# PROBABILISTIC BACK-ENDS FOR ONLINE SPEAKER RECOGNITION AND CLUSTERING

*Alexey Sholokhov*[1*]*, Nikita Kuzmin*[2,3*]*, Kong Aik Lee*[3]*, Eng Siong Chng*[2]

[1]Federal Research Center "Computer Science and Control"
of the Russian Academy of Sciences, Moscow, Russia
[2]Nanyang Technological University, Singapore
[3]Institute for Infocomm Research, A⋆STAR, Singapore

asholokhov@frccsc.ru, s220028@e.ntu.edu.sg, lee_kong_aik@i2r.a-star.edu.sg, aseschng@ntu.edu.sg

## ABSTRACT

This paper focuses on multi-enrollment speaker recognition which naturally occurs in the task of online speaker clustering, and studies the properties of different scoring back-ends in this scenario. First, we show that popular cosine scoring suffers from poor score calibration with a varying number of enrollment utterances. Second, we propose a simple replacement for cosine scoring based on an extremely constrained version of probabilistic linear discriminant analysis (PLDA). The proposed model improves over the cosine scoring for multi-enrollment recognition while keeping the same performance in the case of one-to-one comparisons. Finally, we consider an online speaker clustering task where each step naturally involves multi-enrollment recognition. We propose an online clustering algorithm allowing us to take benefits from the PLDA model such as the ability to handle uncertainty and better score calibration. Our experiments demonstrate the effectiveness of the proposed algorithm.

***Index Terms***— speaker verification, online speaker clustering

## 1. INTRODUCTION

In this paper, we consider a general scenario that we call *online speaker recognition*, where speech segments arrive sequentially, and the speaker recognition system has to identify previously encountered speakers and detect new speakers. At each time, there is a history of previously processed segments and the current segment to be classified.

One application scenario is *household speaker recognition* [1,2]. A household is a small set of family members whose speech data is processed by a shared device such as a smart speaker (*e.g.* Amazon Alexa). First, the device collects speech data from the users to create their profiles (speaker models). Then, at each interaction with a person, the device identifies the user and, optionally, updates (enriches) the corresponding speaker model. The device continuously collects the data of the users to improve its performance by estimating more accurate speaker representations. Also, the recorded speech utterances may belong to unregistered speakers (*e.g.* guests) leading to an open-set identification task. Another related task is low-latency speaker spotting [3], where a previously registered target speaker has to be detected in an audio stream.

Another example is *online speaker diarization* or *clustering* [4–9]. In this case, short speech segments from an audio stream have to be classified with low latency (*e.g.* 1-2 seconds). Unlike household speaker recognition, where all unregistered speakers are not of interest, in the speaker clustering task, there are no speakers registered

---

*Equal contribution.

beforehand, and a new speaker model has to be created for each previously unseen speaker. In the following, we focus on the online speaker clustering task since it is more general, and online speaker recognition can be seen as a special case.

What these scenarios have in common is that speech segments are received *sequentially* in nature and have to be classified on arrival. Specifically, an *open-set identification* problem has to be solved for each new segment. That is, the current segment has to be assigned to either one of the known speakers or a new (unknown) speaker. As a result, the number of segments per speaker continuously increases over time. This requires some way to aggregate information from multiple segments to form a memory-efficient speaker representation. This is usually referred to as *multi-enrollment* (or multi-session) speaker recognition [10–13], that is, when a speaker is represented by multiple speech segments. Moreover, different speakers may be represented by *different* numbers of segments. As shown in [11], this can be a major complicating factor for speaker recognition, since it causes inconsistency in scores from different speaker models. To our best knowledge, this issue has not been studied for modern large-margin speaker embeddings.

Inspired by [11,14], this work focuses on the issues arising from multi-enrollment scoring since it is a core element of online speaker recognition and clustering. We show that popular cosine scoring could have undesirable properties when used for multi-enrollment verification. Then we show that a highly constrained version of PLDA can be a suitable alternative while having better performance and comparable computational complexity. Specifically, we propose a PLDA model with spherical between- and within-covariance matrices as a replacement for cosine scoring back-end. While being *equivalent* to cosine scoring in a special case, this model can naturally handle varying degrees of uncertainty specific to the multi-enrollment scenario.

Further, we propose a probabilistic back-end for online speaker recognition and clustering. It is based on the spherical PLDA model and therefore has several appealing properties compared to cosine scoring. It employs an incremental (online) variant of variational Bayesian inference and provides probabilistic soft decisions for each input observation, based on the history of preceding observations.

Our contributions are summarized as follows:

- We compare scoring back-ends for multi-enrollment verification for modern large-margin embeddings.

- We propose a simple alternative to cosine scoring suitable for multi-enrollment verification.

- We propose a probabilistic back-end for online speaker recognition and clustering.

## 2. BACKGROUND

### 2.1. PLDA

**General formulation.** In this study we focus on a variant of PLDA known as the *two-covariance model* [15]. Let $\mathbf{x}_{i,j} \in \mathbb{R}^d$ denote the $j$th speaker embedding of speaker $i$. Also, let $\mathbf{y}_i$ be the latent speaker identity of speaker $i$. Then, the model is specified by two Gaussian distributions:

$$p(\mathbf{y}_i) = \mathcal{N}(\mathbf{y}_i|\boldsymbol{\mu}, \mathbf{B}), \quad p(\mathbf{x}_{i,j}|\mathbf{y}_i) = \mathcal{N}(\mathbf{x}_{i,j}|\mathbf{y}_i, \mathbf{W}). \quad (1)$$

Here, $\boldsymbol{\mu}$ is a global mean, and $\mathbf{B}, \mathbf{W} \in \mathbb{R}^{d \times d}$ are the between- and within -speaker covariance matrices, respectively.

Being a linear Gaussian model, PLDA allows making inferences about speaker identities in closed form. Given a set of observations (embeddings), one can compare different hypotheses about the partition of this set by computing the corresponding hypothesis likelihoods. This is often referred to as *by-the-book* scoring in the literature [10, 16].

**PLDA with spherical covariances.** Despite being a gold standard for previously popular i-vectors [17], one could recently observe a gradual shift towards replacing PLDA with a simpler parameter-less cosine scoring back-end [18]. As discussed in [14], the high intra-speaker compactness of the large-margin embedding makes the conventional full-rank PLDA model superfluous. It was also observed in [14] that discarding off-diagonal elements in the within-speaker covariance matrix can bring considerable performance gain. Here, we analyze a much more constrained version of the PLDA model, to our knowledge, firstly proposed in [1, 19]. Specifically, we consider PLDA with *spherical covariances*, $\mathbf{B} = \sigma_{\mathrm{B}}^2 \mathbf{I}$, $\mathbf{W} = \sigma_{\mathrm{W}}^2 \mathbf{I}$, where $\sigma_{\mathrm{B}}^2$ and $\sigma_{\mathrm{W}}^2$ are between- and within-speaker variances and $\mathbf{I}$ denotes an identity matrix. In the following text, we will refer to this model as the *spherical PLDA*.

**Relationship with cosine scoring.** As was shown in [18], for length-normalized and centered embeddings, the verification likelihood ratio of the spherical PLDA can be written as a scaled and shifted cosine similarity measure. Since an affine transformation of scores is order-preserving, the two scoring rules are equivalent. This brings up a question about the usefulness of spherical PLDA. As we discuss further, spherical PLDA has several advantages over cosine scoring. For instance, we show that the PLDA by-the-book scoring outperforms different cosine based heuristic scoring methods in multi-enrollment verification.

**Relationship with PSDA.** Another closely related scoring back-end is the so-called probabilistic spherical discriminant analysis (PSDA) recently proposed in [20]. It can be viewed as PLDA model with Gaussian distributions replaced by von Mises-Fisher (VMF) distributions that are defined on the $d - 1$ dimensional unit hypersphere $\mathbb{S}^{d-1}$ [21]:

$$p(\mathbf{y}_i) = \mathcal{V}(\mathbf{y}_i|\boldsymbol{\mu}, b), \quad p(\mathbf{x}_{i,j}|\mathbf{y}_i) = \mathcal{V}(\mathbf{x}_{i,j}|\mathbf{y}_i, w). \quad (2)$$

Here, $\mathcal{V}(\mathbf{y}|\boldsymbol{\mu}, \kappa)$ denotes the density of the VMF distribution with mean direction vector $\boldsymbol{\mu} \in \mathbb{S}^{d-1}$ and scalar concentration $\kappa \geq 0$ parameter. Similar to spherical PLDA, this model is parameterized by the mean direction vector $\boldsymbol{\mu}$ and two scalars: between-speaker, $b$, and within-speaker, $w$, concentrations.

The relation to spherical PLDA follows from the fact that restricting any isotropic Gaussian density to the unit hypersphere gives a VMF density, up to normalization. However, the two models are *not* equivalent, though their behavior is very similar as we show in the experiments.

We use both spherical PLDA and PSDA as a basis for a proposed online speaker clustering algorithm described in Section 3.

### 2.2. Multi-enrollment verification

When available, multiple enrollment utterances may represent various acoustic environments, or channels, that could be useful to better disentangle speaker identity from other irrelevant factors.

The study in [10] analyzes different methods of aggregating information from multiple speech segments for the PLDA scoring. Among them were embedding averaging, score averaging, and by-the-book scoring. Their experiments with i-vectors revealed that embedding averaging systematically outperforms other methods, including by-the-book scoring.

However, these observations were made for previously popular i-vector embeddings and have not been yet confirmed for modern large-margin embeddings. In fact, in our experiments, we observe that by-the-book scoring with spherical PLDA or PSDA outperforms embedding averaging.

## 3. ONLINE PROBABILISTIC SPEAKER CLUSTERING

In this section, we describe the proposed back-end model for online speaker recognition and clustering. The difference between offline (batch) and online settings is that in the former case all the data to be processed is available at once, while in the latter case pieces of data are observed sequentially, in some order.

### 3.1. Online clustering

The general pattern behind many online clustering algorithms is solving a series of successive open-set identification tasks [3,4,6,22]. The basic idea is to compare each new observation to the existing clusters, and either alter the closest cluster or create a new cluster. The generic Algorithm 1 demonstrates this for a single time step $t$.

---
**Algorithm 1** Online clustering (time step $t$)

---
$s_i \leftarrow \mathrm{score}(\mathbf{x}_t, \mathbf{X}_i)$      ▷ Compare $\mathbf{x}_t$ to the existing clusters
$k \leftarrow \arg\max_i s_i$      ▷ Find the most similar cluster
**if** $s_k \geq \tau$ **then** ▷ If the maximal score $s_k$ is above the threshold $\tau$
     $\mathbf{X}_k \leftarrow \{\mathbf{X}_k, \mathbf{x}_t\}$      ▷ Add $\mathbf{x}_t$ to the $k$-th cluster
**else**
     $\mathbf{X}_{K+1} \leftarrow \{\mathbf{x}_t\}$      ▷ Create a new cluster
     $K \leftarrow K + 1$      ▷ Increment the total number of clusters
**end if**

---

First, the observation $\mathbf{x}_t$ is compared to all existing clusters, each represented by a set of observations, $\mathbf{X}_i$. If similarity to the closest cluster is above the threshold, $\tau$, then $\mathbf{x}_t$ is assigned to this cluster. Otherwise, a new cluster is formed.

In this algorithm, clusters are represented by subsets of observations sharing the same label. Therefore, computing similarity to a cluster involves many-to-one comparison, also referred to as multi-enrollment verification in the context of speaker recognition. As discussed in [11], varying cluster sizes may result in miscalibrated scores leading to sub-optimal decisions with a fixed threshold $\tau$.

We aim at addressing this issue and propose an algorithm suitable for online clustering. Specifically, the underlying scoring model should be robust to varying cluster sizes naturally occurring in the online scenario. The proposed algorithm can be seen as a probabilistic extension of the Algorithm 1 constructed upon PLDA or PSDA models. As a result, it benefits from the advantages of PLDA (or PSDA) for multi-enrollment verification.

### 3.2. Model-based clustering

We start with a brief description of a generative model-based clustering [23, 24].

Model-based clustering builds upon a generative model that specifies how a set of data points $\mathbf{X} = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$ is generated from the hidden parameters of $K$ clusters $\mathbf{Y} = \{\mathbf{y}_1, ..., \mathbf{y}_K\}$, given the cluster assignments $\mathbf{Z} = \{\mathbf{z}_1, ..., \mathbf{z}_N\}$. A typical clustering model is given by the following joint distribution: $p(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = p(\mathbf{X}|\mathbf{Z}, \mathbf{Y})p(\mathbf{Y})p(\mathbf{Z})$.

The clustering problem requires finding the most likely partition of the data $\mathbf{Z}_* = \arg\max_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X})$. Our approach is based on the "mean-field" variational Bayesian approximation [25, 26] assuming that the approximate posterior factorizes as $p(\mathbf{Z}, \mathbf{Y}|\mathbf{X}) \approx q(\mathbf{Z})q(\mathbf{Y})$. This assumption leads the algorithm consisting of iterative updates of the factors $q(\mathbf{Z})$ and $q(\mathbf{Y})$.

However, such updates are designed for the conventional clustering setup, where all observations are available *at once*. We modify the standard inference algorithm to make it suitable for *online* clustering, where observations arrive sequentially. This algorithm can be seen as an online version of the VBx [27] with simplified prior on assignments $p(\mathbf{Z})$. It is also similar to the algorithm from [7], where the authors modified the offline variational inference to make it suitable for online processing.

### 3.3. The proposed algorithm

Let us denote the current observation at the time step $t$ as $\mathbf{x}_t$, and use the notation $\mathbf{X}_{1:t} = \{\mathbf{x}_1, ..., \mathbf{x}_t\}$ to denote causal observations.

The algorithm updates posterior distributions of latent identity variables $q(\mathbf{y}_k) \approx p(\mathbf{y}_k|\mathbf{X}_{1:t})$ after receiving a new observation $\mathbf{x}_t$. In general, several update iterations can be done. Our experiments reveal that even a single update can be sufficient for reasonable performance. In this case only posterior for the current data point $q(\mathbf{z}_t)$ needs to be computed, followed by updating each of $q(\mathbf{y}_k)$:

$$q(\mathbf{z}_t) \propto \exp \sum_{k=1}^{K} z_{t,k} \underbrace{\left[ \mathbb{E}_{q(\mathbf{y}_k)}[\log p(\mathbf{x}_t|\mathbf{y}_k)] + \log \pi_k \right]}_{\log \gamma_{t,k}}, \quad (3)$$

$$q(\mathbf{y}_k) \propto \exp \left[ \gamma_{t,k} \log p(\mathbf{x}_t|\mathbf{y}_k) + \log q(\mathbf{y}_k|\mathbf{X}_{1:t-1}) \right]. \quad (4)$$

Here, $\gamma_{t,k}$ is the $k$-th component of the vector of posterior probabilities $q(\mathbf{z}_t)$ over the cluster assignments and $\pi_k$ are the corresponding prior probabilities.

This algorithm continuously updates speaker models defined by $q(\mathbf{y}_k)$. Also, one can obtain speaker labels at each time step $t$ by finding $\arg\max_k \gamma_{t,k}$. For instance, if $\gamma_{t,k} = 0$, then the posterior $q(\mathbf{y}_k)$ stays unchanged. In general, if the soft-assignments $\boldsymbol{\gamma}_t$ were converted into hard decisions, then updating $q(\mathbf{y}_k)$ would be nothing more than the sequential application of the Bayes formula. Also, the algorithm would become very similar to a sequence of multi-enrollment recognition tasks, where predictions are obtained via by-the-book scoring.

These update equations can be used to construct different online recognition and clustering algorithms depending on a particular choice of the underlying generative model defined by $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y})$. In this study, we use two models: spherical PLDA and PSDA. Table 1 demonstrates the update equations for both models.

| PLDA: $q(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{m}_t, \mathbf{S}_t)$ | PSDA: $q(\mathbf{y}) = \mathcal{V}(\mathbf{y}|\mathbf{m}_t, r_t)$ |
|---|---|
| $\boldsymbol{\Lambda}_t = \gamma_t \mathbf{W}^{-1} + \boldsymbol{\Lambda}_{t-1}$ $\boldsymbol{\eta}_t = \gamma_t \mathbf{W}^{-1}\mathbf{x}_t + \boldsymbol{\eta}_{t-1}$ $\mathbf{S}_t = \boldsymbol{\Lambda}_t^{-1}, \mathbf{S}_0 = \mathbf{B}$ $\mathbf{m}_t = \mathbf{S}_t\boldsymbol{\eta}_t, \mathbf{m}_0 = \boldsymbol{\mu}$ | $\boldsymbol{\eta}_t = w\gamma_t\mathbf{x}_t + \boldsymbol{\eta}_{t-1}$ $r_t = \|\boldsymbol{\eta}_t\|, r_0 = b$ $\mathbf{m}_t = \boldsymbol{\eta}_t / r_t, \mathbf{m}_0 = \boldsymbol{\mu}$ |

**Table 1**: Update equations for the full-rank PLDA (1) and PSDA (2) at the time step $t$. The speaker index is omitted for clarity.

To detect new speakers we introduce an extra class corresponding to an unknown speaker. For this class, the posterior for the speaker identity variable is equal to the prior.

Algorithm 2 outlines the time step $t$ of the proposed algorithms.

---
**Algorithm 2** Proposed algorithm (time step $t$)

---
| | |
|---|---|
| $\boldsymbol{\gamma}_t \equiv q(\mathbf{z}_t) \leftarrow$ Eq. (3) | ▷ Cluster membership probabilities |
| $q(\mathbf{y}_k) \leftarrow$ Table 1 | ▷ Update clusters |
| $k \leftarrow \arg\max_i \gamma_{t,i}$ | ▷ Find the most probable cluster |
| **if** $k = K + 1$ **then** | ▷ New class is detected |
| $\quad K \leftarrow K + 1$ | ▷ Increment the total number of clusters |
| **end if** | |

---

The advantage of the proposed algorithm over Algorithm 1 is that it uses soft decisions for updating clusters. This makes the algorithm more robust to classification errors.

As a baseline for our experiments, we use Algorithm 1 with cosine similarity scoring.

## 4. EXPERIMENTS

In this section, we analyze the performance of several back-end scoring models in the multi-enrollment scenario. First, we report results for a rarely investigated speaker verification scenario, *i.e.*, where the number of enrollment and test segments *varies* within an evaluation protocol. Next, we apply the proposed Algorithm 2 for the online speaker diarization task. To support reproducible research, we make the code and evaluation protocols publicly available.

We used open-source speaker embedding extractors in order to make our experiments reproducible. We decided to stick to the following systems: SpeechBrain [28], BUT model [27], and CLOVA [29]. Due to space limitations, we report results only for SpeechBrain, while other results can be found at the project repository[1].

### 4.1. Multi-enrollment verification

In this section, we compare different scoring back-ends in multi-enrollment speaker verification scenario. Specifically, we investigate calibration properties of the verification scores in the case where the number of enrollment and test segment varies within an evaluation protocol.

**Experimental setup.** We created several custom evaluation protocols from the VoxCeleb1 test set [30]. Specifically, we generated four trial lists with configurations $(1, 1)$, $(3, 1)$, $(10, 1)$, and $(3, 3)$, where the notation (#enrollments, #tests) represents the number of enrollment or test segments in a single trial. In addition, we combined all the trial lists to get the *pooled* protocol. The idea behind it is to reveal the robustness of scoring back-ends to the number of enrollment segments. To exclude the effect of utterance duration, the recordings were cropped to 2 seconds before extracting embeddings.

We compared several different scoring variants: cosine similarity with embedding averaging (CSEA) or score averaging (CSSA), PSDA [20], and three versions of PLDA with spherical, diagonal, and full covariance matrices. For PLDA and PSDA by-the-book scoring was used. The VoxCeleb1 dev set [30] was used for training the back-ends. We used two performance metrics: the equal error rate (EER) and the minimum normalized detection cost function (minDCF) with $P_{\text{target}} = 0.01$ [31].

**Results.** Figure 1 demonstrates the distribution of verification scores for different numbers of enrollment segments. One can see
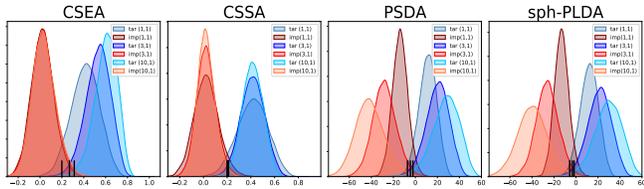
**Fig. 1**: Distributions of target and impostor scores for different numbers of enrollment segments: 1, 3, and 10. Short black vertical lines represent EER thresholds.

considerable distribution shifts for the target scores computed with CSEA. To be precise, a large variation of EER thresholds clearly makes each one sub-optimal for the other protocols. In contrast, the EER thresholds seem to be more stable for other scoring back-ends, even despite large differences in distribution means and variances for PLDA and PSDA. These observations are supported by objective metrics presented in Table 2. Despite low EERs for each pro-

| Back-end | Evaluation protocol | | | | |
|---|---|---|---|---|---|
| | (1, 1) | (3, 1) | (10, 1) | (3, 3) | pooled |
| CSEA | 4.98 | 1.65 | 0.83 | 0.17 | 2.85 / 0.206 |
| CSSA | 4.98 | 1.79 | 1.02 | 0.37 | 2.05 / 0.228 |
| PSDA | 4.85 | 1.55 | 0.78 | 0.13 | 2.08 / 0.172 |
| sph-PLDA | 4.98 | 1.59 | 0.78 | 0.14 | 1.99 / 0.170 |
| diag-PLDA | 4.95 | 1.62 | 0.78 | 0.13 | 1.98 / 0.169 |
| full-PLDA | 4.74 | 1.79 | 1.08 | 0.20 | 2.06 / 0.201 |

**Table 2**: Comparison of the speaker verification performance for different scoring back-ends in terms of EER, %. The last column shows minDCF as well. SpeechBrain embeddings were used.

tocol individually, the performance of CSEA degrades significantly on the pooled protocol. In contrast, CSSA does not suffer from this problem, however, it has higher error rates on the other protocols. Finally, PLDA and PSDA perform the best, overall, handling well all the cases. They also have very similar metrics and distributions of scores. These results are also in line with findings in [1] where spherical PLDA outperformed cosine similarity in the household speaker recognition task. Note that CSEA, CSSA, and sph-PLDA have exactly the same metrics in the $(1, 1)$ protocol because sph-PLDA is equivalent to cosine scoring. Another observation is that models with more parameters, diag- and full-PLDA, have comparable performance to sph-PLDA. This motivates choosing sph-PLDA as a simpler and faster alternative.

It should be noted that, unlike this study, PLDA model studied in [11] was not robust to the number of enrollment utterances. This probably can be explained by the nature of i-vector distribution which is different from the distribution of large-margin embeddings.

### 4.2. Online speaker diarization

In this section, we describe experiments on online speaker diarization. We used the same PLDA and PSDA models as for the previous experiments.

**Experimental setup.** We used two popular datasets of multi-speaker recordings: AMI [32], and VoxConverse [33]. Again, due to space limitations, we report only the results for the first one, while similar observations were made for the VoxConverse.

We used the development/evaluation split for the AMI corpus from [27][2]. The development set was used for tuning the hyper-

parameters of the back-end models, pretrained on the VoxCeleb data.

For AMI, the evaluation was performed on Mix-Headset channel. We extracted embeddings from segments of length 2.0 sec with 1.0 sec overlap within the boundaries obtained by the ground-truth annotation. These embeddings were sequentially processed by several online clustering algorithms, producing the output annotation. We did not use any special heuristics for handling segments with overlapped speakers, thus one speaker was assigned to each segment.

We compared three versions of Algorithm 1: with CSEA, CSSA, and PLDA scoring. Also, we evaluated two versions of the proposed Algorithm 2, with sph-PLDA and PSDA models. All of the algorithms have at least one hyper-parameter (*e.g.* decision threshold) that was tuned on the development split.

**Results.** For the evaluation metrics, we use the diarization error rate (DER) [34] and Jaccard error rate (JER) [35]. The forgiveness collar was set to 0.25, and overlapped speech regions were excluded from evaluation for DER, however, JER is calculated with no forgiveness collar and includes overlapped speech [35]. Table 3 provides the evaluation results. Unlike the previous experiment

| Clustering back-end | DER, % | JER, % |
|---|---|---|
| Algorithm 1 w/ CSEA | 3.63 | 25.20 |
| Algorithm 1 w/ CSSA | 3.67 | 26.33 |
| Algorithm 1 w/ sph-PLDA | 6.58 | 27.49 |
| Algorithm 2 w/ PSDA | 3.34 | 24.47 |
| Algorithm 2 w/ sph-PLDA | 3.32 | 25.21 |

**Table 3**: Online speaker diarization with SpeechBrain embeddings.

on speaker verification, we found that sequential PLDA scoring performs worse than cosine. As was discussed in [36] and [37], this probably can be explained by an inadequate assumption of statistical independence of the enrollment segments, which affects the score calibration. According to the theoretical model, enrollment segments are independent draws from the within-speaker distribution, while in diarization it is clearly not the case because of a shared acoustic environment and recording channel. However, unlike i-vector embeddings considered in [10,36], this effect seems to be less evident for large-margin embeddings. Apparently, PLDA suffers from this effect only in diarization, while yielding adequate score calibration when embeddings are less statistically dependent, as in our previous experiment.

At the same time, the proposed clustering algorithm which uses the same PLDA model delivers lower error rates than Algorithm 1. In the future, we plan to further investigate the properties of this algorithm in other applications such as household speaker recognition, where speech utterances are also processed sequentially.

## 5. CONCLUSION

This paper studies the properties of popular scoring back-ends suitable for large-margin speaker embeddings, with a particular focus on multi-enrollment speaker verification. Our experiments with the state-of-the-art embeddings revealed shortcomings of cosine scoring in the multi-enrollment scenario. To address this, we advocate for using the spherical PLDA that has several attractive properties: absence of numerical instabilities specific to PSDA due to Bessel functions; better performance, comparable computational complexity, and equivalence to cosine scoring in a special case. Also, we introduced a simple online clustering algorithm that uses the advantages of PLDA and PSDA for the multi-enrollment scenario. Empirical evaluation of the online speaker diarization showed superior performance of the proposed algorithm.

# 6. REFERENCES

[1] Sholokhov, A., Liu, X., Sahidullah, M., and Kinnunen, T., "Baselines and protocols for household speaker recognition," in *Odyssey*, 2022, pp. 185–192.

[2] Tan, Z., Yang, Y., Han, E., and Stolcke, A., "Improving speaker identification for shared devices by adapting embeddings to speaker subsets," in *ASRU*, 2021, pp. 1124–1131.

[3] Patino, J. et al., "Low-latency speaker spotting with online diarization and detection," in *Odyssey*, 2018, pp. 140–146.

[4] Liu, D. and Kubala, F., "Online speaker clustering," in *ICASSP*, 2004, pp. 333–336.

[5] Aloni-Lavi, R., Opher, I., and Lapidot, I., "Incremental on-line clustering of speakers' short segments," in *Odyssey*, 2018, pp. 120–127.

[6] Wisniewski, G., Bredin, H., Gelly, G., and Barras, C., "Combining speaker turn embedding and incremental structure prediction for low-latency speaker diarization," in *INTER-SPEECH*, 2017, pp. 3582–3586.

[7] Koshinaka, T., Nagatomo, K., and Shinoda, K., "Online speaker clustering using incremental learning of an ergodic hidden Markov model," *IEICE Trans. Inf. Syst.*, vol. 95-D, no. 10, pp. 2469–2478, 2012.

[8] Zhu, W. and Pelecanos, J. W., "Online speaker diarization using adapted i-vector transforms," in *ICASSP*, 2016, pp. 5045–5049.

[9] Soldi, G., Beaugeant, C., and Evans, N. W. D., "Adaptive and online speaker diarization for meeting data," in *EUSIPCO*, 2015, pp. 2112–2116.

[10] Rajan, P., Afanasyev, A., Hautamaki, V., and Kinnunen, T., "From single to multiple enrollment i-vectors: Practical PLDA scoring variants for speaker verification," *Digital Signal Processing*, vol. 31, pp. 93–101, 2014.

[11] Lee, K., Larcher, A., You, C. H., Ma, B., and Li, H., "Multi-session PLDA scoring of i-vector for partially open-set speaker detection," in *INTERSPEECH*, 2013, pp. 3651–3655.

[12] Soni, M. H. and Panda, A., "LDA-based speaker verification in multi-enrollment scenario using expected vector approach," in *ISCSLP*, 2021, pp. 1–5.

[13] Zeng, C., Wang, X., Cooper, E., Miao, X., and Yamagishi, J., "Attention back-end for automatic speaker verification with multiple enrollment utterances," in *ICASSP*, 2022, pp. 6717–6721.

[14] Wang, Q., Lee, K. A., and Liu, T., "Scoring of large-margin embeddings for speaker verification: Cosine or PLDA?," *arXiv:2204.03965*, 2022.

[15] Brümmer, N. and de Villiers, E., "The speaker partitioning problem," in *Odyssey*, Brno, Czech Republic, 2010, pp. 194–201.

[16] López, J. A. V., Díez, M., Varona, A., and Lleida, E., "Handling recordings acquired simultaneously over multiple channels with PLDA," in *INTERSPEECH*, 2013, pp. 2509–2513.

[17] Dehak, N., Kenny, P., Dehak, R., Dumouchel, P., and Ouellet, P., "Front-end factor analysis for speaker verification," *IEEE Trans. Audio, Speech & Language Processing*, vol. 19, no. 4, pp. 788–798, 2011.

[18] Peng, Z., He, X., Ding, K., Lee, T., and Wan, G., "Unifying cosine and PLDA back-ends for speaker verification," *arXiv:2204.10523*, 2022.

[19] Kuzmin, N., Fedorov, I., and Sholokhov, A., "Magnitude-aware probabilistic speaker embeddings," in *Odyssey*, 2022, pp. 1–8.

[20] Brümmer, N., Swart, A., Mosner, L., Silnova, A., Plchot, O., Stafylakis, T., and Burget, L., "Probabilistic spherical discriminant analysis: An alternative to PLDA for length-normalized embeddings," in *INTERSPEECH*, 2022, pp. 1446–1450.

[21] Mardia, K. V. and Jupp, P. E., *Directional statistics*, vol. 2, Wiley Online Library, 2000.

[22] Mansfield, P. A., Wang, Q., Downey, C., Wan, L., and Lopez-Moreno, I., "Links: A high-dimensional online clustering method," *arXiv:1801.10123*, 2018.

[23] Valente, F. and Wellekens, C., "Variational Bayesian speaker clustering," in *Odyssey*, 2004, pp. 207–214.

[24] Díez, M., Burget, L., Landini, F., and Cernocký, J., "Analysis of speaker diarization based on Bayesian HMM with eigen-voice priors," *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 28, pp. 355–368, 2020.

[25] Corduneanu, A. and Bishop, C., "Variational Bayesian model selection for mixture distributions," in *AISTATS*, 2001, pp. 27–34.

[26] Bishop, C. M., *Pattern Recognition and Machine Learning*, Springer Science+Business Media, LLC, New York, 2006.

[27] Landini, F., Profant, J., Díez, M., and Burget, L., "Bayesian HMM clustering of x-vector sequences (VBx) in speaker diarization: theory, implementation and analysis on standard tasks," *arXiv:2012.14952*, 2020.

[28] Ravanelli, M. et al., "SpeechBrain: A general-purpose speech toolkit," *arXiv:2106.04624*, 2021.

[29] Heo, H. S., Lee, B.-J., Huh, J., and Chung, J. S., "Clova baseline system for the VoxCeleb speaker recognition challenge 2020," *arXiv:2009.14153*, 2020.

[30] Nagrani, A., Chung, J. S., and Zisserman, A., "VoxCeleb: A large-scale speaker identification dataset," in *INTERSPEECH*, 2017, pp. 2616–2620.

[31] Przybocki, M. and Martin, A., "NIST speaker recognition evaluation chronicles," in *Odyssey*, 2004, pp. 15–22.

[32] McCowan, I. et al., "The AMI meeting corpus," *Measuring Behavior 2005, 5th International Conference on Methods and Techniques in Behavioral Research*, pp. 137–140, 2005.

[33] Chung, J. S., Huh, J., Nagrani, A., Afouras, T., and Zisserman, A., "Spot the conversation: Speaker diarisation in the wild," in *INTERSPEECH*, 2020, pp. 299–303.

[34] Galibert, O., "Methodologies for the evaluation of speaker diarization and automatic speech recognition in the presence of overlapping speech," in *INTERSPEECH*, 2013, pp. 1131–1134.

[35] Ryant, N., Church, K., Cieri, C., Cristià, A., Du, J., Ganapathy, S., and Liberman, M., "The second DIHARD diarization challenge: Dataset, task, and baselines," in *INTERSPEECH*, 2019, pp. 978–982.

[36] McCree, A., Sell, G., and Garcia-Romero, D., "Extended variability modeling and unsupervised adaptation for PLDA speaker recognition," in *INTERSPEECH*, 2017, pp. 1552–1556.

[37] Stafylakis, T., Kenny, P., Gupta, V., and Dumouchel, P., "Compensation for inter-frame correlations in speaker diarization and recognition," in *ICASSP*, 2013, pp. 7731–7735.