

Bayesian HMM clustering of x-vector sequences (VBx) in speaker diarization: theory, implementation and analysis on standard tasks

Federico Landini^{a,1,*}, Ján Profant^{b,1}, Mireia Diez^a, Lukáš Burget^a

^a*Brno University of Technology, Faculty of Information Technology, Speech@FIT, Czechia*

^b*Phonexia, Czechia*

Abstract

The recently proposed VBx diarization method uses a Bayesian hidden Markov model to find speaker clusters in a sequence of x-vectors. In this work we perform an extensive comparison of performance of the VBx diarization with other approaches in the literature and we show that VBx achieves superior performance on three of the most popular datasets for evaluating diarization: CALLHOME, AMI and DIHARDII datasets. Further, we present for the first time the derivation and update formulae for the VBx model, focusing on the efficiency and simplicity of this model as compared to the previous and more complex BHMM model working on frame-by-frame standard Cepstral features. Together with this publication, we release the recipe for training the x-vector extractors used in our experiments on both wide and narrowband data, and the VBx recipes that attain state-of-the-art performance on all three datasets. Besides, we point out the lack of a standardized evaluation protocol for AMI dataset and we propose a new protocol for both Beamformed and Mix-Headset audios based on the official AMI partitions and transcriptions.

Keywords: Speaker Diarization, Variational Bayes, HMM, x-vector, DIHARD,

*The work was supported by Czech National Science Foundation (GACR) project “NEUREM3” No. 19-26934X, and European Union’s Horizon 2020 project No. 833635 ROX-ANNE.

*Corresponding author

Email addresses: landini@fit.vutbr.cz (Federico Landini), jan.profant@phonexia.com (Ján Profant), mireia@fit.vutbr.cz (Mireia Diez), burget@fit.vutbr.cz (Lukáš Burget)

¹Equal contribution

1. Introduction

In recent years, speaker diarization works are proliferating. This is due to two main factors: first, new datasets and challenges are providing new benchmarks that bring the interest of the community and foster *healthy* competition. The most relevant examples are the DIHARD series addressing diarization in a wide variety of challenging domains [1, 2], CHIME6 [3] providing a very demanding multi-channel benchmark for diarization, or the recent VoxConverse [4] exploring diarization on several kinds of YouTube videos. Second, the success of the new end-to-end paradigm for speaker recognition is starting to be adopted for diarization tasks. Unlike standard diarization approaches which normally deal with diarization tasks using oracle voice activity detection (VAD), end-to-end diarization systems deal also with the VAD task. The end-to-end approaches have the advantage of being able to cope with overlapped speech [5, 6, 7]. Even if these methods are still limited, e.g. they are restricted to scenarios with a fixed number of speakers, and mainly tested on artificially created short recordings [5, 8], and mostly do not achieve state-of-the-art results [9], this research line is very promising and indeed prolific.

However, due to the several difficulties that end-to-end approaches still have to overcome for diarization tasks, in recent speaker diarization evaluations the best performing systems are based on a more conventional approach, the clustering of x-vectors [10, 11, 12].

In this paper, we show how our current VBx diarization approach, which clusters x-vectors using a Bayesian hidden Markov model (BHMM) [13], combined with a ResNet101 x-vector extractor [14] achieves superior results on CALLHOME [15], AMI [16] and DIHARDII [2] datasets.

Besides establishing new baselines for these representative datasets, we perform a thorough analysis comparing our results to the best numbers found in the literature. In the case of AMI dataset, this proved to be a challenging task.

Most works published on AMI data choose their own partition, references and audio setup for evaluation, making the comparison between works very complicated. We identified works that were presenting superior performance and reproducible setups and we evaluated our system in their respective setups. In this paper, we further provide our own evaluation protocol, which comprises lists for train/dev/eval partition, references and audios. Our setup is based on the official partition of AMI corpus. We believe that this setup could serve as a new standard facilitating a fair comparison of diarization systems on AMI corpus.

The VBx diarization approach has been presented before [17], but the paper did not provide any derivation of update formulae, as it was introduced merely as a special case and simplification of its *big-brother* BHMM with eigen-voice priors [18]. The more complex BHMM model from [18] operates directly on frame-by-frame standard Cepstral features. It is based on a Bayesian HMM where states correspond to speakers and transitions correspond to speaker turns. To robustly model speaker distributions it uses an i-vector like model [19, 20]: the distribution of each speaker is represented by a Gaussian mixture model (GMM). Such GMM is constrained to live in a low-dimensional eigen-voice subspace and each speaker can be therefore robustly represented by a fixed-length i-vector like latent variable.

The VBx used in this paper is based on a similar BHMM model and a similar idea for modeling speaker distributions. However, it is used for directly clustering x-vectors, which allows to use a much simpler probabilistic linear discriminant analysis (PLDA) based model for modeling speaker distributions. The model is essentially the same as the one in [18] but using only a single Gaussian component to model speaker distributions. In fact, when VBx was introduced in [17], it was suggested that the same model and the same inference from the BHMM model working on a frame-basis [18] could be reused just by replacing the GMMs that modeled the speaker-specific distributions by single Gaussians. However, naively re-implementing the algorithm used in [18] is not effective as the design changes made to obtain the VBx model lead to significant

simplification of the inference formulae and derivations. Therefore, in this paper, we present the same derivations and update formulae as in [18], but now for the simpler VBx. This derivation should be much easier to follow for readers interested in this specific model. It also allows us to elaborate on how to make this simplified model much more computationally efficient.

All our code is made publicly available: the recipe for training the x-vector extractor (same architecture for both 8 kHz and 16 kHz), the trained extractors and the pipeline for applying BHMM diarization using agglomerative hierarchical clustering (AHC) as initialization <https://github.com/BUTSpeechFIT/VBx>.

2. VBx Diarization model

This section introduces the VBx model which is used in all the experiments in this paper. The derivation of the inference formulae is also provided. While this derivation is essentially the same as the one that can be found in [18], it only addresses the simple model used in this paper. In fact, the following text is the same as Section II from [18], rewritten and simplified to address only the model considered in this paper. We intentionally reused the text, structure and symbols from the mentioned paper for the following reasons: we want to make the treatment of the model here as self-contained as possible and we would like to facilitate the comparison of the simplified model with the more complex full model proposed in [18].

2.1. Model overview

As described in the previous section, we expect a sequence of x-vectors extracted from consecutive short segments of speech as input to our diarization method, which aims to cluster these x-vectors according to their speaker identity.

Our diarization model assumes that the input sequence of x-vectors is generated by an HMM with speaker-specific state distributions. To facilitate the

discrimination between speakers, the speaker- (or HMM state-) specific distributions are derived from a PLDA [21] model pre-trained on a large number of speaker-labeled x-vectors. More details on how the speaker-specific distributions are derived from the PLDA are given in section 2.3. For now, it is sufficient to note that the speaker distributions will be represented only by a latent vector \mathbf{y}_s of the same dimensionality as the x-vectors.

We use an ergodic HMM with one-to-one correspondence between the HMM states and the speakers, where transitions from any state to any state are possible. Note that our model does not consider any overlapped speech as each speech frame is assumed to be generated from an HMM state corresponding to only one of the S speakers. The transition probabilities can be used to discourage too frequent transitions between speakers in order to reflect speaker turn durations of a natural conversation. More details on setting and learning the transition probabilities can be found in section 2.2.

Let $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ be the sequence of observed x-vectors and $\mathbf{Z} = \{z_1, z_2, \dots, z_T\}$ the corresponding sequence of discrete latent variables defining the hard alignment of x-vectors to HMM states. In our notation, $z_t = s$ indicates that the speaker (HMM state) s is responsible for generating observation \mathbf{x}_t .

To address the speaker diarization (SD) task using our model, the speaker distributions (i.e. the vectors \mathbf{y}_s) and the latent variables z_t are jointly estimated given an input sequence \mathbf{X} . The solution to the SD task is then given by the most likely sequence \mathbf{Z} , which encodes the alignment of speech frames to speakers.

2.2. HMM topology

The HMM topology and transition probabilities model the speaker turn durations. The HMM model is ergodic (transitions between all states are possible). Figure 1 shows an example of the HMM topology for only $S = 3$ speakers. The transition probabilities are set as follows: we transition back to the same speaker/state with probability P_{loop} . This probability is one of the tunable parameters in the model. The remaining probability $(1 - P_{loop})$ is the probability of changing speaker, which corresponds to the transition to the non-emitting node

in Figure 1. From the non-emitting node, we immediately transition to one of the speaker states with probability π_s .² Therefore, the probability of leaving a speaker and entering another speaker s is $(1 - P_{loop})\pi_s$. To summarize, the probability of transitioning from state s' to state s is

$$p(s|s') = (1 - P_{loop})\pi_s + \delta(s = s')P_{loop}, \quad (1)$$

where $\delta(s = s')$ equals 1 if $s = s'$ and is 0 otherwise.

The non-emitting node in Figure 1 is also the initial state of the model. Therefore, the probabilities π_s also control the selection of the initial HMM state (i.e. the state generating the first observation). These probabilities π_s are inferred (jointly with the variables \mathbf{y}_s and z_t) from the input conversation. Thanks to the automatic relevance determination principle [22] stemming from our Bayesian model, zero probabilities will be learned for the π_s corresponding to redundant speakers, which effectively drops such speakers from the HMM model. Typically, we initialize the HMM with a larger number of speakers (see section 3.2) and we make use of this behavior to drop the redundant speakers (i.e. to estimate the number of speakers in the conversation).

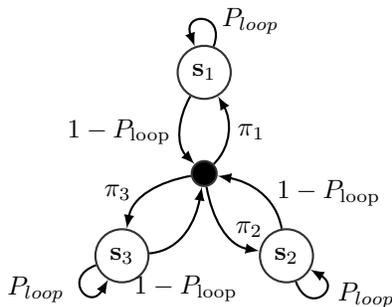


Figure 1: HMM model for 3 speakers with 1 state per speaker, with a dummy non-emitting (initial) state.

²For convenience, we allow to re-enter the same speaker as it leads to simpler update formulae.

2.3. Speaker-specific distributions

The speaker (HMM state) specific distributions are derived from a PLDA which is a standard model used for comparing x-vectors in speaker verification [21]. Here, only a simplified variant of PLDA is considered, which is often referred to as *two-covariance model* [23]. This model assumes that the distribution of x-vectors specific to speaker s is Gaussian $\mathcal{N}(\hat{\mathbf{x}}_t; \hat{\mathbf{m}}_s, \Sigma_w)$, where Σ_w is the within-speaker covariance matrix shared by all speaker models, and $\hat{\mathbf{m}}_s$ is the speaker-specific mean. Speaker means are further assumed to be Gaussian distributed $\mathcal{N}(\hat{\mathbf{m}}_s; \mathbf{m}, \Sigma_b)$, where \mathbf{m} is the global mean and Σ_b is the between-speaker covariance matrix. In general, Σ_w and Σ_b can be full covariance matrices. However, to further simplify and speed up the inference in our model, we assume that the x-vectors are linearly transformed into a space where Σ_b is diagonal and Σ_w is identity. This can be achieved as follows:

Let $\hat{\mathbf{X}}$ be the matrix of original (untransformed) x-vectors that the parameters of the original PLDA model \mathbf{m} , Σ_w and Σ_b were estimated on. The x-vectors that are used as input for the diarization algorithm are obtained as

$$\mathbf{X} = (\hat{\mathbf{X}} - \mathbf{m})\mathbf{E} \quad (2)$$

where \mathbf{E} is the transformation matrix which transforms the x-vectors into the desired space. This matrix can be obtained by solving the standard generalized eigen-value problem

$$\Sigma_b \mathbf{E} = \Sigma_w \mathbf{E} \Phi \quad (3)$$

where \mathbf{E} is the matrix of eigen-vectors and Φ is the diagonal matrix of eigen-values, which is also the between-speaker covariance matrix in the transformed space. Note that the eigen-vectors \mathbf{E} are, in fact, bases of linear discriminant analysis (LDA) estimated directly from the PLDA model parameters. Therefore, if we construct Φ only using R largest eigen-values and assemble \mathbf{E} only using the corresponding eigen-vectors, (2) further performs LDA dimensionality reduction of x-vectors to R -dimensional space. We use R as one of the hyper-parameters of the VBx method. In equation (2), we have also subtracted the global mean from the original x-vectors to have the new set of x-vectors zero-centered.

In summary, the PLDA model compatible with the new set of x-vectors \mathbf{X} suggests that speaker-specific means are distributed as

$$p(\mathbf{m}_s) = \mathcal{N}(\mathbf{m}_s; \mathbf{0}, \Phi). \quad (4)$$

For convenience and for the compatibility with the notation introduced in [18], we further re-parametrize the speaker mean as

$$\mathbf{m}_s = \mathbf{V}\mathbf{y}_s, \quad (5)$$

where diagonal matrix $\mathbf{V} = \Phi^{\frac{1}{2}}$ and \mathbf{y}_s is a standard normal distributed random variable

$$p(\mathbf{y}_s) = \mathcal{N}(\mathbf{y}_s; \mathbf{0}, \mathbf{I}). \quad (6)$$

The speaker-specific distribution of x-vectors is

$$p(\mathbf{x}_t|\mathbf{y}_s) = \mathcal{N}(\mathbf{x}_t; \mathbf{V}\mathbf{y}_s, \mathbf{I}), \quad (7)$$

where \mathbf{I} is identity matrix.

In our diarization model, we use (7) to model the speaker (HMM state) distributions. This distribution is fully defined only in terms of the speaker vector \mathbf{y}_s (and the pre-trained matrix \mathbf{V} shared by all the speakers). The speaker vector \mathbf{y}_s is treated as a latent variable with standard normal prior (6), which is why the BHMM model is called *Bayesian*³. This way, the full PLDA model is incorporated into the BHMM in order to properly model between- and across-speaker variability. Therefore, the model is capable of discriminating between speakers just like PLDA model when used for speaker verification.

2.4. Bayesian HMM

To summarize, our complete model for SD is a Bayesian HMM, which is defined in terms of the state-specific distributions (or so-called output probabilities)

$$p(\mathbf{x}_t|z_t = s) = p(\mathbf{x}_t|s) = p(\mathbf{x}_t|\mathbf{y}_s) \quad (8)$$

³However, unlike other “Fully Bayesian” HMM implementations [24, 25], we do not impose any prior on the transition probabilities.

described in section 2.3 and the transition probabilities

$$p(z_t = s | z_{t-1} = s') = p(s | s') \quad (9)$$

described in section 2.2. By abuse of notation, $p(z_1 | z_0)$ will correspond to the initial state probability $p(z_1 = s) = \pi_s$ in the following formulae.

The complete model can be also defined in terms of the joint probability of the observed and latent random variables (and their factorization) as

$$\begin{aligned} p(\mathbf{X}, \mathbf{Z}, \mathbf{Y}) &= p(\mathbf{X} | \mathbf{Z}, \mathbf{Y}) p(\mathbf{Z}) p(\mathbf{Y}) \\ &= \prod_t p(\mathbf{x}_t | z_t) \prod_t p(z_t | z_{t-1}) \prod_s p(\mathbf{y}_s), \end{aligned} \quad (10)$$

where $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_S\}$ is the set of all the speaker-specific latent variables.

The model assumes that each x-vector sequence corresponding to an input conversation is obtained using the following generative process:

```

for  $s = 1..S$  do
   $\mathbf{y}_s \sim \mathcal{N}(0, \mathbf{I})$ 
for  $t = 1..T$  do
   $z_t \sim p(z_t | z_{t-1})$ 
   $\mathbf{x}_t \sim p(\mathbf{x}_t | z_t)$ 

```

2.5. Diarization inference

The diarization problem consists in finding the assignment of frames to speakers, which is represented by the latent sequence \mathbf{Z} . In order to find the most likely sequence \mathbf{Z} , we need to infer the posterior distribution $p(\mathbf{Z} | \mathbf{X}) = \int p(\mathbf{Z}, \mathbf{Y} | \mathbf{X}) d\mathbf{Y}$. Unfortunately, the evaluation of this integral is intractable, and therefore, we will approximate it using variational Bayes (VB) inference [22], where the distribution $p(\mathbf{Z}, \mathbf{Y} | \mathbf{X})$ is approximated by $q(\mathbf{Z}, \mathbf{Y})$. We use the mean-field approximation [22, 20] assuming that the approximate posterior distribution factorizes as

$$q(\mathbf{Z}, \mathbf{Y}) = q(\mathbf{Z})q(\mathbf{Y}). \quad (11)$$

The particular form of the approximate distributions $q(\mathbf{Z})$ and $q(\mathbf{Y})$ directly follows from the optimization described below.

We search for such $q(\mathbf{Z}, \mathbf{Y})$ that minimizes the Kullback-Leibler divergence $D_{KL}(q(\mathbf{Z}, \mathbf{Y})\|p(\mathbf{Z}, \mathbf{Y}|\mathbf{X}))$, which is equivalent to maximizing the standard VB objective – the evidence lower bound objective (ELBO) [22]

$$\mathcal{L}(q(\mathbf{X}, \mathbf{Y})) = E_{q(\mathbf{Y}, \mathbf{Z})} \left\{ \ln \left(\frac{p(\mathbf{X}, \mathbf{Y}, \mathbf{Z})}{q(\mathbf{Y}, \mathbf{Z})} \right) \right\}. \quad (12)$$

Using the factorization (11), the ELBO can be split into three terms

$$\begin{aligned} \hat{\mathcal{L}}(q(\mathbf{X}, \mathbf{Y})) &= F_A E_{q(\mathbf{Y}, \mathbf{Z})} [\ln p(\mathbf{X}|\mathbf{Y}, \mathbf{Z})] \\ &+ F_B E_{q(\mathbf{Y})} \left[\ln \frac{p(\mathbf{Y})}{q(\mathbf{Y})} \right] + E_{q(\mathbf{Z})} \left[\ln \frac{p(\mathbf{Z})}{q(\mathbf{Z})} \right], \end{aligned} \quad (13)$$

where the first term is the expected log-likelihood of the observed x-vector sequence \mathbf{X} and the second and third terms are Kullback-Leibler divergences $D_{KL}(q(\mathbf{Y})\|p(\mathbf{Y}))$ and $D_{KL}(q(\mathbf{Z})\|p(\mathbf{Z}))$ regularizing the approximate posterior distributions $q(\mathbf{Y})$ and $q(\mathbf{Z})$ towards the priors $p(\mathbf{Y})$ and $p(\mathbf{Z})$. In (13), we modified the ELBO by scaling the first two terms by constant factors F_A and F_B .⁴ The theoretically correct values for these factors leading to the original ELBO (13) are $F_A = F_B = 1$. However, choosing different values gives us finer control over the inference, which can be used to improve diarization performance. For further details on the specific effects these scaling factors have in the inference, we refer the reader to [18].

As described above, we search for the approximate posterior $q(\mathbf{Z}, \mathbf{Y})$ that maximizes the ELBO (13). In the case of the mean-field factorization (11), we proceed iteratively by finding the $q(\mathbf{Y})$ that maximizes the ELBO given fixed $q(\mathbf{Z})$ and vice versa. This section provides all the formulae necessary for implementing these updates or for understanding our open-source Python implementation⁵. In this section, we do not give any details on deriving the

⁴Note that similar scaling factor for the third term would be redundant as only the relative scale of the three factors is relevant for the optimization.

⁵<http://speech.fit.vutbr.cz/software/vb-diarization-eigenvoice-and-hmm-priors>

update formulae. For the readers interested in the derivations, we prepared a technical report [26].

2.5.1. Updating $q(\mathbf{Y})$

Given a fixed $q(\mathbf{Z})$, the distribution over \mathbf{Y} that maximizes the ELBO is

$$q^*(\mathbf{Y}) = \prod_s q^*(\mathbf{y}_s), \quad (14)$$

where the speaker-specific approximate posteriors

$$q^*(\mathbf{y}_s) = \mathcal{N}(\mathbf{y}_s | \boldsymbol{\alpha}_s, \mathbf{L}_s^{-1}) \quad (15)$$

are Gaussians with the mean vector and precision matrix

$$\boldsymbol{\alpha}_s = \frac{F_A}{F_B} \mathbf{L}_s^{-1} \sum_t \gamma_{ts} \boldsymbol{\rho}_t \quad (16)$$

$$\mathbf{L}_s = \mathbf{I} + \frac{F_A}{F_B} \left(\sum_t \gamma_{ts} \right) \boldsymbol{\Phi}. \quad (17)$$

where

$$\boldsymbol{\rho}_t = \mathbf{V}^T \mathbf{x}_t \quad (18)$$

In this update formula, $\gamma_{ts} = q(z_t = s)$ is the marginal approximate posterior derived from the current estimate of the distribution $q(\mathbf{Z})$ (see below), which can be interpreted as the responsibility of speaker s for generating observation \mathbf{x}_t (i.e. defines a soft alignment of \mathbf{x} -vectors to speakers).

If we compare these update formulae to the corresponding ones from the more complex BHMM model in [18], it can be seen that in [18], $\boldsymbol{\Phi}_t$ is a frame-dependent full-matrix computationally expensive to calculate. In contrast, $\boldsymbol{\Phi}$ here does not depend on time frame t and, as pointed out in section 2.3, it is a diagonal matrix. Therefore also matrix \mathbf{L}_s is diagonal, and its inversions and application in (16) become trivial.

2.5.2. Updating $q(\mathbf{Z})$

We never need to infer the complete distribution over all the possible alignments of observations to speaker $q(\mathbf{Z})$. When updating $q(\mathbf{Y})$ using (16) and (17),

we only need the marginals $\gamma_{ts} = q(z_t = s)$. Therefore, when updating $q(\mathbf{Z})$, we can directly search for the responsibilities γ_{ts} that correspond to the distribution $q^*(\mathbf{Z})$ maximizing the ELBO given a fixed $q(\mathbf{Y})$. Similar to the standard HMM training, such responsibilities can be calculated efficiently using a forward-backward algorithm as

$$\gamma_{ts} = \frac{A(t, s)B(t, s)}{\bar{p}(\mathbf{X})} \quad (19)$$

where the forward probability

$$A(t, s) = \bar{p}(\mathbf{x}_t|s) \sum_{s'} A(t-1, s') p(s|s') \quad (20)$$

is recursively evaluated by progressing forward in time for $t=1..T$ starting with $A(0, s) = \pi_s$. Similarly,

$$B(t, s) = \sum_{s'} B(t+1, s') \bar{p}(\mathbf{x}_{t+1}|s') p(s'|s) \quad (21)$$

is the backward probability evaluated using backward recursion for times $t = T..1$ starting with $B(T, s) = 1$.

$$\bar{p}(\mathbf{X}) = \sum_s A(T, s) \quad (22)$$

is the total forward probability and

$$\begin{aligned} \log \bar{p}(\mathbf{x}_t|s) &= F_A \left[\boldsymbol{\alpha}_s^T \boldsymbol{\rho}_t - \frac{1}{2} \text{tr} (\boldsymbol{\Phi} [\mathbf{L}_s^{-1} + \boldsymbol{\alpha}_s \boldsymbol{\alpha}_s^T]) - \frac{D}{2} \ln 2\pi - \frac{1}{2} \mathbf{x}_t^T \mathbf{x}_t \right] \\ &= F_A \left[\boldsymbol{\alpha}_s^T \boldsymbol{\rho}_t - \frac{1}{2} \boldsymbol{\phi}^T [\boldsymbol{\lambda}_s + \boldsymbol{\alpha}_s^2] - \frac{D}{2} \ln 2\pi - \frac{1}{2} \mathbf{x}_t^T \mathbf{x}_t \right] \end{aligned} \quad (23)$$

is the expected log likelihood of observation \mathbf{x}_t given a speaker s taking into account its uncertainty $q(\mathbf{y}_s)$. The second line of (23) corresponds to an efficient evaluation of this term, where vector $\boldsymbol{\phi}$ is the diagonal of the diagonal matrix $\boldsymbol{\Phi}$, vector $\boldsymbol{\lambda}$ is the diagonal of the diagonal matrix \mathbf{L}_s^{-1} and the square in $\boldsymbol{\alpha}_s^2$ is element-wise. Note also that the terms $-\frac{D}{2} \ln 2\pi - \frac{1}{2} \mathbf{x}_t^T \mathbf{x}_t$ in (23) are not only constant over VB iterations, but also constant for different speakers s . As a consequence, contribution of these terms cancels in (19) and therefore does not have to be calculated at all.

2.5.3. Updating π_s

Finally, the speaker priors π_s are updated as maximum likelihood type II estimates [22]: Given fixed $q(\mathbf{Y})$ and $q(\mathbf{Z})$, we search for the values of π_s that maximize the ELBO (13), which gives the following update formula

$$\pi_s \propto \gamma_{1s} + \frac{(1-P_{loop})\pi_s}{\bar{p}(\mathbf{X})} \sum_{t=2}^T \sum_{s'} A(t-1, s') p(\mathbf{x}_t | s) B(t, s) \quad (24)$$

with the constraint $\sum_s \pi_s = 1$. As described in section 2.2, this update tends to drive the π_s corresponding to “redundant speakers” to zero values, which effectively drops them from the model and selects the right number of speakers in the input conversation.

2.5.4. Evaluating the ELBO

The convergence of the iterative VB inference can be monitored by evaluating the ELBO objective. For the Bayesian HMM, the ELBO can be efficiently evaluated (see page 95 of [25]) as

$$\hat{\mathcal{L}} = \ln \bar{p}(\mathbf{X}) + \sum_s \frac{F_B}{2} (R + \ln |\mathbf{L}_s^{-1}| - \text{tr}(\mathbf{L}_s^{-1}) - \boldsymbol{\alpha}_s^T \boldsymbol{\alpha}_s), \quad (25)$$

where R is the dimensionality of the \mathbf{x} -vectors. Note, that since \mathbf{L}_s is a diagonal matrix, $\ln |\mathbf{L}_s^{-1}|$ can be calculated just as the sum of the log of the elements in the diagonal. This way of evaluating the ELBO is very practical as the term $\bar{p}(\mathbf{X})$ from (22) is obtained as a byproduct of “updating $q(\mathbf{Z})$ ” using the forward-backward algorithm. On the other hand, (25) allows to evaluate the ELBO only right after the $q(\mathbf{Z})$ update. It does not allow to monitor the improvement in ELBO obtained from $q(\mathbf{Y})$ or π_s updates, which might be useful for debugging purposes. Therefore, we also provide the derivation formulae for the explicit evaluation of all three ELBO terms from (13) in [26].

The complete VB inference consisting of iterative updates of $q(\mathbf{Y})$, $q(\mathbf{Z})$ and parameters π_s is summarized in the following algorithm:

Initialize all γ_{ts} as described in section 3.2.

repeat

Update $q(\mathbf{y}_s)$ for $s=1..S$ using (15)

for $t = 1..T$ **do**

Calculate $A(t, s)$ for $s=1..S$ using (20)

for $t = T..1$ **do**

Calculate $B(t, s)$ for $s=1..S$ using (21)

Update γ_{ts} for $t=1..T, s=1..S$ using (19)

Update π_s for $s=1..S$ using (24)

Evaluate ELBO $\hat{\mathcal{L}}$ using (25)

until convergence of $\hat{\mathcal{L}}$

3. Experimental setup

3.1. x-vector extractor and PLDA

As described in the previous section, VBx diarization relies on a pre-trained x-vector extractor and a PLDA model. Since we report results on both 16 kHz recordings (DIHARD and AMI) and 8 kHz telephone recordings (CALLHOME), we train two x-vector extractors and the corresponding PLDA models, one for each condition. The complete PyTorch [27] recipe for x-vector extractor and PLDA training is available at <https://github.com/phonexiaresearch/VBx-training-recipe>.

3.1.1. x-vector extractor architecture

Both 8 kHz and 16 kHz x-vector extractors use the same deep neural network architecture based on ResNet101 [14, 28]. In both cases, the neural network inputs are 64 log Mel filter bank features extracted every 10 ms using 25 ms window. The two x-vector extractors differ only in the frequency ranges spanned by the Mel filters, which are 20-7700 Hz and 20-3700 Hz for the 16 kHz and 8 kHz systems, respectively. The x-vector extractor architecture is summarized in Table 1. The first 2D convolutional layer operates on the $64 \times T$ matrix

of log Mel filter bank features, where T is the number of frames in the input segments. For training, we use 4 s segments (i.e. $T = 400$). The following layers are standard ResNet blocks [14]. As in the original x-vector architecture [29], the statistical pooling layer is used to aggregate information over the whole speech segment (i.e. mean and standard deviation of activations is calculated over the time dimension). After the pooling layer, a linear transformation is used to reduce the dimensionality to obtain the 256-dimensional x-vectors.

Table 1: *The structure of the proposed ResNet101 architecture. The first dimension of the input shows the size of the filterbank and the second dimension indicates the number of frames.*

Layer	Structure	Stride	Output
Input	-	-	$64 \times T \times 1$
Conv2D-1	$3 \times 3, 32$	1	$64 \times T \times 32$
ResNetBlock-1	$\begin{bmatrix} 1 \times 1, 32 \\ 3 \times 3, 32 \\ 1 \times 1, 128 \end{bmatrix} \times 3$	1	$64 \times T \times 128$
ResNetBlock-2	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 4$	2	$32 \times T / 2 \times 256$
ResNetBlock-3	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 23$	2	$16 \times T / 4 \times 512$
ResNetBlock-4	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 3$	2	$8 \times T / 8 \times 1024$
Statistics Pooling	-	-	16×1024
Flatten	-	-	16384
Linear	-	-	256

The x-vector extractors are trained using stochastic gradient descent and additive angular margin loss [30] with speaker identities as class labels. We ramp-up the margin during the first two epochs (pass through the training data) and then train the neural network for another epoch with fixed margin $m = 0.2$.

3.1.2. x-vector extractor training data

The 16 kHz x-vector extractor is trained using data from VoxCeleb1 [31] (323 h of speech from 1211 speakers), VoxCeleb2 [32] (2290 h, 5994 speakers) and CN-CELEB [33] (264 h, 973 speakers). The energy-based VAD from Kaldi [34]

toolkit is used to remove silence frames. Speakers with less than 2 recordings are discarded. Further, we drop utterances with less than 4 seconds of speech. This way, about 4% of speech data is discarded. Data augmentation is performed the same way as in the SRE16 Kaldi recipe [35]. This way, we obtain four additional copies of the data with artificially added noise, music or reverberation. Training examples are randomly sampled from the training data. This way we extract about 89 million examples (original and augmented 4s segments), which cover more than 60% of the speech from the training corpora.

To train the 8 kHz x-vector extractor, the same data sets are used as in the 16 kHz case. Additionally, the following data sets were used: Mixer collection (NIST SRE 2004-2010, 3805 h, 4254 speakers), Switchboard (1170 h, 2591 speakers) and DeepMine [36] (688 h, 1858 speakers). Any wide-band data used were downsampled to 8 kHz and passed through a telephone codec. The same data selection and augmentation was used as for the 16 kHz case. Note that about 30% of DeepMine data were discarded as this dataset contains many utterances with less than 4 seconds of speech (mostly phrases for text-dependent speaker verification).

3.1.3. PLDA training

The 8 kHz and 16 kHz PLDA models are trained on the same data as the corresponding x-vector extractors. For this purpose, one x-vector is extracted from each individual recording (e.g. one cut from a YouTube video in the case of the VoxCeleb data). The length of such recordings can range from 4s to several minutes. Note that the PLDA trained on such x-vectors is later used in VBx to operate on x-vectors extracted from much shorter 1.5s segments. This mismatch, however, does not seem to negatively affect diarization performance.

3.2. Diarization pipeline

To perform the diarization, each input recording is first split into speech segments according to the oracle VAD and the segments shorter than 0.1 s are discarded. From these segments, x-vectors are extracted every 0.25 s from over-

lapping sub-segments of 1.5s (or less than 1.5s for the last sub-segments or shorter segments). The x-vectors are centered, whitened and length normalized [37] (which is also done for the PLDA training data).

As described in section 2.2, VBx diarization needs an initial assignment of x-vectors to speakers. For this purpose, the x-vectors are pre-clustered using AHC to obtain the initial speaker labels. The only input to the AHC is the matrix of cosine similarities between all pairs of x-vectors. The threshold used as the stopping criterion for AHC is tuned to under-cluster so that the following VBx has more freedom to search for the optimal results and converge to the right number of speaker models⁶. Nevertheless, the same threshold is used for all our results on all the datasets.

In the final step, the x-vectors are further clustered using the VBx model and the inference described in section 2. For this step, the x-vector dimensionality is further reduced to 128 dimensions (see parameter R in section 2.3). Note that unlike in our previous works [11, 13], we do not perform any adaptation of PLDA models to the target data as this paper aims to present a simple diarization system which performs well for different datasets. Nevertheless, we tune the VBx parameters: F_A , F_B , P_{loop} on the respective DIHARDII, AMI and CALLHOME development sets.

In order to demonstrate the effectiveness of the VBx method, we also report results for baseline systems where only AHC is used to cluster x-vectors. In this case, the stopping threshold is tuned to obtain the best performance on the respective development set.

3.3. Evaluation Datasets

3.3.1. CALLHOME

The 2000 NIST Speaker Recognition Evaluation (LDC2001S97⁷) dataset, usually referred as “CALLHOME”, [38] has been the standard dataset for

⁶Note that the inference in BHMM cannot converge to higher number of speakers than what is suggested by the AHC-based initialization.

⁷<https://catalog.ldc.upenn.edu/LDC2001S97>

diarization in the last decade [39, 40, 41]. In its full form⁸, it consists of 499⁹ recordings of conversational telephone speech in Arabic, English, German, Japanese, Mandarin and Spanish. The number of speakers per recording ranges between 2 and 7, although 87% of the files contain only 2 or 3 speakers. It amounts to around 15 hours of speech after VAD.

Since a development-evaluation split for CALLHOME is not available, we split the dataset into two halves as defined in the Kaldi recipe for CALLHOME¹⁰. We use this split to perform cross-validation to tune parameters i.e. F_A , F_B and P_{loop} .

3.3.2. AMI corpus

When trying to cover the most standard datasets for speaker diarization, we could not leave AMI out [16]. The AMI meeting corpus is a multi-modal data collection of 100 hours of meeting recordings. This corpus was recorded using both close-talking and far-field microphones. It consists of 171 meetings recorded at the University of Edinburgh (U.K.), Idiap (Switzerland), and the TNO Human Factors Research Institute (The Netherlands). The dataset comes with annotations for automatic speech recognition (ASR). AMI has been widely used by the community for diarization purposes. Still, somewhat surprisingly, authors do not use a standard evaluation protocol to report the results on this dataset.

The description of the different evaluation protocols reported in the literature as well as our proposed new protocol for evaluation on AMI can be found in section 4.

⁸Not only the English partition, nor the partition limited to 2 speaker audios sometimes used.

⁹One audio is commonly excluded because its references have formatting errors

¹⁰https://github.com/kaldi-asr/kaldi/blob/master/egs/callhome_diarization/v2/run.sh

3.3.3. DIHARD II

DIHARD II is one of the newest datasets designed for diarization. This dataset was created as an extension of the first DIHARD dataset for the second DIHARD challenge [2], the second of a yearly series of challenges designed to foster research on diarization in hard conditions. One of the main features of this dataset, is that it contains audios from several sources (YouTube, court rooms, meetings, etc.) covering a wide range of numbers of speakers per recording (1 to 10) and large variety of channels and audio conditions. The corpus consists of 192 development and 194 evaluation recordings, containing around 18 h and 17 h of speech, respectively.

3.4. Evaluation protocol

Diarization performance is evaluated in terms of diarization error rate (DER) as defined by NIST [42]:

$$DER = \frac{SER + FA + Miss}{Total_speech} \quad (26)$$

where:

- *SER* stands for speaker error, the amount of time that speech is attributed to incorrect speakers
- *FA* is false alarm, the amount of time that non-speech regions are incorrectly attributed to a speaker (or time when overlapped speech is found in single speaker speech regions)
- *Miss* stands for missed speech, the amount of time that speech is not attributed to any speaker
- *Total_speech* is the total amount of speech, accounting also for speaker overlaps.

Note that, as we use oracle VAD in all our experiments, *FA* error is zero, and the *Miss* (due only to non-handled overlapped speech regions) can be directly calculated as $DER - SER$ as we also report *SER*. We also evaluate the system

in terms of Jaccard error rate (JER), which has been established as secondary metric in the latest diarization challenges [2, 4]. JER is similar to DER, although it weighs every speaker equally, regardless of the amount of speech they produced. All experiments are evaluated using the `dscore` tool¹¹.

For a more thorough analysis of results, for CALLHOME and AMI we consider three setups for evaluation: First, a *forgiving* one in which a 0.25s collar is considered for DER estimation and no overlap is evaluated. This is the standard configuration used for these datasets and allows a comparison of results with previously published works. Second, we consider a similar evaluation using a 0.25s collar but accounting for overlapped speech, as in [4]. This *fair* setup covers a pragmatic scenario where all speech is evaluated, while being flexible on the speaker change points, as no realistic human annotation can achieve frame-precision. Finally, the *full* one, in which no collar is used and overlapped speech is evaluated, which is in line with the setup used in latest diarization challenges [2, 3]. We consider that looking into the numbers for the *fair* and *full* setups is truly relevant on these datasets: very low DER values are already being achieved on these sets with the *forgiving* setup, which suggest that future evaluations could shift into more challenging setups. Also, the hot-topic end-to-end approaches are likely to surpass the performance of current systems on overlapped speech regions, and suitable baselines need to be established.

For the more recent DIHARD II dataset, only the *fair* and *full* evaluation setups are considered.

Note that JER considers no collar and evaluates overlap regions by definition, so it is not affected by these configurations. To avoid confusion (and repetition), we only report this value on the *full* setup.

¹¹<https://github.com/nryant/dscore>

4. AMI evaluation protocols

4.1. Evaluation protocols found in previously published works

As pointed out before, when it comes to evaluation on AMI, there is not a defined evaluation protocol. Different authors evaluate their systems on different audio types (recordings from Mix-Headset microphones, microphone arrays, etc.). Besides, authors use different data partitions (train/dev/eval sets) and use different references. This makes it practically impossible to compare results between sites. We would like to highlight that during our search for baselines in the literature, we found that most works are unaware of this inconsistency of evaluation protocols for AMI, which frequently leads to unfair comparisons.

Based on our literature review we replicated some evaluation protocols from works presenting remarkable performance. We use these protocols to evaluate our approach and fairly compare it with the respective works.

There are two major publicly available recipes for diarization on AMI that are included in Pyannote¹² [43] and Kaldi¹³ [34] toolkits. These two recipes evaluate only on one audio type, which is the independent headset microphone mixed audio (Mix-Headset in AMI). Each recipe also derives their references in different ways from the official ASR transcriptions. Finally, they both use their own partitioning of the data, Pyannote uses the Full-corpus AMI partition and Kaldi claims to use the official Full-corpus-ASR AMI partition¹⁴. Nevertheless, Kaldi partition differs from the Full-corpus-ASR one, as it includes one meeting in the training set (IB4005) which causes speaker overlap between train and dev sets. This meeting is explicitly excluded in the Full-corpus-ASR partition.

All the works that we compare with use different combinations of partitions and references taken mostly from these two recipes. The first six columns in Table 3 summarize this mixture of protocols. Some works use Pyannote partition

¹²https://github.com/pyannote/pyannote-audio/tree/master/tutorials/pipelines/speaker_diarization

¹³<https://github.com/kaldi-asr/kaldi/tree/master/egs/ami/s5c>

¹⁴<http://groups.inf.ed.ac.uk/ami/corpus/datasets.shtml>

[43, 44], while others use the Kaldi one [45, 46, 47, 48], or even a modified version of the Kaldi one, excluding TNO meetings [49, 50]. As for references, AMI corpus comes only with manual and automatic transcriptions for ASR training and there are different ways of deriving diarization references out of them, as will be described later. In previously published works, we find again references from Pyannote [43, 44, 47, 48] and the ones from Kaldi [46], although some works use modified versions of these based on their own ASR forced alignment [45], or simply derive their own [49]. Note also the mix of evaluation setups (better described in section 3.4) with different collars and criteria for including or excluding overlapped speech. Column 5 in the table shows the amount of speech (in seconds) evaluated with each protocol for dev and eval sets. As it can be seen, evaluation can range from around 13000 to up to 52000 seconds of speech for the eval set depending on the evaluation setup used.

We strongly believe a standard evaluation protocol should be established on AMI. Next, we introduce what we believe should be established as this new standard.

4.2. New AMI evaluation protocol

The following evaluation protocol was built after discussions with researchers from different labs, authors of official Kaldi and Pyannote recipes (BUT, CLSP JHU, IRIT).

We propose to use the official Full-corpus-ASR partition¹⁴. This way, we make the scoring of diarization tasks consistent with the scoring of speaker-attributed ASR. As mentioned before, such partition is very similar to the Kaldi partition (in fact, dev and eval sets are the same) but it has no speaker overlap between sets, which makes it suitable for diarization tasks.

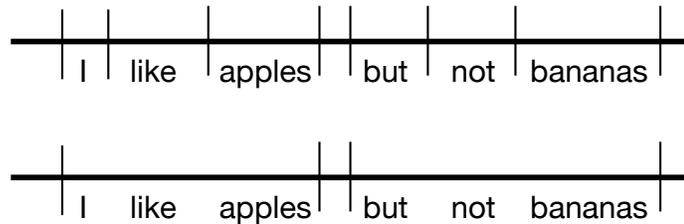
Regarding the references, authors of other works do not make clear how their references were created. Our diarization references are directly derived from the AMI manual annotations, version 1.6.2.¹⁵ These annotations are human tran-

¹⁵<http://groups.inf.ed.ac.uk/ami/download/>

criptions of all the meetings, containing words, vocal sounds and punctuation marks. To generate the references:

- All words are considered as speech and included in the references.
- Sounds of very different nature were annotated as vocal sounds. Some examples are Dutch speech, whistling, yawn, laughter, cough, clicking with tongue, raspberry noise, blowing nose, clapping, etc. Some of these are clearly speech (Dutch speech), while others are clearly noises (blowing nose, clapping, etc.). For sounds such as laughter or whistling it is simply unclear if they should be considered for diarization purposes, as this would depend on the particular application of the system. Besides, we found that several of these vocal sounds are labeled without time annotations, which makes it impossible to add them to the references. We therefore decided to take a well defined, consistent and conservative approach in which all vocal sounds are discarded. This way, only the words that could be recognized by an ASR system are considered in our references. This is also more consistent with the task of speaker-attributed ASR.
- Speaker turns respect precisely the annotations, but adjacent speech segments (words) of the same speaker are merged not to create false “break” points. Consider the following example of an ASR transcription with speech from a single speaker which, as the AMI one, is composed of short segments of speech for each word:

```
starttime="0.86" endtime="0.98" word=I
starttime="0.98" endtime="1.1" word=like
starttime="1.1" endtime="1.40" word=apples
starttime="1.45" endtime="1.55" word=but
starttime="1.55" endtime="1.62" word=not
starttime="1.62" endtime="2.0" word=bananas
```



If adjacent speech segments from the same speaker are not merged, it truly affects the diarization evaluation when collars are considered, as the collar is applied over all VAD borders and these “break points” are considered as one of these borders. These “break points” between adjacent speech segments are common in the references derived with the Kaldi recipe. With our processing, the above transcription results in one speech segment from 0.86 to 1.40 (merged “I like apples”) and another from 1.45 to 2.0 (merged “but not bananas”).

On the other hand, consecutive speech segments from the same speaker separated by pauses (silence) are not merged in any case. Using the above example again, we could think that the pause is too short and maybe it should be discarded and the two segments should be merged into one. But this kind of processing would require some heuristic to determine which is the required pause for considering separate speech segments. We prefer to follow a clean approach keeping the original pauses.

Anyway, in the case of using no collar (which we believe is the best choice

for evaluation) merging the adjacent speech segments has no effect at all. Still, we establish it in case anyone would prefer to use a collar.

In later experiments, we use the same partition and references to evaluate our system with two audio types: AMI Mix-Headset audios and the beamformed microphone array N1, where BeamformIt [51] is applied using the specific setup provided for AMI.

The partition, references and audios are shared in our repository¹⁶. This evaluation protocol will also be adopted in the latest Kaldi and pyannote recipes. Additionally, we also generated an extra version of the references including all (time-labeled) vocal sounds. As mentioned before, we consider these references not well defined, but we understand that some researchers might find them useful.

5. Results

5.1. CALLHOME

We first present in Table 2 the results of our diarization model on CALLHOME data. As a baseline, we also provide the result of a standalone AHC clustering of x-vectors, where its threshold is tuned for optimal performance. We provide these results to illustrate the gain achieved specifically by VBx diarization, which uses AHC as initialization (see section 3.2 for details). Our VBx system achieves 4.42% DER on the *forgiving* evaluation setup, outperforming all systems from previously published works. With the *fair* evaluation setup, considering also the overlapped speech (which our system does not handle), performance drops to 14.21% DER. Note, that the result achieved in [9] is obtained on a subset of the CALLHOME dataset, which makes it not directly comparable with our results. The script provided with their implementation makes random partitions of CALLHOME dataset into dev/eval so it is not possible to replicate this partition to make a fair comparison. Finally, with zero collar (*full*

¹⁶<https://github.com/BUTSpeechFIT/AMI-diarization-setup>

evaluation setup) increases the error to a total 21.77% DER. We believe that future research works should report results using this more challenging evaluation setup.

Table 2: *Diarization performance on CALLHOME. Results marked with * are obtained on a subset of the dataset and are therefore not comparable, see text for more details.*

Evaluation setup			System	SER	DER	JER
Name	Collar	Overlap				
<i>Forgiving</i>	0.25	No	Kaldi (Sell et al. [10])	6.48	–	–
			Zhang et al. [52]	7.60	–	–
			Lin et al. [53]	6.63	–	–
			Pal et al. [49]	6.76	–	–
			Aronowitz et al. [54]	5.10	–	–
			AHC	8.10	–	–
			VBx	4.42	–	–
<i>Fair</i>	0.25	Yes	Horiguchi et al. [9]	–	15.29*	–
			AHC	7.53	17.64	–
			VBx	4.10	14.21	–
<i>Full</i>	0	Yes	AHC	11.06	25.61	35.48
			VBx	7.22	21.77	34.02

5.2. AMI

We first report in Table 3 results for our system using the different evaluation protocols found in the literature for AMI corpus. As explained in section 4.1, the mix of protocols between sites called for running experiments with 3 different data partitions (train/dev/eval sets), 5 sets of references, 2 different types of audio and considering different evaluation setups to be able to compare with all works. We would like to highlight again that the different protocols differ largely in the amount of speech used for evaluation: from only 13309 seconds of speech considered in [49] up to 52317 seconds considered in [43]. When

Table 3: *Diarization performance of the proposed model compared to baselines from the literature. (*) denotes dev and eval pooled results*

Partition	References	Audio type	Evaluation setup		Scored speech dev/eval (s)	System	development		evaluation		
			Name	Collar Overlap			SER	DER	SER	DER	
<i>Pyannote</i>	<i>Pyannote</i>	<i>Mix-Headset</i>	<i>Forgiving</i>	0.25	No	Bredin et al. [43] VBx	–	–	4.6	–	
				–	–	2.14	–	–	2.17		
			<i>Full</i>	0	Yes	Bullock et al. [44] VBx	–	–	–	–	–
<i>Kaldi</i>	<i>Force Aligned</i>	<i>Beamformed mic-array</i>	<i>Forgiving</i>	0.25	No	Sun et al. [45] VBx	16.4	–	15.4	–	
				–	–	1.32	–	–	1.84		
	<i>Kaldi</i>	<i>Mix-Headset</i>	<i>Forgiving</i>	0.25	Yes	Sun et al. [45] VBx	–	19.4	–	17.8	
				–	–	1.26	4.96	1.92	4.67		
	<i>Pyannote</i>	<i>Mix-Headset</i>	<i>Mix-Headset</i>	<i>Forgiving</i>	0.25	No	Maciejewski et al. [46] VBx	–	–	–	(4.8*)
					–	–	2.14	–	–	3.02/(2.58*)	
<i>Pyannote</i>	<i>Pyannote</i>	<i>Mix-Headset</i>	<i>Full</i>	0	Yes	Raj et al. [47] Raj et al. [48] VBx	–	–	10.1	23.6	
				–	–	3.12	22.63	3.56	23.47		
<i>Kaldi no TNO</i>	<i>Work specific</i>	<i>Beamformed mic-array</i>	<i>Forgiving</i>	0.25	No	Raj et al. [48] VBx	7.7	–	5.2	–	
				–	–	4.08	–	–	3.80		
<i>Kaldi no TNO</i>	<i>Work specific</i>	<i>Beamformed mic-array</i>	<i>Forgiving</i>	0.25	No	Pal et al. [49] VBx	5.02	–	4.92	–	
				–	–	6.21	–	–	2.87		
<i>Kaldi no TNO</i>	<i>Work specific</i>	<i>Beamformed mic-array</i>	<i>Forgiving</i>	0.25	No	Pal et al. [49] VBx	4.27	–	4.58	–	
				–	–	4.27	–	–	4.58		

analyzing the results, it can be seen that our VBx system attains the best results on all evaluation protocols, with two exceptions. The first exception is [48], which uses Kaldi partition, Pyannote references and *Full* evaluation setup. This work presents a fusion of 3 different diarization systems dealing with overlapped speech, as compared to our single system with no overlap handling. The same system evaluated only on non overlap regions obtains 7.7% and 5.2% DER on dev and eval respectively, while VBx obtains only 4.08% on dev and 3.8% on eval. The second exception is the eval result from [49], which uses Kaldi “no TNO” partition. While this system has significantly better performance than our system (only) on the eval set (2.87 vs. 4.58 DER), when analyzing the results from the paper, it seems to us that this number is an outlier inconsistent with the other results presented in the work: all systems presented in the paper have consistently similar performance on dev and eval sets, this particular system reduces the error more than a 50% (6.21% DER on dev vs 2.87% on eval). Also, this would not be the system of choice if selected according to the performance on the dev set, for which the system was tuned (also shown in the table, with 5.02% DER on the dev set and 4.92% DER on eval). Out of all the partitions used, this is the one with the smallest eval set, which might result in noisier results.

In Table 3, we have demonstrated that VBx method has superior performance as compared to other published works on AMI dataset. However, to offer a fair comparison, we had to deal with too many protocols. This led us to the proposal of the new evaluation protocol, described in section 4.2 which is, already being adopted by other research labs. In Table 4, we report results obtained with the VBx system with the proposed evaluation protocol. This results can serve as a reference for future works on this corpus. Once again, results are also provided for the baseline standalone AHC when it is tuned for optimal performance (see section 3.2 for details). For the sake of completeness, we provide results of our system when evaluated on beamformed mic-array audios as well as on Mix-Headset audios. The system is evaluated on all evaluation setups presented in 3.4. Performance on Beamformed audios gets as low as 3.9%

DER on the eval set when using the forgiving evaluation. When considering overlapped speech, DER increases to 14.23%. Finally, without any collar, the system achieves 20.84% DER and 26.92% JER. For the Mix-Headset audios, results are consistently better, as expected, achieving 2.10%, 12.53% and 18.99% DER for forgiving, fair and full evaluation setups, respectively.

Table 4: *Diarization performance of the proposed model on AMI with the proposed AMI protocol.*

Audio type	Evaluation setup			System	development			evaluation		
	Name	Collar	Overlap		SER	DER	JER	SER	DER	JER
Beamformed	<i>Forgiving</i>	0.25	No	AHC	6.32		–	7.65		–
				VBx	2.80		–	3.90		–
	<i>Fair</i>	0.25	Yes	AHC	6.43	14.68	–	8.82	18.36	–
				VBx	2.57	10.81	–	4.69	14.23	–
	<i>Full</i>	0	Yes	AHC	8.68	22.14	25.29	10.93	25.48	29.85
				VBx	4.20	17.66	22.26	6.28	20.84	26.92
Mix-Headset	<i>Forgiving</i>	0.25	No	AHC	3.90		–	3.96		–
				VBx	1.56		–	2.10		–
	<i>Fair</i>	0.25	Yes	AHC	4.06	12.31	–	5.05	14.60	–
				VBx	1.43	9.68	–	2.98	12.53	–
	<i>Full</i>	0	Yes	AHC	6.16	19.61	23.90	6.87	21.43	25.50
				VBx	2.88	16.33	20.57	4.43	18.99	24.57

5.3. DIHARD II

Diarization results obtained on the DIHARDII dataset are presented in Table 5. Our VBx system obtains 18.19% DER on the development set and 18.55% DER on the evaluation set. When analyzing the table, it can be seen that, in fact, these are not the best numbers ever published on DIHARDII eval as the system does not overcome the results from [56], nor the ones we attained on the challenge [11] which are, as far as we know, still the best in the literature. However, note that the best results on DIHARDII [11] are obtained with the same VBx method as described in this paper, but including additional adaptation to the DIHARD data and additional steps (overlapped speech handling

Table 5: *Diarization performance of the proposed model on DIHARD II. (*) denotes our own previous results. Results in brackets were obtained when not using the dev set for training nor adaptation purposes, thus they are comparable to our results.*

Evaluation setup			System			development			evaluation			
Name	Collar	Overlap	System	SER	DER	JER	SER	DER	JER	SER	DER	JER
			Landini et al.* [11]	-	17.90 (18.34)	-	-	18.21 (19.14)	-	-	-	-
			Lin et.al. [55]	-	21.36	-	-	18.84	-	-	-	-
<i>Full</i>	0	Yes	Lin et.al. [56]	-	18.76	-	-	18.44 (19.46)	-	-	-	-
			AHC	10.89	21.68	42.28	13.89	23.59	43.93			
			VBx	7.41	18.19	42.53	8.85	18.55	43.91			
			AHC	8.22	14.91	-	10.94	16.67	-			
<i>Fair</i>	0.25	Yes	VBx	5.53	12.23	-	6.55	12.29	-			

and resegmentation). Similarly, in [56] the diarization system is adapted to the DIHARD dev data. As mentioned in section 3.2, in this paper we are aiming to use a generic approach to all datasets, without using the development sets for training or adaptation purposes. If we compare our system to the results of [11, 56] when not performing adaptation on the dev set, which are shown in brackets, our current VBx outperforms both systems.

When evaluating the system with the *fair* setup, DER improves around a 30% on both sets, resulting in 11.75% and 12.41% on dev and eval, respectively.

6. Conclusion

This paper presents a diarization system based on a Bayesian HMM model for clustering x-vectors, also known as VBx. Our VBx diarization achieves state-of-the-art results on CALLHOME, AMI and DIHARDII without performing specific model adaptation to any of the datasets.

Most of the papers, which we have compared our system with, present clustering of new sets of embeddings/x-vectors. All these approaches are complementary with VBx: the VBx clustering can be combined with most of these new embeddings, which shows the further potential of the method.

In the case of AMI dataset, we have presented, for the first time, a fully fair comparison of our system with several works from the literature. One of the major contributions of the paper is the proposed evaluation protocol for AMI, which will be adopted also in future Kaldi and Pyannote recipes, and which we hope will become a new standard.

The analysis of results for CALLHOME and AMI datasets reveals that systems are reaching very low diarization error rates when evaluating with the standard 0.25 s collar and without considering overlapped speech regions. We believe that, as in latest diarization challenges, systems should be tested using more challenging evaluation setups considering also overlapped speech and no collar.

Future work will focus on combining our system with real VAD instead of

the oracle labels, to approach the setup of nowadays largely popular end-to-end techniques.

7. Acknowledgements

We would like to thank Brian Sun [45, 57], Monisankha Pal [49], Matthew Maciejewski [46] and Desh Raj [47, 48] for detailed feedback about their evaluation protocols for AMI. Special thanks to Desh Raj, Paola Liebny García-Perera and Hervé Bredin for their contribution on the new evaluation protocol for AMI.

References

- [1] N. R. et al., DIHARD Corpus. Linguistic Data Consortium. (2018).
- [2] N. R. et. al., The Second DIHARD Diarization Challenge: Dataset, task, and baselines., in: Proceedings of Interspeech, 2019.
- [3] S. Watanabe, M. Mandel, J. Barker, E. Vincent, A. Arora, X. Chang, S. Khudanpur, V. Manohar, D. Povey, D. Raj, D. Snyder, A. S. Subramanian, J. Trmal, B. B. Yair, C. Boeddeker, Z. Ni, Y. Fujita, S. Horiguchi, N. Kanda, T. Yoshioka, N. Ryant, CHiME-6 Challenge: Tackling Multi-speaker Speech Recognition for Unsegmented Recordings (2020). [arXiv:2004.09249](#).
- [4] J. S. Chung, J. Huh, A. Nagrani, T. Afouras, A. Zisserman, Spot the conversation: speaker diarisation in the wild, [arXiv preprint arXiv:2007.01216](#).
- [5] Y. Fujita, S. Watanabe, S. Horiguchi, Y. Xue, K. Nagamatsu, End-to-End Neural Diarization: Reformulating Speaker Diarization as Simple Multi-label Classification (2020). [arXiv:2003.02966](#).
- [6] Z. Huang, S. Watanabe, Y. Fujita, P. García, Y. Shao, D. Povey, S. Khudanpur, Speaker diarization with region proposal network, in: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2020, pp. 6514–6518.

- [7] I. Medennikov, M. Korenevsky, T. Prisyach, Y. Khokhlov, M. Korenevskaya, I. Sorokin, T. Timofeeva, A. Mitrofanov, A. Andrusenko, I. Podluzhny, et al., Target-Speaker Voice Activity Detection: A Novel Approach for Multi-Speaker Diarization in a Dinner Party Scenario, *Interspeech 2020* doi:10.21437/interspeech.2020-1602.
URL <http://dx.doi.org/10.21437/Interspeech.2020-1602>
- [8] K. Kinoshita, M. Delcroix, N. Tawara, Integrating end-to-end neural and clustering-based diarization: Getting the best of both worlds (2020).
arXiv:2010.13366.
- [9] S. Horiguchi, Y. Fujita, S. Watanabe, Y. Xue, K. Nagamatsu, End-to-End Speaker Diarization for an Unknown Number of Speakers with Encoder-Decoder Based Attractors (2020). arXiv:2005.09921.
- [10] G. S. et al., Diarization is Hard: Some Experiences and Lessons Learned for the JHU Team in the Inaugural DIHARD Challenge, in: *Proc. Interspeech, 2018*, pp. 2808–2812. doi:10.21437/Interspeech.2018-1893.
URL http://www.danielpovey.com/files/2018_interspeech_dihard.pdf
- [11] F. Landini, S. Wang, M. Diez, L. Burget, P. Matějka, K. Žmolíková, L. Mošner, A. Silnova, O. Plchot, O. Novotný, et al., BUT System for the Second DIHARD Speech Diarization Challenge, in: *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2020, pp. 6529–6533.
- [12] X. Xiao, N. Kanda, Z. Chen, T. Zhou, T. Yoshioka, S. Chen, Y. Zhao, G. Liu, Y. Wu, J. Wu, S. Liu, J. Li, Y. Gong, Microsoft Speaker Diarization System for the VoxCeleb Speaker Recognition Challenge 2020 (2020).
arXiv:2010.11458.
- [13] M. Diez, L. Burget, F. Landini, S. Wang, H. Černocký, Optimizing Bayesian HMM based x-vector clustering for the second DIHARD speech

- diarization challenge, in: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2020, pp. 6519–6523.
- [14] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [15] A. F. Martin, M. A. Przybocki, Stream-based speaker segmentation using speaker factors and eigenvoices, in: 7th European Conference on Speech Communication and Technology, Eurospeech, Vol. 7, num. 2, 2001, pp. 787–790.
- [16] J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos, W. Kraaij, M. Kronenthal, et al., The AMI meeting corpus: A pre-announcement, in: International workshop on machine learning for multimodal interaction, Springer, 2006, pp. 28–39.
- [17] M. Diez, L. Burget, S. Wang, J. Rohdin, H. Černocký, Bayesian HMM based x-vector clustering for Speaker Diarization, in: Proc. Interspeech 2019, 2019, pp. 346–350. doi:10.21437/Interspeech.2019-2813.
- [18] M. Diez, L. Burget, F. Landini, J. Černocký, Analysis of Speaker Diarization Based on Bayesian HMM With Eigenvoice Priors, IEEE/ACM Transactions on Audio, Speech, and Language Processing 28 (2020) 355–368.
- [19] N. D. et al., Front-End Factor Analysis for Speaker Verification, IEEE Transactions on Audio, Speech, and Language Processing 19 (4) (2011) 788–798. doi:10.1109/TASL.2010.2064307.
- [20] P. Kenny, Bayesian Analysis of Speaker Diarization with Eigenvoice Priors, Tech. rep., Montreal: CRIM (2008).
- [21] P. Kenny, Bayesian Speaker Verification with Heavy-Tailed Priors, in: Proc. Odyssey-10, Brno, Czech Republic, 2010.

- [22] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
- [23] N. Brummer, E. Villiers, The speaker partitioning problem, *Proc. of Odyssey 2010*.
- [24] E. B Fox, E. B Sudderth, M. Jordan, A. S Willsky, The sticky HDP-HMM: Bayesian nonparametric hidden Markov models with persistent states, Technical Report P-2777, MIT LIDS (01 2007).
- [25] M. Beal, *Variational Algorithms for Approximate Bayesian Inference*, Ph.D. thesis, Gatsby Unit, University College London (2003).
- [26] M. Diez, L. Burget, CSL VBx derivations, http://www.fit.vutbr.cz/~mireia/CSL_VBMM_tech_report.pdf, Technical Report, Brno University of Technology (2021).
- [27] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, in: *Advances in neural information processing systems*, 2019, pp. 8026–8037.
- [28] H. Zeinali, S. Wang, A. Silnova, P. Matějka, O. Plchot, BUT system description to VoxCeleb speaker recognition challenge 2019, arXiv preprint arXiv:1910.12592.
- [29] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, S. Khudanpur, X-vectors: Robust DNN embeddings for speaker recognition, in: *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2018, pp. 5329–5333.
- [30] J. Deng, J. Guo, N. Xue, S. Zafeiriou, Arcface: Additive angular margin loss for deep face recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.

- [31] A. Nagrani, J. S. Chung, A. Zisserman, VoxCeleb: a large-scale speaker identification dataset, in: INTERSPEECH, 2017.
- [32] J. S. Chung, A. Nagrani, A. Zisserman, VoxCeleb2: Deep Speaker Recognition, in: INTERSPEECH, 2018.
- [33] Y. Fan, J. Kang, L. Li, K. Li, H. Chen, S. Cheng, P. Zhang, Z. Zhou, Y. Cai, D. Wang, CN-CELEB: a challenging Chinese speaker recognition dataset, in: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2020, pp. 7604–7608.
- [34] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz, et al., The Kaldi speech recognition toolkit, in: IEEE 2011 workshop on automatic speech recognition and understanding, IEEE Signal Processing Society, 2011.
- [35] Kaldi, SRE16 v2, <https://github.com/kaldi-asr/kaldi/tree/master/egs/sre16/v2>, [Downloaded: 2017-12].
- [36] H. Zeinali, H. Sameti, T. Stafylakis, DeepMine Speech Processing Database: Text-Dependent and Independent Speaker Verification and Speech Recognition in Persian and English., in: Odyssey, 2018, pp. 386–392.
- [37] D. Garcia-Romero, C. Espy-Wilson, Analysis of i-vector Length Normalization in Speaker Recognition Systems., 2011, pp. 249–252.
- [38] NIST SRE 2000 Evaluation Plan, https://www.nist.gov/sites/default/files/documents/2017/09/26/spk-2000-plan-v1.0.htm_.pdf.
- [39] S. H. Shum, N. Dehak, R. Dehak, J. R. Glass, Unsupervised Methods for Speaker Diarization: An Integrated and Iterative Approach, IEEE Transactions on Audio, Speech, and Language Processing 21 (10) (2013) 2015–2028. doi:10.1109/TASL.2013.2264673.

- [40] M. Senoussaoui, P. Kenny, T. Stafylakis, P. Dumouchel, A Study of the Cosine Distance-Based Mean Shift for Telephone Speech Diarization, *IEEE/ACM Trans. Audio, Speech and Lang. Proc.* 22 (1) (2014) 217–227. doi:10.1109/TASLP.2013.2285474.
URL <http://dx.doi.org/10.1109/TASLP.2013.2285474>
- [41] D. G.-R. et al., Speaker diarization using deep neural network embeddings, in: *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 4930–4934. doi:10.1109/ICASSP.2017.7953094.
- [42] NIST Rich Transcription Evaluations, <https://www.nist.gov/itl/iad/mig/rich-transcription-evaluation>, version: md-eval-v22.pl.
- [43] H. Bredin, R. Yin, J. M. Coria, G. Gelly, P. Korshunov, M. Lavechin, D. Fustes, H. Titeux, W. Bouaziz, M.-P. Gill, pyannote.audio: neural building blocks for speaker diarization, in: *ICASSP 2020, IEEE International Conference on Acoustics, Speech, and Signal Processing, Barcelona, Spain, 2020*.
- [44] L. Bullock, H. Bredin, L. P. Garcia-Perera, Overlap-aware diarization: re-segmentation using neural end-to-end overlapped speech detection (2019). arXiv:1910.11646.
- [45] G. Sun, C. Zhang, P. Woodland, Combination of Deep Speaker Embeddings for Diarisation (2020). arXiv:2010.12025.
- [46] M. Maciejewski, D. Snyder, V. Manohar, N. Dehak, S. Khudanpur, Characterizing Performance of Speaker Diarization Systems on Far-Field Speech Using Standard Methods, in: *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5244–5248. doi:10.1109/ICASSP.2018.8461546.
- [47] D. Raj, Z. Huang, S. Khudanpur, Multi-class Spectral Clustering with Overlaps for Speaker Diarization (2020). arXiv:2011.02900.

- [48] D. Raj, L. P. Garcia-Perera, Z. Huang, S. Watanabe, D. Povey, A. Stolcke, S. Khudanpur, DOVER-Lap: A Method for Combining Overlap-aware Diarization Outputs (2020). [arXiv:2011.01997](https://arxiv.org/abs/2011.01997).
- [49] M. Pal, M. Kumar, R. Peri, T. J. Park, S. H. Kim, C. Lord, S. Bishop, S. Narayanan, Meta-learning with Latent Space Clustering in Generative Adversarial Network for Speaker Diarization (2020). [arXiv:2007.09635](https://arxiv.org/abs/2007.09635).
- [50] G. Sun, C. Zhang, P. C. Woodland, Speaker diarisation using 2D self-attentive combination of embeddings, CoRR abs/1902.03190. [arXiv:1902.03190](https://arxiv.org/abs/1902.03190).
URL <http://arxiv.org/abs/1902.03190>
- [51] X. Anguera, C. Wooters, J. Hernando, Acoustic Beamforming for Speaker Diarization of Meetings, Audio, Speech, and Language Processing, IEEE Transactions on 15 (2007) 2011 – 2022. [doi:10.1109/TASL.2007.902460](https://doi.org/10.1109/TASL.2007.902460).
- [52] A. Zhang, Q. Wang, Z. Zhu, J. Paisley, C. Wang, Fully supervised speaker diarization, in: ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2019, pp. 6301–6305.
- [53] Q. Lin, R. Yin, M. Li, H. Bredin, C. Barras, LSTM based similarity measurement with spectral clustering for speaker diarization, arXiv preprint [arXiv:1907.10393](https://arxiv.org/abs/1907.10393).
- [54] H. Aronowitz, W. Zhu, M. Suzuki, G. Kurata, R. Hoory, New Advances in Speaker Diarization, in: Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020, ISCA, 2020, pp. 279–283. [doi:10.21437/Interspeech.2020-1879](https://doi.org/10.21437/Interspeech.2020-1879).
URL <https://doi.org/10.21437/Interspeech.2020-1879>
- [55] Q. Lin, W. Cai, L. Yang, J. Wang, J. Zhang, M. Li, DIHARD II is Still Hard: Experimental Results and Discussions from the DKU-LENOVO

Team, in: Proc. Odyssey 2020 The Speaker and Language Recognition Workshop, 2020, pp. 102–109. doi:10.21437/Odyssey.2020-15.

URL <http://dx.doi.org/10.21437/Odyssey.2020-15>

- [56] Q. Lin, Y. Hou, M. Li, Self-Attentive Similarity Measurement Strategies in Speaker Diarization, in: Proc. Interspeech 2020, 2020, pp. 284–288. doi:10.21437/Interspeech.2020-1908.

URL <http://dx.doi.org/10.21437/Interspeech.2020-1908>

- [57] Y. Fathullah, C. Zhang, P. Woodland, Improved Large-Margin Soft-max Loss for Speaker Diarisation, 2020, pp. 7104–7108. doi:10.1109/ICASSP40776.2020.9053373.