## N-CORE: N-View Consistency Regularization for Disentangled Representation Learning in Nonverbal Vocalizations

## Siddhant Bikram Shah<sup>1</sup> Kristina T. Johnson<sup>1,2</sup>

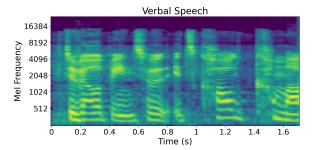
<sup>1</sup>Department of Electrical and Computer Engineering <sup>2</sup>Department of Communication Sciences and Disorders Northeastern University, Boston, USA

### **Abstract**

Nonverbal vocalizations are an essential component of human communication, conveying rich information without linguistic content. However, their computational analysis is hindered by a lack of lexical anchors in the data, compounded by biased and imbalanced data distributions. While disentangled representation learning has shown promise in isolating specific speech features, its application to nonverbal vocalizations remains unexplored. In this paper, we introduce N-CORE, a novel backbone-agnostic framework designed to disentangle intertwined features like emotion and speaker information from nonverbal vocalizations by leveraging N views of audio samples to learn invariance to specific transformations. N-CORE achieves competitive performance compared to state-of-the-art methods for emotion and speaker classification on the VIVAE, ReCANVo, and ReCANVo-Balanced datasets. We further propose an emotion perturbation function that disrupts affective information while preserving speaker information in audio signals for emotion-invariant speaker classification. Our work informs research directions on paralinguistic speech processing, including clinical diagnoses of atypical speech and longitudinal analysis of communicative development. Our code is available at https://github.com/SiddhantBikram/N-CORE.

### 1 Introduction

Nonverbal vocalizations (NVVs) are a fundamental component of human communication, encompassing a diverse range of non-speech sounds like laughter, sighs, cries, and other sounds that convey rich affective information without linguistic content (Cowen et al., 2019). Interpreting these paralinguistic sounds is vital for comprehensive modeling of human communication and developing emotionally intelligent AI systems (Tzirakis et al., 2023). However, their computational analysis presents unique



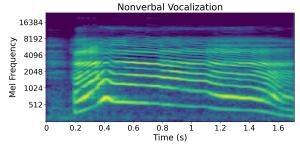


Figure 1: Comparison of mel-spectrograms from verbal (top) and nonverbal vocalizations (bottom). The syllabic structure of word-based speech results in specific temporal variations that are less common in NVVs.

challenges distinct from conventional speech processing.

A primary bottleneck in NVV analysis is the scarcity of annotated data. Unlike speech corpora with millions of hours of recorded content, NVV datasets are typically restricted to only a few hundred hours (Baird et al., 2022; Koudounas et al., 2025), limiting the performance of contemporary data-intensive machine learning (ML) and deep learning (DL) methods. This limitation is exacerbated by substantial biases arising from the intertwined nature of emotion and speaker labels in these datasets, leading to models that learn spurious correlations. Consequently, emotion classification models trained on such corpora are confounded by speaker characteristics like pitch range, vocal timbre, and articulation patterns, and conversely, speaker classification models are confounded by affective characteristics like dynamic

intensity, prosodic contours, and fundamental frequency (F0) variations (Pei et al., 2024). This hinders generalization across demographic groups and emotion categories, especially in low-resource settings.

Several methods have attempted to address these challenges through various representation learning approaches. Foundation models such as HuBERT (Hsu et al., 2021b) and Wav2Vec2 (Baevski et al., 2020) have demonstrated remarkable success in learning generalized speech representations that can be fine-tuned for downstream tasks. These models are predominantly trained on verbal corpora, where canonical phoneme structures and linguistic content serve as strong structural priors. In contrast, NVVs lack the phoneme-based priors these models exploit, causing them to struggle when encoding paralinguistic speech (Lane et al., 2015; Tzirakis et al., 2023). Figure 1 compares melspectrograms of verbal speech and NVVs, highlighting how verbal speech has more complex spectral variability and clear transitions in temporal segmentation as compared to NVVs, which may assist representation learning (Nagamine et al., 2015).

Disentangled representation learning (DRL), the process of separating different informational factors in data, has been extensively explored in the speech domain for tasks including emotion recognition (Yuan et al., 2024; Xi et al., 2022), depression detection (Ravi et al., 2022), and voice conversion (Zuo et al., 2024; Wang et al., 2021a). However, extending DRL methods to NVVs is nontrivial: NVVs lack lexical anchors, and their prosodic characteristics simultaneously contain both speaker and emotion features. Conventional DRL methods on speech data often depend on perturbation strategies that preserve lexical content while altering specific features (Tu et al., 2024; Hsu et al., 2019); however, in the absence of transformation-invariant lexical content, a single perturbation may either disrupt useful information or result in uninformative artifacts persisting in the signal.

In this paper, we investigate DRL for NVVs. Our contributions are summarized as follows:

- We propose N-CORE: N-View COnsistency REgularization (pronounced 'Encore'), a novel framework for supervised DRL of NVVs by using N perturbed views of an audio signal.
- We propose a novel transformation function to perturb affective components in speech sig-

- nals while retaining speaker characteristics. We further examine the validity of an existing speaker perturbation method on NVVs.
- We comprehensively benchmark audio foundation models, domain-specific models, DRL methods, and state-of-the-art frameworks on emotion and speaker classification tasks across the VIVAE, ReCANVo, and ReCANVo-Balanced datasets. To the best of our knowledge, we are the first to study DRL in NVVs.

### 2 Related Work

## 2.1 Machine Learning for Nonverbal Vocalizations and Paralinguistic Speech

Early ML research on NVVs relied predominantly on hand-engineered feature sets. For instance, Schuller et al. (2013) established the ComParE acoustic feature set that captured spectral, prosodic, and voice quality parameters for the paralinguistic analysis of social signals, conflict, and emotion, with application to autism diagnosis. This was further refined by the GeMAPS and eGEMAPS frameworks (Eyben et al., 2015), providing a standardized feature extraction framework for affective computing applications. These approaches have been successfully employed for classifying NVVs (Lefter and Jonker, 2017; Narain et al., 2020), typically using traditional ML classifiers like Support Vector Machines (Cortes and Vapnik, 1995) and Random Forests (Breiman, 2001).

The advent of DL has significantly advanced NVV processing by enabling representation learning directly from raw waveforms and bypassing manual feature engineering. Early approaches used convolutional neural networks (CNNs) like ResNet (He et al., 2016) to process paralinguistic speech for tasks like understanding nonverbal emotion (Hsu et al., 2021a), classifying speakers (Xu et al., 2024), and judging singing voice quality (Xu et al., 2022). More recently, self-supervised learning has revolutionized ML for speech, with models like Wav2Vec2 (Baevski et al., 2020) and HuBERT (Hsu et al., 2021b) achieving state-of-the-art performance on various speech processing benchmarks. Although pretrained primarily on linguistic content, they exhibit strong transferability to paralinguistic tasks (Tzirakis et al., 2023; Shah and Johnson, 2025; Phukan et al., 2025); however, their performance is often limited by domain shifts between pretraining and evaluation data. Addressing this

gap, Koudounas et al. (2025) developed Voc2Vec, which pre-trains the Wav2Vec2 architecture over multiple NVV datasets with a self-supervised learning objective.

In clinical applications, ML approaches have been instrumental in analyzing atypical vocalizations. For instance, Bone et al. (2017) developed a classification framework to identify distinctive acoustic signatures in the vocalizations of children with autism spectrum disorder (ASD). Similarly, Narain et al. (2022) demonstrated that ML methods could effectively classify affective and communicative functions in NVVs from individuals with ASD. Further, these techniques have been applied to speech therapy (Mulfari et al., 2021), automatic speech recognition (Mulfari et al., 2023), and speech conversion (Doshi et al., 2021) for individuals with atypical speech.

## 2.2 Disentangled Representation Learning in Speech

DRL aims to separate distinct informational factors in data, enabling models to extract and manipulate independent semantic dimensions (Wang et al., 2024b). In speech processing, DRL typically focuses on separating speaker characteristics, linguistic content, and emotion from each other (Williams, 2022). This separation is valuable for tasks like voice conversion (Luong and Tran, 2021), speech recognition (Trinh and Braun, 2022), and emotion recognition (Yuan et al., 2024), where isolating specific features leads to improved performance. These methods harness lexical content and phoneme sequences in speech (Hsu et al., 2019) as stable anchors against the transformation of various attributes like emotion or speaker identity, which are conveyed through prosodic modulations (Chu et al., 2006). The application of these techniques to NVVs, which lack explicit lexical anchors and have entangled speaker and emotion information in their prosodic features, remains an unexplored domain, motivating us to investigate DRL in NVVs.

A prominent DRL approach involves a gradient reversal layer (GRL) (Ganin and Lempitsky, 2015), enabling end-to-end training of classifiers invariant to characteristics like domain (Lu et al., 2022) and speaker identity (Oneaţă et al., 2021). Autoencoderbased methods are also widely used to learn disentangled latent spaces by imposing specific constraints on the latent distribution (Yingzhen and Mandt, 2018; Nam et al., 2024). Subsequent frameworks like NANSY (Choi et al., 2021) and Con-

tentVec (Qian et al., 2022) learn speaker-invariant speech representations by encouraging models to learn similar representations for audio pairs with perturbed speaker information; however, a single perturbed view per sample may not expose the model to the spectrum of acoustic and affective variability present in datasets, limiting the robustness of the learned invariances. Further, these methods are limited to speaker-invariant representation learning, as they rely solely on speaker perturbation. To address these gaps, we propose N-CORE, which uses N views of perturbed samples from an audio signal for increased sample diversity. We further propose an emotion perturbation function that selectively alters affective components while preserving speaker information for emotion-invariant representation learning.

## 3 Methodology

In this section, we describe N-CORE, our proposed supervised DRL framework to encode NVVs by isolating either emotion- or speaker-specific information. Our backbone-agnostic framework can utilize any representation learner as the backbone encoder. It applies audio perturbations to suppress either emotion or speaker information while preserving the complementary features. We generate N perturbed views per audio sample to encourage invariance across a broader distribution of irrelevant variations, regulated by a pairwise distance loss for consistency regularization. We use two classification heads-a primary head to predict the target label, and a secondary adversarial head with a GRL mechanism-to simultaneously promote taskrelevant features while discarding task-irrelevant information in the learned representations. We train the model via a composite objective that balances the regularization loss, cross-entropy loss, and gradient reversal loss. Figure 2 presents our audio perturbation functions alongside the overall architecture of the N-CORE framework.

#### 3.1 Problem Formulation

Let X represent an acoustic signal encompassing an NVV with a target positive label  $y^+$  and a negative label  $y^-$ . We aim to learn a representation model R = f(X) that maps X to a learned embedding  $x \in \mathbb{R}^D$ , encapsulating the core components of  $y^+$  from the sample while discarding information that describes  $y^-$ . Specifically, if the learning objective is to predict the emotion label  $y^e$ , x must

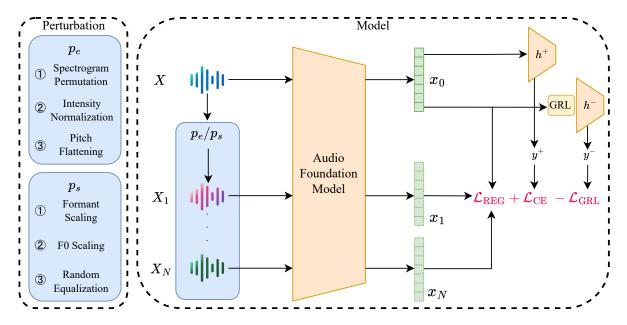


Figure 2: Our proposed framework, N-CORE, predicts label  $y^+$  and disentangles features that inform the label  $y^-$ . Perturbation functions  $p_e$  or  $p_s$  are used to create N views of the original sample X for consistency regularization. Cross-entropy loss is used to predict  $y^+$  with classification head  $h^+$ , and a GRL is used to adversarially disentangle  $y^-$  using classification head  $h^-$ .

retain information pertinent to the underlying emotion expressed in X while remaining uninformative with respect to speaker label  $y^s$ . Conversely, when predicting  $y^s$ , x should encapsulate speaker-specific traits from X while discarding affective content descriptive of  $y^e$ . Achieving such disentanglement is challenging given the inherent entanglement of emotion- and speaker-specific information in the audio signal.

### 3.2 Representation Learner

N-CORE employs a pre-trained audio foundation model (AFM) as its core feature encoder. Our framework is backbone-agnostic, allowing it to use various foundation models like HuBERT (Hsu et al., 2021b) and Wav2Vec2 (Baevski et al., 2020) interchangeably. This model learns a neural embedding x from the raw audio signal X by encoding essential phonetic, prosodic, and stylistic information (Kharitonov et al., 2021), as x = AFM(X).

### 3.3 Feature-Invariant Audio Perturbation

**Emotion Perturbation.** We aim to disrupt affective information in the audio signal while preserving speaker characteristics. The emotion perturbation function  $p_e$  comprises three transformations: 1) We compute the Short-Time Fourier Transform (STFT) of X, resulting in a spectrogram S(X) with  $n_{\rm spec}$  non-overlapping frequency bands. We ran-

domly permute  $\eta_1$  of these bands, retaining the rhythm and energy essential for speaker identification (Quatieri et al., 1994), while distorting content information (Davis and Johnsrude, 2003). 2) We normalize intensity by adjusting the waveform's RMS to a fixed target  $\eta_2$  in order to suppress dynamic intensity correlated with emotion features (Koolagudi and Rao, 2012). 3) We flatten the pitch of the speaker to its mean in their pitch contour  $f_0$ , effectively flattening prosodic variance and the affective content it contains (Mozziconacci, 2002). **Speaker Perturbation.** We adopt the audio transformation pipeline designed by Choi et al. (2021) for the NANSY framework to perturb speaker information while preserving the underlying content information. Similar to ContentVec (Qian et al., 2022), the speaker perturbation function  $p_s$  comprises three transformations: 1) scaling formant frequencies by a factor of  $\rho_1$ ; 2) scaling F0 in every frame by  $\rho_2$ , and 3) applying a random equalizer to account for channel variations.

#### 3.4 N Perturbed Views

Prior work on feature-invariant representation learning (Qian et al., 2022; Tu et al., 2024; Wang et al., 2024a) typically relies on a single perturbed version of each input and then enforces invariance between them. This one-shot strategy constrains the diversity of transformations exposed to the

model, making its learned invariances less robust to unseen distortions.

In contrast, our approach samples N distinct perturbations drawn independently from the original audio signal X. By exposing the model to a spectrum of variations, we increase the range of uninformative factors the encoder is encouraged to ignore, reduce reliance on any single perturbation pattern, and promote the consistent encoding of all views of X into a tight cluster in the representation space. Multiple perturbations are especially crucial in NVVs, which lack lexical anchors that could be preserved after perturbation (Ko et al., 2015). We regularize the pairwise distances among all views by measuring the average squared distance across all unique pairs as a loss function:

$$\mathcal{L}_{REG} = \frac{\sum_{i=0}^{N} \sum_{j=i+1}^{N} ||x_i - x_j||_2^2}{\frac{N(N+1)}{2}}$$
 (1)

where the denominator is the number of unique pairs among the set of N+1 embeddings  $\{x_0, x_1...x_N\}$ , including the unperturbed view  $x_0$ . This loss encourages the model to create the same representation for all views of X, disentangling irrelevant information from task-relevant features.

## 3.5 Classification

We project x to two separate classification heads  $h^+$  and  $h^-$  that use cross-entropy to predict labels  $y^+$  and  $y^-$ , respectively. This step operates solely on the unperturbed view  $x_0$ . Both heads share the same underlying structure: a two-layer multilayer perceptron with ReLU and dropout in between. To enforce invariance to  $y^-$ , we precede  $h^-$  with a GRL that scales embeddings by  $-\alpha$ , encouraging the model to disentangle and suppress features corresponding to  $y^-$  in its learned representations. We obtain losses  $\mathcal{L}_{\text{CE}}$  and  $\mathcal{L}_{\text{GRL}}$  as follows:

$$\mathcal{L}_{CE} = -\sum_{k^{+}=1}^{K^{+}} y_{k^{+}}^{+} \log[h^{+}(x_{0})]_{k^{+}}$$
 (2)

$$\mathcal{L}_{GRL} = -\sum_{k^{-}=1}^{K^{-}} y_{k^{-}}^{-} \log \left[ h^{-} \left( GRL_{\alpha}(x_{0}) \right) \right]_{k^{-}}$$
(3)

### 3.6 Training Objective

Our model is trained by optimizing a composite objective function comprising the three losses obtained from equations 1, 2, and 3, calculated for each dataset sample:

$$\mathcal{L}_{total} = \lambda_{REG} \cdot \mathcal{L}_{REG} + \lambda_{CE} \cdot \mathcal{L}_{CE} - \lambda_{GRL} \cdot \mathcal{L}_{GRL}$$
 (4)

where  $\lambda_{\rm reg}$ ,  $\lambda_{\rm CE}$ , and  $\lambda_{\rm GRL}$  represent scaling factors that regulate the contribution of each loss towards  $\mathcal{L}_{\rm total}$ . The optimizer minimizes  $\mathcal{L}_{\rm total}$  by maximizing the negative term  $\mathcal{L}_{\rm GRL}$ , designed to learn representations that are invariant to the secondary label  $y^-$ .

### 4 Experimental Settings

### 4.1 Datasets

We evaluate our methods on three NVV datasets: Variably Intense Vocalizations of Affect and Emotion (VIVAE) (Holz et al., 2022), Real-World Communicative and Affective Nonverbal Vocalizations (ReCANVo) (Johnson et al., 2023), and ReCANVo-Balanced. Each dataset sample has both an emotion label and a speaker identity label, but has only one class per label type. For each dataset, we evaluate performance on emotion and speaker recognition tasks. We use a train/test split of 80/20 for all datasets. Detailed dataset statistics are presented in Appendix A.2.

**VIVAE.** The VIVAE corpus comprises 1,085 non-speech emotion vocalizations produced by 11 non-professional female actors, 20-39 years old, instructed to express six affective states: achievement/triumph, sexual pleasure, surprise, anger, fear, and physical pain across multiple intensity levels. ReCANVo. The ReCANVo dataset contains 7,077 NVVs collected from eight non- and minimally-speaking individuals, ranging in age from 6-23 years old and diagnosed with various neurodevelopmental disorders, including ASD, cerebral palsy, and genetic neurodevelopmental disorders. Classes with sample counts reaching n≥100 were extracted from this dataset, yielding a derived dataset of 6,551 utterances distributed among seven functions: delighted, dysregulated, frustrated, laughter, request, self-talk, and social. This derived dataset is highly imbalanced with an imbalance factor of 18.66.

**ReCANVo-Balanced.** We use a multi-stage sampling procedure to create a balanced subset from ReCANVo by extracting 100 samples for each emotion class. Within each emotion category, participant diversity was maximized by systematically distributing the sample selection, with the constraint that no single participant would contribute a majority of samples for any given emotion class.

### 4.2 Baselines

We conduct a comprehensive benchmark of audio ML methods on NVVs, organized into two main families based on their foundational architecture. For the Wav2Vec2-based methods, we evaluate Wav2Vec2-Base (Baevski et al., 2020), Voc2Vec (Koudounas et al., 2025), Wav2Vec2-GRL (Ganin and Lempitsky, 2015), and SACE (Dutta and Ganapathy, 2024). For the HuBERT-based methods, we evaluate HuBERT-Base (Hsu et al., 2021b), HuBERT-ER and HuBERT-SID (Yang et al., 2021), and ContentVec (Qian et al., 2022). To demonstrate the flexibility of our proposed backboneagnostic framework, we use N-CORE with both Wav2Vec2 and HuBERT as backbones, denoted as N-CORE<sub>Wav2Vec2</sub> and N-CORE<sub>HuBERT</sub>, respectively. Detailed implementation details are given in Appendix A.1.

### 5 Experimental Results

Tables 1 and 2 present the results for emotion and speaker recognition, respectively. We conduct each experiment on three seeds and report the Mean and Standard Deviation (±) for Accuracy, F1-score (Macro), and Unweighted Average Recall (Macro).

# 5.1 Emotion Classification with Speaker Disentanglement

Foundation Models. In line with previous research on emotion and speaker classification (Wang et al., 2021b), HuBERT consistently achieves the highest performance across all metrics in all datasets compared to the Wav2Vec2 family of models. The Voc2Vec model was trained exclusively on NVVs, allowing it to outperform Wav2Vec2 with the same architecture, demonstrating the advantage of domain-specific pre-training. Further, its self-supervised training objective may enable it to avoid overfitting and classification unfairness (Liu et al., 2021), as demonstrated by the difference in F1-Score and UAR compared to Wav2Vec2. However, despite being specifically designed for NVVs, Voc2Vec underperforms HuBERT on ReCANVo and ReCANVo-Balanced while matching its performance on VIVAE, suggesting that domain-specific pre-training may not solely surpass the representation learning power of a more suitable model.

**Domain-Specific Models.** Notably, neither HuBERT-ER nor HuBERT-SID outperforms the baseline HuBERT model, which may be attributed to the domain shift between the spoken

word datasets used during finetuning and the NVV datasets used for this evaluation. Further, fine-tuning on a smaller corpus limits the generalizability of these models to out-of-distribution data.

GRL-based Models. Both Wav2Vec2-GRL and HuBERT-GRL show improvements in performance across VIVAE and ReCANVo compared to their respective baselines. These results support our hypothesis that using adversarial training to explicitly disentangle speaker information leads to more robust representations that are less sensitive to speaker-specific characteristics and biases. However, they underperform their respective baselines on ReCANVo-Balanced.

**DRL Frameworks.** N-CORE<sub>Wav2Vec2</sub> and N-CORE<sub>HuBERT</sub> outperform all other methods in their respective model families on VIVAE and ReCANVo-Balanced, but fall short for ReCANVo. This may be due to ReCANVo's intertwined speaker and emotion distributions, where models could be relying on speaker characteristics to predict emotions due to a biased sample distribution (see Table 6); thus, N-CORE's superior DRL capabilities may have penalized its performance. ReCANVo-Balanced mitigates this imbalance, and N-CORE outperforms all respective methods here.

# 5.2 Speaker Classification with Affect Disentanglement

**Foundation Models.** For the ReCANVo dataset, Voc2Vec performs worse than HuBERT and Wav2Vec2, despite ReCANVo being a part of its pre-training corpus; however, it surpasses both on ReCANVo-Balanced. Voc2Vec also uses the VI-VAE dataset for pre-training, on which it performs the best, followed by HuBERT and Wav2Vec2, respectively.

Domain-Specific Models. HuBERT-ER shows competitive performance for speaker identification compared to the baseline and even the specialized HuBERT-SID model on ReCANVo, but exhibits a substantial drop on VIVAE, highlighting the importance of task-specific pre-training. However, the model performs poorly on ReCANVo-Balanced, suggesting that it could be exploiting affective information to predict speakers in ReCANVo.

**GRL-based Methods.** On ReCANVo and VI-VAE, both Wav2Vec2-GRL and HuBERT-GRL demonstrate substantial performance gains after disentanglement. On ReCANVo-Balanced, Wav2Vec2-GRL exhibits improvement in perfor-

| Model                      |                   | VIVAE             |                   |                   | ReCANVo                                 |                   | ReCANVo-Balanced  |                   |                   |  |  |  |
|----------------------------|-------------------|-------------------|-------------------|-------------------|-----------------------------------------|-------------------|-------------------|-------------------|-------------------|--|--|--|
|                            | Acc               | F1                | UAR               | Acc               | F1                                      | UAR               | Acc               | F1                | UAR               |  |  |  |
| Wav2Vec2-based Methods     |                   |                   |                   |                   |                                         |                   |                   |                   |                   |  |  |  |
| Wav2Vec2                   | 54.22±2.1         | 53.75±1.2         | 53.82±2.0         | <b>63.95</b> ±0.7 | 51.76±0.8                               | 51.70±1.0         | 28.33±1.4         | 23.40±1.5         | 28.33±1.4         |  |  |  |
| Voc2Vec                    | 59.14±3.6         | 58.69±3.0         | 58.77±3.7         | 61.35±0.8         | 50.48±1.7                               | 48.75±1.5         | 30.29±1.6         | 27.29±2.5         | 30.29±1.6         |  |  |  |
| Wav2Vec2-GRL               | 54.99±1.5         | 54.41±1.7         | 54.43±1.8         | 63.34±0.3         | <b>52.86</b> ±0.6                       | <b>53.06</b> ±1.0 | 27.86±1.5         | 21.18±2.4         | 27.86±1.5         |  |  |  |
| SACE                       | 51.77±3.1         | 50.73±3.2         | 51.42±3.1         | 63.64±0.4         | $\underline{52.11}{\scriptstyle\pm0.5}$ | 51.64±1.4         | 25.95±2.2         | 17.23±3.2         | 25.95±2.2         |  |  |  |
| N-CORE <sub>Wav2Vec2</sub> | <b>59.29</b> ±1.2 | <b>59.26</b> ±1.5 | <b>58.85</b> ±1.3 | 63.54±0.4         | 50.84±2.7                               | 51.12±2.6         | <b>31.43</b> ±2.1 | 28.26±1.7         | <b>31.43</b> ±2.1 |  |  |  |
|                            |                   |                   | HuBl              | ERT-based         | Methods                                 |                   |                   |                   |                   |  |  |  |
| HuBERT                     | 58.83±1.2         | 58.14±1.2         | 58.10±1.1         | 66.21±0.6         | 55.86±2.1                               | 55.06±1.5         | 35.24±1.2         | 31.05±2.7         | <b>35.24</b> ±1.2 |  |  |  |
| HuBERT-ER                  | 57.45±3.7         | 54.84±5.7         | 56.26±4.1         | 65.34±0.4         | 53.19±0.7                               | 53.10±0.6         | 31.90±1.7         | 26.72±1.4         | 31.90±1.7         |  |  |  |
| HuBERT-SID                 | 58.53±2.2         | 57.25±2.3         | 57.73±2.2         | 63.77±0.5         | 55.24±1.0                               | 54.03±0.4         | 29.05±2.7         | 26.69±2.4         | 29.05±2.7         |  |  |  |
| HuBERT-GRL                 | 62.98±3.0         | 62.34±2.8         | 62.12±2.8         | 66.46±0.3         | <b>56.62</b> ±1.1                       | <b>56.05</b> ±0.1 | 34.67±2.0         | 33.05±3.0         | 34.67±2.0         |  |  |  |
| ContentVec                 | 58.68±1.3         | 57.82±1.2         | 55.03±2.7         | 65.93±0.6         | <u>56.28</u> ±2.9                       | 55.03±2.7         | 34.29±0.5         | $32.70 \pm 0.5$   | $34.29 \pm 0.5$   |  |  |  |
| N-CORE <sub>HuBERT</sub>   | <b>65.13</b> ±3.1 | <b>64.37</b> ±3.0 | <b>64.39</b> ±3.1 | <b>66.59</b> ±0.2 | 54.02±1.1                               | 53.56±1.2         | <b>35.24</b> ±2.3 | <b>34.01</b> ±2.4 | <b>35.24</b> ±2.3 |  |  |  |

Table 1: Comparison of model performance on the emotion classification task for VIVAE, ReCANVo, and ReCANVo-Balanced. The results are in the form of Mean ± Standard Deviation. For each model family, the best results are highlighted in **bold** and the second-best results are <u>underlined</u>.

mance after disentanglement, whereas HuBERT-GRL experiences a performance decline relative to its baseline.

**DRL Frameworks.** Across all datasets, N-CORE<sub>Wav2Vec2</sub> and N-CORE<sub>HuBERT</sub> outperform all other methods in their respective model families. Notably, ContentVec outperforms all other methods on ReCANVo despite being trained to be invariant to speakers, indicating that speaker perturbation may not transform all speaker features.

## 5.3 Data Analysis

ReCANVo's data imbalance reflects real-life data distributions, where multi-label data often exhibit inherent biases (Schultheis et al., 2022). In this context, affective vocalizations reflect the idiosyncratic behaviors of individual autistic speakers, and since the vocalizations are not acted, some samples may naturally lie between two emotional categories. These facets limit model performance for emotion classification despite the dataset's relatively large number of samples, with performance deteriorating significantly on ReCANVo-Balanced.

Universally, speaker identification proves more challenging on the VIVAE dataset across all models, with significantly lower performance compared to ReCANVo. This dataset contains acted vocalizations from adults, where emotional expressive-

ness tends to converge on shared cultural templates for what each affective vocalization is expected to sound like. This reduces inter-speaker variability by masking natural speaker-specific cues, making it more difficult for models to distinguish between speakers, especially compared to spontaneous, real-world vocalization datasets like ReCANVo. Disentanglement was particularly effective for speaker classification on VIVAE, suggesting that DRL excels in datasets with homogeneous speaker demographics.

All the models demonstrated remarkably high performance on speaker identification for Re-CANVo, which may be due to the diverse age range of the dataset and the idiosyncratic forms of NVVs across individuals with autism (Pegado et al., 2020), making speaker classification a relatively easier ML task. The competitive performance of all models on the small-scale ReCANVo-Balanced dataset shows that even a relatively small corpus of NVVs can help create effective speaker recognition systems for heterogeneous populations.

### 5.4 Cross-Verification of Perturbation

We conducted a cross-verification experiment to validate the efficacy of our affect and speaker perturbation functions by applying each to both classification tasks in VIVAE using N-CORE<sub>HuBERT</sub>, and our results are presented in Table 3. Applying

| Model                      |                   | VIVAE             |                                |           | ReCANVo                                    |                                         | ReCANVo-Balanced  |                   |                             |  |  |  |
|----------------------------|-------------------|-------------------|--------------------------------|-----------|--------------------------------------------|-----------------------------------------|-------------------|-------------------|-----------------------------|--|--|--|
|                            | Acc               | F1                | UAR                            | Acc       | F1                                         | UAR                                     | Acc               | F1                | UAR                         |  |  |  |
| Wav2Vec2-based Methods     |                   |                   |                                |           |                                            |                                         |                   |                   |                             |  |  |  |
| Wav2Vec2                   | 44.85±1.0         | 39.88±2.0         | 42.89±2.0                      | 92.91±0.8 | 91.22±0.8                                  | 91.53±0.3                               | 63.57±7.1         | 59.31±10.8        | 60.91±9.4                   |  |  |  |
| Voc2Vec                    | 65.75±2.6         | <u>64.16</u> ±2.8 | 64.36±2.7                      | 91.30±0.1 | 89.88±0.4                                  | 89.74±0.9                               | 75.95±2.9         | 75.86±2.6         | 76.25±2.0                   |  |  |  |
| Wav2Vec-GRL                | 57.30±1.5         | 54.83±0.8         | 55.12±1.0                      | 93.95±0.1 | 92.52±0.1                                  | 92.55±0.5                               | 68.33±0.8         | 67.60±2.0         | 67.63±2.0                   |  |  |  |
| SACE                       | 43.78±0.9         | 36.79±2.7         | $40.98 {\scriptstyle \pm 0.6}$ | 93.06±0.6 | 91.54±0.6                                  | 91.83±0.4                               | 62.38±1.7         | 55.06±1.7         | 59.32±2.0                   |  |  |  |
| N-CORE <sub>Wav2Vec2</sub> | <b>75.42</b> ±3.6 | <b>74.12</b> ±3.1 | <b>73.72</b> ±3.3              | 95.25±0.4 | <b>94.07</b> ±0.5                          | <b>94.02</b> ±0.5                       | <b>85.24</b> ±2.0 | <b>84.98</b> ±3.0 | <b>85.03</b> ±3.0           |  |  |  |
|                            |                   |                   | HuB                            | ERT-based | Methods                                    |                                         |                   |                   |                             |  |  |  |
| HuBERT                     | 60.83±1.6         | 57.50±2.1         | 58.80±1.7                      | 94.38±0.2 | 93.13±0.5                                  | 93.54±0.7                               | 80.71±3.0         | 80.20±3.1         | 80.33±3.0                   |  |  |  |
| HuBERT-ER                  | 47.16±5.8         | 39.83±8.1         | 44.70±5.5                      | 94.51±0.4 | 93.47±0.6                                  | 93.65±0.6                               | 65.71±3.2         | 62.24±4.6         | 63.65±3.3                   |  |  |  |
| HuBERT-SID                 | 64.52±2.8         | 63.30±3.3         | 63.30±3.3                      | 94.43±3.3 | 93.23±0.5                                  | 93.56±0.5                               | 76.90±2.4         | 76.77±2.9         | 77.31±3.0                   |  |  |  |
| HuBERT-GRL                 | 71.43±1.9         | 69.00±1.1         | 69.81±1.5                      | 95.12±0.2 | 93.95±0.3                                  | 94.21±0.2                               | 76.90±4.7         | 77.06±3.5         | 77.20±2.9                   |  |  |  |
| ContentVec                 | 67.43±2.5         | 65.10±2.7         | 65.55±2.6                      | 95.17±0.2 | $\underline{94.13} {\scriptstyle \pm 0.2}$ | $\underline{94.48}{\scriptstyle\pm0.1}$ | 77.14±1.5         | 76.63±1.9         | $76.74{\scriptstyle\pm1.8}$ |  |  |  |
| N-CORE <sub>HuBERT</sub>   | 77.57±3.4         | <b>76.32</b> ±3.2 | <b>76.20</b> ±3.4              | 95.27±0.3 | 94.22±0.4                                  | <b>94.52</b> ±0.2                       | <b>83.10</b> ±2.4 | <b>82.62</b> ±2.7 | <b>82.91</b> ±2.3           |  |  |  |

Table 2: Comparison of model performance on the speaker classification task for VIVAE, ReCANVo, and ReCANVo-Balanced. The results are in the form of Mean ± Standard Deviation. For each model family, the best results are highlighted in **bold** and the second-best results are underlined.

speaker perturbation  $p_s$  to speaker classification or emotion perturbation  $p_e$  to emotion classification significantly degrades performance, indicating the successful disruption of cues that the respective perturbation function targets. Conversely, applying the inverse pairing leads to improved performance, indicating that the model learns to become invariant to the perturbed features, and that the respective transformations do not disrupt features informative to the classification task. This experiment validates our proposed transformation pipeline  $p_e$ , and proves the applicability of  $p_e$  and  $p_s$  to NVVs.

| Task    | Perturbation                           | Performance           |                       |                       |  |  |  |  |
|---------|----------------------------------------|-----------------------|-----------------------|-----------------------|--|--|--|--|
|         |                                        | Acc.                  | F1                    | UAR                   |  |  |  |  |
| Speaker | $p_e \ p_s$                            | <b>75.12</b> 70.97    | <b>74.25</b> 66.16    | <b>74.54</b> 68.21    |  |  |  |  |
| Emotion | $egin{array}{c} p_e \ p_s \end{array}$ | 61.29<br><b>64.06</b> | 61.18<br><b>63.01</b> | 60.64<br><b>63.52</b> |  |  |  |  |

Table 3: Cross-Verification of signal perturbation efficacy using N-CORE<sub>HuBERT</sub> on VIVAE. The best results are highlighted in **bold**.

### 5.5 Optimal number of perturbations

To identify the optimal number of perturbations (N) for N-CORE<sub>HuBERT</sub>, we evaluated the model's classification accuracy on VIVAE while varying N from 1 to 7, with results presented in Figure 3. We find that N=5 leads to the best result for

this dataset; however, this may vary with dataset size and the distribution of multi-labeled samples. We hypothesize that the drop in performance when N>5 is an early sign of overfitting to the perturbation process. With too many augmented views, the model may begin to learn idiosyncrasies of the specific transformations rather than the core, invariant features, leading to a drop in performance. This may also be the reason we observe consistent performance gains as N increases from 2 to 5.

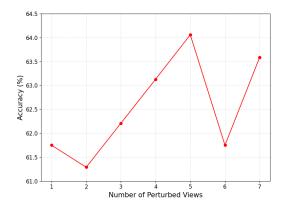


Figure 3: Accuracy vs. number of perturbed views with N-CORE<sub>HuBERT</sub> for emotion classification on VIVAE. The y-axis is limited from 61.0 to 64.5 for clarity.

### 5.6 Ablation Study

We conduct a systematic ablation study on N-CORE<sub>HuBERT</sub> to evaluate the individual contribu-

tion of its components, and the results are presented in Table 4. We observe a consistent progression in performance across all metrics as we sequentially add GRL, regularization loss, and especially the N perturbed views.

|          | Compo | nent | Performance |       |                         |       |  |  |
|----------|-------|------|-------------|-------|-------------------------|-------|--|--|
| HB       | GRL   | RL   | NV          | Acc.  | F1                      | UAR   |  |  |
| 1        |       |      |             | 58.06 | 56.51                   | 56.81 |  |  |
| /        | ✓     |      |             | 59.91 | 56.51<br>59.16<br>61.03 | 59.20 |  |  |
| ✓        | ✓     | ✓    |             | 61.75 | 61.03                   | 60.98 |  |  |
| <b>✓</b> | ✓     | ✓    | ✓           | 64.06 | 63.01                   | 63.52 |  |  |

Table 4: Ablation studies conducted on N-CORE<sub>HuBERT</sub> for emotion recognition in VIVAE. The abbreviations HB, GRL, RL, and NV refer to HuBERT, Gradient Reversal Layers, Regularization Loss, and N-Views, respectively. The final row corresponds to the entire framework. The best results are highlighted in **bold**.

### **5.7** Disentanglement Training

N-CORE<sub>HuBERT</sub>'s DRL optimization for emotion classification on VIVAE is illustrated through the loss and accuracy curves presented in Figure 4 and Figure 5, respectively. Figure 4 shows the emotion classification loss decreasing and stabilizing over epochs, while the adversarial speaker classification loss increases, as intended with the use of a GRL. Concurrently, Figure 5 shows that the emotion classification accuracy consistently improves until stabilization, whereas the speaker classification accuracy rapidly drops to random chance. These trends infer N-CORE's success in learning representations that are discriminative for emotion while simultaneously becoming invariant to speaker characteristics over the training period.

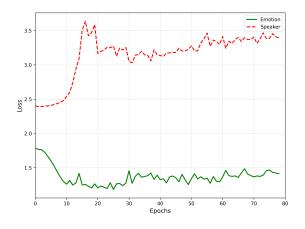


Figure 4: Loss vs. Number of Epochs for emotion classification on VIVAE by N-CORE $_{\text{HuBERT}}$ .

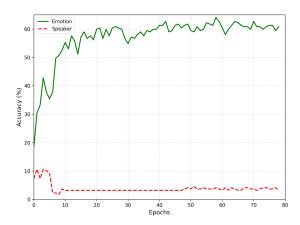


Figure 5: Accuracy vs. Number of Epochs for emotion classification on VIVAE by N-CORE<sub>HuBERT</sub>.

### 6 Conclusion

In this paper, we investigate DRL specifically for NVVs. We proposed N-CORE, a novel backboneagnostic disentanglement method using N-views of perturbed audio signals to disentangle relevant features from uninformative ones. Our experiments demonstrate that multi-view perturbation enhances performance compared to traditional single-view approaches, with N-CORE achieving competitive performance on both emotion and speaker classification tasks for VIVAE and ReCANVo-Balanced datasets. We further propose a signal transformation pipeline that perturbs emotions in speech signals while preserving speaker information. Further, we validate our emotion perturbation technique and a previously proposed speaker transformation, finding that both are generalizable to NVVs.

Our work further establishes that DRL is indeed achievable for NVVs and applies to both typical and atypical paralinguistic speech. This opens several promising directions for future research and applications, including privacy-preserving encoding of NVVs, disentangled voice conversion for NVVs, and the clinical analysis of vocalizations from non- and minimally-speaking individuals. N-CORE further empowers longitudinal studies of communicative development through NVVs that remain invariant to changes in speaker characteristics over time. The backbone-agnostic design of N-CORE allows it to scale with advances in DL, potentially benefiting from larger foundation models as they become available. Our work is an important step toward more inclusive and accurate computational models of human paralinguistic communication.

### **Ethical Considerations**

Potential Risks. We acknowledge the privacy implications of technologies that can separate speaker characteristics from communicative content. While our work demonstrates benefits for privacy-preserving representations by removing identifying speaker information from emotion-specific embeddings, this same capability could potentially be misused for unauthorized voice anonymization or modification. We emphasize that any deployment of these technologies should adhere to strict privacy protocols and informed consent requirements, particularly when working with data from vulnerable populations such as non- and minimally-speaking individuals.

**Biases.** Our experimental results highlight how dataset imbalances can significantly affect model performance. Demographic limitations of training data may introduce biases that could impact the equitable performance of these systems across different populations. We urge caution in applying these models to populations not well-represented in the training data.

**Reproducibility Statement.** We include implementation details and hyperparameter settings for all models in Appendix A.1. The code for N-CORE is available at https://github.com/SiddhantBikram/N-CORE.

## Limitations

Our study primarily focuses on disentangling emotion and speaker features. NVVs, however, convey a rich spectrum of paralinguistic information, including varying levels of intensity, different communicative intents beyond broad affective categories, and other subtle cues, which N-CORE does not explicitly disentangle. The generalizability of our findings is also constrained by the two datasets and one derived dataset we use; while diverse, they do not encompass the full variability of NVVs across different cultures, age ranges, real-world acoustic environments, or clinical populations. The general challenge of limited annotated NVV data also impacts the scale at which models can be trained and validated.

N-CORE's performance was comparatively lower for emotion recognition on the highly imbalanced ReCANVo dataset. This suggests that in scenarios with extreme data imbalances or where speaker and affective cues are deeply convoluted, our model's strong disentanglement capabilities

might not directly translate to optimal performance for classification.

#### References

- Alexei Baevski, 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.
- Alice Baird, Panagiotis Tzirakis, Jeffrey A Brooks, Chris B Gregory, Björn Schuller, Anton Batliner, Dacher Keltner, and Alan Cowen. 2022. The acii 2022 affective vocal bursts workshop & competition. In 2022 10th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW), pages 1–5. IEEE.
- Daniel Bone, Chi-Chun Lee, Theodora Chaspari, James Gibson, and Shrikanth Narayanan. 2017. Signal processing and machine learning for mental health research and clinical applications [perspectives]. *IEEE Signal Processing Magazine*, 34(5):196–195.
- Leo Breiman. 2001. Random forests. *Machine learning*, 45:5–32.
- Hyeong-Seok Choi, Juheon Lee, Wansoo Kim, Jie Lee, Hoon Heo, and Kyogu Lee. 2021. Neural analysis and synthesis: Reconstructing speech from self-supervised representations. *Advances in Neural Information Processing Systems*, 34:16251–16265.
- Min Chu, Yong Zhao, and Eric Chang. 2006. Modeling stylized invariance and local variability of prosody in text-to-speech synthesis. *Speech Communication*, 48(6):716–726.
- Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20:273–297.
- Alan S Cowen, Petri Laukka, Hillary Anger Elfenbein, Runjing Liu, and Dacher Keltner. 2019. The primacy of categories in the recognition of 12 emotions in speech prosody across two cultures. *Nature human behaviour*, 3(4):369–382.
- Matthew H Davis and Ingrid S Johnsrude. 2003. Hierarchical processing in spoken language comprehension. *Journal of Neuroscience*, 23(8):3423–3431.
- Rohan Doshi, Youzheng Chen, Liyang Jiang, Xia Zhang, Fadi Biadsy, Bhuvana Ramabhadran, Fang Chu, Andrew Rosenberg, and Pedro J Moreno. 2021. Extending parrotron: An end-to-end, speech conversion and speech recognition model for atypical speech. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6988–6992. IEEE.
- Soumya Dutta and Sriram Ganapathy. 2024. Zero shot audio to audio emotion transfer with speaker disentanglement. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10371–10375. IEEE.

- Florian Eyben, Klaus R Scherer, Björn W Schuller, Johan Sundberg, Elisabeth André, Carlos Busso, Laurence Y Devillers, Julien Epps, Petri Laukka, Shrikanth S Narayanan, and 1 others. 2015. The geneva minimalistic acoustic parameter set (gemaps) for voice research and affective computing. *IEEE transactions on affective computing*, 7(2):190–202.
- Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on* computer vision and pattern recognition, pages 770– 778.
- Natalie Holz, Pauline Larrouy-Maestri, and David Poeppel. 2022. The variably intense vocalizations of affect and emotion (vivae) corpus prompts new perspective on nonspeech perception. *Emotion*, 22(1):213.
- Jia-Hao Hsu, Ming-Hsiang Su, Chung-Hsien Wu, and Yi-Hsuan Chen. 2021a. Speech emotion recognition considering nonverbal vocalization in affective conversations. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:1675–1686.
- Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021b. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM transactions on audio, speech, and language processing*, 29:3451–3460.
- Wei-Ning Hsu, Yu Zhang, Ron J Weiss, Yu-An Chung, Yuxuan Wang, Yonghui Wu, and James Glass. 2019. Disentangling correlated speaker and noise for speech synthesis via data augmentation and adversarial factorization. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5901–5905. IEEE.
- Kristina T Johnson, Jaya Narain, Thomas Quatieri, Pattie Maes, and Rosalind W Picard. 2023. Recanvo: A database of real-world communicative and affective nonverbal vocalizations. *Scientific Data*, 10(1):523.
- Eugene Kharitonov, Ann Lee, Adam Polyak, Yossi Adi, Jade Copet, Kushal Lakhotia, Tu-Anh Nguyen, Morgane Rivière, Abdelrahman Mohamed, Emmanuel Dupoux, and 1 others. 2021. Text-free prosodyaware generative spoken language modeling. *arXiv* preprint arXiv:2109.03264.
- Tom Ko, Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur. 2015. Audio augmentation for speech recognition. In *Interspeech*, volume 2015, page 3586.
- Shashidhar G Koolagudi and K Sreenivasa Rao. 2012. Emotion recognition from speech: a review. *International journal of speech technology*, 15:99–117.

- Alkis Koudounas, Moreno La Quatra, Sabato Marco Siniscalchi, and Elena Baralis. 2025. voc2vec: A foundation model for non-verbal vocalization. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.
- Nicholas D Lane, Petko Georgiev, and Lorena Qendro. 2015. Deepear: robust smartphone audio sensing in unconstrained acoustic environments using deep learning. In *Proceedings of the 2015 ACM international joint conference on pervasive and ubiquitous computing*, pages 283–294.
- Iulia Lefter and Catholijn M Jonker. 2017. Aggression recognition using overlapping speech. In 2017 Seventh International Conference on Affective Computing and Intelligent Interaction (ACII), pages 299–304. IEEE.
- Hong Liu, Jeff Z HaoChen, Adrien Gaidon, and Tengyu Ma. 2021. Self-supervised learning is more robust to dataset imbalance. *arXiv preprint arXiv:2110.05025*.
- Cheng Lu, Yuan Zong, Wenming Zheng, Yang Li, Chuangao Tang, and Björn W Schuller. 2022. Domain invariant feature learning for speaker-independent speech emotion recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:2217–2230.
- Manh Luong and Viet Anh Tran. 2021. Many-tomany voice conversion based feature disentanglement using variational autoencoder. *arXiv* preprint arXiv:2107.06642.
- Sylvie Mozziconacci. 2002. Prosody and emotions. In *Speech prosody*, volume 2002, pages 1–9.
- Davide Mulfari, Lorenzo Carnevale, and Massimo Villari. 2023. Toward a lightweight asr solution for atypical speech on the edge. *Future Generation Computer Systems*, 149:455–463.
- Davide Mulfari, Gabriele Meoni, Marco Marini, and Luca Fanucci. 2021. Machine learning assistive application for users with speech disorders. *Applied Soft Computing*, 103:107147.
- Tasha Nagamine, Michael L Seltzer, and Nima Mesgarani. 2015. Exploring how deep neural networks form phonemic categories. In *Interspeech*, pages 1912–1916.
- KiHyun Nam, Hee-Soo Heo, Jee-weon Jung, and Joon Son Chung. 2024. Disentangled representation learning for environment-agnostic speaker recognition. *arXiv preprint arXiv:2406.14559*.
- Jaya Narain, Kristina T Johnson, Craig Ferguson, Amanda O'Brien, Tanya Talkar, Yue Zhang Weninger, Peter Wofford, Thomas Quatieri, Rosalind Picard, and Pattie Maes. 2020. Personalized modeling of real-world vocalizations from nonverbal individuals. In *Proceedings of the 2020 International Conference on Multimodal Interaction*, pages 665–669.

- Jaya Narain, Kristina T Johnson, Thomas F Quatieri, Rosalind W Picard, and Pattie Maes. 2022. Modeling real-world affective and communicative nonverbal vocalizations from minimally speaking individuals. *IEEE Transactions on Affective Computing*, 13(4):2238–2253.
- Dan Oneaţă, Adriana Stan, and Horia Cucu. 2021. Speaker disentanglement in video-to-speech conversion. In 2021 29th European Signal Processing Conference (EUSIPCO), pages 46–50. IEEE.
- Felipe Pegado, Michelle HA Hendriks, Steffie Amelynck, Nicky Daniels, Jean Steyaert, Bart Boets, and Hans Op de Beeck. 2020. Adults with high functioning autism display idiosyncratic behavioral patterns, neural representations and connectivity of the 'voice area' while judging the appropriateness of emotional vocal reactions. *Cortex*, 125:90–108.
- Guanxiong Pei, Haiying Li, Yandi Lu, Yanlei Wang, Shizhen Hua, and Taihao Li. 2024. Affective computing: Recent advances, challenges, and future trends. *Intelligent Computing*, 3:0076.
- Orchid Chetia Phukan, Mohd Mujtaba Akhtar, Swarup Ranjan Behera, Sishir Kalita, Arun Balaji Buduru, Rajesh Sharma, SR Mahadeva Prasanna, and 1 others. 2025. Strong alone, stronger together: Synergizing modality-binding foundation models with optimal transport for non-verbal emotion recognition. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.
- Kaizhi Qian, Yang Zhang, Heting Gao, Junrui Ni, Cheng-I Lai, David Cox, Mark Hasegawa-Johnson, and Shiyu Chang. 2022. Contentvec: An improved self-supervised speech representation by disentangling speakers. In *International conference on machine learning*, pages 18003–18017. PMLR.
- Thomas F Quatieri, CR Jankowski, and Douglas A Reynolds. 1994. Energy onset times for speaker identification. *IEEE Signal Processing Letters*, 1(11):160–162.
- Vijay Ravi, Jinhan Wang, Jonathan Flint, and Abeer Alwan. 2022. A step towards preserving speakers' identity while detecting depression via speaker disentanglement. In *Interspeech*, volume 2022, page 3338.
- Björn Schuller, Stefan Steidl, Anton Batliner, Alessandro Vinciarelli, Klaus Scherer, Fabien Ringeval, Mohamed Chetouani, Felix Weninger, Florian Eyben, Erik Marchi, and 1 others. 2013. The interspeech 2013 computational paralinguistics challenge: Social signals, conflict, emotion, autism. In *Proceedings INTERSPEECH 2013, 14th Annual Conference of the International Speech Communication Association, Lyon, France.*
- Erik Schultheis, Marek Wydmuch, Rohit Babbar, and Krzysztof Dembczynski. 2022. On missing labels, long-tails and propensities in extreme multi-label

- classification. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1547–1557.
- Siddhant Bikram Shah and Kristina T Johnson. 2025. Multi-feature audio fusion for nonverbal vocalization classification. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.
- Viet Anh Trinh and Sebastian Braun. 2022. Unsupervised speech enhancement with speech recognition embedding and disentanglement losses. In *ICASSP* 2022-2022 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 391–395. IEEE.
- Youzhi Tu, Man-Wai Mak, and Jen-Tzung Chien. 2024. Contrastive self-supervised speaker embedding with sequential disentanglement. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Panagiotis Tzirakis, Alice Baird, Jeffrey Brooks, Christopher Gagne, Lauren Kim, Michael Opara, Christopher Gregory, Jacob Metrick, Garrett Boseck, Vineet Tiruvadi, and 1 others. 2023. Large-scale nonverbal vocalization detection using transformers. In ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1–5. IEEE.
- Huimeng Wang, Zengrui Jin, Mengzhe Geng, Shujie Hu, Guinan Li, Tianzi Wang, Haoning Xu, and Xunying Liu. 2024a. Enhancing pre-trained asr system fine-tuning for dysarthric speech recognition using adversarial data augmentation. In *ICASSP* 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 12311–12315. IEEE.
- Jie Wang, Jingbei Li, Xintao Zhao, Zhiyong Wu, Shiyin Kang, and Helen Meng. 2021a. Adversarially learning disentangled speech representations for robust multi-factor voice conversion. *arXiv preprint arXiv:2102.00184*.
- Xin Wang, Hong Chen, Si'ao Tang, Zihao Wu, and Wenwu Zhu. 2024b. Disentangled representation learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Yingzhi Wang, Abdelmoumene Boumadane, and Abdelwahab Heba. 2021b. A fine-tuned wav2vec 2.0/hubert benchmark for speech emotion recognition, speaker verification and spoken language understanding. arXiv preprint arXiv:2111.02735.
- Jennifer Williams. 2022. *Learning disentangled speech representations*. Ph.D. thesis, University of Edinburgh.
- Yu-Xuan Xi, Yan Song, Li-Rong Dai, Ian McLoughlin, and Lin Liu. 2022. Frontend attributes disentanglement for speech emotion recognition. In *ICASSP* 2022-2022 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7712–7716. IEEE.

Anfeng Xu, Kevin Huang, Tiantian Feng, Helen Tager-Flusberg, and Shrikanth Narayanan. 2024. Audiovisual child-adult speaker classification in dyadic interactions. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8090–8094. IEEE.

Yanze Xu, Weiqing Wang, Huahua Cui, Mingyang Xu, and Ming Li. 2022. Paralinguistic singing attribute recognition using supervised machine learning for describing the classical tenor solo singing voice in vocal pedagogy. *EURASIP Journal on Audio, Speech, and Music Processing*, 2022(1):8.

Shu-wen Yang, Po-Han Chi, Yung-Sung Chuang, Cheng-I Jeff Lai, Kushal Lakhotia, Yist Y Lin, Andy T Liu, Jiatong Shi, Xuankai Chang, Guan-Ting Lin, and 1 others. 2021. Superb: Speech processing universal performance benchmark. *arXiv* preprint *arXiv*:2105.01051.

Li Yingzhen and Stephan Mandt. 2018. Disentangled sequential autoencoder. In *International Conference on Machine Learning*, pages 5670–5679. PMLR.

Zhichen Yuan, CL Philip Chen, Shuzhen Li, and Tong Zhang. 2024. Disentanglement network: Disentangle the emotional features from acoustic features for speech emotion recognition. In *ICASSP* 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 11686–11690. IEEE.

Lishi Zuo, Man-Wai Mak, and Youzhi Tu. 2024. Promoting independence of depression and speaker features for speaker disentanglement in speech-based depression detection. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10191–10195. IEEE.

### A Appendix

### A.1 Implementation Details

We conducted all our experiments on Python 3.9.21 and PyTorch 2.6.0 with NVIDIA V100 and H200 GPUs. We set the physical batch size to 16, using gradient accumulation steps when necessary. We trained each model for 100 epochs with an early stopping patience of 20 while monitoring validation accuracy to report the best model for each run. By default, we used the training hyperparameters described by the authors of each tested method. When unspecified, we used a learning rate of  $10^{-5}$  with the AdamW optimizer, which was used for N-CORE. We use  $N=5, \alpha=1$ ,  $\lambda_{CE} = 1, \lambda_{REG} = 0.005, \lambda_{GRL} = 0.01, \eta_1 = 20,$  $\eta_2 = 0.05, \, \rho_1 \in [0.7, 1.4], \, \text{and} \, \rho_2 \in [0.5, 2.0]$ for all experiments on N-CORE. We used a linear scheduler with  $0.1 \times$  the number of training steps

as warmup steps. We conduct each experiment on three seeds and report the Mean and Standard Deviation. We set the three experimental seeds to 42, 100, and 510.

We implemented HuBERT<sup>1</sup>, Wav2Vec2<sup>2</sup>, and Voc2Vec2<sup>3</sup>, HuBERT-ER<sup>4</sup>, HuBERT-SID<sup>5</sup>, and ContentVec<sup>6</sup> through the HuggingFace library. We implemented GRL<sup>7</sup> using its PyTorch implementation on GitHub. We implemented SACE<sup>8</sup> using the code released by the authors.

#### A.2 Dataset Distribution

Detailed dataset statistics for VIVAE, ReCANVo, and ReCANVo-Balanced are presented in Tables 5, 6, and 7.

### A.3 TSNE Plots

We use TSNE plots to compare HuBERT and N-CORE<sub>HuBERT</sub> on the testing sets of VIVAE in Figures 6 and 7, and ReCANVo in Figures 8 and 9. Representations from N-CORE<sub>HuBERT</sub> were generated solely using the HuBERT backbone.

https://huggingface.co/facebook/hubert-base-ls960

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/facebook/wav2vec2-base-960h

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/alkiskoudounas/voc2vec

<sup>&</sup>lt;sup>4</sup>https://huggingface.co/superb/hubert-base-superb-er

<sup>&</sup>lt;sup>5</sup>https://huggingface.co/superb/hubert-base-superb-sid

<sup>&</sup>lt;sup>6</sup>https://huggingface.co/lengyue233/content-vec-best

<sup>&</sup>lt;sup>7</sup>https://github.com/tadeephuy/GradientReversal

<sup>8</sup>https://github.com/iiscleap/ZEST/

| Label       | S01 | S02 | S03 | S04 | S05 | S06 | S07 | S08 | S09 | S10 | S11 | Total |
|-------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-------|
| achievement | 16  | 11  | 12  | 18  | 20  | 12  | 17  | 16  | 18  | 14  | 7   | 161   |
| anger       | 12  | 18  | 15  | 18  | 18  | 20  | 14  | 19  | 17  | 16  | 7   | 174   |
| fear        | 16  | 17  | 14  | 18  | 19  | 19  | 17  | 18  | 17  | 13  | 8   | 176   |
| pain        | 17  | 20  | 21  | 17  | 19  | 20  | 18  | 14  | 19  | 12  | 8   | 185   |
| pleasure    | 19  | 19  | 20  | 17  | 15  | 19  | 20  | 20  | 18  | 18  | 17  | 202   |
| surprise    | 13  | 16  | 19  | 20  | 20  | 21  | 17  | 21  | 19  | 14  | 7   | 187   |
| Total       | 93  | 101 | 101 | 108 | 111 | 111 | 103 | 108 | 108 | 87  | 54  | 1085  |

Table 5: Data distribution of the VIVAE dataset.

| Label        | P01  | P02 | P03 | P05 | P06 | P08  | P11 | P16 | Total |
|--------------|------|-----|-----|-----|-----|------|-----|-----|-------|
| delighted    | 357  | 43  | 25  | 235 | 227 | 39   | 207 | 139 | 1272  |
| dysregulated | 212  | 0   | 302 | 116 | 5   | 13   | 22  | 34  | 704   |
| frustrated   | 150  | 56  | 47  | 283 | 30  | 781  | 27  | 162 | 1536  |
| request      | 130  | 13  | 61  | 6   | 124 | 44   | 22  | 19  | 419   |
| self-talk    | 564  | 34  | 55  | 286 | 56  | 503  | 33  | 354 | 1885  |
| social       | 182  | 247 | 0   | 0   | 1   | 93   | 52  | 59  | 634   |
| laughter     | 0    | 38  | 8   | 13  | 0   | 42   | 0   | 0   | 101   |
| Total        | 1595 | 431 | 498 | 939 | 443 | 1515 | 363 | 767 | 6551  |

Table 6: Data distribution of the ReCANVo dataset.

| Label        | P01 | P02 | P03 | P05 | P06 | P08 | P11 | P16 | Total |
|--------------|-----|-----|-----|-----|-----|-----|-----|-----|-------|
| delighted    | 13  | 13  | 13  | 13  | 12  | 12  | 12  | 12  | 100   |
| dysregulated | 17  | 0   | 17  | 16  | 5   | 13  | 16  | 16  | 100   |
| frustrated   | 12  | 12  | 13  | 13  | 12  | 13  | 13  | 12  | 100   |
| request      | 14  | 13  | 13  | 6   | 14  | 13  | 13  | 14  | 100   |
| self-talk    | 13  | 12  | 13  | 12  | 12  | 13  | 13  | 12  | 100   |
| social       | 20  | 20  | 0   | 0   | 1   | 20  | 20  | 19  | 100   |
| laughter     | 0   | 38  | 8   | 13  | 0   | 41  | 0   | 0   | 100   |
| Total        | 89  | 108 | 77  | 73  | 56  | 125 | 87  | 85  | 700   |

Table 7: Data distribution of the ReCANVo-Balanced dataset.

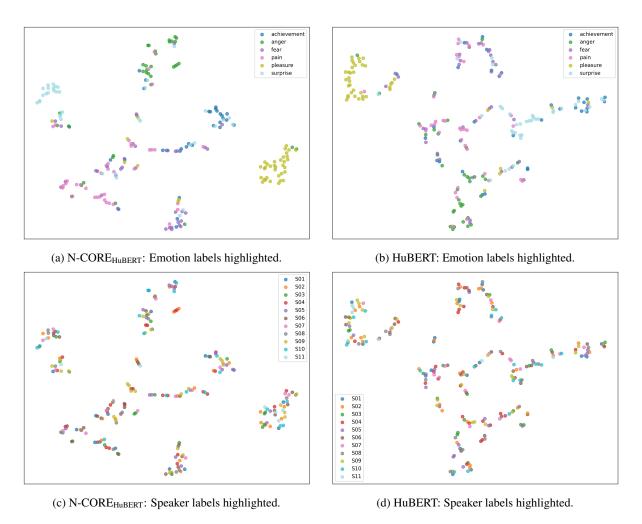


Figure 6: TSNE plots for emotion classification on VIVAE.

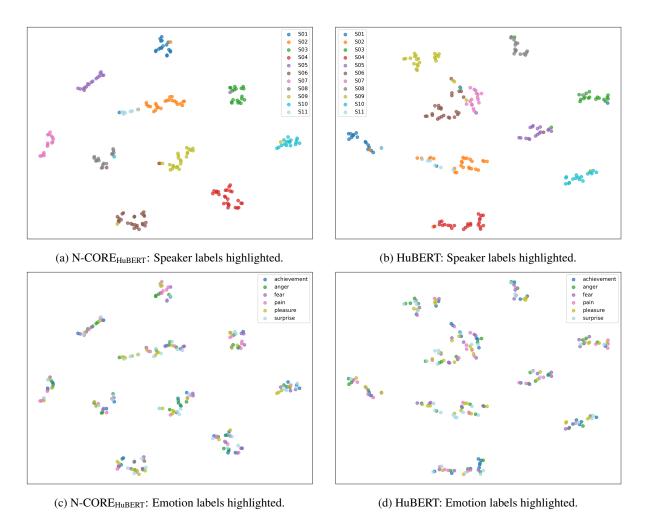


Figure 7: TSNE plots for speaker classification on VIVAE.

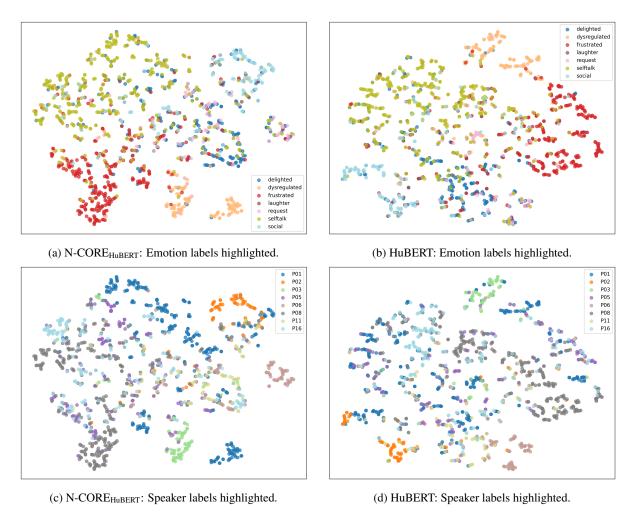


Figure 8: TSNE plots for emotion classification on ReCANVo.

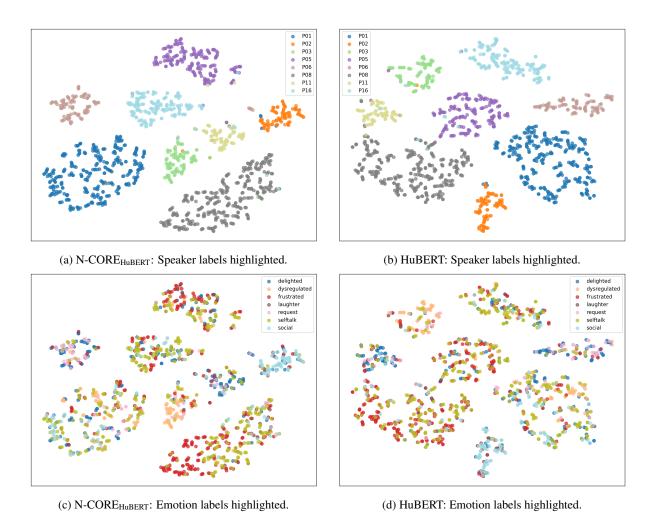


Figure 9: TSNE plots for speaker classification on ReCANVo.