# 腾讯说话人识别反欺骗技术进展与应用探索

## Shanshan Zhang

参与人：Xinyue Ma, Ji Gao, Dan Su, Sheng Huang

2020/11/21

# 说话人Anti-Spoofing问题

2019年《网络音视频信息管理规定》指出，明确禁止任何组织和个人滥用AI技术制作发布传播虚假新闻信息。

# ASVspoof 比赛介绍

- logical access (LA)：TTS/VC
- physical access (PA)：Replay

# TTS vs VC

# 真假语音对比



真实语音

合成语音

# 音频特点分析



真实语音
- 细节丰富
- 韵律自然



伪造语音
- 细节不自然—高频部分
- 韵律相对单一

# Feature extraction



MFCC/LFCC/IMFCC特征在ASV2019 eval测试集上的性能对比



Mel Filterbank

Linear Filterbank

iMel Filterbank

Frequency (kHz) →

# *Modeling*

- Light CNN with Max-Feature-Map(MFM) activation



Table 1: *LCNN architecture*

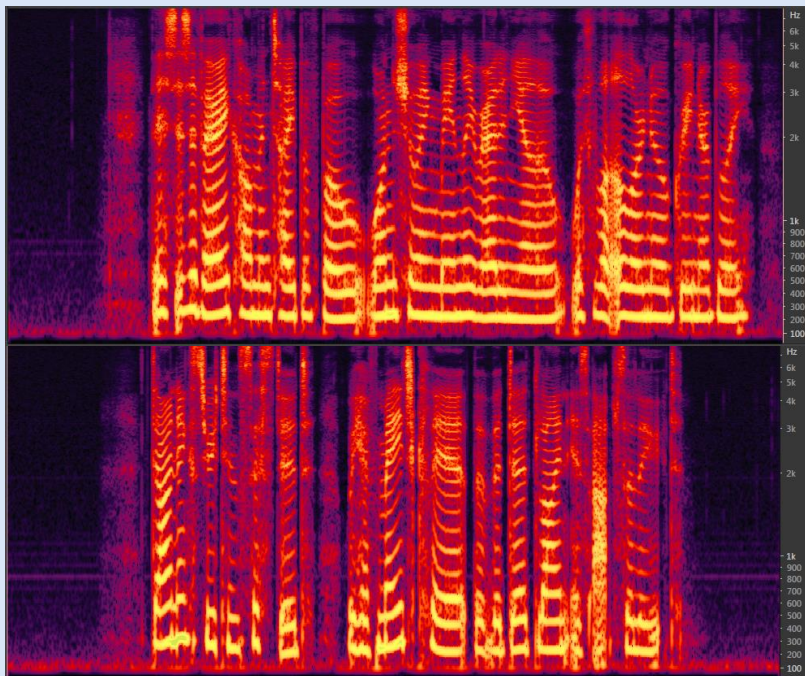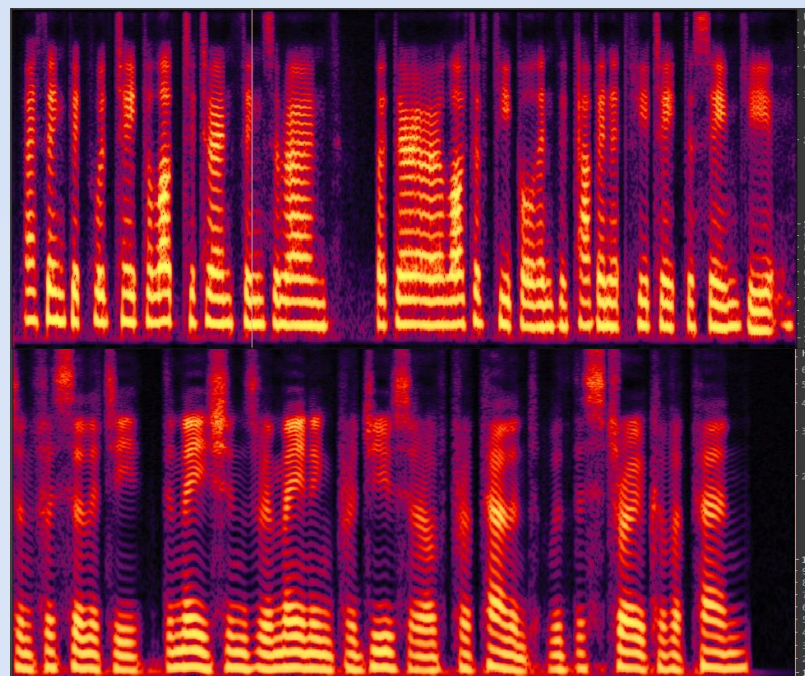| Type | Filter / Stride | Output | Params |
|------|-----------------|--------|--------|
| Conv_1 | 5 × 5 / 1 × 1 | 863 × 600 × 64 | 1.6K |
| MFM_2 | – | 864 × 600 × 32 | – |
| MaxPool_3 | 2 × 2 / 2 × 2 | 431 × 300 × 32 | – |
| Conv_4 | 1 × 1 / 1 × 1 | 431 × 300 × 64 | 2.1K |
| MFM_5 | – | 431 × 300 × 32 | – |
| BatchNorm_6 | – | 431 × 300 × 32 | – |
| Conv_7 | 3 × 3 / 1 × 1 | 431 × 300 × 96 | 27.7K |
| MFM_8 | – | 431 × 300 × 48 | – |
| MaxPool_9 | 2 × 2 / 2 × 2 | 215 × 150 × 48 | – |
| BatchNorm_10 | – | 215 × 150 × 48 | – |
| Conv_11 | 1 × 1 / 1 × 1 | 215 × 150 × 96 | 4.7K |
| MFM_12 | – | 215 × 150 × 48 | – |
| BatchNorm_13 | – | 215 × 150 × 48 | – |
| Conv_14 | 3 × 3 / 1 × 1 | 215 × 150 × 128 | 55.4K |
| MFM_15 | – | 215 × 150 × 64 | – |
| MaxPool_16 | 2 × 2 / 2 × 2 | 107 × 75 × 64 | – |
| Conv_17 | 1 × 1 / 1 × 1 | 107 × 75 × 128 | 8.3K |
| MFM_18 | – | 107 × 75 × 64 | – |
| BatchNorm_19 | – | 107 × 75 × 64 | – |
| Conv_20 | 3 × 3 / 1 × 1 | 107 × 75 × 64 | 36.9K |
| MFM_21 | – | 107 × 75 × 32 | – |
| BatchNorm_22 | – | 107 × 75 × 32 | – |
| Conv_23 | 1 × 1 / 1 × 1 | 107 × 75 × 64 | 2.1K |
| MFM_24 | – | 107 × 75 × 32 | – |
| BatchNorm_25 | – | 107 × 75 × 32 | – |
| Conv_26 | 3 × 3 / 1 × 1 | 107 × 75 × 64 | 18.5K |
| MFM_27 | – | 107 × 75 × 32 | – |
| MaxPool_28 | 2 × 2 / 2 × 2 | 53 × 37 × 32 | – |
| FC_29 | – | 160 | 10.2 MM |
| MFM_30 | – | 80 | – |
| BatchNorm_31 | – | 80 | – |
| FC_32 | – | 2 | 64 |

Refer to:  [1] Wu X, He R, Sun Z, et al. A light cnn for deep face representation with noisy labels[J]. IEEE Transactions on Information Forensics and Security, 2018, 13(11): 2884-2896.
[2] Galina Lavrentyeva, Sergey Novoselov, Andzhukaev Tseren, Marina Volkova, Artem Gorlanov, and Alexandr Kozlov, "STC Antispoofing Systems for the ASVspoof2019 Challenge," Interspeech, 09 2019.

# *Loss function*

- ## Two class classification



- ## Softmax

$$L_i = -\log\left(\frac{e^{\|\boldsymbol{W}_{y_i}\|\|\boldsymbol{x}_i\|\cos(\theta_{y_i})}}{\sum_j e^{\|\boldsymbol{W}_j\|\|\boldsymbol{x}_i\|\cos(\theta_j)}}\right)$$

Introduce margin m

- ## Large margin softmax

$$L_i = -\log\left(\frac{e^{\|\boldsymbol{W}_{y_i}\|\|\boldsymbol{x}_i\|\psi(\theta_{y_i})}}{e^{\|\boldsymbol{W}_{y_i}\|\|\boldsymbol{x}_i\|\psi(\theta_{y_i})} + \sum_{j\neq y_i} e^{\|\boldsymbol{W}_j\|\|\boldsymbol{x}_i\|\cos(\theta_j)}}\right)$$
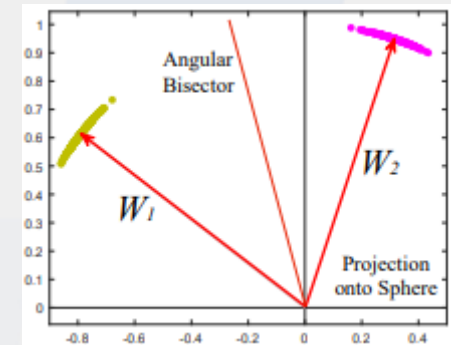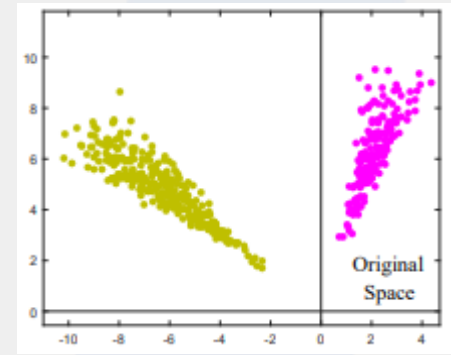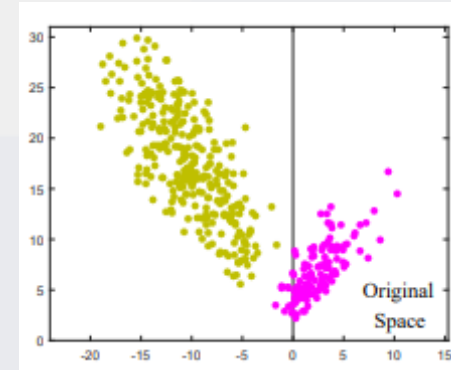
in which we generally require

$$\psi(\theta) = \begin{cases} \cos(m\theta), & 0 \le \theta \le \frac{\pi}{m} \\ \mathcal{D}(\theta), & \frac{\pi}{m} < \theta \le \pi \end{cases}$$

Normalize the weights
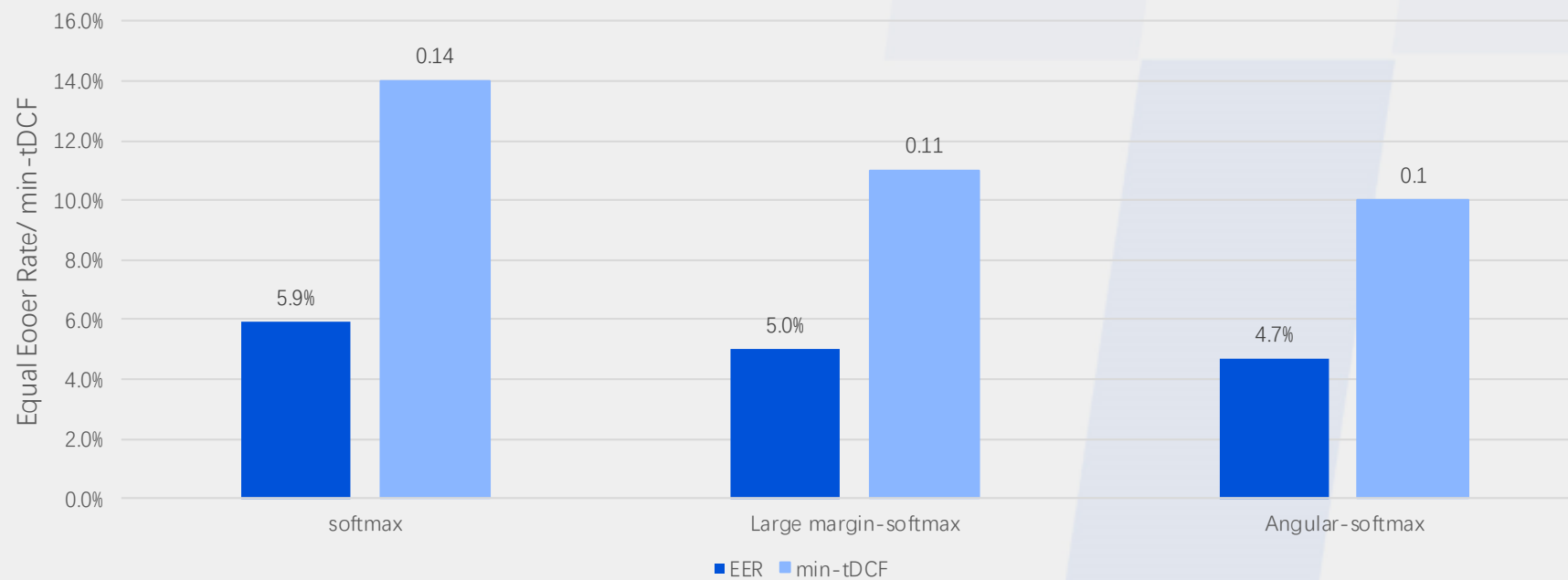and zero the biases

- ## Angular softmax

$$L_{\mathrm{ang}} = \frac{1}{N}\sum_i -\log\left(\frac{e^{\|\boldsymbol{x}_i\|\psi(\theta_{y_i,i})}}{e^{\|\boldsymbol{x}_i\|\psi(\theta_{y_i,i})} + \sum_{j\neq y_i} e^{\|\boldsymbol{x}_i\|\cos(\theta_{j,i})}}\right)$$

in which we define $\psi(\theta_{y_i,i}) = (-1)^k \cos(m\theta_{y_i,i}) - 2k$

Refer to:  [1] large-Margin Softmax Loss for Convolutional Neural Networks. Weiyang Liu, Yandong Wen, Zhiding Yu, Meng Yang. ICML 2016
[2] SphereFace: Deep Hypersphere Embedding for Face Recognition. Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, Le Song. ICML 2017
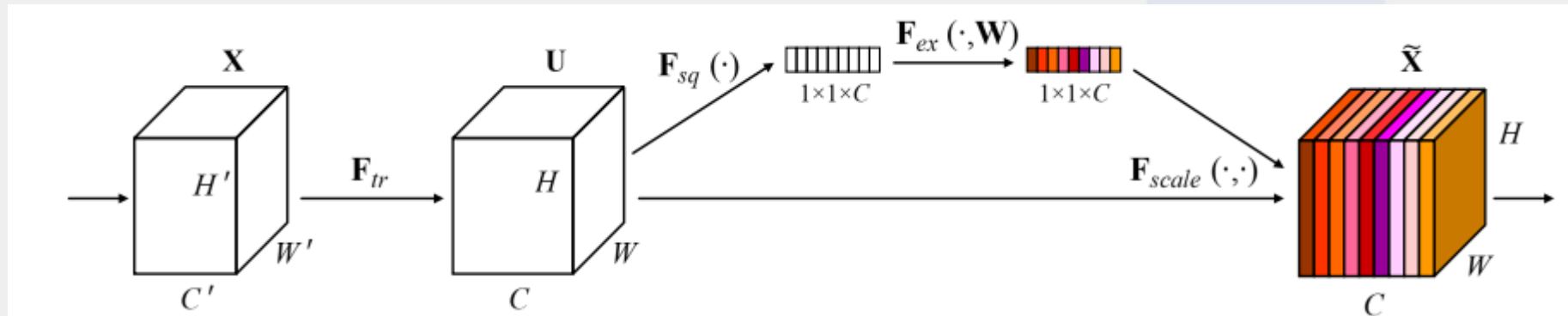
# Loss function



基于不同Loss function的LCNN网络在ASVSpoof2019 eval集上的性能对比

# *Attention modules*

- Squeeze-and-Excitation (SE)
  - Squeeze: global average pooling in channel
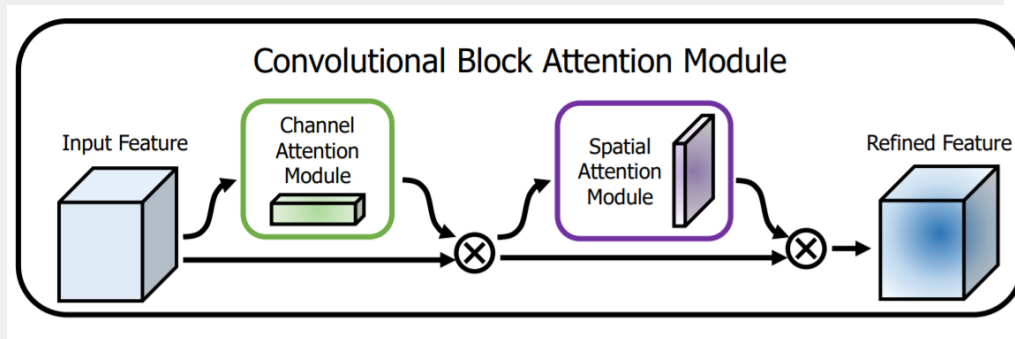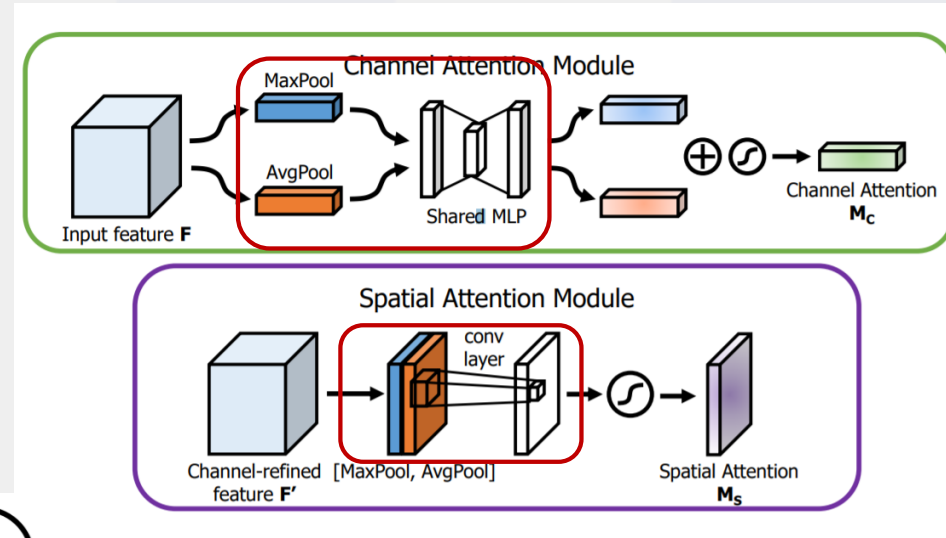  - Excitation: produces modulation weights

# Attention modules

- Convolutional Block Attention Module (CBAM)
  - channel attention module

$$\mathbf{M_c}(\mathbf{F}) = \sigma(MLP(AvgPool(\mathbf{F})) + MLP(MaxPool(\mathbf{F})))$$
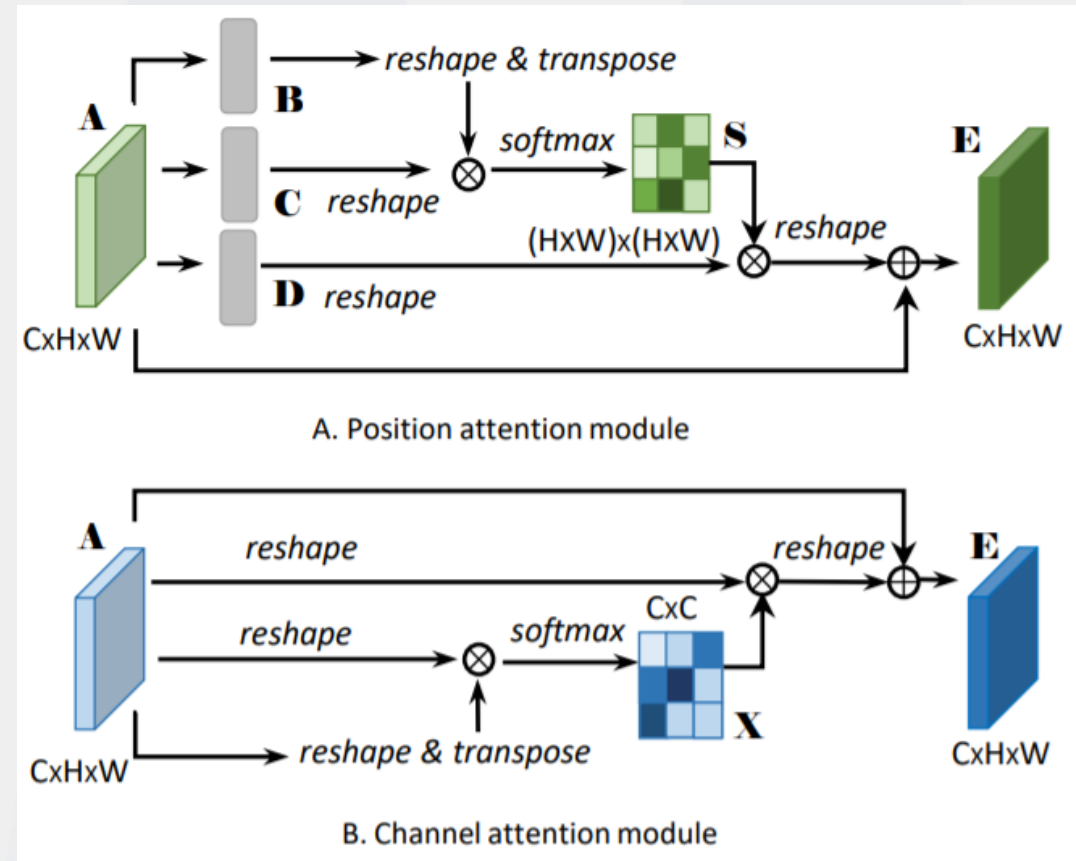$$= \sigma(\mathbf{W_1}(\mathbf{W_0}(\mathbf{F_{avg}^c})) + \mathbf{W_1}(\mathbf{W_0}(\mathbf{F_{max}^c}))),$$
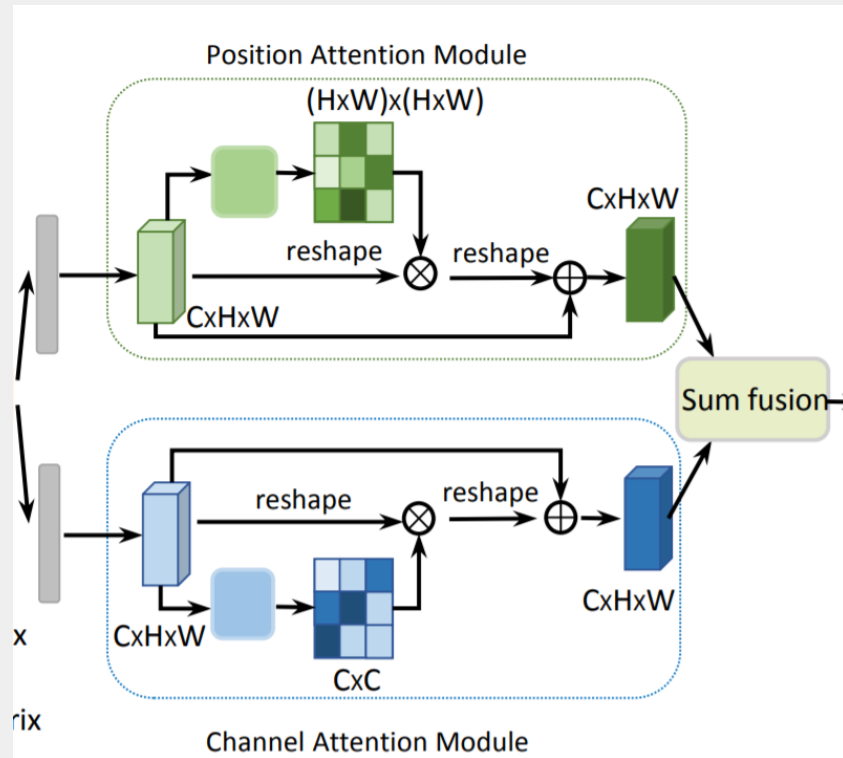
  - spatial attention module

$$\mathbf{M_s}(\mathbf{F}) = \sigma(f^{7\times7}([AvgPool(\mathbf{F}); MaxPool(\mathbf{F})]))$$
$$= \sigma(f^{7\times7}([\mathbf{F_{avg}^s}; \mathbf{F_{max}^s}])),$$

Refer to: Sanghyun Woo, JongChan Park, Joon-Young Lee, and Inso Kweon, "CBAM: Convolutional Block Attention Module," arXiv: Computer Vision and Pattern Recognition, pp. 3–19, 2018.

# *Attention modules*

- Dual Attention (DA)
  - Channel 和position 并行关系

Refer to:Jun Fu, Jing Liu, Haijie Tian, Zhiwei Fang, and HanqingLu, "Dual Attention Network for Scene Segmentation," 2019 IEEE/CVF Conference on Computer Vision andPatternRecognition(CVPR),pp.3141–3149,2018.

# Result on ASVSpoof 2019

- Attention modules
  - Squeeze-and-Excitation (SE)
  - Convolutional Block Attention Module (CBAM)
  - Dual Attention （DA)

| | min-tDCF | EER （%) | min-tDCF | EER （%) |
|---|---|---|---|---|
| Baseline-GMM | 0.066 | 2.71 | 0.021 | 8.09 |
| LCNN | 0.008 | 0.27 | 0.101 | 4.74 |
| LCNN_SE | 0.006 | 0.20 | 0.137 | 6.06 |
| LCNN_CBAM | 0.028 | 0.93 | 0.094 | 3.67 |
| LCNN_DA | 0.025 | 0.78 | **0.078** | **2.76** |

注：测试数据为 说话人Anti-Spoofing比赛 ASVSpoof2019 dev和eval测试集，评价指标为EER和min-tDCF

# *Deepfake Detection Challenge*

- 与图像团队参加由Facebook主办的DFDC百万美金比赛
  作为变声检测方案提供方

| | EER(%) | min-tDCF |
|---|---|---|
| Dev集 | 1.16 | 0.03 |

| | 提交结果 | |
|---|---|---|
| 变声检测 | 0.67 | |
| 变脸检测 | 0.36 | |
| 两者融合 | 0.35 | 对总名次提升22名 |

Kaggle-DFDC竞赛获得银牌

注：从提供的dev数据来看，变声数据比例大概占全部Fake数据10%。
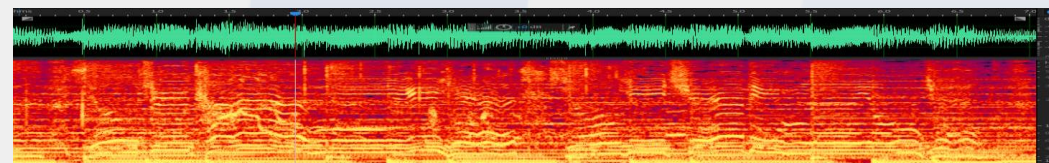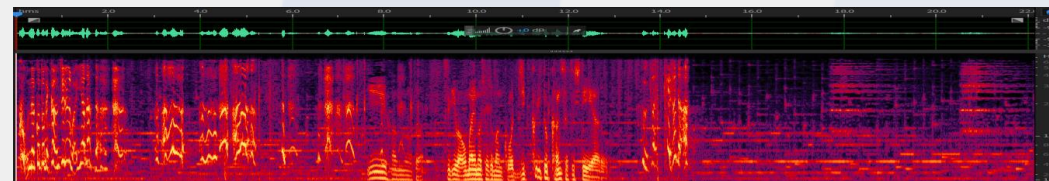
# 业务落地—数据分析

- **学术公开集数据**
  - 单一纯净音频

- **业务场景数据**
  - 多种音频内容
  - 信噪比低

# 业务落地一数据分布

- 不同场景音频类型分布差异大

场景1数据比例

32.50%

17.65%

26.41%

11.26%

5.38%

■语音 ■歌曲 ■音乐 ■干扰音 ■其他

场景2数据比例

3.25%

6.17%

7.33%

8.64%

74.53%

■语音 ■歌曲 ■音乐 ■干扰音 ■其他
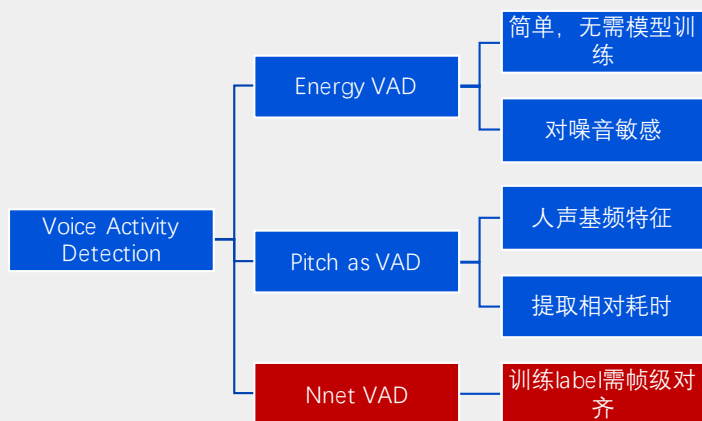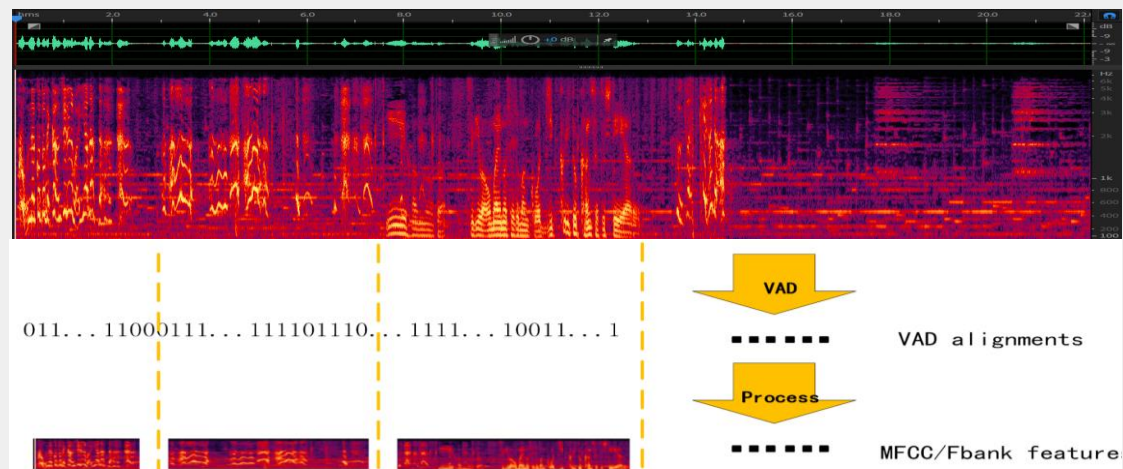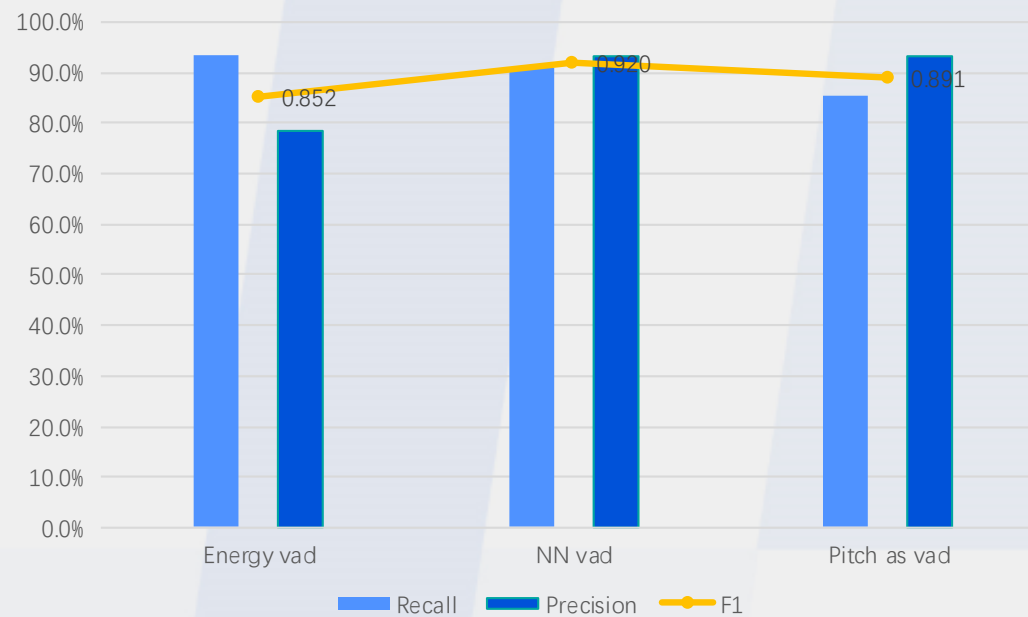
# 业务落地一问题定义

- **学术公开集任务**
  - 该音频中说话人是否为 spoofed
- **学术公开集LA spoofing 范围**
  - TTS
  - VC

- **业务场景任务**
  - 该音频中是否出现spoofed说话人
- **业务场景中spoofing 范围**
  - TTS
  - VC
  - 鬼畜拼接

# Voice Activity Detection (VAD)



011...11000111...111101110...1111...10011...1

VAD

‒ ‒ ‒ ‒ ‒ ‒   VAD alignments

Process

‒ ‒ ‒ ‒ ‒ ‒   MFCC/Fbank feature

Voice Activity Detection

- Energy VAD
  - 简单，无需模型训练
  - 对噪音敏感
- Pitch as VAD
  - 人声基频特征
  - 提取相对耗时
- Nnet VAD
  - 训练label需帧级对齐

## 不同VAD方式在某业务测试集上性能对比



0.852　　　0.920　　　0.891

Recall　　Precision　　F1

# Data Augmentation

- 加入真实背景音乐、噪声数据

Data Augmentation

Add background music/noise

resampling rate

frequency masking

time masking

Data augmentation对ASV Spoofing 系统（LCNN_CBAM）在实际场景的性能影响

Recall/False Alarm

80.02%    80.01%

3.64%    1.71%

without data augmentation    with data augmentation

合成音    非合成音

不同音频类型误报分析
Without data augmentation

语音+背景音乐    9.10%
语音+背景噪音    2.96%
语音+背景歌曲    5.26%
纯净语音    1.65%

0.00%  2.00%  4.00%  6.00%  8.00%  10.00%
False Alarm

不同音频类型误报分析
With data augmentation

语音+背景音乐    2.09%
语音+背景噪音    1.78%
语音+背景歌曲    2.50%
纯净语音    1.32%

0.00%  2.00%  4.00%  6.00%  8.00%  10.00%
False Alarm

# 特定人的*antispoofing* 识别

# 音频内容理解



| | 语音 | 音乐 | 语音 | 其他 | 音乐 | 音乐 |
|---|---|---|---|---|---|---|
| **分类** | 语音 | 音乐 | 语音 | 其他 | 音乐 | 音乐 |
| **子类标签** | 语种 | 歌曲 | 语种<br>说话人，真假语音 | 干扰音 | 纯音乐 | 歌曲 |
| **内容** | "今天天气真好。" | | " Hello, welcome to Tencent! " | | | |
| **业务应用** | | | | | | |

音频

*Thanks*