# Reliable Visualization for Deep Speaker Recognition

*Pengqi Li[1], Lantian Li[2], Askar Hamdulla[1], Dong Wang[2]*

[1]Information Science and Engineering Institute, Xinjiang University
[2]Center for Speech and Language Technologies, Tsinghua University

{lipq,lilt}@cslt.org, askar@xju.edu.cn, wangdong99@mails.tsinghua.edu.cn

## Abstract

In spite of the impressive success of convolutional neural networks (CNNs) in speaker recognition, our understanding to CNNs' internal functions is still limited. A major obstacle is that some popular visualization tools are difficult to apply, for example those producing saliency maps. The reason is that speaker information does not show clear spatial patterns in the temporal-frequency space, which makes it hard to interpret the visualization results, and hence hard to confirm the reliability of a visualization tool.

In this paper, we conduct an extensive analysis on three popular visualization methods based on CAM: Grad-CAM, Score-CAM and Layer-CAM, to investigate their reliability for speaker recognition tasks. Experiments conducted on a state-of-the-art ResNet34SE model show that the Layer-CAM algorithm can produce reliable visualization, and thus can be used as a promising tool to explain CNN-based speaker models. The source code and examples are available in our project page: *http://project.cslt.org/*.

**Index Terms**: visual explanations, convolutional neural networks, speaker recognition, class activation maps

## 1. Introduction

Deep convolutional neural networks (CNNs) have attained remarkable success in computer vision [1, 2, 3]. Besides the unprecedented performance on a broad range of tasks, a special reason is that there are multiple visualization tools that can be used to explain the decisions of the model [4, 5, 6]. Some of the most representative visualization tools include guided backpropagation [7], deconvolution [8], CAM and its variants [6, 9, 10, 11]. These tools can produce saliency maps that identify the important regions in an image that lead to the model's decision. Importantly, humans can easily interpret a saliency map of an image, and hence judge the quality of a visualization tool.

Recently, CNN models have been widely adopted in speaker recognition and achieved fairly good performance [12]. However, how the models obtain such performance is hard to explain. A major obstacle is that the established visualization tools cannot be used directly. Basically, this is because humans cannot 'see' speech, which makes interpretation and quality judgement for saliency maps on speech signals quite difficult [13]. Recently, some researchers have recognized the difficulty, and provide some solutions [14, 13, 15]. Nearly all the research focused on visualizing phone classes.

So far very few research in literature reports visualization for speaker recognition. This is not surprising as speaker traits spread among nearly all frequency bands and temporal segments, making the question 'where is important' difficult to answer. In contrast, phone information is much more localized, at least in the temporal axis. Among the few exceptions, [16] employed Grad-CAM [17] to compare the behavior of ResNet and Res2Net under noisy corruption. They found that the saliency map produced by Grad-CAM is more stable with Res2Net compared with ResNet, thus explaining the advantage of Res2Net. In [18], the authors used Grad-CAM to analyze genuine speech and spoof speech, and found that CNN models look into high-frequency components to identify spoof speech. All the mentioned studies employed the visualization tools directly by assuming that they are correct. Unfortunately, *so far we have no idea if any of the visualization tools are reliable when applied to speaker recognition, which makes the conclusions obtained from visualization not fully convincing.*

In this paper, we focus on the popular CAM-based algorithms [6], and try to answer the question *if these algorithms, or any of them, can be reliably applied to speaker recognition tasks*. Three CAM algorithms will be investigated: Grad-CAM++ [9], Score-CAM [10] and Layer-CAM [11]. The main idea of these algorithms is to generate a saliency map by combing the activation maps (channels) of a convolutional layer.

Our investigation starts from a deletion and insertion experiment, as suggested in [19]. In the deletion process, the most relevant regions (MoRRs) in Mel spectrograms are gradually masked by setting the values to 0 according to a saliency map. In the insertion process, the MoRRs are gradually unmasked (exposed), starting from a totally-masked Mel spectrogram. For each insertion and deletion config, we examine the accuracy of a speaker recognition model. This results in a deletion curve and an insertion curve, by which we can tell if a saliency map really identifies the salient regions, and compare quality of different saliency maps, hence different visualization algorithms.

We conducted experiments with a ResNet34SE x-vector model. The results show that all the three CAM algorithms outperform random masking and time-aligned masking, demonstrating that they are effective. Among the three CAMs, Layer-CAM shows superiority in the deletion/insertion test, and produces more localized patterns. However, no clear temporal-frequency (T-F) patterns are detected.

Further more, we conducted the deletion/insertion experiment on multi-speaker speech. This time, Layer-CAM demonstrates surprisingly good performance in distinguishing target speakers and interfering speakers, and the other two CAMs largely fail. This clearly shows that only Layer-CAM is a valid visualization tool for speaker recognition.

## 2. Related work

Understanding the behavior of deep speaker recognition systems by visualization is a common practice. A popular approach

is to visualize the distributions of frame-level or utterance-level representations [20, 21, 22, 23, 24, 25], via manifold learning algorithms [26] such as PCA and t-SNE [27]. Some researchers conduct visualization based on scores of trials. For instance, [28] visualizes the relation of different systems by *multidimensional scaling* (MDS), where the distance between a pair of systems is derived from the scores produced by each of them. Recently, the authors proposed a novel config-performance (C-P) map tool that visualizes the performance of an ASV system in a 2-dimensional map [29]. All these visualization methods can (partly) show the behavior of a system, but provide little *understanding* for the behavior. CAM, or other saliency-map algorithms, is supposed to offer better explanations for a CNN model, hence better understanding.

The saliency-map algorithms can be categorized into three classes: gradient-based approach [4, 7, 8, 30, 31], perturbation-based approach [8, 19, 32, 33] and CAM-based approach [6, 9, 10, 11, 17]. The gradient-based approach is subject to low quality and noisy interference [34], and the perturbation-based approach usually needs additional regularizations [32] and is time-costing. In comparison, the CAM-based approach often creates more clear and localized saliency maps, thus being adopted widely by the computer vision community.

Considering the high quality and localization capability, this paper focuses on CAMs, with the goal of identifying a CAM algorithm that can be reliably used to visualize speaker recognition models. To the best of our knowledge, this is the first work towards this direction.

# 3. Methodology

A class activation map (CAM) is a saliency map that shows the important regions used by the CNN to identify a particular class. In this section, we will firstly revisit three CAM algorithms that we will experiment with, and then present the normalization process designed to attain suitable T-F masks.

## 3.1. Revisit CAMs

### 3.1.1. Grad-CAM and Grad-CAM++

We start from Grad-CAM [17]. Let $f$ denote the speaker classifier instantiated by a CNN, and $\theta$ represents its parameters. For a given input $x$ from class $c$, the prediction score (posterior probability) for the target class can be computed by a forward pass:

$$y^c = f_c(x; \theta). \tag{1}$$

Then for the $k$-th activation map (i.e., the $k$-th channel) $A^k$ of a convolutional layer, the gradient of $y^c$ with respect to $A_{ij}^k$ is computed and the values at all the locations are averaged to obtain the weight for $A^k$ with respect to class $c$:

$$w_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}, \tag{2}$$

where $Z$ is a constant corresponding to the number of points in the map. Grad-CAM then produces the saliency map for class $c$ by linearly combining $A^k$ with weight $w_k^c$ and followed by $relu(\cdot)$:

$$S^c = relu(\sum_k w_k^c \cdot A^k)). \tag{3}$$

Grad-CAM++ [9] is a derived version of Grad-CAM, where the weight $w_k^c$ for $A^k$ is computed by:

$$w_k^c = \frac{1}{Z} \sum_i \sum_j \alpha_{ij}^{kc} \cdot relu(\frac{\partial y^c}{\partial A_{ij}^k}). \tag{4}$$

Specifically, Grad-CAM++ focuses on positive gradients only, and weights $\frac{\partial y^c}{\partial A_{ij}^k}$ by $\alpha_{ij}^{kc}$. This change helps identify multiple occurrences of the same class.

### 3.1.2. Score-CAM

Gradients of a deep neural network can be noisy and vanished. Score-CAM is a gradient-free algorithm [10], which computes the weights $w_k^c$ for activation map $A^k$ by forward activation rather than backward gradient. Specifically, it firstly forwards $x$ through the CNN to generate activation map $\{A^k\}$, and then use $A^k$ to mask $x$:

$$\hat{x}_k = x \circ \{\text{Norm}(\text{Upsampling}(A^k))\} \tag{5}$$

where Upsampling($\cdot$) stretches $A^k$ to meet the size of $x$, and Norm($\cdot$) performs a min-max normalization. $\hat{x}_k$ is then passed through the CNN again, and the generated posterior probability $f_c(\hat{x}_k)$ is used as the weight $w_k^c$.

### 3.1.3. Layer-CAM

Layer-CAM [11] is gradient-based, and reweights $A^k$ in a point-wised way. The weight for activation map $A^k$ for class $c$ at location $(i, j)$ is defined as the gradient at that location:

$$w_{ij}^{kc} = relu(\frac{\partial y^c}{\partial A_{ij}^k}). \tag{6}$$

Note that only positive gradients are considered, as in Grad-CAM++. The saliency map is produced as follows:

$$S_{ij}^c = relu\{\sum_k w_{ij}^{kc} \cdot A_{ij}^k\}. \tag{7}$$

It was demonstrated that the point-wised weighting is beneficial to produce more fine-grained and localized saliency maps [11].

## 3.2. Saliency normalization

The values in the saliency map $S^c$ may vary in a large range, and the range could be very different with different CAMs. To make the maps comparable in visualization, and to use them for speaker localization (Ref. Section 4), we firstly re-arrange the saliency values into the interval [0,1] by min-max normalization:

$$\hat{S}^c = \frac{S^c - \min S^c}{\max S^c - \min S^c}. \tag{8}$$

Furthermore, a scale function is applied to redistribute the saliency values. According to [11], we choose $tanh$ as the scale function:

$$\hat{S}_n^c = tanh(\frac{\gamma * \hat{S}^c}{\max \hat{S}^c}), \tag{9}$$

where $\gamma$ is the scale coefficient which was set to be 5 in our experiments.

# 4. Experiments

In this section, we quantitatively evaluate the reliability of three CAM algorithms: Grad-CAM++, Score-CAM and Layer-CAM, using a well-trained deep speaker model.

## 4.1. Speaker model

The development set of VoxCeleb2 [35] was used to train the speaker model, which contains 5,994 speakers in total. No data augmentation was used. The structure of the model is ResNet34 with squeeze-and-excitation (SE) layers [36], shown in Table 1. The model was trained by the Adam optimizer, following the voxceleb/v2 recipe of the Sunine toolkit [1].

Table 1: *The topology of ResNet34SE model.*

| Layer | Module | Output |
|---|---|---|
| Input | – | $80 \times 200 \times 1$ |
| Conv2D | $3 \times 3 \times 32$, Stride 1 | $80 \times 200 \times 32$ |
| ResNetBlock1 | $\begin{bmatrix} 3 \times 3 \times 32 \\ 3 \times 3 \times 32 \\ \text{SE Layer} \end{bmatrix} \times 3$, Stride 1 | $80 \times 200 \times 32$ |
| ResNetBlock2 | $\begin{bmatrix} 3 \times 3 \times 64 \\ 3 \times 3 \times 64 \\ \text{SE Layer} \end{bmatrix} \times 4$, Stride 2 | $40 \times 100 \times 64$ |
| ResNetBlock3 | $\begin{bmatrix} 3 \times 3 \times 128 \\ 3 \times 3 \times 128 \\ \text{SE Layer} \end{bmatrix} \times 6$, Stride 2 | $20 \times 50 \times 128$ |
| ResNetBlock4 | $\begin{bmatrix} 3 \times 3 \times 256 \\ 3 \times 3 \times 256 \\ \text{SE Layer} \end{bmatrix} \times 3$, Stride 2 | $10 \times 25 \times 256$ |
| Pooling | TSP [37] | $20 \times 256$ |
| Flatten | – | 5120 |
| Dense | – | 256 |
| Dense | AM-Softmax [38] | 5994 |

## 4.2. Single-speaker experiment

In this section, we probe the behavior of the three CAM algorithms using single-speaker utterances, i.e., only a single target speaker exists in an utterance. Our purpose is to examine if the CAM algorithms can detect salient components in Mel spectrograms. We randomly choose 2,000 utterances of 200 speakers (10 utterances per speaker) from the training data to perform the evaluation. The saliency maps are generated from ResNet-Block4 (S4), by Grad-CAM++, Score-CAM and Layer-CAM respectively.

First of all, we qualitatively compare the saliency maps produced by the three CAM-based methods. Two examples are shown in Figure 1. It can be observed that all the saliency maps clearly separate speech and non-speech segments, demonstrating the basic capacity of the CAMs. Another observation is that Grad-CAM++ and Score-CAM tend to regard all the speech segments being important, while Layer-CAM produces more selective and localized patterns. Nevertheless, no clear T-F patterns are found in any of the saliency maps. More examples are provided in our project page[2].

We further conduct the deletion and insertion experiment suggested in [19]. In the deletion experiment, we monitor the Top-1 accuracy of the CNN model as more and more important regions of the input are masked, and in the insertion experiment, we monitor the Top-1 accuracy as more and more important regions are unmasked. The Area Under the Curve
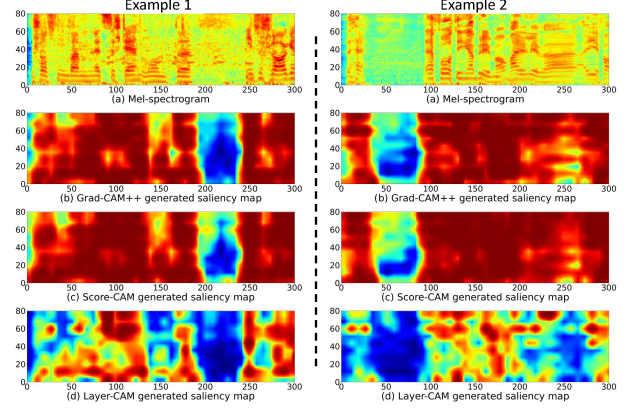


Figure 1: *The saliency maps of two speaker utterances generated by GradCAM++, Score-CAM and Layer-CAM. The deeper the color, the more important the region.*
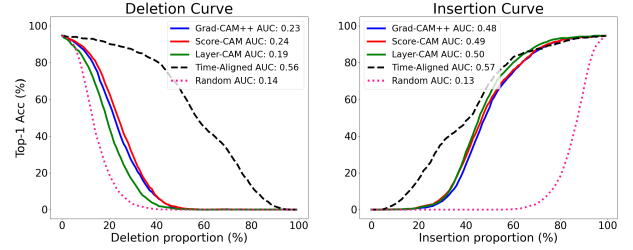


Figure 2: *Deletion and insertion curves of three CAM algorithms with single-speaker speech.*

(AUC) value is used to measure the quality of the saliency maps. More the saliency map is accurate, more the AUC lower in the deletion curve and higher in the insertion curve. Since we have known that all the CAM algorithms are capable of detecting non-speech segment, we focus on the speech segment in this experiment. For that purpose, voice activity detection (VAD) has been firstly employed to remove non-speech segments. Moreover, we show the performance with random (un)masking and left-to-right time-aligned (un)masking as reference (un)masking methods.

Figure 2 shows the results. It can be observed that the three CAMs are comparable in this deletion/insertion test, though Layer-CAM is slightly better. The comparison with random masking and time-aligned masking is also interesting: it shows that the three CAM algorithms indeed find salient regions. For example, in the insertion experiment, the curves of CAM algorithms clearly are much higher than that of the random masking, indicating that the regions exposed earlier by CAMs are indeed more important than random regions. And in the deletion experiment, the curves of CAM algorithms are much lower than that of the time-aligned masking, showing that the regions deleted in the early stage by CAMs are more important than real speech with the same amount of T-F bins. Note that the quick accuracy increase with time-aligned masking in the insertion curve is understandable, as it inserts entire frames which are valuable for speaker recognition when the utterance is short. Similarly, the quick accuracy drop with random masking in the deletion curve is not surprising, as random noise is very harmful for speaker recognition.

---

[1] https://gitlab.com/csltstu/sunine
[2] http://project.cslt.org/

Table 2: *Top-1 Acc (%) on target-speaker localization and recognition task with three CAMs. A represents target speaker; B and C represent non-target speakers. G-CAM: Grad-CAM++; S-CAM: Score-CAM; L-CAM: Layer-CAM. S1∼S4 represents ResNetBlock1∼ResNetBlock4. '+' denotes element-wised average of multiple saliency maps.*

| Cases | A-B | | | A-B-A | | | B-A-B | | | A-B-C | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Original | 49.15% | | | 83.55% | | | 15.35% | | | 30.30% | | |
| Settings | G-CAM | S-CAM | L-CAM | G-CAM | S-CAM | L-CAM | G-CAM | S-CAM | L-CAM | G-CAM | S-CAM | L-CAM |
| S1 | 43.00% | 34.00% | 6.75% | 75.55% | 62.50% | 8.90% | 15.15% | 12.10% | 4.15% | 22.40% | 17.45% | 3.90% |
| S2 | 46.60% | 46.60% | 61.85% | 79.90% | 79.25% | 85.00% | 15.70% | 16.00% | 35.20% | 26.85% | 26.85% | 45.15% |
| S3 | 48.45% | 48.60% | 49.40% | 82.65% | 82.35% | 80.15% | 15.75% | 16.00% | 20.20% | 29.20% | 29.55% | 31.05% |
| S4 | 49.20% | 48.25% | 53.20% | 82.10% | 82.65% | 82.90% | 17.20% | 16.15% | 24.05% | 30.10% | 29.20% | 34.65% |
| S4+S3 | 48.65% | 48.15% | 51.15% | 82.50% | 82.25% | 82.60% | 16.50% | 16.15% | 21.90% | 29.40% | 29.30% | 33.55% |
| S4+S3+S2 | 48.55% | 48.40% | 59.85% | 82.20% | 82.00% | 87.30% | 16.10% | 16.20% | 28.65% | 29.65% | 29.15% | 42.75% |
| S4+S3+S2+S1 | 47.70% | 47.50% | **71.55%** | 81.50% | 80.65% | **92.20%** | 16.10% | 16.10% | **44.60%** | 27.95% | 27.45% | **58.90%** |

### 4.3. Multi-speaker experiment

In the multi-speaker experiment, we concatenate an utterance of the target speaker with one or two utterances of other interfering speakers, and draw the saliency map. Figure 3 shows a 'B-A-B' test example. A denotes the target speaker while B denotes the interfering speaker. This time, Layer-CAM shows surprisingly good performance: it can accurately locate the segments of the target speaker, and mask non-target speakers almost perfectly. In comparison, Grad-CAM++ and Score-CAM are very weak in detecting non-target speakers. Moreover, Figure 4 shows the results of the deletion and insertion curves with 2,000 multi-speaker utterances in the B-A-B form. It can be seen that Layer-CAM gains much better AUCs than the other two CAMs.
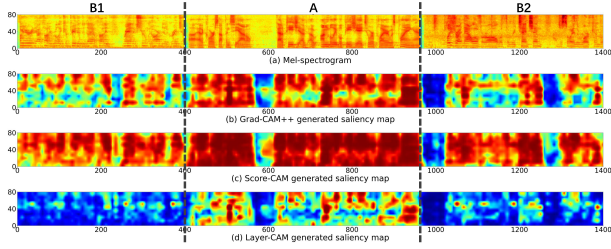


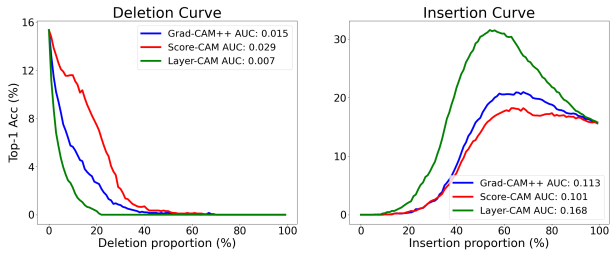Figure 3: *Saliency maps on a 'B-A-B' test example.*



Figure 4: *Deletion and insertion curves of three CAM algorithms with multi-speaker speech.*

### 4.4. Localization and recognition

Since Layer-CAM can localize target speakers, we can use it as a tool to perform localization and recognition, i.e., firstly identify where the target speaker resides and then perform speaker recognition with the located segments only. We assume this is better than using the entire utterance.

To test the hypothesis, we select 100 speakers from the training set, each with 20 utterances. We randomly concatenate the utterances in forms A-B, B-A-B, A-B-C, and A-B-A, where A denotes target speakers and B/C denotes interfering speakers. We use saliency maps produced at the layers of different ResNet blocks (S1-S4) to mask the input utterance by simple element-wised multiplication. Saliency maps of different layers can be combined as well. Top-1 Acc (%) are reported in Table 2.

Firstly, we observe that neither Grad-CAM++ nor Score-CAM can achieve performance better than the baseline (using the whole utterance). Layer-CAM, in contrast, delivers remarkable performance improvement, and this is the case for the saliency maps at all layers. This provides a very strong evidence that Layer-CAM can identify the important speaker-discriminative regions, while the other two algorithms cannot. This furthermore suggests that Layer-CAM is the only valid visualization tool among the three variants.

Secondly, we find that although saliency maps at all layers produced by Layer-CAM are informative, the one from S2 seems the most discriminative. More investigation is required here, but one possibility is that the saliency map of S2 is more conservative and retains more regions when compared to the ones obtained from higher layers.

Finally, we find that for Layer-CAM, aggregating saliency maps from different layers can improve performance. This observation seems consistent with the feature aggregation technique [39, 40]. Note that there is no such a trend with the other two CAMs. If we believe features are truly complimentary, then it is another evidence that only Layer-CAM can be used as a reliable visualization tool.

## 5. Conclusion

The ultimate goal of our study is to identify a reliable visualization tool. We focused on three CAM-based algorithms: Grad-CAM, Score-CAM and Layer-CAM. Experiments conducted with a state-of-the-art ResNet34SE model showed that although all the three algorithms can identify important regions in single-speaker utterances, Layer-CAM can localize target speakers in multi-speaker utterances. We therefore conclude that Layer-CAM is a reliable visualization tool for speaker recognition, and is the only one among the three CAM variants. In future, we will use the same protocol to test other visualization tools. Furthermore, the localization and recognition experiments conducted here suggests that integrating saliency maps may improve speaker recognition. This deserves more investigation.

# 6. References

[1] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," in *International conference on machine learning*, 2009, pp. 609–616.

[2] J. Murphy, "An overview of convolutional neural network architectures for deep learning," *Microway Inc*, 2016.

[3] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai *et al.*, "Recent advances in convolutional neural networks," *Pattern Recognition*, vol. 77, pp. 354–377, 2018.

[4] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[5] M. T. Ribeiro, S. Singh, and C. Guestrin, "" why should i trust you?" explaining the predictions of any classifier," in *ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.

[6] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *IEEE conference on computer vision and pattern recognition*, 2016, pp. 2921–2929.

[7] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," *arXiv preprint arXiv:1412.6806*, 2014.

[8] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.

[9] A. Chattopadhay, A. Sarkar, P. Howlader, and V. N. Balasubramanian, "Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks," in *2018 IEEE WACV*. IEEE, 2018, pp. 839–847.

[10] H. Wang, Z. Wang, M. Du, F. Yang, Z. Zhang, S. Ding, P. Mardziel, and X. Hu, "Score-cam: Score-weighted visual explanations for convolutional neural networks," in *IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 24–25.

[11] P.-T. Jiang, C.-B. Zhang, Q. Hou, M.-M. Cheng, and Y. Wei, "Layercam: Exploring hierarchical class activation maps for localization," *IEEE Transactions on Image Processing*, vol. 30, pp. 5875–5888, 2021.

[12] Z. Bai and X.-L. Zhang, "Speaker recognition based on deep learning: An overview," *Neural Networks*, vol. 140, pp. 65–99, 2021.

[13] A. Krug and S. Stober, "Introspection for convolutional automatic speech recognition," in *2018 EMNLP Workshop*, 2018, pp. 187–199.

[14] V. A. Trinh, "Identifying, evaluating and applying importance maps for speech," Ph.D. dissertation, City University of New York, 2022.

[15] V. A. Trinh and M. I. Mandel, "Large scale evaluation of importance maps in automatic speech recognition," *arXiv preprint arXiv:2005.10929*, 2020.

[16] T. Zhou, Y. Zhao, and J. Wu, "Resnext and res2net structures for speaker verification," in *2021 IEEE SLT*. IEEE, 2021, pp. 301–307.

[17] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *IEEE international conference on computer vision*, 2017, pp. 618–626.

[18] I. Himawan, S. Madikeri, P. Motlicek, M. Cernak, S. Sridharan, and C. Fookes, "Voice presentation attack detection using convolutional neural networks," in *Handbook of Biometric Anti-Spoofing*. Springer, 2019, pp. 391–415.

[19] V. Petsiuk, A. Das, and K. Saenko, "Rise: Randomized input sampling for explanation of black-box models," *arXiv preprint arXiv:1806.07421*, 2018.

[20] L. Li, Y. Chen, Y. Shi, Z. Tang, and D. Wang, "Deep speaker feature learning for text-independent speaker verification," *arXiv preprint arXiv:1705.03670*, 2017.

[21] S. Shon, H. Tang, and J. Glass, "Frame-level speaker embeddings for text-independent speaker recognition and analysis of end-to-end model," in *2018 SLT*. IEEE, 2018, pp. 1007–1013.

[22] Y. Jia, Y. Zhang, R. Weiss, Q. Wang, J. Shen, F. Ren, P. Nguyen, R. Pang, I. Lopez Moreno, Y. Wu *et al.*, "Transfer learning from speaker verification to multispeaker text-to-speech synthesis," *Advances in neural information processing systems*, vol. 31, 2018.

[23] L. Xu, R. K. Das, E. Yılmaz, J. Yang, and H. Li, "Generative x-vectors for text-independent speaker verification," in *2018 IEEE SLT*. IEEE, 2018, pp. 1014–1020.

[24] G. Bhattacharya, M. J. Alam, and P. Kenny, "Deep speaker recognition: Modular or monolithic?" in *INTERSPEECH*, 2019, pp. 1143–1147.

[25] Y. Shi, Q. Huang, and T. Hain, "H-vectors: Utterance-level speaker embedding using a hierarchical attention model," in *2020 IEEE ICASSP*. IEEE, 2020, pp. 7579–7583.

[26] K. Q. Weinberger and L. K. Saul, "Unsupervised learning of image manifolds by semidefinite programming," *International journal of computer vision*, vol. 70, no. 1, pp. 77–90, 2006.

[27] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of machine learning research*, vol. 9, no. 11, 2008.

[28] T. Kinnunen, A. Nautsch, M. Sahidullah, N. Evans, X. Wang, M. Todisco, H. Delgado, J. Yamagishi, and K. A. Lee, "Visualizing classifier adjacency relations: A case study in speaker verification and voice anti-spoofing," *arXiv preprint arXiv:2106.06362*, 2021.

[29] L. Li, D. Wang, W. Du, and D. Wang, "Cp map: A novel evaluation toolkit for speaker verification," *arXiv preprint arXiv:2203.02942*, 2022.

[30] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *International conference on machine learning*. PMLR, 2017, pp. 3319–3328.

[31] J. Adebayo, J. Gilmer, I. Goodfellow, and B. Kim, "Local explanation methods for deep neural networks lack sensitivity to parameter values," *arXiv preprint arXiv:1810.03307*, 2018.

[32] R. C. Fong and A. Vedaldi, "Interpretable explanations of black boxes by meaningful perturbation," in *IEEE international conference on computer vision*, 2017, pp. 3429–3437.

[33] C.-H. Chang, E. Creager, A. Goldenberg, and D. Duvenaud, "Explaining image classifiers by counterfactual generation," *arXiv preprint arXiv:1807.08024*, 2018.

[34] D. Omeiza, S. Speakman, C. Cintas, and K. Weldermariam, "Smooth grad-cam++: An enhanced inference level visualization technique for deep convolutional neural network models," *arXiv preprint arXiv:1908.01224*, 2019.

[35] A. Nagrani, J. S. Chung, W. Xie, and A. Zisserman, "Voxceleb: Large-scale speaker verification in the wild," *Computer Speech & Language*, vol. 60, p. 101027, 2020.

[36] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.

[37] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust dnn embeddings for speaker recognition," in *2018 IEEE ICASSP*. IEEE, 2018, pp. 5329–5333.

[38] F. Wang, J. Cheng, W. Liu, and H. Liu, "Additive margin softmax for face verification," *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 926–930, 2018.

[39] B. Desplanques, J. Thienpondt, and K. Demuynck, "Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification," *arXiv preprint arXiv:2005.07143*, 2020.

[40] Y. Tu and M.-W. Mak, "Aggregating frame-level information in the spectral domain with self-attention for speaker embedding," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2022.