# ResMax: Detecting Voice Spoofing Attacks with Residual Network and Max Feature Map

Il-Youp Kwak[1], Sungsu Kwag[2], Junhee Lee[2], Jun Ho Huh[2], Choong-Hoon Lee[2], Youngbae Jeon[3], Jeonghwan Hwang[3], Ji Won Yoon[3]

Chung-Ang University[1], Samsung Research[2], Korea University[3]

## ABSTRACT

The "2019 Automatic Speaker Verification Spoofing And Countermeasures Challenge" (ASVspoof) competition aimed to facilitate the design of highly accurate voice spoofing attack detection systems. the competition did not emphasize model complexity and latency requirements; such constraints are strict and integral in real-world deployment. Hence, most of the top performing solutions from the competition all used an ensemble approach, and combined multiple complex deep learning models to maximize detection accuracy – this kind of approach would sit uneasily with real-world deployment constraints. To design a lightweight system, we combined the notions of skip connection (from ResNet) and max feature map (from Light CNN), and evaluated the accuracy of the system using the ASVspoof 2019 dataset. With an optimized constant Q transform (CQT) feature, our single model achieved a replay attack equal error rate (EER) of 0.37% on the evaluation set, surpassing the top ensemble system from the competition that achieved an EER of 0.39%.

## OBJECTIVES

### Why Voice Spoofing Detection?

6-year-old orders $170 dollhouse, cookies with Amazon's Alexa

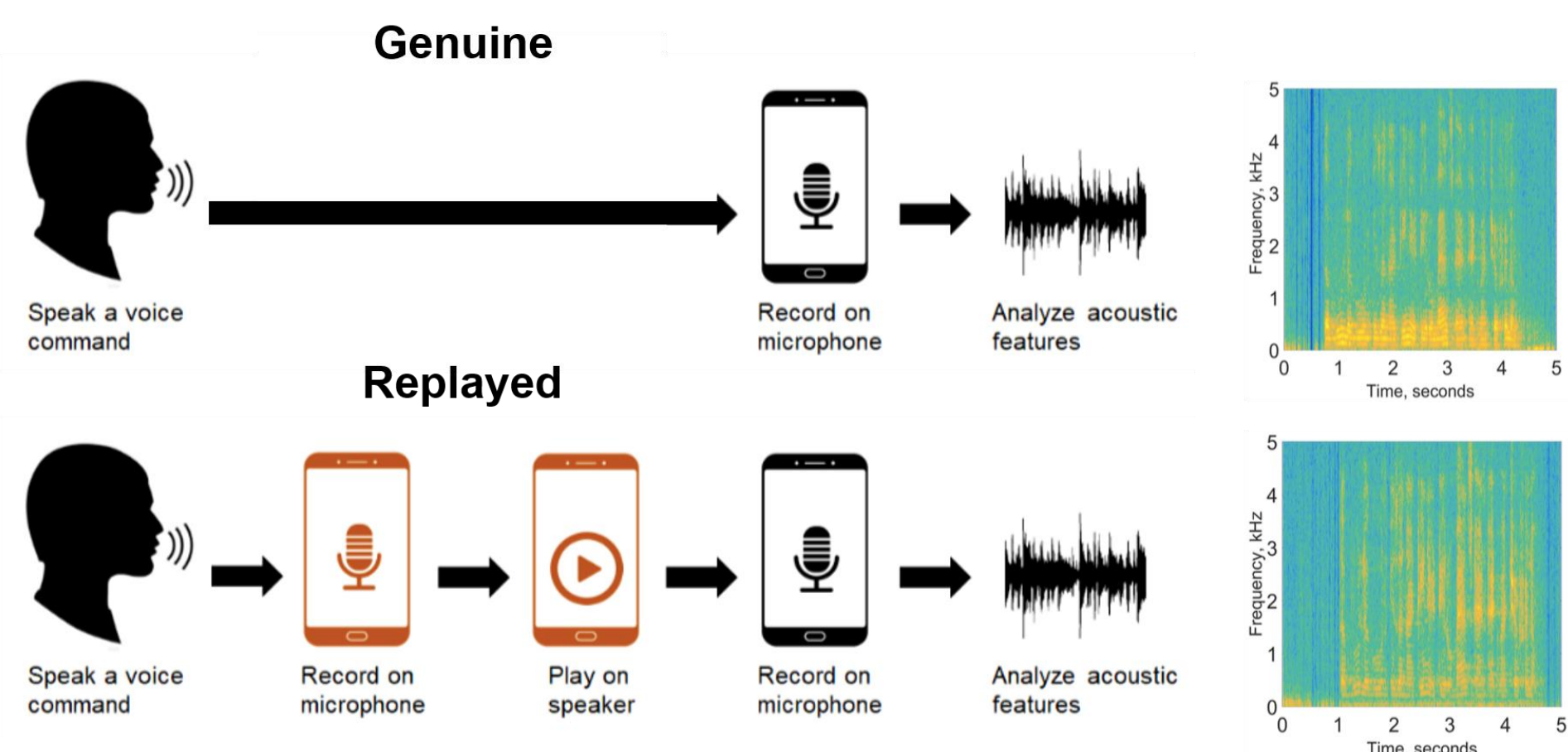TV anchor says live on-air 'Alexa, order me a dollhouse' – guess what happens next

Story on accidental order begets story on accidental order begets accidental order

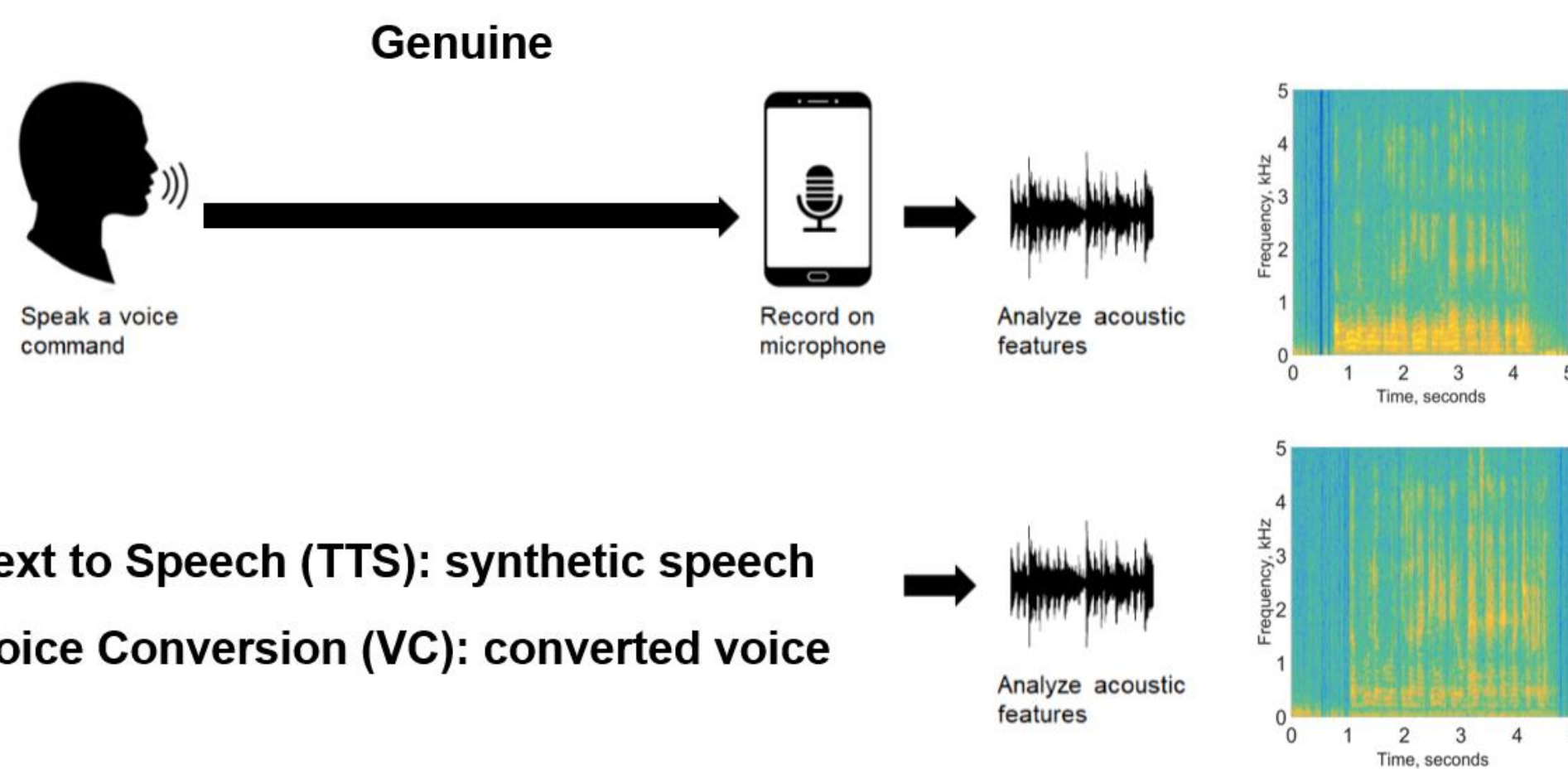By Shaun Nichols in San Francisco 7 Jan 2017 at 00:58    344🗨    SHARE ▼

A San Diego TV station sparked complaints this week – after an on-air report about a girl who ordered a dollhouse via her parents' Amazon Echo caused Echoes in viewers' homes to also attempt to order dollhouses.

### Voice Spoofing Detection, Physical Access (PA)

**Genuine**

Speak a voice command → Record on microphone → Analyze acoustic features

**Replayed**

Speak a voice command → Record on microphone → Play on speaker → Record on microphone → Analyze acoustic features

### Voice Spoofing Detection, Logical Access (LA)

**Genuine**

Speak a voice command → Record on microphone → Analyze acoustic features

**Text to Speech (TTS):** synthetic speech
**Voice Conversion (VC):** converted voice

Analyze acoustic features

## METHODS

### Spoofing Classification

- **Classifying Human or Speaker**  $f : \mathbf{x} \xrightarrow{f_\theta} \mathbb{R}_{[0,1]}$
- **Dimension for x is (n_freq, n_time)**
- **Data:** $(\mathbf{x}, y)$    $\mathbf{x} \in \mathbb{R}^{(n_f, n_t)}, y \in \{0, 1\}$

    $m$ training examples $\{(\mathbf{x}^{(1)}, y^{(1)}), \ldots, (\mathbf{x}^{(m)}, y^{(m)})\}$

    $X = [\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(m)}]$    $Y = [y^{(1)}, \ldots, y^{(m)}]$

- Given $\mathbf{x} \in \mathbb{R}^{(n_f, n_t)}$, want $\hat{y} = P(y = 1 | \mathbf{x}) = f_\theta(\mathbf{x}) \in \mathbb{R}^{[0,1]}$
- **Objective is to minimize Cost, C($\theta$), w.r.t $\theta$ :**

    $C(\theta) = \sum_{i=1}^{m} L(f_\theta(\mathbf{x}_i), y_i)$

### What would be a good model( $f_\theta(\mathbf{x})$ ) ?

$f_\theta(\mathbf{x})$

A speech sample → Voice Liveness Detection model → Score

### Automatic Speaker Verification Spoofing And Countermeasures Challenge (ASVspoof 2015, 2017 and 2019)

Wu et al. (2015) ASVspoof 2015: the First Automatic Speaker Verification Spoofing and Countermeasures Challenge
Kinnunen et al. (2017) The asvspoof 2017 challenge: Assessing the limits of replay spoofing attack detection
Todisco et al (2019) ASVspoof 2019: Future Horizons in Spoofed and Fake Audio Detection

### Deep-learning based methods and ensemble solutions are dominating in voice liveness detection challenge

Top 2017 PA scenario

Top 2019 LA scenario

Top 2019 PA scenario

Grey    Used Neural Networks
**Bold**    Used Ensemble

Kinnunen et al. (2017) The asvspoof 2017 challenge: Assessing the limits of replay spoofing attack detection
Todisco et al (2019) ASVspoof 2019: Future Horizons in Spoofed and Fake Audio Detection

### How to develop **well performing light-weighted** model?

LCNN architecture    ResNet architecture

TABLE I
ENSEMBLE SOLUTIONS FROM ASVSPOOF 2019 AND THE LIST OF MODELS USED.

| Model | Data | All models used |
|---|---|---|
| T10 [13] | PA | LFCC ResNet, GD gram ResNet, Joint gram ResNet |
| T44 [12] | PA | logspec SENet34, CQCC ResNet, logspec SENet50 |
| T45 [6] | LA | LFCC LCNN, LFCC-CMVN LCNN, CQT LCNN |
| | PA | CQT LCNN, LFCC LCNN, DCT LCNN |
| T50 [14] | LA | CQT-CGCNN, CQT ResNet18, CQT ResNet18 Vec |
| T60 [15] | PA | FFT-CNN, FFT-CRNN, IMFCC-GMM, SVм-iVec |

## light-weighted ResMax Architecture

MFM

$h(x) = \max(x^1, x^2)$

$H \times W \times C$    $H \times W \times C$

(a) MFM

ResMax(f, k, l, m)

| | |
|---|---|
| Conv(k) | $H \times W \times 2C$ |
| MFM | $H \times W \times C$ |
| Conv(k = 1) | $H \times W \times 2C$ |
| MFM | $H \times W \times C$ |
| + | |
| MaxPool(2,2) | $\frac{H}{2} \times \frac{W}{2} \times C$ |
| BN | |

Feature

| |
|---|
| ResMax(16,3,1,1) |
| ResMax(16,5,1,1) |
| ResMax(24,3,1,1) |
| ResMax(32,3,0,0) |
| ResMax(32,3,0,1) |
| ResMax(48,3,0,0) |
| ResMax(48,3,0,1) |
| ResMax(64,3,0,0) |
| ResMax(64,3,0,1) |
| Score |

(b) ResMax Block    (c) Model Architecture

## RESULTS

### High performance of ResMax

**LA**

| # | Model | t-DCF (Dev) | EER (Dev) | t-DCF (Eval) | EER (Eval) | #Mo | # Params |
|---|---|---|---|---|---|---|---|
| 1 | T05 | 0.0000 | 0.000 | 0.0069 | 0.22 | 6 | - |
| 2 | T45 | 0.0000 | 0.000 | 0.0510 | 1.86 | 5 | 1484K |
| 3 | CQT-1_100-ResMax | 0.0179 | 0.56 | 0.0600 | 2.19 | 1 | 262K |
| 4 | T60 | 0.0 | 0.0 | 0.0755 | 2.64 | 4 | - |
| 5 | T24 | - | - | 0.0953 | 3.45 | - | - |
| 6 | T50 | 0.027 | 0.90 | 0.1118 | 3.56 | - | - |
| | T45 (FFT-LCNN) | 0.0009 | 0.040 | 0.1028 | 4.53 | 1 | 371K |
| | T45 (LFCC-LCNN) | 0.0043 | 0.157 | 0.1000 | 5.06 | 1 | 371K |

**PA**

| # | Model | t-DCF (Dev) | EER (Dev) | t-DCF (Eval) | EER (Eval) | #Mo | # Params |
|---|---|---|---|---|---|---|---|
| 1 | CQT-1_120-ResMax | 0.0066 | 0.23 | 0.0091 | 0.37 | 1 | 262K |
| 2 | T28 | - | - | 0.0096 | 0.39 | - | - |
| 3 | T45 | 0.0001 | 0.0154 | 0.0122 | 0.54 | 3 | 1113K |
| 4 | T44 | 0.003 | 0.129 | 0.0161 | 0.59 | 5 | 5811K |
| 5 | T10 | 0.0064 | 0.24 | 0.0168 | 0.66 | 6 | 1330K |
| 6 | T24 | - | - | 0.0215 | 0.77 | - | - |
| | T28 | - | - | - | 0.50 | 1 | - |
| | T45 (CQT-LCNN) | 0.0197 | 0.800 | 0.0295 | 1.23 | 1 | 371K |
| | T44 (logspec-SENet) | 0.015 | 0.575 | 0.0360 | 1.29 | 1 | 1344K |

### Non-speech part have information?

amplitude

(a) Original sound    (b) Processed sound

Model: With processed sound / With original sound / With processed sound / With original sound

Dev    Eval
Average EER

(c) Performance comparison

Fig. 3.  The non-speech part remover suggested and tested. The ResMax model worked better without the non-speech part remover. The barplot indicates averaged EER with one standard deviation error bar.

### The longer you listen the better the performance

Dev    Eval
Model: 9 sec / 6 sec / 3 sec / 9 sec / 6 sec / 3 sec
Average EER

Fig. 4.  The 9-second model performed best for both development and evaluation sets in LA and PA data. The barplot indicates averaged EER with one standard deviation error bar.

## Performance depend on replay device quality

TABLE III
DETECTION PERFORMANCE ON THE ASVSPOOF2019 PHYSICAL ACCESS EVALUATION SETS IN VARIOUS ENVIRONMENTS. THE **A**, **B**, **C** REPRESENT THE CLASSES OF EACH FACTOR WHICH IS WELL DESCRIBED IN [5]. ALL NUMERICAL VALUES REPRESENT THE AVERAGE OF EER.

| | Factors | A | B | C |
|---|---|---|---|---|
| **Verification Env.** | Room size (S) | 0.0047 | 0.0044 | 0.0041 |
| | T60 (R) | 0.0055 | 0.0029 | 0.0038 |
| | Talker-to-ASV distance | 0.0059 | 0.0036 | 0.0042 |
| **Recording Env.** | Attacker-to-talker distance (D_a) | 0.0051 | 0.0036 | 0.0041 |
| | Replay Device Quality (Q) | 0.0067 | 0.0036 | 0.0009 |

## High performance on best performing TTS, VC systems

| ID | Type | Description | EER |
|---|---|---|---|
| A07 | TTS | vocoder + GAN | 0.0022 |
| A08 | TTS | neural waveform | 0.0388 |
| A09 | TTS | vocoder | 0.0003 |
| A10 | TTS | neural waveform | 0.0045 |
| A11 | TTS | griffin lim | 0.0039 |
| A12 | TTS | neural waveform | 0.0002 |
| A13 | TTS,VC | waveform concatenation & filtering | 0.0051 |
| A14 | TTS,VC | vocoder | 0.0012 |
| A15 | TTS,VC | neural waveform | 0.0030 |
| A16 | TTS | waveform concatenation | 0.0039 |
| A17 | VC | waveform filtering | 0.0561 |
| A18 | VC | vocoder | 0.0225 |
| A19 | VC | spectral filtering | 0.0317 |

Fig. 5.  The averaged EER for 13 attack types in evaluation set. The barplot indicate averaged EER with one standard deviation error bar.

**Known attacks are A16,A19 and 4 others**

## CONCLUSIONS

Existing voice spoofing attack detection solutions have been designed without considering real-world model complexity and detection latency requirements, and often consist of multiple heavy and complex deep learning models. Such solutions would not be considered suitable given the tight model size and latency requirements. In comparison, our CQT-1 120-ResMax model used only a single deep learning model with far fewer model parameters to outperform the top performing PA solution from evaluation set, achieving an EER of 0.37% compared to the current best competition EER of 0.39%, which is an ensemble solution. As for the LA set, we rank third with an EER of 2.19%, just behind the second best ensemble solution that achieved an EER of 1.86% in the evaluation set. Among the single model systems, although CQT-1 120-ResMax used the least number of parameters, it demonstrated significant superiority in detection accuracy.

## ACKNOWLEDGMENT

## CONTACTS

For any questions, feel free to ask me via ikwak2@cau.ac.kr