, W.-N. Hsu, P. Tomasello, and A. Mohamed, “Efficient speech representation learning with low-bit quantization,” *arXiv*, 2022.
- [19] S.-w. Yang *et al.*, “SUPERB: Speech processing universal performance benchmark,” in *Interspeech*, 2021.
- [20] H.-S. Tsai *et al.*, “SUPERB-SG: Enhanced speech processing universal PERFORMANCE benchmark for semantic and generative capabilities,” in *ACL*, 2022.
- [21] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” in *NeurIPS*, 2020.
- [22] A. Babu *et al.*, “Xls-r: Self-supervised cross-lingual speech representation learning at scale,” *Interspeech*, 2022.
- [23] V. Pratap *et al.*, “Scaling speech technology to 1,000+ languages,” *arXiv*, 2023.
- [24] A. Baevski, W.-N. Hsu, Q. Xu, A. Babu, J. Gu, and M. Auli, “data2vec: A general framework for self-supervised learning in speech, vision and language,” in *ICML*, 2022.
- [25] A. Gulati *et al.*, “Conformer: Convolution-augmented transformer for speech recognition,” *Interspeech*, 2020.
- [26] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv*, 2018.
- [27] C. Wang, A. Wu, J. Gu, and J. Pino, “Covost 2 and massively multilingual speech translation,” in *Interspeech*, 2021.
- [28] A. Conneau, M. Ma, S. Khanuja, Y. Zhang, V. Axelrod, S. Dalmia, J. Riesa, C. Rivera, and A. Bapna, “Fleurs: Few-shot learning evaluation of universal representations of speech,” in *SLT*, 2023.
- [29] J. Shi *et al.*, “MI-superb: Multilingual speech universal performance benchmark,” *Interspeech*, 2023.
- [30] T.-h. Feng *et al.*, “Superb@ slt 2022: Challenge on generalization and efficiency of self-supervised speech representation learning,” in *SLT*, 2022.
- [31] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, “Unsupervised cross-lingual representation learning for speech recognition,” *Interspeech*, 2021.
- [32] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli, “fairseq: A fast, extensible toolkit for sequence modeling,” in *NAACL-HLT*, 2019.
- [33] H.-J. Chang, A. H. Liu, and J. Glass, “Self-supervised Fine-tuning for Improved Content Representations by Speaker-invariant Clustering,” in *Interspeech*, 2023.
- [34] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *ICLR*, 2015.
- [35] J. Zhao, H. Yang, E. Shareghi, and G. Haffari, “M-adapter: Modality adaptation for end-to-end speech-to-text translation,” *Interspeech*, 2022.
- [36] M. R. Costa-jussà *et al.*, “No language left behind: Scaling human-centered machine translation,” *arXiv*, 2022.
- [37] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *NIPS*, 2017.
- [38] S. Watanabe *et al.*, “Espnet: End-to-end speech processing toolkit,” *Interspeech*, 2018.

	fra	deu	spa	cat	pes	ita	rus	por	cmn	arb	cym	est	ind	jpn	khk	lvs	nld	slv	swe	tam	tur
w2vb2 XX-Large	39.4	36.0	38.5	36.3	24.9	36.4	51.1	48.8	17.8	40.4	51.1	24.8	44.6	16.4	6.4	26.6	37.8	35.1	37.1	4.6	31.4
w2vb2 X-Large	38.3	34.3	37.9	35.1	23.8	35.3	48.9	47.0	14.9	39.5	49.4	22.4	43.1	13.0	5.1	24.8	36.6	32.3	36.6	4.0	29.3
w2vb2 Large	35.5	30.9	35.8	31.9	21.9	32.6	45.0	43.3	13.5	34.5	44.0	17.3	38.8	9.5	3.6	21.4	32.0	27.4	29.2	2.9	23.4
CoLLD Large ₄₀	37.0	33.0	36.7	34.8	23.3	34.0	47.9	45.5	17.1	37.7	47.6	21.5	43.0	14.9	6.3	24.1	34.1	32.5	30.3	3.7	27.4

Fig. 4. Complete BLEU scores on the CoVoST 2 X-Eng S2T task. w2vb2 denotes w2v-BERT 2.0.

	arb	asm	bel	ben	bul	cat	ces	cmn	cym	dan	deu	ell	est	fin	fra	hin	hun	ind	ita	jpn	kat	khk	lit	lvs	mar	mlt	nld	pan	pes	pol	por	ron	rus	slk
w2vb2 XX-Large	27.8	15.9	14.6	20.5	28.7	36.0	29.8	16.7	23.2	33.2	34.7	22.4	26.4	24.2	31.4	23.2	22.9	27.0	23.5	13.9	16.3	12.4	20.9	25.3	19.0	36.1	24.8	21.2	25.4	21.2	34.6	31.3	26.2	29.2
w2vb2 X-Large	25.2	13.2	13.2	17.8	28.0	34.6	28.0	13.9	20.0	31.2	33.1	22.1	24.7	22.0	30.0	21.1	20.8	24.8	22.4	10.3	13.0	9.3	19.2	23.6	17.3	34.7	24.6	19.2	22.7	19.8	33.7	29.7	24.8	27.2
w2vb2 Large	23.0	9.2	13.0	14.0	24.8	32.7	24.6	11.9	15.3	27.4	30.8	18.4	20.1	18.4	27.6	16.6	14.3	23.0	20.7	8.5	10.6	6.5	15.1	18.8	12.3	30.6	21.2	14.6	19.9	16.9	31.8	26.2	22.4	24.4
CoLLD Large ₄₀	27.5	13.1	13.8	19.4	26.1	34.2	26.4	15.8	21.7	28.3	31.6	21.8	24.7	22.6	27.3	20.9	18.5	23.9	21.4	12.1	15.1	11.7	17.8	22.6	15.7	33.4	21.8	18.4	23.3	18.1	32.3	27.2	23.3	26.9
	slv	spa	swe	swh	tam	tha	tur	ukr	urd	uzn	vie	ckb	gle	kir	lug	yue	afr	amh	azj	bos	glg	guj	heb	hrv	hye	isl	jav	kaz	khm	kan	kor	ltz	lin	lao
w2vb2 XX-Large	23.4	22.3	35.2	25.9	15.7	16.7	22.7	29.2	20.3	20.2	16.8	19.3	7.1	14.9	16.2	10.4	38.2	14.7	14.9	31.6	31.6	24.2	26.5	28.1	23.6	18.0	17.9	19.8	17.6	19.6	16.1	31.0	8.8	18.7
w2vb2 X-Large	20.9	21.8	33.5	22.9	13.5	14.0	21.0	26.1	17.6	16.2	14.2	14.8	5.5	12.8	13.9	4.8	36.7	9.1	13.2	29.2	30.4	20.7	21.3	27.1	18.8	15.1	16.0	16.9	12.6	18.0	10.5	27.1	7.6	15.4
w2vb2 Large	18.2	20.2	28.2	20.8	10.6	12.2	16.2	23.7	14.6	13.8	9.1	12.1	3.8	10.8	11.3	2.4	34.5	8.3	11.5	27.0	27.5	17.0	16.4	25.0	15.5	11.4	13.5	14.6	11.4	14.7	9.4	19.9	5.9	12.7
CoLLD Large ₄₀	20.0	21.0	28.2	22.9	12.3	15.6	20.7	26.6	19.2	17.3	14.0	14.2	7.0	14.3	14.0	5.1	34.0	12.1	13.5	28.1	28.6	21.9	21.8	25.6	22.9	19.5	17.1	19.5	15.3	18.8	14.3	17.1	3.4	15.5
	mkd	zlm	npi	pbt	snd	sna	som	srp	tel	tgk	yor	mal	mya	tgl	ceb	ibo	kam	kea	luo	nso	nya	gaz	ory	umb	xho	zul	nob	mri	oci	hau	ast	ful	wol	
w2vb2 XX-Large	32.5	25.9	21.9	14.4	9.4	3.8	13.2	33.7	17.8	24.0	11.6	20.5	13.0	16.6	6.9	2.7	2.2	29.1	0.8	2.3	15.1	0.6	18.6	0.9	12.0	12.7	33.2	1.1	20.1	11.8	26.0	0.8	5.1	
w2vb2 X-Large	30.6	24.3	18.2	9.2	6.4	2.6	9.8	31.6	13.7	21.2	9.6	17.4	8.7	15.4	6.1	1.1	2.0	24.4	0.7	2.4	13.6	0.3	15.4	0.5	8.2	6.8	30.1	1.1	18.4	8.5	25.3	0.6	2.7	
w2vb2 Large	28.8	21.7	14.6	7.5	5.2	3.0	8.8	28.5	11.0	19.2	8.9	14.6	6.8	10.1	4.7	1.1	2.0	21.0	1.1	1.7	11.7	0.3	11.2	0.7	5.3	5.2	28.5	0.7	13.0	6.7	21.9	0.7	3.1	
CoLLD Large ₄₀	29.9	24.0	19.0	9.2	5.1	2.3	10.9	30.8	15.4	22.0	10.9	18.2	11.3	12.7	4.7	1.1	2.0	22.0	0.7	2.0	14.2	0.3	14.9	0.6	4.4	4.7	30.4	0.4	11.1	8.2	21.8	0.6	2.3	

Fig. 5. Complete BLEU scores on the FLEURS-101 X-Eng S2T task. w2vb2 denotes w2v-BERT 2.0. Underlined languages indicate unseen languages in X-Eng fine-tuning data.

Table 4. The l^{th} student layer to the \hat{l}^{th} teacher layer mapping for CoLLD derived from Eq. 1 when distilling from a 1B teacher.

Architecture	L^S	L^T	(l, \hat{l})
Large ₁₂	12	40	(1, 1), (2, 5), (3, 8), (4, 12), (5, 15), (6, 19), (7, 22), (8, 26), (9, 29), (10, 33), (11, 36), (12, 40)
Large ₄₀	40	40	(1, 1), (2, 2), (3, 3), ..., (40, 40)

6. APPENDIX

6.1. Knowledge Distillation Details

Here, we offer details about the knowledge distillation implementation. In Table 4, we show the student-to-teacher layer mapping in our distillation experiments. Next, the L2 regression loss for an utterance can be expressed as

$$\mathcal{L}_{\ell_2} = \frac{1}{DL^S|\mathcal{T}|} \sum_{l=1}^{L^S} \sum_{t \in \mathcal{T}} \left\| \mathbf{z}_t^l - \mathbf{h}_t^{\hat{l}} \right\|_2^2, \quad (4)$$

where D is the dimension of the representations \mathbf{z} and \mathbf{h} , and $|\mathcal{T}|$ is the number of masked time steps. For contrastive learning, the loss function is

$$\mathcal{L}_{\text{Contrastive}} = -\frac{1}{L^S|\mathcal{T}|} \sum_{l=1}^{L^S} \sum_{t \in \mathcal{T}} \log \frac{\exp(\cos(\mathbf{z}_t^l, \mathbf{h}_t^{\hat{l}}) / \tau)}{\sum_{\mathbf{h}' \in \mathcal{H}_t^{\hat{l}}} \exp(\cos(\mathbf{z}_t^l, \mathbf{h}') / \tau)}. \quad (5)$$

Finally, the losses of all utterances within a mini-batch are averaged to obtain the total loss function for optimization.

6.2. S2T Fine-tuning Details

This section offers implementation details of X-Eng S2T fine-tuning. Some fine-tuning hyperparameters for different model architectures

Table 5. X-Eng S2T fine-tuning hyperparameters for different model architectures.

Model	Learning Rate	Batch Size Per GPU	Gradient Accumulation	GPUs
XX-Large	2×10^{-4}	30 sec	2	64
X-Large	1×10^{-4}	60 sec	2	32
Large ₁₂	1×10^{-4}	60 sec	2	32
Large ₄₀	1×10^{-4}	48 sec	2	40

are shown in Table 5. First, the maximum length of an input utterance is 30 seconds, and the maximum number of output tokens is 113. Second, the input frames are randomly masked similar to the distillation process, but with a mask length of 5 and a masking probability of 0.02. Next, layer dropping of probability 0.1 is applied to both w2v-BERT 2.0 and NLLB models. Moreover, the NLLB transformer model is pre-trained with machine translation tasks, which take text as input, so we add a length adaptor [35] after the speech encoder to match the sequence length between speech and text. The adaptor begins with a 1-D CNN layer (kernel size = stride = 8) and a gated linear unit, followed by a single Conformer encoder layer with a convolution kernel size of 31. After this adaptor, the utterance length is reduced by a factor of eight to match the text modality.

6.3. Complete X-Eng S2T Results

In Fig. 4 and 5, we show the BLEU scores of several models of all languages in the CoVoST 2 and FLEURS evaluation sets. The details of different languages in the fine-tuning dataset can be found in Table 35 of [2]. Most unseen languages in the FLEURS testing sets have low BLEU scores. However, some unseen languages like ast (Asturian) and ltz (Luxembourgish) have high BLEU scores. We suspect high-resource languages in the same language family cause this phenomenon.