PRE-TRAINING OF SPEAKER EMBEDDINGS FOR LOW-LATENCY SPEAKER CHANGE DETECTION IN BROADCAST NEWS

Leda Sari^{*}, Samuel Thomas[†], Mark Hasegawa-Johnson^{*}, Michael Picheny[†]

* Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, USA † IBM Research AI, Yorktown Heights, USA

ABSTRACT

In this work, we investigate pre-training of neural network based speaker embeddings for low-latency speaker change detection. Our proposed system takes two speech segments, generates embeddings using shared Siamese layers and then classifies the concatenated embeddings depending on whether they are spoken by the same speaker. We investigate gender classification, contrastive loss and triplet loss based pre-training of the embedding layers and also joint training of the embedding layers along with a same/different classifier. Training is performed on 2-second single speaker segments based on ground truth speaker segmentation of broadcast news data. However, during test, we use the detection system in a practical low-latency setting for annotating automatic closed captions. In contrast to training, test pairs are now created around automatic speech recognition (ASR) based segmentation boundaries. The ASR segments are often shorter than 2 seconds causing duration mismatch during testing. In our experiments, although the baseline i-vector based classifier performs well, the proposed triplet loss based pre-training followed by joint training provides 7-50% relative F-measure improvement in matched and mismatched conditions. In addition, the degradation in performance is less severe for network based embeddings as compared to using i-vectors in the variable duration test conditions.

Index Terms— Speaker change detection, sequence embedding, Siamese networks

1. INTRODUCTION

Speaker change detection is the task of finding the time instances in audio recordings when a different speaker starts to speak. One general approach to this problem is to use a distance-based method. These methods extract features using sliding windows, compare feature representations of consecutive windows using a distance measure and then threshold the distance [1]. On the other hand, modelbased approaches fit a model to the features of individual segments and their concatenation, and choose the hypothesis with a higher score; this score can be the Bayesian information criterion (BIC) [2] or Gaussian likelihood score [3].

Among the most commonly used features to represent speaker characteristics of a speech segment are i-vectors [4]. Although i-vectors have been successfully used in speaker verification applications, reliability of these vectors depend on segment duration [5, 6]. In order to solve this problem, short speech segments are often clustered using BIC, Gaussian divergence [7] or x-means [8, 9] prior to computing the i-vector. However, these clustering methods are mainly designed for offline processing and cannot be used in low-latency applications [8].

Recently, neural network based speaker embeddings have been used as an alternative [10, 11] or as complementary features [12] to i-vectors. Studies have shown that network based embeddings can achieve better performance than using BIC based approaches on mel-frequency cepstral coefficients (MFCCs) [10, 11] or filterbank coefficients [13]. In [12], network embeddings are used in a speaker classification task with a probabilistic linear discriminant analysis backend and have been shown to achieve better performance than ivectors especially when the inputs are short (10s). These embedding networks are trained using multiclass cross-entropy for speaker classification using a large number of speakers [12, 13], using contrastive loss on two inputs processed in a Siamese architecture [14] or using triplet loss [10]. In order to map variable length sequences to fixed dimensional embeddings, long short-term memory (LSTM) [15] layers are usually employed.

In addition to generating embeddings, neural networks have also been used in end-to-end speaker change detection systems [16, 17, 18, 11] where the change decision is made at the end of a network instead of thresholding a distance measure. These systems can be classified into cases where the problem is reduced to taking two speech segments as input and comparing them [16, 11] or deciding if there is a change point within a given single speech segment [17, 18]. The networks that compare two segments usually have a Siamese structure where the initial few layers processing the two inputs share their weights. A similar structure is also usually used in embedding generating networks where the training objective consists of comparing the features extracted from the shared Siamese layers.

In this work, we combine learning speaker embeddings with an end-to-end approach for speaker change detection. We feed the Siamese embeddings of two segments into a fully-connected classifier for speaker change detection and use embedding learning as a pre-training method for these Siamese layers. We investigate three pre-training mechanisms - gender classification, contrastive loss and triplet loss with both cosine and Euclidean distance measures. After pre-training, we either freeze the embedding layers and train the classifier alone or we jointly update them. In order to handle variable length segments, we use bidirectional LSTM (BLSTM) layers in the Siamese part. In earlier studies, similar loss functions have been used for embedding generation and then distance based change detection is applied [10, 13] or the change detection is performed by a network without explicitly training the intermediate speaker embeddings [11]. In this paper, we propose to use pre-training to get better intermediate features for an end-to-end speaker change detection system.

Our second objective is to operate in the low-latency regime which limits our input segment duration; we choose 2s as the maximum segment duration as we consider 2s to be an acceptable latency for consumer applications. As observed in [8], large improvements in the speed of speaker diarization can be obtained if automatic speech recognition (ASR) based decision boundaries are used. Based on this observation, our test setup considers comparisons only



Fig. 1. Overview of the change detection network

around segment boundaries. During training, since manually annotated speaker boundaries are available, our training segments are extracted based on those boundaries. During testing, we use (a) matched condition to the training setup where the decisions are made around ground truth segments, (b) mismatched condition where the decision boundaries around which we take our two inputs are based on ASR segments. Especially this latter test condition introduces an additional mismatch with shorter segment durations. Using this setting for change detection also allows us to annotate speech recognition outputs with speaker change detections for richer closed captioning in real time.

In [10], although the triplet loss training can handle variable durations, the authors do not report the effect of duration mismatch between training and test segments and their change detection is based on distance thresholding. The end-to-end approach of [11] has a similar architecture to our network except that the authors take the absolute difference after Siamese layers rather than concatenating them, they do not employ pre-training, and their tests include only fixed-length segment comparison on synthetically concatenated speech segments. Both of these studies use a sliding window approach; conversely, as discussed above, we use ASR based segment boundaries in our tests for low-latency applications. In effect, one of our contributions is the application of ASR boundaries, proposed in [3], to a low-latency end-to-end neural speaker change detector. In summary, our main contributions are (1) incorporating existing speaker embedding learning methods into an end-to-end system as pre-training, (2) comparison of speaker change detection based on ivectors and neural network embeddings under mismatched test segment duration and (3) evaluation of change points based on ASR segment boundaries for low-latency applications.

The rest of the paper is organized as follows: In Section 2 we introduce our change detection system and describe the pre-training of the Siamese layers in the network. In Section 3, we present experimental setups and results, then we conclude the paper in Section 4.

2. SPEAKER CHANGE DETECTION

Figure 1 gives an overview of the speaker change detection network. Two speech segments are first passed through shared Siamese layers to generate their corresponding embeddings. The embeddings are then concatenated and fed into the classifier which is a fully-connected network. Therefore, in our end-to-end approach, same/different classification decisions are obtained at the output of a network rather than thresholding the distance between the embeddings. After training, speaker change times are determined by evaluating the network around segment boundaries and outputting the beginning time of the right segment.

The Siamese network consists of three BLSTM layers followed by two fully connected tanh layers. Transition from the BLSTM to the fully connected layers is achieved by concatenating the forward and backward average activations over time from the last BLSTM layer. Thus, we can map variable length sequences into fixed dimensional embeddings. These embedding layers are pretrained either for gender classification or using contrastive divergence loss or triplet loss. After pretraining, we either freeze the Siamese layers and learn a same/different classifier on the embeddings or update the parameters of the Siamese network along with the classifier to achieve better classification accuracy.

2.1. Pre-training with Gender Classification

In the first approach, we train a gender classification network using binary cross-entropy objective. This pretraining can be considered as a simplified version of using multi-class speaker classification to generate the embeddings. To get the gender probability $y^{(m)}$ for the *m*-th sample, we pass the embeddings through an affine layer followed by sigmoid nonlinearity. If the gender label is represented as with binary labels *g*, we can write the binary cross-entropy objective \mathcal{L}_x as

$$\mathcal{L}_x = \sum_{m=1}^{M} g^{(m)} \log(y^{(m)}) + (1 - g^{(m)}) \log(1 - y^{(m)})$$
 (1)

2.2. Pre-training with Contrastive Loss

In the second approach, we use a pair of inputs and use contrastive loss as our objective. In this setup, we try to minimize the distance between the embeddings of two inputs if they are uttered by the same speaker and we try to maximize the distance otherwise. Let s_l and s_r denote the speakers of the left and right segments, respectively. Also let δ to be the indicator function and Δ_c be a margin parameter. For a set of M training segments, the contrastive divergence loss \mathcal{L}_c that we want to minimize is

$$\mathcal{L}_{c} = \sum_{m=1}^{M} \delta[s_{l}^{(m)} = s_{r}^{(m)}] d(x_{l}^{(m)}, x_{r}^{(m)}) + \delta[s_{l}^{(m)} \neq s_{r}^{(m)}] \max(0, \Delta_{c} - d(x_{l}^{(m)}, x_{r}^{(m)}))$$
(2)

where d(., .) is a distance measure such as Euclidean or cosine distance and x_l and x_r are the embeddings obtained from the network.

2.3. Pre-training with Triplet Loss

In the third approach, we use triplets and we try to minimize the triplet loss. In this setup, we have an anchor segment, a positive segment uttered by the same speaker as the anchor and a negative segment uttered by a different speaker. The aim is to find embeddings such that for a given triplet, the anchor and the negative sample are separated more than the anchor and positive sample with a margin Δ_{tri} . If we denote the embeddings for the anchor, positive and negative samples using x_a, x_p and x_n , respectively, then the triplet loss \mathcal{L}_{tri} that we want to minimize is written as

$$\mathcal{L}_{tri} = \sum_{m=1}^{M} \max(0, \Delta_{tri} + d(x_a^{(m)}, x_p^{(m)}) - d(x_a^{(m)}, x_n^{(m)}))$$
(3)

As in (2), d(., .) denotes the distance measure.

3. EXPERIMENTS

3.1. Dataset

For our speaker change detection experiments we use 144 hours of audio from the Hub4 acoustic training data set collected between May 1996 and January 1997 [19, 20]. This data set (BN-144) which covers 288 TV shows is manually transcribed and annotated with speaker information.

In the experiments, we trained embeddings networks using a fixed set of pairs or triplets of inputs and did not resample the lists. All training segments are 2s and are extracted from speech segments greater than 2s in length after silence removal around manually labeled change points. To get our pairs, from each speaker we sample n_s segments; for each segment, we sample n_p positive samples, we sample n_p speakers from the remaining speaker set and sample one segment from each of them to get the negative pairs. The dataset has more segments spoken by male speakers therefore the resulting pairs mainly have male-male comparisons. To achieve a balance, we subsampled the initial list of pairs to get equal number of same gender and different gender comparisons and at the end we have 500000 pairs. For triplets, we followed an approach similar to [10] where we picked n random samples from each of S speakers, for each speaker we generated all possible pairs from those n segments, which gave the anchor-positive pairs, and then for each pair we randomly chose a single negative segment from the remaining (S-1)n segments. This resulted in 329000 triplets.

3.2. Training setups

Our 100 dimensional i-vectors [4] are extracted after using the maximum likelihood criterion to train a 2048 component, 40-dimensional diagonal covariance based UBM followed by i-vector extraction matrices. To train these systems, vectors of 9 consecutive PLP features are spliced together and projected down to 40 dimensions using LDA. The i-vector extraction systems are trained on the BN-144 dataset described above.

The embedding network consists of three BLSTM layers followed by two fully connected tanh layers. Transition from the BLSTM to the fully connected layers is achieved by concatenating the forward and backward average activations over time from the last BLSTM layer. In contrastive loss and triplet loss training with Euclidean distance, we also added an L2-normalization layer after the second fully connected layer so that the norms of the embeddings and therefore the distance between the embeddings become bounded. Inputs of the Siamese layers are 19-dimensional PLP features, appended with their deltas and delta-deltas.

Same/different classifier network is a 3-layer fully connected network with rectified linear unit nonlinearity. Its input is the concatenation of either PLPs, i-vectors or the Siamese embeddings of the left and right segments. Since our embeddings are 32 dimensional, the input of the classifier is 64 dimensional. The layers of the classifier have 64 units. Final layer has a single output node with sigmoid nonlinearity.

3.3. Test setups

Ten audio files are taken as test data which are not used in training. However, 117 out of 174 test speakers also have data in the training set. Audio pairs used in testing are determined based on speech segment boundaries. These boundaries are determined from either the segments obtained from an ASR system or the segments based on the ground truth speaker labeling with inter-speaker silences removed. The ASR segments are produced by the IBM Cloud Speech-To-Text (STT) service [21]. Although we trained our networks using 2s clips determined from long enough ground truth segments with silences



Fig. 2. Extraction of test segments around boundaries for fixed and variable length cases

Table 1. Validation data same/different classification accuracy (%) depending on the pre-training method, distance measure used in the objective and whether embedding layers are frozen. S: Siamese, C: classifier

Feat.	Net.	Pre-training	Distance	Freeze	Accu.	
				Embed.		
PLP	С	-	-	-	52.2	
i-vector	С	-	-	-	86.6	
PLP	S+C	Gender	-	Yes	76.9	
PLP	S+C	Gender	-	No	78.1	
PLP	S+C	Contrastive	Cosine	Yes	76.7	
PLP	S+C	Contrastive	Cosine	No	87.3	
PLP	S+C	Contrastive	Euclidean	Yes	77.4	
PLP	S+C	Contrastive	Euclidean	No	87.5	
PLP	S+C	Triplet	Cosine	Yes	84.6	
PLP	S+C	Triplet	Cosine	No	87.9	
PLP	S+C	Triplet	Euclidean	Yes	82.7	
PLP	S+C	Triplet	Euclidean	No	89.0	

removed, during test time we do not always have segments that are at least 2s. Therefore, as shown in Figure 2, for each segment type we performed two types of tests. In the first case, we still used fixed length (2s) data. For each boundary, we took 2s of data before and after the boundary, regardless of whether or not we have multiple speakers within that 2s. In the second case, we took variable length clips depending on the lengths of the segments. If they are shorter than 2s, we took the whole segment, otherwise we took 2s of data from the beginning (end) of the segment after (before) the boundary. To limit latency, we did not use longer segments. Note that especially in the variable length segment case, there is a significant amount of mismatch between training and testing because we do not have 2s clips as input and also the left and right clip possibly have different durations. For the variable length case, average segment duration for the ASR based test segments is 1.49s whereas for the ground truth based test segments it is 1.54s.

3.4. Results

Table 1 shows the same/different classification accuracy on the validation data. The first two rows show the performance of the classifier (C) when its input is concatenated PLP features or i-vectors of the two segments. The rest of the table summarizes the performance of Siamese layers followed by a classifier (S+C). In these rows, we use pre-training with different objectives, which are binary cross entropy for gender classification, contrastive or triplet loss. The distance measure used in the objective is cosine or Euclidean, and the embedding layers are frozen or updated during classifier training. As shown in the table, fine-tuning the embedding layers during classifier training improves the accuracy in all cases. If we freeze the Siamese embedding layers after pre-training, that is, use the Siamese network as a fixed feature extractor, we see that using the triplet loss as the objective leads to better accuracy. This implies that the triplet



(c) Variable length-Ground truth (

(d) 2s-Ground truth

Fig. 3. Precision-recall curves for i-vector and Siamese network based change detection depending on the type of the test segments

loss leads to more dicriminative features as compared to contrastive loss or gender classification. If we check the mean and variance of the Euclidean distance, we also observed that after training with the triplet loss, the average distance between different-pairs get much higher than that of the similar-pairs and the variances of the distances decrease for both same- and different-pairs.

Although the single sigmoid layer at the end of the classifier implies a default threshold of 0.5 for change decision, we can threshold the output to make change decisions when the network output has a higher value, i.e., more confident about a change. By varying this threshold, we obtained the precision-recall curves for the test data which are shown in Figures 3(a)-(d). They show the results obtained for ASR boundary based variable length and 2s segments, and also for ground truth boundary based variable length and 2s segments. In each figure, we include the curves for i-vector classification (ivector) and Siamese network setups pre-trained with triplet loss (Tri) depending on the distance measure, which is cosine (Cos) or Euclidean (Eucl) distance, and depending on whether we freeze (F) or jointly train (T) the embedding layers.

In all four figures, Euclidean distance based triplet loss pretraining followed by fine tuning during classifier training achieves the best performance. Joint training of the Siamese layers with the classifier achieves better performance than freezing them irrespective of the distance measure and the segment type. In the tests with variable length segments (Figures 3a and 3c), frozen embeddings perform as well as i-vectors, jointly trained embeddings perform better. For 2s segments (Figures 3b and 3d), differences among systems are smaller, but jointly trained embeddings are still the best.

For the default threshold of 0.5 imposed by the sigmoid output, we report the F-measures for four types of test segments depending on the classifier model in Table 2. In each column, the best F-measure is highlighted. As seen from the table, joint training of embedding layers with the classifier after Euclidean distance based triplet loss pre-training achieves the best performance for all segment types. For threshold 0.5, performance of the frozen embed-

 Table 2. F-measures on the test data for threshold 0.5 depending on the model and the segments used for testing

	ASR b	oundary	Ground truth boundary		
	Variable	2-second	Variable	2-second	
i-vector	0.3150	0.4902	0.5036	0.6109	
Tri-Cos-F	0.3626	0.4752	0.5130	0.5820	
Tri-Cos-T	0.4354	0.5014	0.5821	0.5999	
Tri-Eucl-F	0.3332	0.4591	0.4722	0.5736	
Tri-Eucl-T	0.4746	0.5323	0.6141	0.6511	

ding layers is slightly worse than i-vectors in tests with 2s segments. However, the reduction in performance is less severe for Siamese network based architectures compared to the i-vector based setup when the test segments have variable duration rather than fixed 2s duration. This results from the fact that i-vector estimation becomes harder for shorter segments whereas the BLSTM based Siamese embeddings are more robust against changes in input duration. The relative improvements in F-measures as compared to i-vectors are 50.7% in the highly mismatched condition (Variable length-ASR), and 6.6% in the matched condition (2s-Ground truth).

Table 2 reports the results without any threshold tuning at 0.5. If the threshold were to be tuned, we observed that having a threshold around 0.7 achieves a better F-measure in all cases. The conclusions from Table 2 are still applicable for this new threshold.

We also combined the i-vector and embedding based systems by averaging their sigmoid outputs. In the variable length case, low performance of i-vector deteriorates the F-measure in the combined system as compared to the Siamese network. However, in the 2s conditions, where i-vectors perform reasonably well, the combined system achieves better F-measure than the individual systems, showing the complementary nature of the two systems. The combined i-vector and Tri-Eucl-T system has an F-measure of 0.5599 and 0.6783 for 2s ASR based and ground truth based test segments, respectively, for a threshold of 0.5. These are around 5% relatively higher than Tri-Eucl-T on 2s segments.

4. CONCLUSIONS

In this work, we presented an end-to-end speaker change detection setup that consists of Siamese layers for speaker embedding generation and a classifier that makes same/different decisions. We investigated gender classification, contrastive loss and triplet loss based pre-training of the embedding layers. Since our objective is to operate in the low-latency regime, we trained our networks using 2s speech segments extracted around ground truth speaker boundaries and during test time we evaluated them on 2s or shorter segments. For low-latency applications, test comparisons should be made around ASR based segment boundaries. This framework also allows us to annotate automatic closed captions with speaker change information. However, these segments tend to be shorter than 2s, causing a mismatch in duration. In our experiments, we compared neural network embeddings with i-vectors and PLP on a broadcast news dataset. We showed that jointly trained neural network embeddings perform substantially better than i-vectors or frozen embeddings. Experimental results showed that pre-training using triplet loss with Euclidean distance followed by joint training of the classifier achieved higher F-measure (7-50% relative improvements) in both matched and mismatched (duration and boundary) cases than i-vectors. Therefore, lower dimensional jointly trained network embeddings provide a compact representation that is robust to the mismatch between training and testing conditions.

5. REFERENCES

- Matthew A Siegler, Uday Jain, Bhiksha Raj, and Richard M Stern, "Automatic segmentation, classification and clustering of broadcast news audio," in *Proc. DARPA speech recognition workshop*, 1997, vol. 1997.
- [2] Scott Chen, Ponani Gopalakrishnan, et al., "Speaker, environment and channel change detection and clustering via the bayesian information criterion," in *Proc. DARPA broadcast news transcription and understanding workshop*. Virginia, USA, 1998, vol. 8, pp. 127–132.
- [3] Daben Liu and Francis Kubala, "Fast speaker change detection for broadcast news transcription and indexing," in *Sixth European Conference on Speech Communication and Technology*, 1999.
- [4] Najim Dehak, Patrick J Kenny, Réda Dehak, Pierre Dumouchel, and Pierre Ouellet, "Front-end factor analysis for speaker verification," *IEEE Transactions on Audio, Speech,* and Language Processing, vol. 19, no. 4, pp. 788–798, 2011.
- [5] Achintya Kumar Sarkar, Driss Matrouf, Pierre Michel Bousquet, and Jean-François Bonastre, "Study of the effect of ivector modeling on short and mismatch utterance duration for speaker verification," in *Proc. Interspeech*, 2012, pp. 2662– 2665.
- [6] Ahilan Kanagasundaram, Robbie Vogt, David B Dean, Sridha Sridharan, and Michael W Mason, "I-vector based speaker recognition on short utterances," in *Proceedings of the 12th Annual Conference of the International Speech Communication Association*. International Speech Communication Association (ISCA), 2011, pp. 2341–2344.
- [7] Claude Barras, Xuan Zhu, Sylvain Meignier, and J-L Gauvain, "Multistage speaker diarization of broadcast news," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 5, pp. 1505–1512, 2006.
- [8] Dimitrios Dimitriadis and Petr Fousek, "Developing on-line speaker diarization system," in *Proc. Interspeech*, 2017, pp. 2739–2743.
- [9] Dan Pelleg, Andrew W Moore, et al., "X-means: Extending k-means with efficient estimation of the number of clusters.," in *ICML*, 2000, vol. 1, pp. 727–734.
- [10] Hervé Bredin, "Tristounet: triplet loss for speaker turn embedding," in *IEEE International Conference on Acoustics, Speech* and Signal Processing (ICASSP). IEEE, 2017, pp. 5430–5434.
- [11] Arindam Jati and Panayiotis Georgiou, "An unsupervised neural prediction framework for learning speaker embeddings using recurrent neural networks," *Proc. Interspeech*, pp. 1131– 1135, 2018.
- [12] David Snyder, Daniel Garcia-Romero, Daniel Povey, and Sanjeev Khudanpur, "Deep neural network embeddings for textindependent speaker verification," in *Proc. Interspeech*, 2017, pp. 999–1003.
- [13] Renyu Wang, Mingliang Gu, Lantian Li, Mingxing Xu, and Thoms Fang Zheng, "Speaker segmentation using deep speaker vectors for fast speaker change scenarios," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).* IEEE, 2017, pp. 5420–5424.

- [14] Neil Zeghidour, Gabriel Synnaeve, Nicolas Usunier, and Emmanuel Dupoux, "Joint learning of speaker and phonetic similarities with siamese networks.," in *Proc. Interspeech*, 2016, pp. 1295–1299.
- [15] Sepp Hochreiter and Jürgen Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [16] Sree Harsha Yella, Andreas Stolcke, and Malcolm Slaney, "Artificial neural network features for speaker diarization," in *Spoken Language Technology Workshop (SLT)*. IEEE, 2014, pp. 402–406.
- [17] Ruiqing Yin, Hervé Bredin, and Claude Barras, "Speaker change detection in broadcast tv using bidirectional long shortterm memory networks," in *Proc. Interspeech*. ISCA, 2017.
- [18] Marek Hrúz and Zbyněk Zajíc, "Convolutional neural network for speaker change detection in telephone speaker diarization system," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).* IEEE, 2017, pp. 4945–4949.
- [19] "1996 English Broadcast News Speech (HUB4)," https:// catalog.ldc.upenn.edu/LDC97S44.
- [20] "1997 English Broadcast News Speech (HUB4)," https:// catalog.ldc.upenn.edu/LDC98S71.
- [21] "IBM Cloud Speech-to-Text," ttps://www.ibm.com/ watson/services/speech-to-text/.