

DELULU: Discriminative Embedding Learning Using Latent Units for Speaker-Aware Self-Trained Speech Foundational Model

Anonymous ACL submission

Abstract

Self-supervised speech models have achieved remarkable success on content-driven tasks, yet they remain limited in capturing speaker-discriminative features critical for verification, diarization, and profiling applications. We introduce DELULU, a speaker-aware self-trained foundational model that addresses this limitation by incorporating speaker-informed structure into pseudo-label generation. DELULU leverages frame-level embeddings from ReDimNet, a state-of-the-art speaker verification model, to guide k-means clustering during pre-training, introducing a speaker-discriminative inductive bias that aligns representation learning with speaker identity. DELULU significantly outperforms prior SSL models across a range of speaker-centric tasks, achieving up to **62% relative improvement** in equal error rate (EER) for speaker verification and consistent gains on zero-shot profiling tasks including gender, age, accent, and speaker counting—notably surpassing even its teacher model on zero-shot evaluations. Our findings demonstrate that **DELULU is a strong universal encoder for speaker-aware speech processing**, enabling superior performance without task-specific fine-tuning.

1 Introduction

Speaker information is essential for a wide range of speech-related applications, including speaker verification, diarization, and personalized speech generation (Reynolds et al., 2000; Anguera et al., 2012; Casanova et al., 2022). Despite the recent success of self-supervised learning (SSL) in speech representation (Baevski et al., 2020; Hsu et al., 2021; Chen et al., 2022), existing models still struggle to capture speaker-specific characteristics effectively. The lack of robust speaker-aware representations poses a fundamental limitation for building systems that rely heavily on identity cues (Qian et al., 2022; Zhang et al., 2023). Although self-supervised models have achieved strong results across a variety

of speech and audio tasks (Baevski et al., 2020; Hsu et al., 2021; Chen et al., 2022; Waheed et al., 2024), their performance on speaker-related applications remains limited. The key bottleneck lies in their reliance on pseudo-labels generated through acoustic-only clustering, which are insufficiently aligned with speaker-discriminative structure. In HuBERT (Hsu et al., 2021), the k-means clustering step relies on shallow acoustic features that prioritize phonetic similarity, often suppressing speaker-specific information such as voice quality, prosody, and speaking style. Although WavLM (Chen et al., 2022) introduces additional context modeling and denoising objectives, it inherits the same clustering mechanism and, as a result, continues to struggle with learning robust speaker representations. We introduce DELULU, a speaker-aware self-trained foundational model that explicitly incorporates speaker-informed structure into the pretraining process. DELULU leverages ReDimNet (Yakovlev et al., 2024), a state-of-the-art speaker verification network, to guide k-means clustering with frame-level speaker embeddings rather than purely acoustic features. While ReDimNet itself is trained in a supervised manner and provides indirect supervision through pseudo-label generation, DELULU remains fundamentally self-trained in the classical sense (Yarowsky, 1995; Amini et al., 2025): it learns from its own predictions on unlabeled data, with the teacher serving only to initialize the target structure. Crucially, the pseudo-label signals DELULU derives from ReDimNet are misaligned with ReDimNet’s original training objectives—ReDimNet optimizes for utterance-level speaker discrimination, whereas DELULU’s clustering operates at the frame level for masked prediction—making the direct influence of ReDimNet’s supervision minimal. This design introduces a speaker-aware inductive bias into the learning pipeline while preserving the scalability and generality of self-training, enabling the model

to capture speaker-relevant information more effectively without compromising general acoustic modeling. A stronger speaker-oriented foundation model should yield representations that better encode speaker identity, leading to improved performance on forensic and identity-centric applications both in zero-shot settings and after task-specific fine-tuning. We show that this holds true: DELULU outperforms existing SSL models across speaker verification, profiling (age, gender, accent), and speaker counting tasks, with particularly large gains in zero-shot evaluation. These improvements translate directly to downstream fine-tuning, confirming that introducing speaker-aware structure at the clustering stage produces more transferable representations. Beyond performance, DELULU illustrates a broader principle: self-trained speech models can be strengthened by integrating targeted external signals to guide pseudo-label formation, without requiring direct task-aligned supervision. This framework provides a scalable path toward speaker-aware foundation models and can be generalized to other architectures and domains beyond speaker modeling. Our main contributions are as follows:

- We introduce DELULU, a speaker-aware self-trained speech model that addresses the fundamental limitation of speaker discriminability in existing SSL approaches by integrating external speaker-informed structure into the pseudo-label generation process through ReDimNet-guided k-means clustering.
- We achieve state-of-the-art performance across a wide suite of speaker-centric benchmarks, including speaker verification, profiling (age, gender, accent), and speaker counting, demonstrating up to 62% relative improvement in EER.
- We establish DELULU as a strong universal encoder for speaker tasks, outperforming prior SSL models by a large margin in zero-shot settings and showing competitive results in upstream performance comparisons—notably surpassing even its teacher model, ReDimNet.

2 Related Work

2.1 Self-Supervised Speech Representation Learning

Self-supervised learning (SSL) has revolutionized speech processing by enabling models to learn

rich representations from large-scale unlabeled audio data (Baevski et al., 2020; Hsu et al., 2021; Chen et al., 2022). These models typically employ pretext tasks such as contrastive predictive coding or masked language modeling to capture phonetic and acoustic structures without explicit labels. Wav2vec 2.0 (Baevski et al., 2020) introduced a contrastive loss over quantized latent representations, achieving strong performance on speech recognition benchmarks. HuBERT (Hsu et al., 2021) advanced this by using offline k-means clustering on MFCC features to generate pseudo-labels, focusing on discrete unit discovery that aligns well with phonetic content. WavLM (Chen et al., 2022) further incorporated denoising objectives and utterance mixing to improve robustness to noise and overlapping speech, making it suitable for a broader range of tasks including speaker identification. Despite these advances, standard SSL models often underperform on speaker-centric tasks due to their emphasis on content over speaker identity (Chen et al., 2022). Recent efforts have explored scaling these models to multilingual settings (Pratap et al., 2024) or enhancing them with multi-task learning (Hu et al., 2024), but speaker discriminability remains a challenge.

2.2 Supervised Speaker Representation Learning

Traditional speaker verification systems rely on the fundamental assumption that the human voice is a unique biometric trait (Singh and Raj, 2025). This expected uniqueness underpins decades of research into supervised speaker modeling pipelines. The x-vector system (Snyder et al., 2018) uses time-delay neural networks (TDNNs) with statistical pooling to produce fixed-dimensional utterance embeddings, achieving robust performance on benchmarks like VoxCeleb (Nagrani et al., 2017). Advancements include ECAPA-TDNN (Desplanques et al., 2020), which incorporates channel and context-dependent attention to better capture speaker variability. More recently, ReDimNet (Yakovlev et al., 2024) introduced a reshaped dimensionality network that optimizes for both local and global speaker features, setting new state-of-the-art results on speaker verification tasks. However, these supervised embeddings are often overspecialized for identity discrimination, limiting their generalization to broader speaker profiling tasks such as age and gender estimation or forensic analysis (Baali et al., 2024). For instance,

184	phonetic biases in attention mechanisms can con-	3 Architecture	233
185	found speaker-specific cues with content-related	Based on the architecture depicted in Figure 1,	234
186	artifacts, reducing transferability to non-identity	DELULU adopts a masked training design that	235
187	attributes (Baali et al., 2024). These supervised	integrates self-supervised learning with speaker-	236
188	approaches excel in controlled settings but require	discriminative guidance through an external teacher	237
189	large labeled datasets and struggle with domain	model.	238
190	shifts (Desplanques et al., 2020). Integrating their	3.1 DELULU Encoder	239
191	strengths into SSL pipelines offers a promising	The DELULU encoder follows the wav2vec 2.0	240
192	direction to enhance unsupervised representation	architecture (Baevski et al., 2020), consisting	241
193	learning, enabling foundational models that pro-	of a convolutional feature extractor, a Trans-	242
194	duce versatile, speaker-aware representations suit-	former encoder, a projection layer, and a	243
195	able for diverse downstream applications.	code embedding layer. The convolutional en-	244
196	2.3 Self-Supervised Learning for Speaker	coder is composed of seven 512-channel layers	245
197	Tasks	with strides [4, 2, 2, 2, 2, 2] and kernel widths	246
198	A growing body of work adapts SSL specif-	[10, 3, 3, 3, 3, 2, 2], differing from HuBERT (Hsu	247
199	ically for speaker-related tasks such as veri-	et al., 2021) by adjusting the stride pattern to en-	248
200	fication and diarization, often leveraging con-	sure temporal alignment with the teacher model’s	249
201	trastive or distillation-based objectives to learn	frame-level outputs. For 16 kHz input speech, this	250
202	speaker-discriminative representations from unlabeled	produces a latent feature sequence with a 16 ms	251
203	data (Lepage and Dehak, 2025). Early ef-	frame rate ($256\times$ down-sampling factor). These	252
204	forts applied contrastive learning frameworks from	latent features are then passed through a stack of	253
205	computer vision, such as SimCLR adapted for	Transformer blocks that model long-range depen-	254
206	speech (Jiang et al., 2020), which minimizes dis-	dencies, after which a projection layer reduces di-	255
207	tances between augmented views of the same ut-	dimensionality and prepares features for clustering	256
208	terance while maximizing separations from oth-	and loss computation. Overall, the DELULU en-	257
209	ers. Self-distillation approaches, inspired by	coder serves as the student model that learns to	258
210	DINO (Caron et al., 2021; Ashihara et al., 2024),	predict clustered representations.	259
211	have shown particular promise. More recently,	3.2 Teacher-Guided Clustering	260
212	CoLMbo (Baali et al., 2025b) proposed a speaker	Unlike conventional self-supervised speech models	261
213	language model that integrates prompt-conditioned	that rely purely on acoustic clustering, DELULU	262
214	speaker encoders to generate descriptive speaker	leverages ReDimNet (Yakovlev et al., 2024), a	263
215	profiles. This approach showed strong zero-shot	state-of-the-art speaker verification model, to guide	264
216	generalization on demographic traits such as di-	the pseudo-label generation process. Instead of us-	265
217	allect and age, demonstrating the potential of us-	ing pooled utterance-level embeddings, we extract	266
218	ing large-scale speaker models beyond classifica-	<i>prepooled frame-level features</i> from ReDimNet to	267
219	tion tasks. Meanwhile, CAARMA (Baali et al.,	preserve temporal resolution and ensure alignment	268
220	2025c) introduced an augmentation framework	with DELULU’s encoder outputs. These features	269
221	where HuBERT’s hidden layers serve as discrim-	are then used for k-means clustering with $k = 256$	270
222	inators during adversarial training, improving the	clusters. This approach ensures that the discrete	271
223	speaker discriminability of the learned embeddings.	targets encode speaker-specific characteristics such	272
224	These methods suggest that integrating architec-	as voice quality, prosody, and speaking style, rather	273
225	tural and training-level bias can enhance SSL’s	than prioritizing only phonetic content.	274
226	ability to encode speaker traits but they do not	3.3 Training	275
227	intervene directly at the pseudo-label generation	DELULU is trained using a masked prediction loss.	276
228	step as DELULU does.	Masked Prediction Loss: Following the masked	277
229	DELULU builds on these by directly guiding	language modeling paradigm, a portion of the input	278
230	SSL clustering with supervised speaker features,	time steps are randomly masked. The model is	279
231	bridging the gap between general acoustic model-	trained to predict the cluster assignments (derived	280
232	ing and targeted speaker discrimination.		

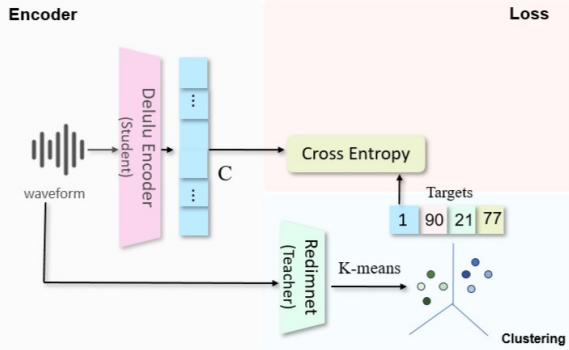


Figure 1: Illustration of DELULU architecture. The student encoder is trained with a masked prediction objective, where frame-level features from the teacher (ReDimNet) are clustered into discrete speaker-aware targets for cross-entropy training.

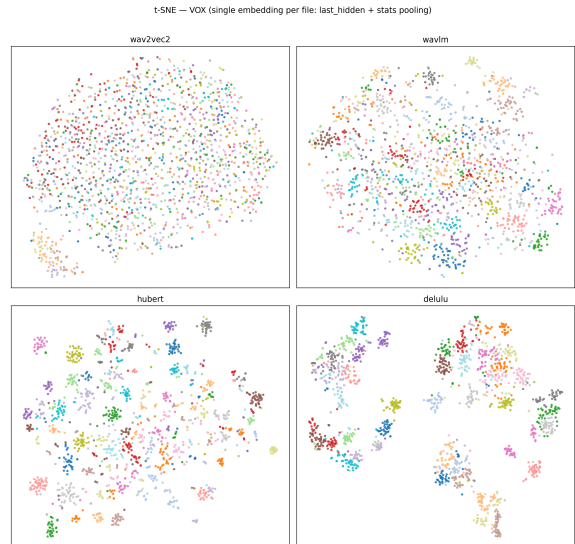


Figure 2: t-SNE visualization of embeddings from 40 VoxO1 speakers. DELULU yields compact, well-separated clusters with lower intra-speaker variability and higher inter-speaker separation than baseline SSL models.

from ReDimNet-guided k-means) at the masked positions. Given the model’s output C_i at masked position i and the corresponding target cluster ID, the cross-entropy loss is computed as:

$$\mathcal{L}_{\text{mask}} = - \sum_{i \in \mathcal{M}} \log P(c_i^* | C_i) \quad (1)$$

where \mathcal{M} denotes the set of masked positions, C_i represents the logits over the $k = 256$ possible clusters, and c_i^* is the ground-truth cluster assignment obtained from the teacher model.

4 Experimental Setup

We build our system on top of the official torchaudio self-supervised learning recipes. For clustering, we adopt the MiniBatchKMeans algorithm from scikit-learn, with a mini-batch size of 10,000 frames. Initialization is performed with k-means++ using 20 random starts to ensure stability, and cluster assignments are set to $k = 256$.

Pre-training is conducted on 960 hours of LibriSpeech audio using 4 NVIDIA H100 GPUs, with each GPU processing 87.5 seconds of audio per batch. The model is trained for a total of 400k updates. Optimization follows the AdamW optimizer ($\beta_1 = 0.9, \beta_2 = 0.98$) with an initial learning rate of $5e-4$. A linear warmup is applied for the first 32k steps, followed by a polynomial decay back to zero. To regularize training, we apply a weight decay of 0.01 and clip gradients at a maximum norm of 10.0.

5 Ablation Study

To quantify the contribution of each component in DELULU’s design, we perform systematic ab-

lation experiments on the VoxCeleb1-O dataset, evaluating upstream speaker verification performance using EER, with all models trained for 400k updates on LibriSpeech 960h under identical optimization settings for fair comparison. Examining clustering feature sources, MFCC features with $k = 100$ yield a baseline EER of 37.73%, while higher-level acoustic representations from a pretrained HuBERT model with $k = 500$ improve it to 34.05%, indicating modest gains from context-aware features; in contrast, frame-level ReDimNet embeddings (dimension 2304) with $k = 256$ reduce EER to 13.53%, a 60% relative improvement over HuBERT, confirming that external speaker-discriminative supervision during clustering drives DELULU’s performance. For cluster size using ReDimNet features, we evaluate $k \in \{256, 500, 1024\}$, obtaining EERs of 13.53%, 14.16%, and 14.16% respectively, showing $k = 256$ offers the optimal trade-off between granularity and generalization, as 256 discrete speaker-aware units suffice for identity variation at this scale. We also varied the DELULU encoder’s stride, to find the optimal stride. Since the training paradigm requires frame-level synchrony between the teacher and student, the RedimNet stride was modified to be identical to that of the encoder. A stride of 16ms achieves 13.52% EER, while both lower (15ms or less) and higher (20ms or greater) strides re-

Features	Clusters	Stride	EER (%)	Rel. Impr.
MFCC	100	20	37.73	-
HuBERT (Stage 2)	500	20	34.05	9.8%
ReDimNet	256	16	13.53	60.2%
ReDimNet	500	16	14.16	58.4%
ReDimNet	1024	16	14.16	58.4%
ReDimNet	256	15	14.16	58.4%

Table 1: Ablation study on VoxCeleb1-O upstream speaker verification (EER %). Lower is better. Results demonstrate that ReDimNet-guided clustering is the primary factor in DELULU’s performance, with optimal cluster size $k=256$ and stride=16 for temporal alignment.

Method	Vox1-O	LibriSpeech	Avg.
DELULU-Utt	34.62	11.18	22.90
DELULU-Frame	13.53	15.95	14.74
<i>Rel. Impr.</i>	<i>60.9%</i>	<i>-42.7%</i>	<i>35.6%</i>

Table 2: Ablation study comparing frame-level versus utterance-level pseudo-labeling (EER %). Lower is better. DELULU-Frame demonstrates superior cross-domain generalization on VoxO1 despite comparable in-domain performance, validating the benefit of fine-grained temporal supervision.

sulted in EERs greater than 14%, hence we chose 16ms as our final stride. Overall, these ablations highlight speaker-discriminative clustering as the key factor, with ReDimNet-guided pseudo-labels providing a 60% relative improvement over HuBERT’s acoustic-only approach, optimal $k = 256$ balancing discriminability and generalization, and stride=16 ensuring effective supervision, validating the introduction of external speaker supervision into pseudo-label generation as a powerful strategy for speaker-aware representations.

A central contribution of DELULU is the leveraging the power of frame-level pseudo-labeling, which enables masked prediction during pre-training and provides fine-grained temporal supervision. To quantify this benefit, we compare against an utterance-level variant (DELULU-Utt) where DELULU features are pooled, with pseudo-labels derived from ReDimNet utterance embeddings rather than frame-level representations. As shown in Table 6, while DELULU-Utt achieves competitive in-domain performance on LibriSpeech (11.18% EER), it suffers a substantial degradation on the out-of-domain VoxCeleb1-O benchmark (34.62% EER). In contrast, DELULU-Frame maintains consistent performance across

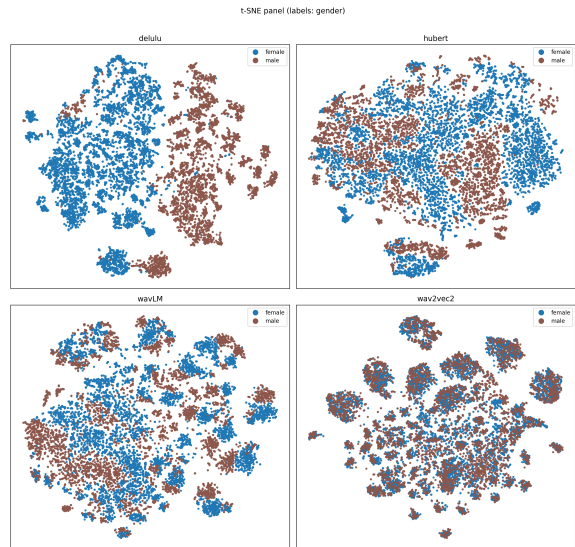


Figure 3: t-SNE visualization of speaker embeddings by gender on the EARS dataset. DELULU shows clear speaker clusters and strong gender separation.

both evaluations (13.53% and 15.95% EER), yielding a 35.6% relative improvement in average EER. These results demonstrate that frame-level supervision not only facilitates masked prediction but also promotes more robust speaker representations that generalize beyond the training distribution. Notably, while utterance-level clustering aligns pseudo-labels directly with the downstream verification objective—where decisions are inherently made at the utterance level—this apparent advantage leads to overfitting on in-domain data. By contrast, frame-level supervision forces the model to learn temporally localized speaker characteristics, yielding representations that transfer more effectively to unseen acoustic conditions.

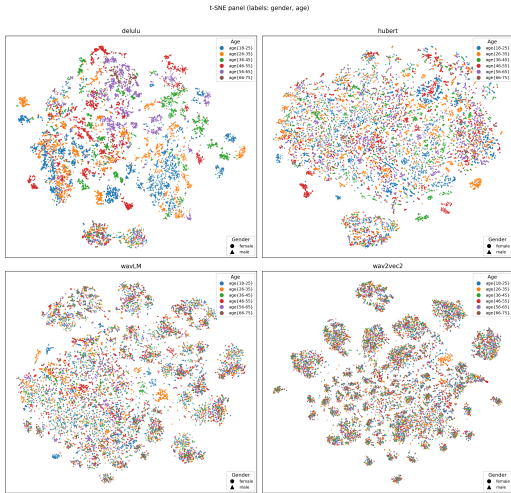


Figure 4: t-SNE visualization of EARS embeddings by gender and age groups (Table 7). DELULU clusters primarily by speaker identity, with demographic traits forming secondary structure and no apparent bias.

Dataset	WavLM	Hubert	DELULU	Wav2Vec2
VoxO1	35.93	34.05	13.53	43.17
SITW	44.00	42.60	25.40	42.20

Table 3: EER (%) across upstream speaker verification benchmarks (\downarrow better)

6 Evaluation

We evaluate DELULU across multiple speaker-centric benchmarks described below to assess its effectiveness in capturing speaker-discriminative information. To verify our hypothesis that, as a more speaker-aware foundation model, DELULU may be expected to result in representations that are naturally better organized by speaker characteristics, and deliver better performance both in zero-shot “upstream” settings where they are used as-is, and in “downstream” settings where the model is further fine-tuned to the task, we perform the evaluations (comparing it to other baseline models) in both settings, with a primary focus on speaker verification for the downstream task. We also visualize the derived representations for further confirmation of our hypothesis. We report on these tests below.

For fair comparison, all baseline models (wav2vec 2.0, HuBERT, and WavLM) are pre-trained on the same 960 hours of LibriSpeech data, ensuring that differences in performance arise from architectural and training objective choices rather than data scale.

Features	Clusters	EER (%)	Rel. Impr.
MFCC	100	13.00	–
HuBERT (Stage 2)	500	7.45	–
ReDimNet	256	5.63	24.5%
ReDimNet	500	6.28	15.7%
ReDimNet	1024	5.99	19.6%

Table 4: EER (%) across downstream speaker verification (VoxO1) (\downarrow better)

6.1 Upstream Evaluation: Speaker Verification

We first evaluate the models in an upstream setting, where representations are extracted directly from the pretrained encoders without any fine-tuning. This zero-shot evaluation reveals the inherent speaker-discriminative properties encoded during self-supervised pretraining.

Protocol. We evaluate on two widely-used speaker verification benchmarks: VoxO1 (Nagrani et al., 2017) and SITW (McLaren et al., 2016) (see Appendix B). For each model, we extract utterance-level embeddings by mean pooling over the temporal dimension of final layer representations. Verification trials are scored using cosine similarity, and we report EER, where lower is better.

Results. Table 3 shows that DELULU achieves substantial improvements over all baseline SSL models. On VoxO1, DELULU obtains an EER of 13.53%, representing a 62% relative improvement over HuBERT (34.05%) and a 60% improvement over WavLM (35.93%). Similar gains are observed on SITW, where DELULU achieves 25.40% EER compared to HuBERT’s 42.60% and WavLM’s 44.00%. These results demonstrate that ReDimNet guided clustering successfully embeds strong speaker-discriminative structure into the learned representations.

6.2 Downstream Evaluation: Speaker Verification

To assess transferability beyond upstream frozen representations, we evaluate DELULU on a downstream speaker verification task using VoxCeleb1-O. Since our objective is to evaluate DELULU’s transferability rather than optimize for maximum downstream performance, we employ a minimal fine-tuning architecture that adds only a simple classification head on top of the frozen encoder, as described in Appendix C and then compute EER

Task	WavLM	Hubert	DELULU	Wav2Vec2
Spoof Detection	51.44	53.51	57.20	52.88
Speaker Count	62.71	64.83	67.13	64.20
Accent Detection	77.76	62.86	78.38	58.60
HowFarSpk	69.37	70.58	73.36	57.84
Gender Detection	95.75	93.97	96.18	92.73
Age Estimation	32.69	29.43	36.00	31.99

Table 5: Zero-shot Macro-F1 score (%) across profiling tasks (\uparrow better).

Table 6: Teacher vs. Student comparison on zero-shot profiling tasks. We report accuracy (%) at both frame-level and utterance-level evaluation. ReDimNet (teacher) was trained on VoxBlink2 (>100K speakers), while DELULU (student) was trained only on LibriSpeech (~1K speakers).

Evaluation	Task	ReDimNet	DELULU
Frame-Level	Speaker Count	47.00	51.00
	HowFarAreYou	50.72	55.35
	Spoof Detection	89.50	91.50
Utterance-Level	Speaker Count	29.50	65.50
	HowFarAreYou	51.00	73.58
	Spoof Detection	90.50	92.08

to measure verification accuracy.

As shown in Table 4, DELULU achieves 5.63% EER with $k = 256$ clusters, representing a 24.5% relative improvement over HuBERT Stage 2 (7.45%) and a 56.7% gain over the MFCC baseline (13.00%). Larger vocabularies ($k = 500, 1024$) yield marginal improvements, confirming that $k = 256$ provides the optimal balance between discriminability and generalization. These results demonstrate that the speaker-discriminative structure learned during pretraining transfers effectively to supervised fine-tuning, further validating DELULU as a strong foundational encoder for speaker verification.

6.3 Demographic Subgroup Analysis

To understand how speaker-discriminative information is distributed across demographic groups following the SVERITAS (Baali et al., 2025a) benchmark, we analyze model performance on the EARS (Richter et al., 2024) dataset (details in Appendix D), stratified by gender and age subgroups.

Protocol. Following the upstream protocol, we extract mean-pooled embeddings and compute EER for speaker verification trials within each demographic subgroup. This analysis reveals whether models learn speaker representations uniformly

across age and gender categories or exhibit demographic biases.

Results. Appendix D Table 7 presents EER across gender and age subgroups. DELULU consistently outperforms baselines across all demographics, with particularly strong improvements for middle-aged speakers (36-55 years). For male speakers aged 36-45, DELULU achieves 24.53% EER compared to HuBERT’s 39.47%. The consistent improvements across subgroups indicate that DELULU’s speaker-discriminative representations generalize effectively across demographic variations without introducing systematic biases.

6.4 Zero-Shot Speaker Profiling Tasks

Beyond verification, we assess DELULU’s capacity to encode diverse speaker attributes through zero-shot evaluation on multiple profiling tasks from the DynamicSUPERB benchmark (Huang et al., 2024). We evaluate six speaker-related tasks, including age, gender, accent, speaker counting, and spoof detection (details in Appendix E).

Protocol. We extract layer-wise representations from each model and train K-Nearest Neighbors (KNN) classifiers without any fine-tuning. All evaluations are conducted in a fully zero-shot setting; no supervised training is used for adaptation of the models. Given that these models are trained without explicit supervision, we expect performance to be poor, possibly approaching random levels; however, all SSL-based models achieve meaningful results, with DELULU consistently outperforming other models. For each task, we perform k-fold cross-validation and report Macro-F1 scores averaged across folds, along with standard deviations to indicate variability and robustness.

Results. Table 5 summarizes the zero-shot results across six profiling tasks. DELULU achieves the highest performance on all tasks, considerably outperforming other SSL models in zero-shot settings.

The second-best model varies by task—for instance, WavLM performs competitively on gender classification, while HuBERT fares better on accent detection. On spoof detection, DELULU achieves 57.20% F1 compared to HuBERT’s 53.51% and WavLM’s 51.44%; on accent detection, it reaches 78.38% F1, surpassing HuBERT (62.86%) and wav2vec 2.0 (58.60%). Even on more challenging tasks such as age estimation (36.00% F1), DELULU exhibits clear gains despite the absence of task supervision. These findings demonstrate that the speaker-guided clustering in DELULU enables robust and discriminative representations that generalize across diverse zero-shot profiling scenarios.

6.5 Teacher vs. Student Evaluation.

To evaluate the robustness of DELULU’s representations under different evaluation protocols, we compare frame-level and utterance-level performance against its teacher model, ReDimNet, on three profiling tasks: speaker counting, speaker distance estimation (HowFarAreYou), and spoof detection. For utterance-level evaluation, we apply mean pooling over frame embeddings, while for frame-level evaluation, we classify each frame independently using KNN and aggregate predictions via majority voting. Notably, ReDimNet was trained on VoxBlink2 (Lin et al., 2024), a large-scale audio-visual dataset comprising over 100,000 speakers with diverse acoustic conditions, whereas DELULU was trained solely on LibriSpeech (Panayotov et al., 2015), a significantly smaller corpus of approximately 1,000 speakers. Despite this substantial difference in training data scale and diversity, Table 6 shows that DELULU consistently outperforms ReDimNet on all three tasks under both evaluation protocols. On speaker counting, DELULU achieves 51.00% and 65.50% accuracy at the frame and utterance levels respectively, compared to ReDimNet’s 47.00% and 29.50%. Similarly, on HowFarAreYou, DELULU reaches 55.35% (frame) and 73.58% (utterance) versus ReDimNet’s 50.72% and 51.00%. For spoof detection, DELULU maintains a consistent edge at both frame level (91.50% vs. 89.50%) and utterance level (92.08% vs. 90.50%). These results demonstrate that DELULU’s speaker-guided distillation enables superior generalization to zero-shot profiling tasks, even when trained on way smaller data compared to its teacher.

6.6 Representation Analysis

We analyze the learned embedding space using t-SNE visualizations to assess whether DELULU produces structured and discriminative speaker representations. In Figure 2, we expect embeddings from the same speaker to form compact clusters, distinct from those of other speakers. Indeed, DELULU’s representations exhibit clear and well-separated speaker clusters, indicating strong identity preservation and reduced intra-speaker variability compared to baseline models. In Figure 3, we expect representations to separate naturally by gender. The visualization confirms this behavior, showing distinct groupings corresponding to male and female speakers, while maintaining tight clustering within each group. Finally, in Figure 4, we examine whether age-related patterns emerge within the embedding space. As expected, DELULU produces smooth transitions across age groups, suggesting that the model implicitly encodes demographic cues such as vocal maturity and pitch characteristics alongside dominant speaker identity. These results demonstrate that DELULU learns a representation space that reflects both speaker individuality and meaningful demographic structure, confirming the model’s capacity for fine-grained, interpretable speaker encoding.

7 Conclusion

We introduce DELULU, a speaker-aware self-trained foundational model that enhances speaker representation learning by guiding pseudo-label generation with frame-level embeddings from ReDimNet. Across benchmarks in speaker verification, profiling, and counting, DELULU consistently outperforms existing SSL models, establishing it as a strong universal encoder for speaker-centric tasks. Since the core idea lies in incorporating task-relevant structure into pseudo-label generation without direct supervision, we expect this strategy to generalize to other foundational model architectures and downstream applications. Future work will explore replacing clustering with distillation-based objectives to further simplify the pretraining pipeline.

Limitations

While DELULU demonstrates substantial improvements in speaker verification, profiling, and zero-shot speaker-aware tasks, several limitations warrant further investigation.

608 Extending DELULU to large-scale, multi-
 609 domain datasets with higher speaker and environ-
 610 mental variability remains an open challenge. Also,
 611 incorporating an external model (ReDimNet) into
 612 the clustering process introduces additional com-
 613 putational overhead and memory requirements, po-
 614 tentially limiting accessibility for research groups
 615 with restricted computational resources. Finally,
 616 while DELULU excels in identity-aware modeling,
 617 its downstream adaptability to non-speaker-centric
 618 applications, such as emotion or intent recognition,
 619 has yet to be comprehensively explored. Future
 620 work will focus on addressing these limitations
 621 through more efficient training strategies, broader
 622 data coverage, and cross-domain generalization
 623 studies.

624 Ethics Statement

625 The development of DELULU is guided by a
 626 strong commitment to ethical responsibility and
 627 privacy preservation. As a speaker-aware foun-
 628 dational model, DELULU possesses the potential
 629 to be misused for surveillance, impersonation, or
 630 identity inference without consent. We emphasize
 631 that its use must strictly comply with ethical re-
 632 search standards, data protection laws, and insti-
 633 tutional review protocols. All datasets employed
 634 in this study (LibriSpeech, VoxCeleb, Dynamic-
 635 Superb, and SVeritas) consist of publicly avail-
 636 able and consented speech samples intended solely
 637 for research purposes. No personally identifiable
 638 or private data was used or generated. Further-
 639 more, while DELULU improves representation
 640 learning for speaker-related tasks, it does not per-
 641 form speaker identification on unconsented audio,
 642 nor is it intended for forensic or monitoring appli-
 643 cations. In line with the ACL Ethics Policy, we
 644 advocate for transparent deployment of DELULU,
 645 ensuring that its advancements contribute positively
 646 to responsible and equitable speech technology re-
 647 search.

648 References

649 Massih-Reza Amini, Vasilii Feofanov, Loic Pauletto,
 650 Lies Hadjadj, Emilie Devijver, and Yury Maximov.
 651 2025. Self-training: A survey. *Neurocomputing*,
 652 616:128904.

653 Xavier Anguera, Simon Bozonnet, Nicholas Evans,
 654 Corinne Fredouille, Gerald Friedland, and Oriol
 655 Vinyals. 2012. Speaker diarization: A review of

recent research. *IEEE Transactions on audio, speech,
 and language processing*, 20(2):356–370. 656 657

Takanori Ashihara, Marc Delcroix, Takafumi Moriya,
 Kohei Matsuura, Taichi Asami, and Yusuke Ijima.
 2024. What do self-supervised speech and speaker
 models learn? new findings from a cross model layer-
 wise analysis. In *ICASSP 2024-2024 IEEE Interna-
 tional Conference on Acoustics, Speech and Signal
 Processing (ICASSP)*, pages 10166–10170. IEEE. 658 659 660 661 662 663 664

Massa Baali, Abdulhamid Aldoobi, Hira Dharmyal, Rita
 Singh, and Bhiksha Raj. 2024. PDAF: A phonetic de-
 biasing attention framework for speaker verification.
 In *SLT*. 665 666 667 668

Massa Baali, Sarthak Bisht, Francisco Teixeira,
 Kateryna Shapovalenko, Rita Singh, and Bhiksha Raj.
 2025a. SVeritas: Benchmark for robust speaker ver-
 ification under diverse conditions. In *Findings of the
 Association for Computational Linguistics: EMNLP*. 669 670 671 672 673

Massa Baali, Shuo Han, Syed Abdul Hannan, Purusot-
 tam Samal, Karanveer Singh, Soham Deshmukh, Rita
 Singh, and Bhiksha Raj. 2025b. CoLMbo: Speaker
 language model for descriptive profiling. In *ASRU*.
 IEEE. 674 675 676 677 678

Massa Baali, Xiang Li, Hao Chen, Syed Abdul Hannan,
 Rita Singh, and Bhiksha Raj. 2025c. CAARMA:
 Class augmentation with adversarial mixup regular-
 ization. In *Findings of the Association for Computa-
 tional Linguistics: EMNLP*. 679 680 681 682 683

Alexei Baeovski, Yuhao Zhou, Abdelrahman Mohamed,
 and Michael Auli. 2020. wav2vec 2.0: A framework
 for self-supervised learning of speech representations.
Advances in neural information processing systems,
 33:12449–12460. 684 685 686 687 688

Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé
 Jégou, Julien Mairal, Piotr Bojanowski, and Armand
 Joulin. 2021. Emerging properties in self-supervised
 vision transformers. In *Proceedings of the IEEE/CVF
 international conference on computer vision*, pages
 9650–9660. 689 690 691 692 693 694

Edresson Casanova, Julian Weber, Christopher D
 Shulby, Arnaldo Candido Junior, Eren Gölge, and
 Moacir A Ponti. 2022. Yourtts: Towards zero-shot
 multi-speaker tts and zero-shot voice conversion for
 everyone. In *International conference on machine
 learning*, pages 2709–2720. PMLR. 695 696 697 698 699 700

Sanyuan Chen, Chengyi Wang, Zhengyang Chen,
 Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki
 Kanda, Takuya Yoshioka, Xiong Xiao, and 1 oth-
 ers. 2022. Wavlm: Large-scale self-supervised pre-
 training for full stack speech processing. *IEEE
 Journal of Selected Topics in Signal Processing*,
 16(6):1505–1518. 701 702 703 704 705 706 707

Brecht Desplanques, Jenthe Thienpondt, and Kris De-
 muynck. 2020. ECAPA-TDNN: Emphasized chan-
 nel attention, propagation and aggregation in tdnn
 based speaker verification. *Interspeech*. 708 709 710 711

712	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai,	self-supervised speech representation by disentangling speakers. In <i>International conference on machine learning</i> , pages 18003–18017. PMLR.	767
713	Kushal Lakhota, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. <i>IEEE/ACM transactions on audio, speech, and language processing</i> , 29:3451–3460.		768
714			769
715			
716		Douglas A Reynolds, Thomas F Quatieri, and Robert B Dunn. 2000. Speaker verification using adapted gaussian mixture models. <i>Digital signal processing</i> , 10(1-3):19–41.	770
717			771
718	Shujie Hu, Long Zhou, Shujie Liu, Sanyuan Chen, Lingwei Meng, Hongkun Hao, Jing Pan, Xunying Liu, Jinyu Li, Sunit Sivasankaran, and 1 others. 2024. Wavllm: Towards robust and adaptive speech large language model. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 4552–4572.		772
719			773
720		Julius Richter, Yi-Chiao Wu, Steven Krenn, Simon Welker, Bunlong Lay, Shinjii Watanabe, Alexander Richard, and Timo Gerkmann. 2024. EARS: An anechoic fullband speech dataset benchmarked for speech enhancement and dereverberation. In <i>Inter-speech</i> .	774
721			775
722			776
723			777
724			778
725	Chien-yu Huang, Ke-Han Lu, Shih-Heng Wang, Chi-Yuan Hsiao, Chun-Yi Kuan, Haibin Wu, Siddhant Arora, Kai-Wei Chang, Jiatong Shi, Yifan Peng, and 1 others. 2024. Dynamic-superb: Towards a dynamic, collaborative, and comprehensive instruction-tuning benchmark for speech. In <i>ICASSP</i> , pages 12136–12140. IEEE.		779
726			
727		Rita Singh and Bhiksha Raj. 2025. Human voice is unique. <i>arXiv preprint arXiv:2506.18182</i> .	780
728			781
729		D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur. 2018. X-vectors: Robust dnn embeddings for speaker recognition. In <i>ICASSP</i> , pages 5329–5333. IEEE.	782
730			783
731			784
732			785
733	Dongwei Jiang, Wubo Li, Miao Cao, Wei Zou, and Xiangang Li. 2020. Speech simclr: Combining contrastive and reconstruction objective for self-supervised speech representation learning. <i>arXiv preprint arXiv:2010.13991</i> .		786
734			787
735		Abdul Waheed, Hanin Atwany, Bhiksha Raj, and Rita Singh. 2024. What do speech foundation models not learn about speech? <i>arXiv preprint arXiv:2410.12948</i> .	788
736			789
737	Theo Lepage and Reda Dehak. 2025. Self-supervised frameworks for speaker verification via bootstrapped positive sampling. <i>arXiv preprint arXiv:2501.17772</i> .		
738			
739		Feng Wang, Jian Cheng, Weiyang Liu, and Haijun Liu. 2018. Additive margin softmax for face verification. <i>IEEE Signal Processing Letters</i> , 25(7):926–930.	790
740	Yuke Lin, Ming Cheng, Fulin Zhang, Yingying Gao, Shilei Zhang, and Ming Li. 2024. Voxblink2: A 100k+ speaker recognition corpus and the open-set speaker-identification benchmark. <i>arXiv preprint arXiv:2407.11510</i> .		791
741			792
742		Ivan Yakovlev, Rostislav Makarov, Andrei Balykin, Pavel Malov, Anton Okhotnikov, and Nikita Torgashov. 2024. Reshape dimensions network for speaker recognition. In <i>Proc. Interspeech 2024</i> , pages 3235–3239.	793
743			794
744			795
745			796
746	Mitchell McLaren, Luciana Ferrer, Diego Castan, and Aaron Lawson. 2016. The speakers in the wild (sitw) speaker recognition database. In <i>Interspeech</i> , pages 818–822.		797
747			
748			
749	Arsha Nagrani, Joon Son Chung, and Andrew Senior. 2017. Voxceleb: A large-scale speaker identification dataset. <i>Interspeech 2017</i> , page 2616.		798
750			799
751			800
752	Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: an asr corpus based on public domain audio books. In <i>2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)</i> , pages 5206–5210. IEEE.		801
753			
754			
755			
756			
757			
758	Vineel Pratap, Andros Tjandra, Bowen Shi, Paden Tomasello, Arun Babu, Sayani Kundu, Ali Elkahky, Zhaoheng Ni, Apoorv Vyas, Maryam Fazel-Zarandi, and 1 others. 2024. Scaling speech technology to 1,000+ languages. <i>Journal of Machine Learning Research</i> , 25(97):1–52.		802
759			803
760			804
761			805
762			
763			
764	Kaizhi Qian, Yang Zhang, Heting Gao, Junrui Ni, Cheng-I Lai, David Cox, Mark Hasegawa-Johnson, and Shiyu Chang. 2022. Contentvec: An improved		
765			
766			

806	A Experimental Setup		
807	Computational complexity	Compared to the	
808		HuBERT baseline, DELULU introduces additional	
809		computational cost during the preprocessing stage	
810		due to RedimNet’s convolutional feature extrac-	
811		tion. However, this overhead is performed offline	
812		and does not affect training or inference. Empir-	
813		ically, on an NVIDIA H100 GPU, HuBERT re-	
814		quires approximately 14.5 ms per batch during in-	
815		ference, while RedimNet requires approximately	
816		51.1 ms under identical hardware and batch con-	
817		ditions. Since preprocessing is handled once and	
818		cached, this difference does not impact overall	
819		training efficiency or deployment latency in our	
820		setup.	
821	B Upstream Evaluation		
822	SV Benchmarks		
823	• VoxCeleb1-O (Nagrani et al., 2017):	The original	
824		VoxCeleb1 test set containing 37,720 verifica-	
825		tion trials across 40 speakers.	
826	• SITW (McLaren et al., 2016):	Speakers in the	
827		Wild, a challenging dataset with 6,445 trials fea-	
828		ture diverse acoustic conditions and speaking	
829		styles. LibriSpeech	
830	C Downstream Model		
831	Model Architecture.	Our downstream model	
832		consists of the pretrained SSL encoder (frozen)	
833		followed by temporal mean pooling and an L2 nor-	
834		malization layer. A single linear embedding head	
835		projects the pooled representations to speaker logits	
836		for classification. This minimal design isolates the	
837		quality of the pretrained representations by limiting	
838		the capacity of the downstream classifier.	
839	Training Setup.	We train the embedding head	
840		on VoxCeleb development set using cross-entropy	
841		loss with the AM-Softmax objective (Wang et al.,	
842		2018). The SSL encoders remain frozen through-	
843		out training, ensuring that performance differences	
844		reflect the quality of pretrained features rather than	
845		fine-tuning capacity. We train for 30 epochs with a	
846		batch size of 32 and learning rate of 1e-3.	
847	D Demographic Subgroup Analysis		
848	Description of the EARS Dataset.	The EARS	
849		(Expressive Anechoic Recordings of Speech)	
850		dataset contains speakers spanning ages 18-75,	
		with balanced gender representation. We evalu-	851
		ate on 102 speakers (59 female, 43 male) across six	852
		age brackets.	853
	E Zero-Shot		854
	Profiling Tasks.		855
	• Spoof Detection:	Binary classification task	856
		detecting synthesized or manipulated speech	857
		versus genuine human speech.	858
	• Speaker Counting:	Predicting the number of	859
		unique speakers in an audio segment.	860
	• Accent Detection:	Classifying the regional	861
		accent of the speaker. The dataset includes	862
		nine distinct distance classes.	863
	• HowFarSpk:	Determining the spatial dis-	864
		tance between speaker and microphone. The	865
		dataset includes three distinct distance classes.	866
	• Gender Recognition:	Binary classification of	867
		speaker gender.	868
	• Age Classification:	Multi-class prediction of	869
		speaker age group.	870

Category	Subgroup	WavLM	Hubert	DELULU	Wav2Vec2
Gender	Female (59 spks)	41.69	41.06	28.11	45.60
	Male (43 spks)	40.76	41.47	28.54	44.70
Age	F (18–25), 13 spks	41.93	39.67	31.01	45.03
	F (26–35), 13 spks	40.60	41.61	30.90	43.90
	F (36–45), 7 spks	43.37	44.26	29.38	47.40
	F (46–55), 14 spks	42.17	40.48	29.56	45.78
	F (56–65), 10 spks	42.71	42.02	31.28	47.78
	F (66–75), 2 spks	43.66	42.25	40.12	48.32
	M (18–25), 14 spks	41.55	43.06	35.01	45.87
	M (26–35), 10 spks	40.51	40.60	30.90	42.97
	M (36–45), 10 spks	39.40	39.47	24.53	45.35
	M (46–55), 4 spks	41.27	42.55	29.63	44.25
	M (56–65), 5 spks	42.53	43.32	31.37	47.53

Table 7: EER (%) across demographic subgroups in EARS dataset on upstream speaker verification (\downarrow better)