

SPEECHNAS: TOWARDS BETTER TRADE-OFF BETWEEN LATENCY AND ACCURACY FOR LARGE-SCALE SPEAKER VERIFICATION

Wentao Zhu*, Tianlong Kong*, Shun Lu, Jixiang Li,
Dawei Zhang, Feng Deng, Xiaorui Wang, Sen Yang, Ji Liu

Kuaishou Technology

ABSTRACT

Recently, x-vector [1] has been a successful and popular approach for speaker verification, which employs a time delay neural network (TDNN) and statistics pooling to extract speaker characterizing embedding from variable-length utterances. Improvement upon the x-vector has been an active research area, and enormous neural networks have been elaborately designed based on the x-vector, *e.g.*, extended TDNN (E-TDNN) [2], factorized TDNN (F-TDNN) [3], and densely connected TDNN (D-TDNN) [4]. In this work, we try to identify the optimal architectures from a TDNN based search space employing neural architecture search (NAS), named SpeechNAS. Leveraging the recent advances in the speaker recognition, such as high-order statistics pooling, multi-branch mechanism, D-TDNN and angular additive margin softmax (AAM) loss with a minimum hyper-spherical energy (MHE), SpeechNAS automatically discovers five network architectures, from SpeechNAS-1 to SpeechNAS-5, of various numbers of parameters and GFLOPs on the large-scale text-independent speaker recognition dataset VoxCeleb1. Our derived best neural network achieves an equal error rate (EER) of 1.02% on the standard test set of VoxCeleb1, which surpasses previous TDNN based state-of-the-art approaches by a large margin.

Index Terms— speaker verification, speaker recognition, SpeechNAS, neural architecture search, TDNN

1. INTRODUCTION

There are numerous measurements and signals, such as fingerprint, face, iris and voice, being investigated for biometric recognition systems [5]. Among these most popular measurements, voice has been one of the most compelling biometrics, because 1) the microphone system has been one of the most widely adopted intelligent agent to extract the speech signal in various hardwares, and 2) the speech sample can be widely accepted and does not considered threatening by users. Most importantly, the speaker recognition area has been well studied for over fifty years, and there is a rich scientific basis and extensive development over the area.

*Equal contributions. <https://github.com/wentaozhu/speechnas.git>

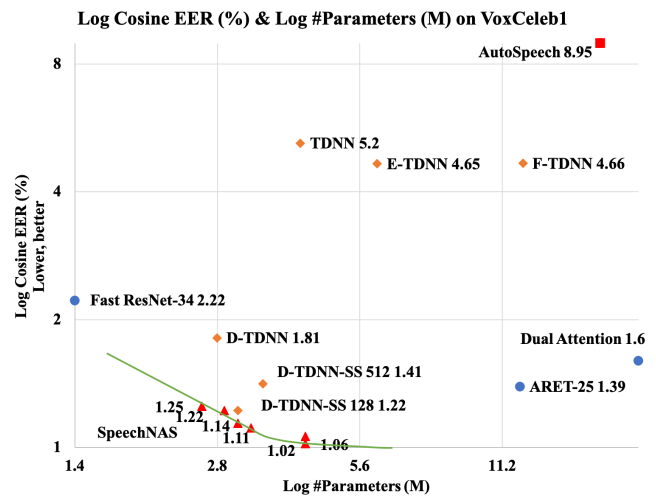


Fig. 1. SpeechNAS achieves better cosine EERs with lower number of parameters.

Deep neural networks have been widely adopted to extract speaker representations [1]. Recently, time delay neural network (TDNN) with x-vectors [1] has been a paradigm for speaker verification. Compared with vanilla TDNN, extended TDNN (E-TDNN) [2] interleaves feed-forward neural network (FNN) layers between the TDNN layer for multi-speaker conversations. Factorized TDNN (F-TDNN) [3] further reduces the number of parameters of the TDNN x-vector by factorizing the weight matrix of each TDNN layer into the product of two low-rank matrices. Densely connected TDNN (D-TDNN) [4] adopts bottleneck layers and dense connectivity, and it can further integrate channel-wise selection mechanism to achieve the state-of-the-art accuracy for TDNN-based speaker verification. The accuracy of TDNN based speaker verification can be further improved leveraging the advanced neural architecture search (NAS).

On the other hand, studies have been extensively conducted to enhance the loss functions [6]. Triplet loss [7] selects appropriate training samples and performs well for speaker verification. Speaker identity subspace loss [8] learns a space where it measures the similarity of speakers by Eu-

clidean distance. Center loss [9, 10] attempts to reduce the intra-speaker variability by constraining features close to the center. Instead, angular distance [11] focuses on cosine similarity and normalizes the features and the weights of the output layer before softmax. Generalized end-to-end loss [12] minimizes a scaled cosine score between the features and the estimated speaker centers. Various angular and large margin based loss functions [13, 14, 15, 16] have been investigated in the speaker verification. Liu *et al.* [6] conducts a comprehensive investigation on the loss functions and combines an additive margin softmax loss and a minimum hyper-spherical energy (MHE) [17] to achieve a desired accuracy. We also employ an advanced large margin based loss to train the candidate architectures for the large scale speaker verification.

In this work, we attempt to construct optimal TDNN-based architectures leveraging the recent advanced study of network designs and loss functions. Our search space consists of multi-branch mechanism, densely connected TDNN (D-TDNN), and channel selection. We conduct the number of branches search, the dimension number of D-TDNN search and the dimension number of channel selection search employing Bayesian optimization. The candidate architecture is trained by a hybrid loss between AAM and a MHE criterion.

Our contributions can be summarized as follows: 1) We design a neural architecture search (NAS) based large scale speaker verification system, SpeechNAS, to identify the optimal architectures leveraging TDNN variants and an AAM related hybrid loss. 2) We conduct a comprehensive comparison to the state-of-the-art speaker verification approaches considering a variety of metrics, including the number of parameters, GFLOPs, latency, the equal error rate (EER) [1], and the minimum of detection cost function (DCF) [1] with target probabilities set to 0.01 and 0.001. 3) Our NAS based speaker verification derives five architectures, from SpeechNAS-1 to SpeechNAS-5, with various numbers of parameters and FLOPs. The SpeechNAS-5 achieves much better performance than previous TDNN based state-of-the-art approaches on the VoxCeleb1 test set as shown in Fig. 1.

2. RELATED WORKS

The recent deep learning based speaker verification approaches can be primarily categorized into two main aspects: advanced network structure constructions [1, 2, 3, 4, 18] and effective loss function designs [6, 19, 20, 21].

Recently, enormous neural nets have been elaborately designed for speaker recognition, such as TDNN [1], E-TDNN [2], F-TDNN [3], D-TDNN [4]. ECAPA-TDNN [18] employs Res2Net [22], SE block [23] and channel-dependent attention, which outperforms TDNN based methods. Zhou *et al.* [24] integrates the phonetic information into the attention based ResNet and improves the speaker verification accuracy significantly. Based on the syllables obtained from the pre-trained HMM models, the SCL [25] directly improves

the discriminative power of the learned frame-level features during training stage. DNN-SAT [26] investigates the use of embeddings for speaker-adaptive training with a small amount of adaptation data per speaker. LSTM [27, 28] can be employed and TDNN-LSTM [29, 27] trained with four different data augmentation methods outperforms the baselines of both i-vector [30] and x-vector [1].

Various loss functions have been studied for speaker verification. Wang *et al.* [20] jointly optimizes classification and clustering with a large margin softmax loss and a large margin Gaussian mixture loss. Logistic affinity loss [19] instead optimizes an end-to-end speaker verification model by building a learnable decision boundary to distinguish the similar pairs and dissimilar pairs. The quartet loss [21] explicitly computes a pair-wise distance loss in the embedding space and increases the gap between the similarity score distributions between the same class pairs and different class pairs. Self-adaptive soft voice activity detection (VAD) [31] incorporates a deep neural network based VAD into a deep speaker embedding system to reduce the domain mismatch. Jung *et al.* [32] applies a teacher-student learning framework to short utterance compensation.

There are few works for NAS based speaker verification. AutoSpeech [33] identifies the optimal operation combination in a neural cell and then derives a ConvNet by stacking the neural cell for multiple times. Auto-Vector [34] searches various choices of temporal context windows based on the x-vector and validates the performance on a private dataset. Concurrent to our work, EfficientTDNN [35] searches different depths, kernels, and widths by progressive once-for-all strategy [36] for speaker recognition in the wild. We construct the search space based on multi-branch, advanced D-TDNN block and channel-wise selection, and employ Bayesian optimization in the search.

3. METHOD

The overall framework of SpeechNAS is illustrated in Fig 2, which consists of supernet construction and training, Bayesian optimization search and candidate architecture retraining.

3.1. Search Space

Let \mathcal{A} be an NAS search space represented by a directed acyclic graph (DAG), and a sub-graph $a \in \mathcal{A}$ is a network architecture denoted as $\mathcal{N}(a, w)$ which is parameterized by weights w . The weight sharing strategy [37] encodes the search space \mathcal{A} in a supernet $\mathcal{N}(\mathcal{A}, W)$, and all the candidate architectures share the weights W of the supernet. Differentiable NAS [38] further relaxes the discrete search space \mathcal{A} to a continuous one $\mathcal{A}(\theta)$ and jointly optimizes the shared weights W and architecture distribution parameter θ . One shot NAS [37] decouples the supernet training and architecture search, which yields a better accuracy. We also conduct

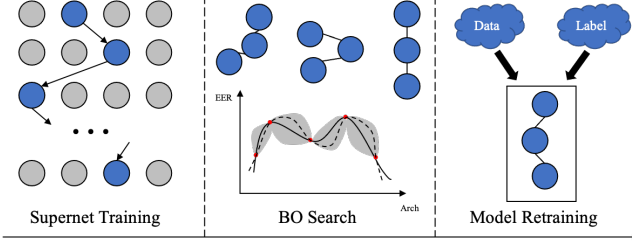


Fig. 2. The overview of SpeechNAS workflow.

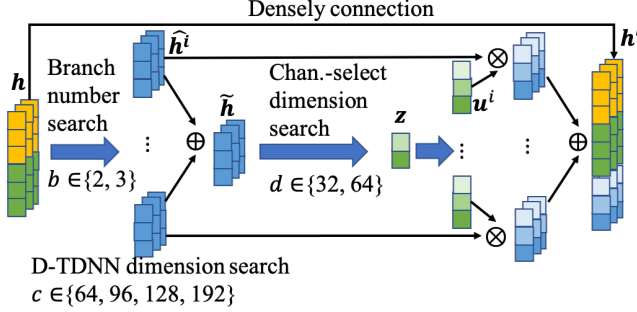


Fig. 3. The illustration of the SpeechNAS search space.

two sequential steps for supernet training and architecture search as

$$\begin{aligned} W_{\mathcal{A}} &= \operatorname{argmin}_W \mathcal{L}_{\text{train}}(\mathcal{N}(\mathcal{A}, W)), \\ a^* &= \operatorname{argmin}_{a \in \mathcal{A}} \text{EER}_{\text{val}}(\mathcal{N}(a, W_{\mathcal{A}(a)})), \end{aligned} \quad (1)$$

where $W_{\mathcal{A}}$ is the shared weights after supernet training and a^* is the searched optimal architecture with the best EER on the validation set.

Inspired by the superior performance of D-TDNN [4], we construct the supernet based on D-TDNN blocks, and define the search space \mathcal{A} consisting of the number of branches, the feature dimension of each D-TDNN block, and the dimension of channel-wise selection as illustrated in Fig. 3.

Let $b \in \mathcal{B}$ be the number of branches where \mathcal{B} is the search branch numbers. We conduct the D-TDNN dimension number search and use $\text{TDNN}_i(\cdot)$ to obtain the i -th branch feature for the t -th frame \hat{h}_t^i . Then we search the channel selection dimension numbers and construct a feed-forward net (FNN) layer f . We obtain the D-TDNN feature z

$$\begin{aligned} \tilde{h}_t &= \sum_{i=1}^b \hat{h}_t^i = \sum_{i=1}^b (\text{TDNN}_i(h))_t, \quad \mu = \frac{1}{T} \sum_{t=1}^T \tilde{h}_t, \\ \sigma &= \sqrt{\frac{1}{T} \sum_{t=1}^T \tilde{h}_t \odot \tilde{h}_t - \mu \odot \mu}, \quad s = \frac{1}{T} \sum_{t=1}^T \left(\frac{\tilde{h}_t - \mu}{\sigma} \right)^3, \\ k &= \frac{1}{T} \sum_{t=1}^T \left(\frac{\tilde{h}_t - \mu}{\sigma} \right)^4, \quad z = f([\mu, \sigma, s, k]), \end{aligned} \quad (2)$$

where h is the feature from the last block and T is the total number of frames. Then we construct a FNN layer f_i^t for the i -th branch. The feature h_t^i can be derived

$$t^i = f_i^t(z), \quad u^i = \text{Softmax}(t^i), \quad h_t^i = [h_t, \sum_{i=1}^b u^i \odot \hat{h}_t^i] \quad (3)$$

where the $\text{Softmax}(\cdot)$ is conducted along the branch dimension and $[\cdot, \cdot]$ is the concatenation of the two features.

Specifically, we conduct the number of branches b search and define two options, 2 or 3 branches. We expect the branches to learn varied features, and define the dilation rate [39] of (1, 3), (1, 3, 5) for the two branch options. Then we conduct D-TDNN feature dimension number c search and define four options of 64, 96, 128 and 192. Lastly, we also employ the channel-wise selection based on statistics-and-selection to enhance the feature representational power of SpeechNAS, and conduct the channel selection feature dimension number d search. We define two options of 32 and 64 for channel selection feature number search. For each block, there are $2 \times 4 \times 2 = 16$ different candidates. We stack the component into 18 layers to construct the base feature extractor. The search space consists of 16^{18} different candidate architectures, which requires an efficient supernet training strategy and an advanced search algorithm to find optimal architectures.

3.2. Loss Function

The large search space leads to slow training of the supernet. To accelerate the training speed of supernet, we employ two different loss functions for the supernet training and the candidate network retraining, respectively. For supernet training, we employ cross entropy loss directly as a proxy loss function to accelerate the supernet training

$$\mathcal{L}_{\text{train}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{w_{y_i}^T g_i}}{\sum_{j=1}^C e^{w_j^T g_i}}, \quad (4)$$

where N is the batch size in the stochastic gradient descent (SGD), y_i is the ground truth speaker for the sample, w_{y_i} is the related weights for speaker y_i in the last linear layer, g_i is the feature vector, and C is the total number of speakers. For the supernet training, we utilize single path a and uniform sampling $U(\mathcal{A})$ strategy to reduce the co-adaptation between node weights [37] as

$$W_{\mathcal{A}} = \operatorname{argmin}_W \mathbb{E}_{a \sim U(\mathcal{A})} [\mathcal{L}_{\text{train}}(\mathcal{N}(a, W(a)))]. \quad (5)$$

For the candidate architecture retraining, we employ an additive angular margin softmax (AAM) loss with a minimum hyper-spherical energy (MHE) criterion inspired by the comprehensive investigation of loss functions in the speaker

verification [6].

$$\mathcal{L}_{\text{retr}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s \cdot \cos(\theta_{y_i} + m)}}{e^{s \cdot \cos(\theta_{y_i} + m)} + \sum_{j=1, j \neq y_i}^C e^{s \cdot \cos \theta_j}} + \frac{\lambda}{N(C-1)} \sum_{i=1}^N \sum_{j=1, j \neq y_i}^C \frac{1}{\|\tilde{\mathbf{w}}_{y_i} - \tilde{\mathbf{w}}_j\|^2}, \quad (6)$$

where s is a scale factor, θ_j is the angle between \mathbf{w}_j and \mathbf{g}_i , m is the additive angular margin to enhance the discriminative power and robustness of feature extractor, λ is the trade-off between the AAM loss and MHE regularization, $\tilde{\mathbf{w}}_j$ is the L_2 normalized weights. MHE loss enlarges the inter-class feature separability.

3.3. Search Algorithm

Bayesian optimization (BO) [40] iterates between fitting probabilistic surrogate models and determining which configuration to evaluate next by maximizing an acquisition function. We employ a random exploration in the initialization. Gaussian process (GP) with a Hamming kernel k is utilized as the surrogate function, and the generative model of Gaussian process can be defined as

$$p|\mathbf{a} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}), \quad \mathbf{EER}|p, \sigma^2 \sim \mathcal{N}(p, \sigma^2 \mathbf{I}), \quad (7)$$

where variables \mathbf{p} are jointly Gaussian, \mathbf{a} are a set of observed architectures, $K_{ij} = k(a_i, a_j)$, and \mathbf{EER} are the evaluated EER values for these architectures with weight sharing. The parameters of GP, $\boldsymbol{\mu}$ and σ , can be estimated by maximizing the marginal log-likelihood. For the acquisition function α , we use probability of feasibility (PoF) [41].

The overall algorithm is described in Algorithm 1, which consists of supernet training, architecture search and candidate architecture retraining.

4. EXPERIMENT

4.1. Dataset

Standard training set For the standard training without augmentation, we follow the same dataset preparation in D-TDNN [4] for VoxCeleb1 [45] and VoxCeleb2 [46]. The two versions have their own explicit train and test splits and consist of 7,323 speakers with over one million utterances and more than 2,000 hours audio data. The training samples are generated by following the Kaldi toolkit [47] recipe. We use the whole VoxCeleb2 and the training set of VoxCeleb1 as our training set, and validate our method on the test set of VoxCeleb1.

For the preprocessing, we extract 30-dimensional Mel-frequency cepstrum coefficients (MFCCs) over a 25 ms long window every 10 ms. To remove silent frames, cepstral mean

Algorithm 1: The SpeechNAS algorithm

Input: Dataset $\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{val}} = \{(\mathbf{x}_i, y_i) | i = 1, \dots, n\}$, search space \mathcal{A} , the number of epochs n_1 and the number of candidates n_2 in BO, and hyper-parameters in the training

Output: Optimal architectures \mathbf{a}^* of low EERs on the validation set

```

/* Supernet training */
1 Construct a supernet based on the search space  $\mathcal{A}$ ;
2 Train the supernet using equation (5) and the loss function in equation (4) with SGD;
/* Architecture search */
3 Randomly explore  $n_2$  candidates  $\mathbf{a}_0$ ; // Initial
4 Evaluate EERs for  $\mathbf{a}_0$  with weight sharing;
5 Add  $\mathbf{a}_0$  and EERs into queue  $Q$ ;
6 Learn GP based on equation (7);
7 for  $i = 1, 2, \dots, n_1$  do
8   Select new architectures  $\mathbf{a}_i$  by optimizing acquisition function  $\alpha$ 
    $\mathbf{a}_i = \operatorname{argmax}_{\mathbf{a}} \alpha(\mathbf{a}_{i-1}; Q)$ ;
9   Evaluate EERs for  $\mathbf{a}_i$ ;
10  Append  $\mathbf{a}_i$  and the EERs to queue  $Q$ ;
11  Update GP based on equation (7) using data in  $Q$ ;
12 end
/* Candidate networks retraining */
13 for  $a$  in  $\mathbf{a}_{n_1}$  do
14   Train network  $a$  with SGD and equation (6);
15   Save the best trained model and evaluate the EER;
16 end
17 return Optimal architectures  $\mathbf{a}^*$  of low EERs and trained models;

```

normalization (CMN) is applied over a 3 seconds long sliding window and energy based VAD is utilized. We randomly split the spectrograms into 200 to 400 frames.

Generalization on augmented training set We further validate the *generalization ability* of our explored optimal network on an augmented dataset. We only augment the training set of VoxCeleb2 and follow the same procedure of ECAPA-TDNN [18] for the data augmentation, which generates a total of six extra samples for each utterance. We utilize the Kaldi recipe [47] in combination with the publicly available MUSAN dataset (babble and noise) [48] and the RIR dataset (reverb) [49] for the first three augmentations. We generate the remaining three augmentations using the open-source SoX [50] (tempo up and tempo down) and FFmpeg (alternating opus or aac compression) libraries.

We extract 80-dimensional MFCCs from a 25 ms window with a 10 ms frame shift as the input features on the augmented training set. We utilize the CMN to normalize two second random crops of the MFCCs feature vectors. We

Table 1. The structure of searched optimal networks.

Layer	SpeechNAS-3		SpeechNAS-4		SpeechNAS-5		
	<i>c</i>	<i>d</i>	<i>b</i>	<i>c</i>	<i>b</i>	<i>c</i>	<i>d</i>
1	96	64	3	128	3	128	64
2	128	32	3	64	3	192	32
3	96	32	3	128	3	192	64
4	96	32	3	192	3	192	64
5	96	32	2	192	2	192	32
6	64	64	2	64	2	192	64
7	64	64	3	64	3	192	32
8	96	64	2	64	3	128	32
9	64	64	3	64	3	192	64
10	64	32	3	64	2	128	64
11	64	32	3	128	3	128	64
12	96	64	3	192	2	192	64
13	96	64	3	64	2	192	32
14	96	32	2	192	2	128	64
15	128	32	2	128	2	192	32
16	128	32	3	128	2	192	64
17	96	64	2	128	3	192	32
18	96	64	3	192	2	192	64

apply SpecAugment [51] on the log Mel spectrogram of the samples for the augmentation. In the frequency domain, we randomly mask zero to five frames in the time domain and zero to ten channels in the frequency domain.

4.2. Implementation Detail

We implement the whole framework based on PyTorch. For the training, we use SGD with momentum of 0.95 and the initial learning rate of 0.01. The weight decay is set to 5×10^{-4} , and batch size is set to 128. We use 26 epochs for both the supernet training and candidate network retraining. The learning rate is adapted to 0.001 at the epoch of 14 and 0.0001 at the epoch of 20.

The scaling factor s in the AAM loss is set to be 30 and the margin m is set to be 0.2. The strength of the MHE criterion λ is set to be 0.01. These hyper-parameters are set basically according to those in D-TDNN [4] firstly and tuned a little due to the variant in our method and re-implementation. We use the reported performance directly from published papers for previous state-of-the-art approaches in Table 2. We use the public implementations¹ to calculate the latency, GFLOPs and the number of parameters (M) of TDNN and F-TDNN. We implement E-TDNN and D-TDNN, and reproduce the results of E-TDNN and D-TDNN. For D-TDNN-SS, we achieve an EER of 1.24% compared to reported 1.22%.

In the search phase, we set n_1 to be 100 and n_2 to be 64. To train a well-performed Gaussian process model, we randomly evaluate 1,200 architectures in the initialization step.

¹<https://github.com/cvqluu/TDNN>; <https://github.com/cvqluu/Factorized-TDNN>

For both the supernet training and candidate architecture re-training, we use one single NVIDIA GEFORCE RTX 2080 Ti graphics card. We use eight 2080 Ti graphics cards to accelerate the search phase. We utilize the OpenBox library [52] to implement the Bayesian optimization in the search.

4.3. Result

We conduct a comprehensive investigation of our SpeechNAS and adopt six metrics for evaluation, including the number of parameters (M), GFLOPs, Latency (ms), the EER with cosine similarity scoring [2] and the minimum of detection cost function (DCF) [2] with target probabilities set to 0.01 and 0.001.

The structure of searched optimal networks are shown in Table 1. From Table 1, we observe that the first four layers of all the three searched networks have large numbers of parameters, probably because the large first few layers encourage various feature exploration. Our SpeechNAS automatically design complicated network structures which are difficult to be manually designed.

The effect of search spaces We try three different search space design strategies, **(1)** both the D-TDNN feature and channel selection dimension number search $\{\{2\}, \{64, 96, 128\}, \{32, 64\}\}$, *i.e.*, without the number of branches search, **(2)** the number of branches search and the D-TDNN feature dimension number search $\{\{2, 3\}, \{64, 128, 192\}, \{32\}\}$, *i.e.*, without the channel selection feature dimension number search, **(3)** the full search space $\{\{2, 3\}, \{128, 192\}, \{32, 64\}\}$, *i.e.*, the number of branches search and the feature dimension number search for each D-TDNN block and channel selection. For **search space (1)**, we obtain three optimal candidate architectures with various numbers of parameters, GFLOPs, latency and EERs, named SpeechNAS-1, SpeechNAS-2, and SpeechNAS-3. For **search space (2)**, we obtain SpeechNAS-4. For the full **search space (3)**, we obtain SpeechNAS-5. The comprehensive performance comparison is listed in Table 2.

From Table 2, SpeechNAS-5 yields a better EER than SpeechNAS-4 and SpeechNAS-3, which shows that the full **search space (3)** offers the best EER. SpeechNAS-4 yields a better EER than SpeechNAS-3, which demonstrates the accuracy gain from search on branch numbers is larger than that from DTDNN feature dimensions in the SpeechNAS. For the model complexity, the number of parameters in SpeechNAS-4 is a little larger than SpeechNAS-3, and the number of parameters in SpeechNAS-5 is the largest. Although the full **search space (3)** offers the best EER, the model complexity of obtained architecture SpeechNAS-5 is the largest. From the latency perspective, the SpeechNAS-5 only increases a little compared to SpeechNAS-3.

The effect of loss functions For the supernet training, we investigate using cross entropy in equation (4) and AAM with MHE in equation (6) in Table 3. Because the AAM with

Table 2. Comparisons to state-of-the-art approaches on the `VoxCeleb1` test set. The **bold** font denotes the best result. \star denotes training using augmented training set. Latency is measured by averaging the inference time running 1,000 times on one NVIDIA GEFORCE RTX 2080 Ti graphics card with batch size of 128. \ast denotes that batch size of 24 is used for TDNN because of GPU memory size limitation.

Model	Embedding Size	Parameters (M)	GFLOPs	Latency (ms)	Cosine EER (%)	DCF _{0.01}	DCF _{0.001}
Dual Attention [42]	512	21.7	-	-	1.60	-	-
ARET-25 [43]	512	12.2	2.9	-	1.39	0.20	-
Fast ResNet-34 [44]	-	1.4	0.45	-	2.22	-	-
TDNN [1]	512	4.2	5.34	146 \ast	5.20	0.44	0.60
E-TDNN [2]	512	6.1	0.91	52	4.65	0.43	0.53
F-TDNN [3]	512	12.4	2.29	115	4.66	0.41	0.57
D-TDNN [4]	512	2.8	-	-	1.81	0.20	0.28
D-TDNN-SS [4]	512	3.5	0.56	71	1.41	0.19	0.24
D-TDNN-SS [4]	128	3.1	0.55	70	1.22	0.13	0.20
AutoSpeech [33]	2048	18	-	-	8.95	-	-
Space (1): SpeechNAS-1	128	2.6	0.44	66	1.25	0.07	0.18
Space (1): SpeechNAS-2	128	2.9	0.49	68	1.22	0.07	0.17
Space (1): SpeechNAS-3	128	3.1	0.44	66	1.14	0.06	0.17
Space (2): SpeechNAS-4	128	3.3	0.60	62	1.11	0.06	0.24
Space (3): SpeechNAS-5	128	4.3	0.76	71	1.06	0.06	0.12
Space (3): SpeechNAS-5 \star	128	4.3	0.77	72	1.02	0.05	0.17

Table 3. Cosine EERs (%) of loss functions for supernets.

Model	Searched arch.s	Retrain
Space (2) w/ cross entropy	1.92 - 2.06	1.11
Space (2) w/ AAM + MHE	1.54 - 1.61	1.13
Space (3) w/ cross entropy	1.72 - 1.80	1.06
Space (3) w/ AAM + MHE	1.29 - 1.35	1.09

MHE induces a better EER, the EERs of searched architectures using AAM + MHE without retraining are better than supernet trained with cross entropy. However, the retraining cannot yield better EERs for supernet training with AAM + MHE, probably because cross entropy loss enables a more adequate training for weights in the supernet which leads to a more accurate performance estimator than AAM + MHE in the weight sharing. Because training time of supernet with AAM + MHE typically is more than twice than that of cross entropy loss, we use cross entropy loss to train the supernet.

Comparison to state-of-the-art approaches We conduct a comprehensive comparison of SpeechNAS to previous TDNN based state-of-the-art methods in Table 2. From Table 2 and Fig. 1, our SpeechNAS achieves much better accuracy with a low latency which is friendly to be deployed to data intensive applications. Retraining the SpeechNAS-5 on an augmented dataset yields further EER improvement, which demonstrates a good generalization of SpeechNAS. We notice that ECAPA-TDNN [18] (embedding size 512) of 6.2M parameters, which is based on Res2Net [22] and SE-ResNet [23], achieves 1.01% EER. Our SpeechNAS is based

on TDNN framework and achieves a comparable accuracy with much less number of parameters of 4.3M.

5. CONCLUSION

In this work, we introduce a new neural architecture search based speaker verification framework, called SpeechNAS. Specifically, we investigate the optimal TDNN based network architectures leveraging the recent advances of the multiple branches design, D-TDNN block, channel-wise selection and the additive angular margin softmax loss (AAM) with a minimum hyper-spherical energy (MHE) criterion in the signal processing. We construct a supernet based on the designed search space and employ single path and uniform sampling with cross entropy loss to efficiently train the supernet. For the search, we utilize the Bayesian optimization with Gaussian process to find the best performed candidate architectures. The candidate architectures are retrained with a hybrid loss of AAM and MHE. Ablation studies and comprehensive experimental results on a large-scale speaker verification dataset, `VoxCeleb1`, demonstrate the effectiveness of each component of our SpeechNAS and that our SpeechNAS achieves a much better accuracy than other TDNN based state-of-the-art variants with a reasonable model complexity. Future NAS related work can be conducted to further improve the speaker verification accuracy by designing advanced search spaces with more effective components, blocks and/or attention mechanisms.

6. REFERENCES

- [1] David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev 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.
- [2] David Snyder, Daniel Garcia-Romero, Gregory Sell, Alan McCree, Daniel Povey, and Sanjeev Khudanpur, “Speaker recognition for multi-speaker conversations using x-vectors,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 5796–5800.
- [3] Jesús Villalba, Nanxin Chen, David Snyder, Daniel Garcia-Romero, Alan McCree, Gregory Sell, Jonas Borgstrom, Fred Richardson, Suwon Shon, François Grondin, et al., “State-of-the-art speaker recognition for telephone and video speech: The jhu-mit submission for nist sre18.,” in *Proc. Interspeech 2019*, 2019, pp. 1488–1492.
- [4] Ya-Qi Yu and Wu-Jun Li, “Densely connected time delay neural network for speaker verification,” in *Proc. Interspeech 2020*, 2020, pp. 921–925.
- [5] Frédéric Bimbot, Jean-François Bonastre, Corinne Fredouille, Guillaume Gravier, Ivan Magrin-Chagnolleau, Sylvain Meignier, Teva Merlin, Javier Ortega-García, Dijana Petrovska-Delacretaz, and Douglas A Reynolds, “A tutorial on text-independent speaker verification,” *EURASIP Journal on Advances in Signal Processing*, vol. 2004, no. 4, pp. 1–22, 2004.
- [6] Yi Liu, Liang He, and Jia Liu, “Large margin softmax loss for speaker verification,” in *Proc. Interspeech 2019*, 2019, pp. 2873–2877.
- [7] Chunlei Zhang and Kazuhito Koishida, “End-to-end text-independent speaker verification with triplet loss on short utterances.,” in *Interspeech*, 2017, pp. 1487–1491.
- [8] Ruifang Ji, Xinyuan Cai, and Bo Xu, “An end-to-end text-independent speaker identification system on short utterances.,” in *Interspeech*, 2018, pp. 3628–3632.
- [9] Na Li, Deyi Tuo, Dan Su, Zhifeng Li, Dong Yu, and A Tencent, “Deep discriminative embeddings for duration robust speaker verification.,” in *Interspeech*, 2018, pp. 2262–2266.
- [10] Sarthak Yadav and Atul Rai, “Learning discriminative features for speaker identification and verification.,” in *Interspeech*, 2018, pp. 2237–2241.
- [11] Chunlei Zhang, Fahimeh Bahmaninezhad, Shivesh Ranjan, Harishchandra Dubey, Wei Xia, and John HL Hansen, “Utdcrss systems for 2018 nist speaker recognition evaluation,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 5776–5780.
- [12] Li Wan, Quan Wang, Alan Papir, and Ignacio Lopez Moreno, “Generalized end-to-end loss for speaker verification,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4879–4883.
- [13] Zili Huang, Shuai Wang, and Kai Yu, “Angular softmax for short-duration text-independent speaker verification,” *Proc. Interspeech 2018*, pp. 3623–3627, 2018.
- [14] Gautam Bhattacharya, Jahangir Alam, and Patrick Kenny, “Adapting end-to-end neural speaker verification to new languages and recording conditions with adversarial training,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6041–6045.
- [15] Weidi Xie, Arsha Nagrani, Joon Son Chung, and Andrew Senior, “Utterance-level aggregation for speaker recognition in the wild,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 5791–5795.
- [16] Rongjin Li, Na Li, Deyi Tuo, Meng Yu, Dan Su, and Dong Yu, “Boundary discriminative large margin cosine loss for text-independent speaker verification,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6321–6325.
- [17] Weiyang Liu, Rongmei Lin, Zhen Liu, Lixin Liu, Zhiding Yu, Bo Dai, and Le Song, “Learning towards minimum hyperspherical energy,” in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 6225–6236.
- [18] Brecht Desplanques, Jenthe Thienpondt, and Kris Demuynck, “Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification,” in *Proc. Interspeech 2020*, 2020, pp. 3830–3834.
- [19] Junyi Peng, Rongzhi Gu, and Yuexian Zou, “Logistic similarity metric learning via affinity matrix for text-independent speaker verification,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 704–709.
- [20] Zhiming Wang, Kaisheng Yao, Shuo Fang, and Xiaolong Li, “Joint optimization of classification and clustering for deep speaker embedding,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 284–290.
- [21] Hira Dharmyal, Tianyan Zhou, Bhiksha Raj, and Rita Singh, “Optimizing neural network embeddings using a pair-wise loss for text-independent speaker verification,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 742–748.
- [22] Shanghua Gao, Ming-Ming Cheng, Kai Zhao, Xin-Yu Zhang, Ming-Hsuan Yang, and Philip HS Torr, “Res2net: A new multi-scale backbone architecture,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [23] Jie Hu, Li Shen, and Gang Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7132–7141.
- [24] Tianyan Zhou, Yong Zhao, Jinyu Li, Yifan Gong, and Jian Wu, “Cnn with phonetic attention for text-independent speaker verification,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 718–725.

- [25] Junyi Peng, Yuexian Zou, Na Li, Deyi Tuo, Dan Su, Meng Yu, Chunlei Zhang, and Dong Yu, "Syllable-dependent discriminative learning for small footprint text-dependent speaker verification," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 350–357.
- [26] Joanna Rownicka, Peter Bell, and Steve Renals, "Embeddings for dnn speaker adaptive training," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 479–486.
- [27] Wentao Zhu, Cuiling Lan, Junliang Xing, Wenjun Zeng, Yanghao Li, Li Shen, and Xiaohui Xie, "Co-occurrence feature learning for skeleton based action recognition using regularized deep lstm networks," in *Proceedings of the AAAI conference on artificial intelligence*, 2016, vol. 30.
- [28] Yufang Huang, Wentao Zhu, Deyi Xiong, Yiye Zhang, Changjian Hu, and Feiyu Xu, "Cycle-consistent adversarial autoencoders for unsupervised text style transfer," in *Proceedings of the 28th International Conference on Computational Linguistics*, 2020, pp. 2213–2223.
- [29] Chien-Lin Huang, "Exploring effective data augmentation with tdnn-lstm neural network embedding for speaker recognition," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 291–295.
- [30] 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, 2010.
- [31] Youngmoon Jung, Yeunju Choi, and Hoirin Kim, "Self-adaptive soft voice activity detection using deep neural networks for robust speaker verification," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 365–372.
- [32] Jee-weon Jung, Hee-Soo Heo, Hye-jin Shim, and Ha-Jin Yu, "Short utterance compensation in speaker verification via cosine-based teacher-student learning of speaker embeddings," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 335–341.
- [33] Shaojin Ding, Tianlong Chen, Xinyu Gong, Weiwei Zha, and Zhangyang Wang, "Autospeech: Neural architecture search for speaker recognition," in *Proc. Interspeech 2020*, 2020, pp. 916–920.
- [34] Xiaoyang Qu, Jianzong Wang, and Jing Xiao, "Evolutionary algorithm enhanced neural architecture search for text-independent speaker verification," *Proc. Interspeech 2020*, pp. 961–965, 2020.
- [35] Rui Wang, Zhihua Wei, Haoran Duan, Shouling Ji, and Zhen Hong, "Efficienttdnn: Efficient architecture search for speaker recognition in the wild," *arXiv preprint arXiv:2103.13581*, 2021.
- [36] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han, "Once-for-all: Train one network and specialize it for efficient deployment," in *International Conference on Learning Representations*, 2019.
- [37] Zichao Guo, Xiangyu Zhang, Haoyuan Mu, Wen Heng, Zechun Liu, Yichen Wei, and Jian Sun, "Single path one-shot neural architecture search with uniform sampling," in *European Conference on Computer Vision*. Springer, 2020, pp. 544–560.
- [38] Hanxiao Liu, Karen Simonyan, and Yiming Yang, "Darts: Differentiable architecture search," in *International Conference on Learning Representations*, 2018.
- [39] Fisher Yu, Vladlen Koltun, and Thomas Funkhouser, "Dilated residual networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 472–480.
- [40] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas, "Taking the human out of the loop: A review of bayesian optimization," *Proceedings of the IEEE*, vol. 104, no. 1, pp. 148–175, 2015.
- [41] Jacob R Gardner, Matt J Kusner, Zhixiang Eddie Xu, Kilian Q Weinberger, and John P Cunningham, "Bayesian optimization with inequality constraints," in *ICML*, 2014, vol. 2014, pp. 937–945.
- [42] Jingyu Li and Tan Lee, "Text-independent speaker verification with dual attention network," in *Proc. Interspeech 2020*, 2020, pp. 956–960.
- [43] Ruiteng Zhang, Jianguo Wei, Wenhuan Lu, Longbiao Wang, Meng Liu, Lin Zhang, Jiayu Jin, and Junhai Xu, "Aret: Aggregated residual extended time-delay neural networks for speaker verification," in *Proc. Interspeech 2020*, 2020, pp. 946–950.
- [44] Joon Son Chung, Jaesung Huh, Seongkyu Mun, Minjae Lee, Hee-Soo Heo, Soyeon Choe, Chiheon Ham, Sunghwan Jung, Bong-Jin Lee, and Icksang Han, "In defence of metric learning for speaker recognition," in *Proc. Interspeech 2020*, 2020, pp. 2977–2981.
- [45] Arsha Nagrani, Joon Son Chung, and Andrew Zisserman, "Voxceleb: A large-scale speaker identification dataset," in *Proc. Interspeech 2017*, 2017, pp. 2616–2620.
- [46] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman, "Voxceleb2: Deep speaker recognition," in *Proc. Interspeech 2018*, 2018, pp. 1086–1090.
- [47] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, et al., "The kaldi speech recognition toolkit," in *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding (ASRU)*. IEEE Signal Processing Society, 2011.
- [48] David Snyder, Guoguo Chen, and Daniel Povey, "Musan: A music, speech, and noise corpus," *arXiv preprint arXiv:1510.08484*, 2015.
- [49] Tom Ko, Vijayaditya Peddinti, Daniel Povey, Michael L Seltzer, and Sanjeev Khudanpur, "A study on data augmentation of reverberant speech for robust speech recognition," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 5220–5224.
- [50] Benjamin Barras, "Sox: Sound exchange," Tech. Rep., 2012.
- [51] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le, "SpecAugment: A simple data augmentation method for automatic speech recognition," in *Proc. Interspeech 2019*, 2019, pp. 2613–2617.
- [52] Yang Li et al., "Openbox: A generalized black-box optimization service," in *KDD-2021*, 2021.