

DNN-based SRE Systems in Multi-Language Conditions — Technical Report

Ondřej Novotný, Pavel Matějka, Ondřej Glembek, Oldřich Plchot, František Grézl, Lukáš Burget, and Jan “Honza” Černocký

Abstract

This work studies the usage of the (currently state-of-the-art) Deep Neural Networks (DNN) i-vector/PLDA-based speaker recognition systems in multi-language (especially non-English) conditions. On the “Language Pack” of the PRISM set, we evaluate the systems’ performance using NIST’s standard metrics. We study the use of multi-lingual DNN in place of the original English DNN on these multi-language conditions. We show that not only the gain from using DNNs vanishes, but also the DNN-based systems tend to produce de-calibrated scores under the studied conditions. This work gives suggestions for directions of future research rather than any particular solutions.

1 Introduction

During the last decade, neural networks have experienced a renaissance as a powerful machine learning tool. Deep Neural Networks (DNN) have been also successfully applied to the field of speech processing. After their great success in automatic speech recognition (ASR) [1], DNNs were also found very useful in other fields of speech processing such as speaker [2, 3, 4] or language recognition [5, 6, 7]. In speech recognition, DNNs are often directly trained for the “target” task of frame-by-frame classification of speech sounds (e.g. phones). Similarly, a DNN directly trained for frame-by-frame classification of languages was successfully used for language recognition in [7]. However, this system provided competitive performance only for speech utterances of short durations.

In the field of speaker recognition, DNNs are usually used in more elaborate and indirect way: One approach is to use DNNs for extracting frame-by-frame speech features. Such features are than used in the usual way (e.g. input to i-vector based system [8]). These features can be directly derived from the DNN output posterior probabilities [9] and combined with the conventional features (PLP or MFCC) [10]. More commonly, however, bottleneck (BN) DNNs are trained for a specific task, where the features are taken from a narrow hidden layer compressing the relevant information into low dimensional feature vectors [6, 5, 11]. Alternatively, standard DNN (with no bottleneck) can be used,

where the high-dimensional outputs of one of the hidden layers can be converted to features using a dimensionality reduction technique such as PCA [12].

In [13], we analyzed various DNN approaches to speaker recognition (as was similarly studied e.g. in [14, 15]). We used two different DNN’s (a monolingual—trained on the Fisher English data corpus—and a multi-lingual—trained on 11 languages of the Babel data collection). The rest of the system was trained on the PRISM set, i.e. mainly on the English data. We reported our results only on the NIST SRE 2010 telephone condition (i.e. only on English speech) via the Equal Error Rates (EERs) and the minimum DCF NIST metrics.

However, when tested on non-English test sets, we observed that the benefit of using the DNNs performance of the systems degraded dramatically. We used the “lan” Language Pack of the PRISM set (described later in the paper), and its Chinese subset—the “chn” pack in comparison with the originally used NIST SRE 2010 telephone condition. Not only we saw performance degradation in terms of EER and the minimum DCFs, but more so in terms of the actual DCFs, i.e. the systems produce heavily de-calibrated scores.

Our hypothesis was that when we use the DNN trained for the target language, the error rates would decrease. To match the sre10, “lan”, and “chn” test conditions, we chose two DNNs: i) the Fisher English, and the ii) Multilingual DNN. However, it turned out that, apart from the Fisher English being optimal for the NIST SRE 2010 test, there was no clear correlation between the test language and the DNN training language.

This paper merely analyzes the problems that emerged when applying the current state-of-the-art SRE systems to non-English domains, and rather provides directions for future research.

2 Theoretical Background

2.1 i-vector Systems

The i-vectors [8] provide an elegant way of reducing large-dimensional input data to a small-dimensional feature vector while retaining most of the relevant information. The main principle is that the utterance-dependent Gaussian Mixture Model (GMM) supervector of concatenated mean vectors \mathbf{s} is modeled as

$$\mathbf{s} = \mathbf{m} + \mathbf{T}\mathbf{w}, \tag{1}$$

where $\mathbf{m} = [\boldsymbol{\mu}^{(1)'}, \dots, \boldsymbol{\mu}^{(C)'}]'$ is the Universal Background Model (UBM) GMM mean supervector (of C components), $\mathbf{T} = [\mathbf{T}^{(1)'}, \dots, \mathbf{T}^{(C)'}]'$ is a low-rank matrix representing M bases spanning subspace with important variability in the mean supervector space, and \mathbf{w} is a latent variable of size M with standard normal distribution.

The i-vector $\boldsymbol{\phi}$ is the Maximum a Posteriori (MAP) point estimate of the variable \mathbf{w} . It maps most of the relevant information from a variable-length observation \mathcal{X} to a fixed- (small-) dimensional vector. $\mathbf{L}_{\mathcal{X}}$ is the precision of the posterior distribution.

The closed-form solution for computing the i-vector can be expressed as a function of the *zero-* and *first-order statistics*: $\mathbf{n}_{\mathcal{X}} = [N_{\mathcal{X}}^{(1)}, \dots, N_{\mathcal{X}}^{(C)}]'$ and $\mathbf{f}_{\mathcal{X}} = [\mathbf{f}_{\mathcal{X}}^{(1)'}, \dots, \mathbf{f}_{\mathcal{X}}^{(C)'}]'$, where

$$N_{\mathcal{X}}^{(c)} = \sum \gamma_t^{(c)} \quad (2)$$

$$\mathbf{f}_{\mathcal{X}}^{(c)} = \sum_t \gamma_t^{(c)} \mathbf{o}_t, \quad (3)$$

where $\gamma_t^{(c)}$ is the posterior (or occupation) probability of frame t being generated by the mixture component c . The tuple $\gamma_t = (\gamma_t^{(1)}, \dots, \gamma_t^{(C)})$ is usually referred to as *frame alignment*. Note that this variable can be computed either using the GMM UBM or using a completely different model [2, 14, 15]. We will refer to this approach as a *DNN alignment* approach later in this paper. The i-vector is then expressed as

$$\phi_{\mathcal{X}} = \mathbf{L}_{\mathcal{X}}^{-1} \bar{\mathbf{T}}' \bar{\mathbf{f}}_{\mathcal{X}} \quad (4)$$

where $\mathbf{L}_{\mathcal{X}}$ is the precision matrix of the posterior distribution, computed as:

$$\mathbf{L}_{\mathcal{X}} = \mathbf{I} + \sum_{c=1}^C N_{\mathcal{X}}^{(c)} \bar{\mathbf{T}}^{(c)'} \bar{\mathbf{T}}^{(c)}, \quad (5)$$

with c being the GMM UBM component index, and the ‘bar’ symbols denote normalized variables:

$$\bar{\mathbf{f}}_{\mathcal{X}}^{(c)} = \Sigma^{(c)-\frac{1}{2}} \left(\mathbf{f}_{\mathcal{X}}^{(c)} - N_{\mathcal{X}}^{(c)} \boldsymbol{\mu}^{(c)} \right) \quad (6)$$

$$\bar{\mathbf{T}}^{(c)} = \Sigma^{(c)-\frac{1}{2}} \mathbf{T}^{(c)}, \quad (7)$$

where $\Sigma^{(c)-\frac{1}{2}}$ is a symmetrical decomposition (such as Cholesky decomposition) of an inverse of the GMM UBM covariance matrix $\Sigma^{(c)}$.

2.2 Stacked Bottleneck Features (SBN)

Bottleneck Neural-Network (BN-NN) refers to such topology of a NN, one of whose hidden layers has significantly lower dimensionality than the surrounding layers. A bottleneck feature vector is generally understood as a by-product of forwarding a primary input feature vector through the BN-NN and reading off the vector of values at the bottleneck layer. We have used a cascade of two such NNs for our experiments. The output of the first network is *stacked* in time, defining context-dependent input features for the second NN, hence the term Stacked Bottleneck Features.

The NN input features are 24 log Mel-scale filter bank outputs augmented with fundamental frequency features from 4 different f_0 estimators (Kaldi, Snack¹, and two other according to [16] and [17]). Together, we have 13 f_0

¹<http://kaldi.sourceforge.net>, www.speech.kth.se/snack/

related features, see [18] for more details. The conversation-side based mean subtraction is applied on the whole feature vector. 11 frames of log filter bank outputs and fundamental frequency features are stacked together. Hamming window followed by DCT consisting of 0^{th} to 5^{th} base are applied on the time trajectory of each parameter resulting in $(24 + 13) \times 6 = 222$ coefficients on the first stage NN input.

The configuration for the first NN is $222 \times D_H \times D_H \times D_{BN} \times D_H \times K$, where K is the number of targets. The dimensionality of the bottleneck layer, D_{BN} was fixed to 80. This was shown as optimal in [6]. The dimensionality of other hidden layers was set to 1500. The bottleneck outputs from the first NN are sampled at times $t-10$, $t-5$, t , $t+5$ and $t+10$, where t is the index of the current frame. The resulting 400-dimensional features are input to the second stage NN with the same topology as first stage. The 80 bottleneck outputs from the second NN (referred as SBN) are taken as features for the conventional GMM/UBM i-vector based SID system.

We experimented with English and multilingual BN features. In the case of multilingual training, we adopted training scheme with block-softmax, which divides the output layer into parts according to individual languages. During training, only the part of the output layer is activated that corresponds to the language that the given target belongs to. See [19, 20] for detailed description.

2.3 DNN Alignment

The true frame alignment is a hidden variable in GMM modeling. Traditionally, it is computed using the GMM UBM (as used in the “baseline” and “SBN” experiments further in the paper). However, it was shown that DNNs (as well as other models, e.g. [13, 15, 14]) can be used directly for posterior computation [2].

For completeness, we report the performance of the DNN alignment systems, where the posteriors of the SBN-NNs from the previous section were used. In other words, we show the utility of the trained DNNs as both feature- and posterior-extractors.

Note that the output activation function of the Multilingual SBN is a block-softmax, giving a set of posterior probabilities (one set per training language). Therefore, we cannot utilize the Multilingual SBN for this purpose in a straightforward way.

Note also that the *normalization GMM UBM* (i.e. the $\boldsymbol{\mu}^{(c)}$ and $\boldsymbol{\Sigma}^{(c)}$ parameters) should be computed via the same alignment as used in eq. (2) and (3), i.e. the DNN alignment.

Table 1: Comparison of the systems under the PRISM “lan” and “chn”, and the SRE2010-condition 5 (tel-tel) tests. We expected (without result) the Multilang SBN to perform best in the “lan” and “chn” conditions.

Test set	System	DCF _{new} ^{min}		DCF _{old} ^{min}		EER [%]	
		male	female	male	female	male	female
chn	Baseline	0.1834	0.3019	0.0621	0.0894	1.44	2.27
	English SBN	0.1491	0.2251	0.0418	0.0838	1.00	1.99
	Multilang SBN	0.2121	0.1907	0.0439	0.0670	1.16	1.93
	English DNN	0.1373	0.3621	0.0616	0.1192	1.29	3.05
lan	Baseline	0.2979	0.9836	0.1021	0.2007	2.60	5.05
	English SBN	0.2963	0.9848	0.0979	0.2305	2.45	4.93
	Multilang SBN	0.4008	0.9854	0.0898	0.2997	2.16	5.03
	English DNN	0.2963	0.9463	0.0914	0.2228	2.70	5.68
sre10	Baseline	0.3577	0.3387	0.0967	0.1013	1.84	1.94
	English SBN	0.1295	0.1679	0.0387	0.0471	1.17	1.11
	Multilang SBN	0.1280	0.1696	0.0416	0.0544	1.21	1.16
	English DNN	0.1200	0.2212	0.0352	0.0449	0.71	0.93

3 Experiments

3.1 DNN Training Data

For training the Multilingual neural networks, the IARPA Babel Program data² were mainly used. This data set simulates the scenario of what one could collect in a limited time from a completely new language. It consists mainly of conversational telephone speech (CTS), but scripted recordings, as well as far field recordings, are present. We used 11 languages to train our multilingual SBN feature extractor. The *language list* (as referred to later in this paragraph) consists of Cantonese, Assamese, Bengali, Pashtu, Turkish, Tagalog, Vietnamese, Haiti, Lao, Tamil, and Zulu. More details about the characteristics of the languages can be found in [21]. The phone-state target labels were obtained using forced-alignment with our BABEL ASR system [22], with $471 + 141 + 147 + 216 + 126 + 252 + 303 + 99 + 411 + 102 + 219 = 2487$ phone states, respectfully to the *language list*.

For the English DNN variant, we have used a selection of 250 hours of data derived from the Fisher English Part 1 and 2 with 2423 tied tri-phone states.

²Collected by Appen, <http://www.appenbutlerhill.com>

Table 2: Analysis of the actual DCF’s under the PRISM “lan” and “chn”, and the SRE2010-condition 5 (tel-tel) tests. Note the system de-calibration on the “lan” and “chn” conditions. Also note that de-calibration is more emphasized for the female conditions. (Due to the dynamic range of the values, we prefer to report a table of numbers rather than a graph plot.)

Test	System	DCF _{new}				DCF _{old}			
		actual		min		actual		min	
		male	female	male	female	male	female	male	female
chn	Baseline	5.7461	16.0798	0.1834	0.3019	0.1206	0.2785	0.0621	0.0894
	English SBN	1.5201	10.4024	0.1491	0.2251	0.0515	0.1857	0.0418	0.0838
	Multilang SBN	3.9156	12.3843	0.2121	0.1907	0.0863	0.2189	0.0439	0.0670
	English DNN	10.2419	46.4058	0.1373	0.3621	0.1856	0.6857	0.0616	0.1192
lan	Baseline	3.5369	14.0482	0.2979	0.9836	0.1142	0.2812	0.1021	0.2007
	English SBN	2.1503	24.4566	0.2963	0.9848	0.0702	0.3476	0.0979	0.2305
	Multilang SBN	5.2089	38.1320	0.4008	0.9854	0.1121	0.4855	0.0898	0.2997
	English DNN	6.6261	36.8887	0.2963	0.9463	0.1427	0.5451	0.0914	0.2228
sre10	Baseline	0.4323	0.3442	0.3577	0.3387	0.1587	0.2171	0.0967	0.1013
	English SBN	0.1472	0.1750	0.1295	0.1679	0.0976	0.1098	0.0387	0.0471
	Multilang SBN	0.1530	0.1921	0.1280	0.1696	0.1171	0.1339	0.0416	0.0544
	English DNN	0.1234	0.2286	0.1200	0.2212	0.0800	0.1204	0.0352	0.0449

3.2 Test Set and Evaluation Metric

We report our results on the “Language Set” pack of the PRISM set [23], referred to as “lan” later in the results. It was crafted from the NIST SRE 2005–2008 datasets by selecting 500 speakers for which there exists at least one session in a language other than English. Additional 300 speakers (that appear only in English conversations) were added from the NIST SRE 2010. The trials were created as a Cartesian product of all sessions sessions, resulting in 3590/130880 male, and 6304/297683 female target/non-target trials, respectively. Note that half of the trials are still English.

Moreover, results on the Chinese subset of the “lan” condition, referred to as “chn” are reported. The set comprises of 1027/59004 male, and 1555/113405 female target/non-target trials, respectively.

To provide a contrastive view, we also report the results on the NIST SRE 2010 data extended core condition (telephone-telephone, “condition-5”), referred to as “sre10”, with 3465/175873 male, and 3704/233077 female target/non-target trials, respectively.

The detection cost function (DCF) is used as a primary evaluation metric. We report two numbers: DCF_{old}^{min} and DCF_{new}^{min} , corresponding to the primary

evaluation metric for the NIST speaker recognition evaluation in 2008 and 2010, respectively. We also report their *actual* variants DCF_{old}^{act} and DCF_{new}^{act} . Equal Error Rate (EER) is also reported. For more details, see the evaluation plans of NIST SRE ³.

3.3 System Description

Voice Activity Detection (VAD) was performed using Neural Network speech/non-speech classifier. The NN was trained on Czech CTS data where we artificially added noise with different levels of SNR to 30% of the database. The NN had two hidden layers each comprising of 300 neurons. We used a vectorized block of 31 frames of 15 Mel filter bank energies as input features. For the *interview data*, we removed the interviewer based on the ASR transcripts provided by NIST.

As the baseline features, we used 19 MFCC coefficients + energy augmented with their delta and double delta coefficients, resulting in 60-dimensional feature vectors. The analysis window was 20 ms long with the shift of 10 ms. First, we removed silence frames according to VAD, after which we applied short-time (300 frames) cepstral mean and variance normalization.

The PRISM set [23] was chosen as the base training dataset platform. It contains the following telephone data: NIST SRE 2004, 2005, 2006, 2008, 2010 Switchboard II Phases 2 and 3, Switchboard Cellular Parts 1 and 2, Fisher English Parts 1 and 2 giving 9670 female speakers. We have not included any noisy or reverberated data.

A gender-independent UBM was represented as a full or diagonal covariance 2048-component GMM. It was trained on a subset of PRISM, giving 15602 files equally distributed between telephone and microphone condition, and male and female portions. The variance flooring was used in each iteration of EM algorithm during the UBM training.

Gender-independent i-vector extractor was trained (in 10 iterations of a joint Expectation Maximization and Minimum Divergence steps) using the entire PRISM set. The results are reported with 600-dimensional i-vectors.

Gender-independent LDA and PLDA was trained on the same data as the i-vector extractor.

3.4 Results

Tab. 1 shows the overall results of all systems in terms of (calibration insensitive) DCF_{old}^{min} , DCF_{new}^{min} , and EER. For the “sre10” test, the best performing system is the DNN-alignment with the DNN trained on the Fisher English data, as expected. However, when looking at the “lan” condition, there is no gain from switching from the Baseline system to English DNN (and only a negligible gain in switching to English SBN).

³www.itl.nist.gov/iad/mig/tests/sre/

Our hypothesis was that this behavior would be fixed by using a more general DNN, such as the Multilingual DNN (only in the SBN variant, as explained in Sec. 2.3), since the test comprises of numerous languages.

Looking at the “chn” condition, we again expected the Multilingual SBN to significantly outperform the English DNN (or SBN), but with no result.

Our initial hypothesis was that the English training corpus is the largest, and therefore had to provide best phoneme accuracy and thus a better acoustic space clustering. However, it was observed in many cases (e.g. in [24]) that better phoneme accuracy does not necessarily imply better SRE performance. Therefore, we leave this question open for future research.

Let us also note that the UBM/i-vector/PLDA training data are identical—i.e., mainly English—across the different systems. Our hypothesis is that even if the DNN matches the target language, the acoustic space clustering does not correspond to the observed data. Therefore, the first-order statistics (3) for the i-vector extractor computation are “warped”, and the i-vector extractor captures a different “total” variability than is in fact used for the test. One of the possible indications for this hypothesis is the fact that the performance on the “sre10” condition does not vary dramatically across different systems. Similar hypothesis holds for the PLDA/LDA modeling, where the within/across variabilities are modeled using these “warped” i-vectors.

Tab. 2 shows the overall performance summary in terms of the actual vs. the minimum DCF values, i.e., it directly shows the calibration loss. We see that the “sre10” condition is well calibrated, i.e., the actual values are close enough to the minimum counterparts. However, looking at the “chn” and “lan” tests, and especially at the new DCF metric, the calibration losses are extremely high. This effect is even more pronounced for the female part of the tests. All this behavior indicates a heavy language-dependent score modality.

4 Conclusions

In this work, we have studied the behavior of the DNN techniques in SRE i-vector/PLDA systems, currently considered to be state-of-the-art, as evaluated on the most common NIST SRE English test sets, such as the NIST SRE 2010, condition 5. We have shown that when applied to non-English test sets, these techniques stop being effective and are susceptible to de-calibration of the scores produced by the traditional i-vector/PLDA systems.

This work suggests that we focus on the analysis of the DNN acoustic space clustering with regard to multiple languages and other types of variability.

References

- [1] Geoffrey Hinton, Li Deng, Dong Yu, Abdel rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara Sainath George Dahl, and Brian Kingsbury, “Deep neural networks for acoustic modeling in speech recognition,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 82–97, November 2012.

- [2] Y. Lei, N. Scheffer, L. Ferrer, and M. McLaren, “A novel scheme for speaker recognition using a phonetically-aware deep neural network,” in *ICASSP*, 2014.
- [3] Y. Lei, L. Ferrer, M. McLaren, and N. Scheffer, “Comparative study on the use of senone-based deep neural networks for speaker recognition,” *Submitted to IEEE Trans. ASLP*, 2014.
- [4] Garcia-Romero D., Zhang X., McCree A., and Povey D., “Improving speaker recognition performance in the domain adaptation challenge using deep neural networks,” in *SLT*, 2014.
- [5] Y. Song et al, “i-vector representation based on bottle neck feature for language identification,” in *IEEE Electronics Letters*, 2013.
- [6] Pavel Matějka et al., “Neural network bottleneck features for language identification,” in *IEEE Odyssey: The Speaker and Language Recognition Workshop*, Joensuu, Finland, 2014.
- [7] Ignacio Lopez-Moreno, Javier Gonzalez-Dominguez, and Oldřich Plchot, “Automatic language identification using deep neural networks,” in *ICASSP 2014*, Florence, Italy, 2014.
- [8] N. Dehak, P. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, “Front-end factor analysis for speaker verification,” *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. PP, no. 99, 2010.
- [9] Mireia Diez, Amparo Varona, Mikel Penagarikano, Luis Javier Rodriguez-Fuentes, and Germán Bordel, “Using phone log-likelihood ratios as features for speaker recognition,” in *Interspeech 2013*, Lyon, France, 2013.
- [10] Jeff Ma et al., “Improvements in language identification on the RATS noisy speech corpus,” in *Interspeech 2013*, Lyon, France, 2013.
- [11] Najim Dehak Fred Richardson, Douglas A. Reynolds, “A unified deep neural network for speaker and language recognition,” in *Interspeech*, 2015.
- [12] Yuan Liu, Yanmin Qian, Nanxin Chen, Tianfan Fu, Ya Zhang, and Kai Yu, “Deep feature for text-dependent speaker verification,” *Speech Communication*, vol. 73, pp. 1–13, October 2015.
- [13] Pavel Matějka, Ondřej Glembek, Ondřej Novotný, Oldřich Plchot, František Grézl, Lukáš Burget, and Jan Černocký, “Analysis of DNN approaches to speaker identification,” in *Proceedings of the 41th IEEE International Conference on Acoustics, Speech and Signal Processing*. 2016, pp. 5100–5104, IEEE Signal Processing Society.
- [14] Yao Tian, Meng Cai, Liang He, and Jia Liu, “Investigation of bottleneck features and multilingual deep neural networks,” in *Interspeech*, 2015.

- [15] Sandro Cumani, Olda Plchot, and Pietro Laface, “Comparison of hybrid dnn-gmm architectures for speaker recognition,” in *ICASSP. 2016*, IEEE Signal Processing Society.
- [16] Kornel Laskowski and Jens Edlund, “A Snack implementation and Tcl/Tk interface to the fundamental frequency variation spectrum algorithm,” in *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC’10)*, Valletta, Malta, may 2010.
- [17] David Talkin, “A robust algorithm for pitch tracking (RAPT),” in *Speech Coding and Synthesis*, W. B. Kleijn and K. Paliwal, Eds., New York, 1995, Elsevier.
- [18] Martin Karafiát, František Grézl, Karel Veselý, Mirko Hannemann, Igor Szóke, and Jan Černocký, “BUT 2014 Babel system: Analysis of adaptation in NN based systems,” in *Interspeech 2014*, 2014, pp. 3002–3006.
- [19] Karel Veselý, Martin Karafiát, František Grézl, Miloš Janda, and Ekaterina Egorova, “The language-independent bottleneck features,” in *Proceedings of IEEE 2012 Workshop on Spoken Language Technology*, 2012, pp. 336–341.
- [20] Radek Fér, Pavel Matějka, František Grézl, Oldřich Plchot, and Jan Černocký, “Multilingual bottleneck features for language recognition,” *Interspeech 2015*, 2015.
- [21] M. Harper, “The BABEL program and low resource speech technology,” in *ASRU 2013*, Dec 2013.
- [22] Martin Karafiát, František Grézl, Mirko Hannemann, Karel Veselý, and Jan ”Honza” Černocký, “BUT BABEL System for Spontaneous Cantonese,” in *Interspeech 2013*, Lyon, France, 2013, pp. 2589–2593.
- [23] L. Ferrer, H. Bratt, L. Burget, H. Cernocky, O. Glembek, M. Graciarena, A. Lawson, Y. Lei, P. Matejka, O. Plchot, et al., “Promoting robustness for speaker modeling in the community: the prism evaluation set,” <https://code.google.com/p/prism-set/>, 2012.
- [24] Alicia Lozano-Diez, Anna Silnova, Pavel Matejka, Ondrej Glembek, Oldrich Plchot, Jan Pesan, Lukas Burget, and Joaquin Gonzalez-Rodriguez, “Analysis and optimization of bottleneck features for speaker recognition,” in *Odyssey 2016: The Speaker and Language Recognition Workshop*, Bilbao, Spain, June 21-24 2016, pp. 352–357.