

Self-supervised Speaker Verification Employing a Novel Clustering Algorithm



Abderrahim Fathan, Jahangir Alam
Computer Research Institute of Montreal (CRIM), Montreal (Quebec), Canada



Summary

- We propose a **new general-purpose clustering algorithm** called **CAMSAT**.
- CAMSAT combines the benefits of **mutual information (MI) maximization** of IMSAT clustering framework & the regularization benefit of AUGMIX (mix of augmentations at the predictions level).
- Using a thorough comparative analysis via clustering metrics, **CAMSAT allows us to outperform all other clustering algorithms** for speaker clustering.
- We achieve **better speaker verification performance (SV)** than all **SOTA SV** baselines.
- Benefits: better **generalizability, robustness, and stability** of clustering under **data shift** for **large-scale datasets** or/and a **high number of clusters**.
- We perform an ablation study to **analyze the contribution of the different components** of our proposed framework.

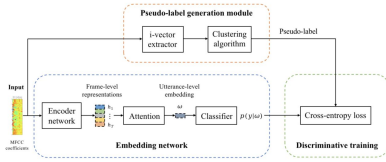
Introduction & Motivation

- **Labeled datasets** can be **expensive** and **time-consuming** to obtain.
- **Label noise** can significantly impact performance of self-supervised learning-based models.
- There are generally 2 groups of methods to learn from noisy data:
 - Noise-robust algorithms: learn directly from noisy labels
 - Label-cleansing approaches: remove or correct mislabeled data.
- How can we **improve the performance of deep clustering models**?
- How can we **improve the estimation of the ground truth number of clusters**?

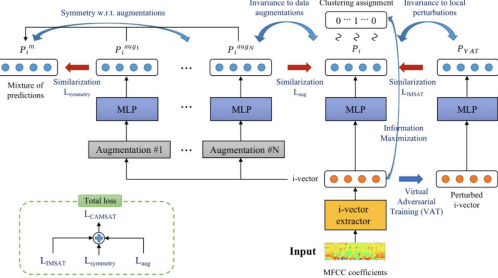
Proposed Approach

Clustering-based Self-supervised Systems

- The general process for training our clustering generated pseudo-label-based self-supervised speaker embedding networks:



Proposed CAMSAT Approach



CAMSAT | Math Equations

- We minimize the following objective:

$$L_{CAMSAT} = L_{Aug} + L_{IMSAT} + L_{Symmetry} = L_{Aug} + R_{SAT}(\theta, T_{VAT}) + \lambda(H(Y|X) - \mu H(Y)) + L_{Symmetry}$$

where:

$$L_{Aug} = \frac{1}{N} \sum_{i=1}^N KL(p_i^{aug+i} || p_i)$$

$$L_{Symmetry} = \frac{1}{N} \sum_{i=1}^N \sum_{j \in \{1, \dots, J\}} KL(p_i^{aug+j} || p_i^j)$$

$$p_i^j = \frac{1}{|J|+1} \sum_{j \in \{1, \dots, |J|\}} \alpha_j p_i^{aug+j} + p_i$$

$$R_{SAT}(\theta; T) = \frac{1}{N} \sum_{n=1}^N R_{SAT}(\theta; x_n, T(x_n))$$

$$R_{SAT}(\theta; x, T(x)) = - \sum_{c=1}^C \sum_{y_c=0}^1 p_{\theta}(y_c|x) \log p_{\theta}(y_c|T(x))$$

- L_{Aug} forces the predicted representations of augmented samples to be close to those of the original data points.
- R_{SAT} loss term allows the following:
 - Representations of the augmented samples are pushed close to those of the original samples.
 - Regularize the complexity of the network against local perturbations using Virtual Adversarial Training (VAT).
- $L_{Symmetry}$ allows the following:
 - Produce **consistent feature representations** that are labeled identically irrespective of transformations, generalize better, and generate compact clusters.
 - Mixing augmentations generates further **diverse transformations** at the latent predictions level.
 - Induce robustness and reduce the memorization of used augmentations.
 - Enforce **representation smoothness** and symmetry w.r.t. data augmentations.
 - Bootstrapping mixed predictions can be considered a **majority vote clustering** method.
 - Conduct **entropy minimization** implicitly.

Experiments & Results

Experimental Setup

- **VoxCeleb Corpus**
 - Trained using the VoxCeleb2 development set: 5994 speakers (1,092,009 utterances in total)
 - Evaluated on the VoxCeleb1-0 test set: 40 speakers (4,874 utterances in total)
- **Acoustic features:** 40-dimensional MFCCs.
- **Data Augmentation**
 - Waveform-level data augmentations: **additive noise** and **room impulse response (RIR)** simulation.
 - **SpecAugment** is applied on the fly.
 - Both **time** and **frequency maskings** are performed.
 - Additional augmentation for clustering (on i-vectors):
 - **Gaussian noise**
 - **Input Masking** (5-10% of input)
- **Clustering Evaluation Metrics**
 - **7 Supervised metrics:** Unsupervised clustering accuracy (ACC), Normalized Mutual Information (NMI), Adjusted MI (AMI), Completeness, Homogeneity, Purity, Fowlkes-Mallows index (FMI).
 - **3 Unsupervised metrics:** Silhouette score, Calinski-Harabasz score (CHS), Davies-Bouldin score (DBS).

Results

Results | Clustering & SV performances

Model	Clustering Metrics							Speaker Verification EER (%)
	ACC	AMI	NMI	No. of clusters	Completeness	Homogeneity	FMI	
Supervised (Our Labels)	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.00
IMSAT (1000)	0.82	0.62	0.57	1000	0.79	0.78	0.62	0.20
IMSAT (500)	0.80	0.62	0.56	500	0.78	0.77	0.61	0.21
IMSAT (250)	0.78	0.61	0.55	250	0.76	0.75	0.59	0.22
IMSAT (100)	0.75	0.59	0.53	100	0.73	0.72	0.56	0.23
IMSAT (50)	0.72	0.57	0.51	50	0.70	0.69	0.53	0.24
IMSAT (25)	0.68	0.54	0.48	25	0.66	0.65	0.49	0.25
IMSAT (10)	0.64	0.51	0.45	10	0.62	0.61	0.45	0.26
IMSAT (5)	0.60	0.48	0.42	5	0.58	0.57	0.41	0.27
IMSAT (2)	0.56	0.44	0.38	2	0.54	0.53	0.37	0.28
IMSAT (1)	0.52	0.41	0.35	1	0.50	0.49	0.33	0.29
IMSAT (0.5)	0.48	0.38	0.32	0.5	0.46	0.45	0.29	0.30
IMSAT (0.2)	0.44	0.34	0.28	0.2	0.42	0.41	0.27	0.31
IMSAT (0.1)	0.40	0.30	0.24	0.1	0.38	0.37	0.23	0.32
IMSAT (0.05)	0.36	0.26	0.20	0.05	0.34	0.33	0.19	0.33
IMSAT (0.01)	0.32	0.22	0.16	0.01	0.30	0.29	0.15	0.34
IMSAT (0.005)	0.28	0.18	0.12	0.005	0.26	0.25	0.11	0.35
IMSAT (0.001)	0.24	0.14	0.08	0.001	0.22	0.21	0.07	0.36
IMSAT (0.0005)	0.20	0.10	0.04	0.0005	0.18	0.17	0.03	0.37
IMSAT (0.0001)	0.16	0.06	0.00	0.0001	0.14	0.13	0.00	0.38
IMSAT (0.00005)	0.12	0.02	0.00	0.00005	0.10	0.09	0.00	0.39
IMSAT (0.00001)	0.08	0.00	0.00	0.00001	0.06	0.05	0.00	0.40
IMSAT (0.000005)	0.04	0.00	0.00	0.000005	0.02	0.01	0.00	0.41
IMSAT (0.000001)	0.00	0.00	0.00	0.000001	0.00	0.00	0.00	0.42
IMSAT (0.0000005)	0.00	0.00	0.00	0.0000005	0.00	0.00	0.00	0.43
IMSAT (0.0000001)	0.00	0.00	0.00	0.0000001	0.00	0.00	0.00	0.44
IMSAT (0.00000005)	0.00	0.00	0.00	0.00000005	0.00	0.00	0.00	0.45
IMSAT (0.00000001)	0.00	0.00	0.00	0.00000001	0.00	0.00	0.00	0.46
IMSAT (0.000000005)	0.00	0.00	0.00	0.000000005	0.00	0.00	0.00	0.47
IMSAT (0.000000001)	0.00	0.00	0.00	0.000000001	0.00	0.00	0.00	0.48
IMSAT (0.0000000005)	0.00	0.00	0.00	0.0000000005	0.00	0.00	0.00	0.49
IMSAT (0.0000000001)	0.00	0.00	0.00	0.0000000001	0.00	0.00	0.00	0.50
IMSAT (0.00000000005)	0.00	0.00	0.00	0.00000000005	0.00	0.00	0.00	0.51
IMSAT (0.00000000001)	0.00	0.00	0.00	0.00000000001	0.00	0.00	0.00	0.52
IMSAT (0.000000000005)	0.00	0.00	0.00	0.000000000005	0.00	0.00	0.00	0.53
IMSAT (0.000000000001)	0.00	0.00	0.00	0.000000000001	0.00	0.00	0.00	0.54
IMSAT (0.0000000000005)	0.00	0.00	0.00	0.0000000000005	0.00	0.00	0.00	0.55
IMSAT (0.0000000000001)	0.00	0.00	0.00	0.0000000000001	0.00	0.00	0.00	0.56
IMSAT (0.00000000000005)	0.00	0.00	0.00	0.00000000000005	0.00	0.00	0.00	0.57
IMSAT (0.00000000000001)	0.00	0.00	0.00	0.00000000000001	0.00	0.00	0.00	0.58
IMSAT (0.000000000000005)	0.00	0.00	0.00	0.000000000000005	0.00	0.00	0.00	0.59
IMSAT (0.000000000000001)	0.00	0.00	0.00	0.000000000000001	0.00	0.00	0.00	0.60
IMSAT (0.0000000000000005)	0.00	0.00	0.00	0.0000000000000005	0.00	0.00	0.00	0.61
IMSAT (0.0000000000000001)	0.00	0.00	0.00	0.0000000000000001	0.00	0.00	0.00	0.62
IMSAT (0.00000000000000005)	0.00	0.00	0.00	0.00000000000000005	0.00	0.00	0.00	0.63
IMSAT (0.00000000000000001)	0.00	0.00	0.00	0.00000000000000001	0.00	0.00	0.00	0.64
IMSAT (0.000000000000000005)	0.00	0.00	0.00	0.000000000000000005	0.00	0.00	0.00	0.65
IMSAT (0.000000000000000001)	0.00	0.00	0.00	0.000000000000000001	0.00	0.00	0.00	0.66
IMSAT (0.0000000000000000005)	0.00	0.00	0.00	0.0000000000000000005	0.00	0.00	0.00	0.67
IMSAT (0.0000000000000000001)	0.00	0.00	0.00	0.0000000000000000001	0.00	0.00	0.00	0.68
IMSAT (0.00000000000000000005)	0.00	0.00	0.00	0.00000000000000000005	0.00	0.00	0.00	0.69
IMSAT (0.00000000000000000001)	0.00	0.00	0.00	0.00000000000000000001	0.00	0.00	0.00	0.70
IMSAT (0.000000000000000000005)	0.00	0.00	0.00	0.000000000000000000005	0.00	0.00	0.00	0.71
IMSAT (0.000000000000000000001)	0.00	0.00	0.00	0.000000000000000000001	0.00	0.00	0.00	0.72
IMSAT (0.0000000000000000000005)	0.00	0.00	0.00	0.0000000000000000000005	0.00	0.00	0.00	0.73
IMSAT (0.0000000000000000000001)	0.00	0.00	0.00	0.0000000000000000000001	0.00	0.00	0.00	0.74
IMSAT (0.00000000000000000000005)	0.00	0.00	0.00	0.00000000000000000000005	0.00	0.00	0.00	0.75
IMSAT (0.00000000000000000000001)	0.00	0.00	0.00	0.00000000000000000000001	0.00	0.00	0.00	0.76
IMSAT (0.000000000000000000000005)	0.00	0.00	0.00	0.000000000000000000000005	0.00	0.00	0.00	0.77
IMSAT (0.000000000000000000000001)	0.00	0.00	0.00	0.000000000000000000000001	0.00	0.00	0.00	0.78
IMSAT (0.0000000000000000000000005)	0.00	0.00	0.00	0.0000000000000000000000005	0.00	0.00	0.00	0.79
IMSAT (0.0000000000000000000000001)	0.00	0.00	0.00	0.0000000000000000000000001	0.00	0.00	0.00	0.80
IMSAT (0.00000000000000000000000005)	0.00	0.00	0.00	0.00000000000000000000000005	0.00	0.00	0.00	0.81
IMSAT (0.00000000000000000000000001)	0.00	0.00	0.00	0.00000000000000000000000001	0.00	0.00	0.00	0.82
IMSAT (0.000000000000000000000000005)	0.00	0.00	0.00	0.000000000000000000000000005	0.00	0.00	0.00	0.83
IMSAT (0.000000000000000000000000001)	0.00	0.00	0.00	0.000000000000000000000000001	0.00	0.00	0.00	0.84
IMSAT (0.0000000000000000000000000005)	0.00	0.00	0.00	0.0000000000000000000000000005	0.00	0.00	0.00	0.85
IMSAT (0.0000000000000000000000000001)	0.00	0.00	0.00	0.0000000000000000000000000001	0.00	0.00	0.00	0.86
IMSAT (0.00000000000000000000000000005)	0.00	0.00	0.00	0.00000000000000000000000000005	0.00	0.00	0.00	0.87
IMSAT (0.00000000000000000000000000001)	0.00	0.00	0.00	0.00000000000000000000000000001	0.00	0.00	0.00	0.88
IMSAT (0.000000000000000000000000000005)	0.00	0.00	0.00	0.000000000000000000000000000005	0.00	0.00	0.00	0.89
IMSAT (0.000000000000000000000000000001)	0.00	0.00	0.00	0.000000000000000000000000000001	0.00	0.00	0.00	0.90
IMSAT (0.0000000000000000000000000000005)	0.00	0.00	0.00	0.0000000000000000000000000000005	0.00	0.00	0.00	0.91
IMSAT (0.0000000000000000000000000000001)	0.00	0.00	0.00	0.0000000000000000000000000000001	0.00	0.00	0.00	0.92
IMSAT (0.00000000000000000000000000000005)	0.00	0.00	0.00	0.00000000000000000000000000000005	0.00	0.00	0.00	0.93
IMSAT (0.00000000000000000000000000000001)	0.00	0.00	0.00	0.00000000000000000000000000000001	0.00	0.00	0.00	0.94
IMSAT (0.000000000000000000000000000000005)	0.00	0.						