

# Audio Adversarial Examples: Targeted Attacks on Speech-to-Text

Nicholas Carlini, David Wagner  
University of California, Berkeley  
SYSSEC

# Related works

- **Attacks on speech recognition**
  - **Side-Channel** → **Dolphin Attack, light command**
  - **Mangled Voice** → **Hidden Voice Command**

# Related works

- Dolphin Attack

frequency sound

**Commands**  
**Without Background Noise**

accurate meaning.

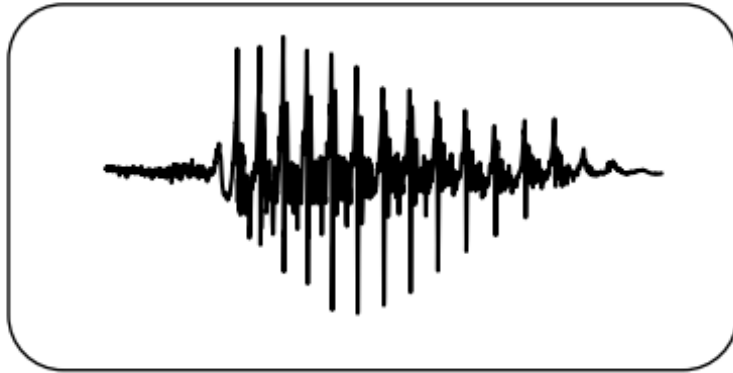
# Introduction

- Attack on speech recognition using **adversarial example**
- Overcome the limitations
  - Need a high-performance device? **No**
  - Victims notice the attack? **No**
- Target : DeepSpeech, an **automatic speech recognition system**

**ASR**  
a technology that converts  
speech to text in real time

# Neural Networks for Automatic Speech Recognition (ASR)

audio sample



HELLO



Training data  
pairs of (audio, text)  
of **variable lengths**  
with **no alignment**

# Connectionist Temporal Classification, CTC

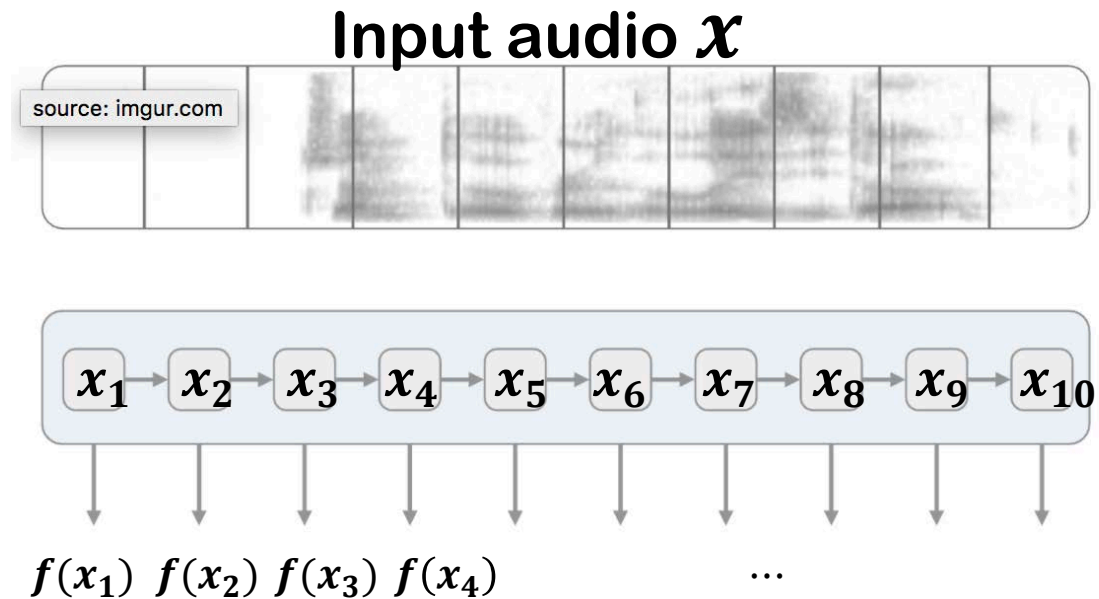
## CTC Loss

A differentiable measure of **distance** from **NN output** to the **true target phrase**.

## Training objective

Minimize CTC Loss between the training audio and corresponding text.

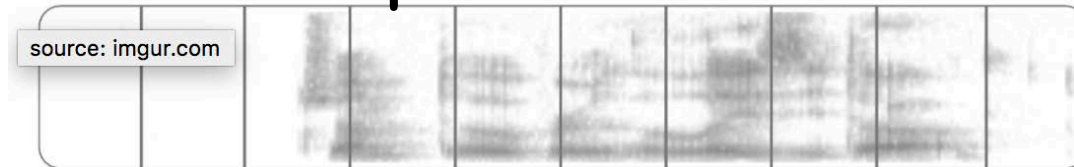
# Connectionist Temporal Classification, CTC



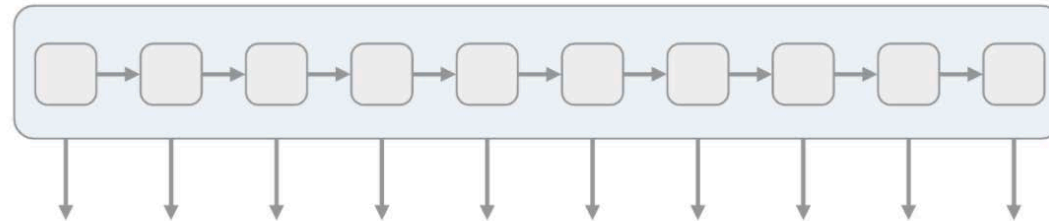
a probability distribution  
over all the possible labels  $f(x)$

# Connectionist Temporal Classification, CTC

Input audio  $x$



$x = \text{Hello}$



h	h	h	h	h	h	h	h	h	h
e	e	e	e	e	e	e	e	e	e
l	l	l	l	l	l	l	l	l	l
o	o	o	o	o	o	o	o	o	o
$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$

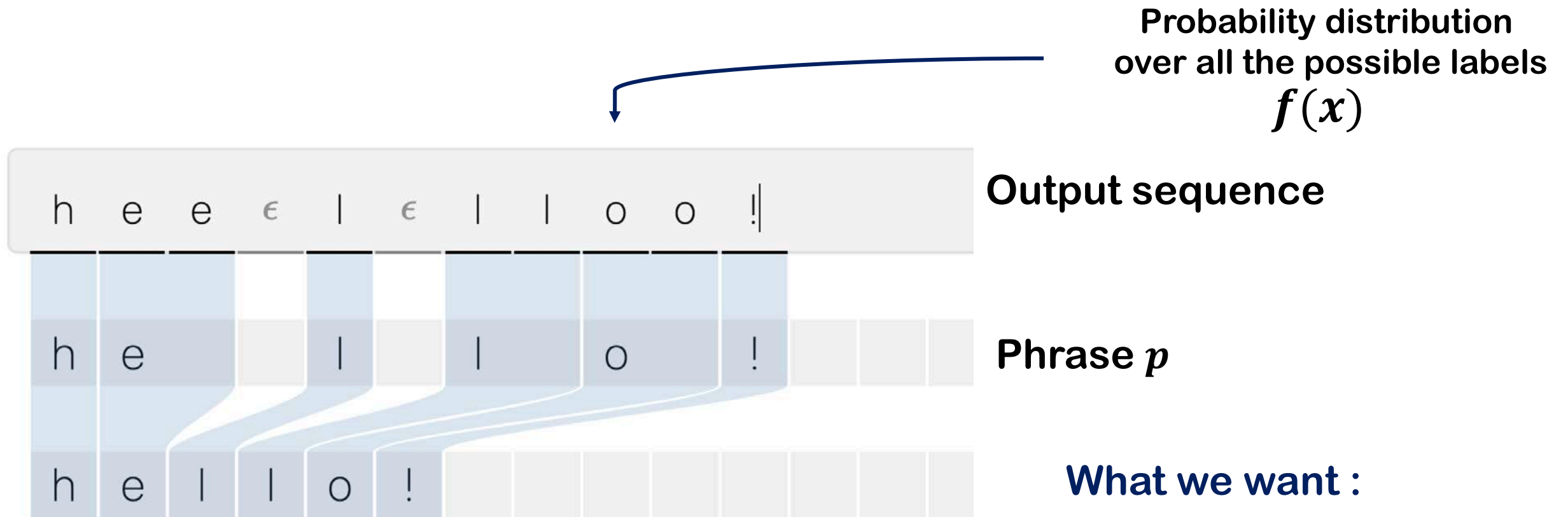
Possible labels  
 $\{h, e, l, o\} + \{\epsilon\}$

a probability distribution  
over all the possible labels

$f(x)$



# Connectionist Temporal Classification, CTC

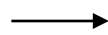


Maximize the probability that we will get  $f(x)$  that collapses into  $p$

**Goal** Minimize  $CTC - Loss(f(x), p) = -\log Pr(p|f(x))$

# Targeted Adversarial Examples

- Given an input  $x$ , classified as  $f(x) = L$
- Find the AE  $x'$  close to  $x$  so that  $f(x') = T$  [for  $T \neq L$ ]



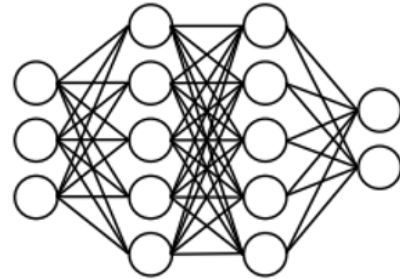
cat



dog

# Targeted Adversarial Examples

- Make the machine recognize my voice as that of Professor Yongdae Kim.



- Make the machine recognize the Professor Yongdae Kim's speech as the target phrase.

# This Paper/ Overview

## Purpose

Propose an **targeted white-box adversarial attack** method on **speech recognition system**.

## Contribution

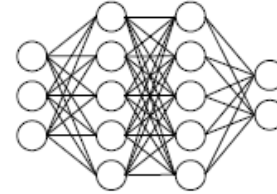
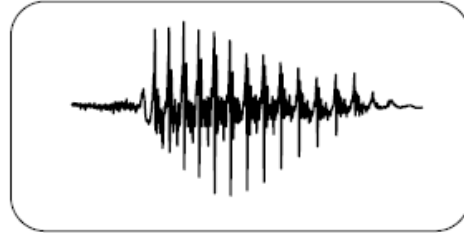
The attack works with **100 % success rate**, regardless of the **desired transcription** or **initial source audio sample**.

# Attack Model



“We have class tomorrow”

**Speech**



**Text**  
DeepSpeech transcribe as

“We have class tomorrow”

# Algorithm

- **Formulation**

given an original audio sample  $x$ ,  
find an adversarial example  $x'$

where

minimize  
such that

$$d(x, x')$$
$$f(x') = T$$

ASR model

Target phrase  
that adversary has chosen

**What is the distance metric?**

**Magnitude of perturbation (in dB)  
relative to the source audio signal**

# Algorithm

- **Formulation**

given  $x$ , find  $x'$  where

minimize  $d(x, x') + I(x', T)$

What is the loss function  $I$ ?

**CTC Loss!**

- $I(x', T)$  is a loss function **measuring how close  $f(x')$  is to target  $T$** 
  - $I(x', T)$  is small, if  $f(x') = T$
  - $I(x', T)$  is large, if  $f(x') \neq T$

# Algorithm

- **Formulation**

given  $x$ , find  $x'$  where

minimize  $\|\delta\|_2^2 + c \cdot \text{CTC-Loss}(x + \delta, T)$

such that  $\text{dB}_x(\delta) < \tau$

- **Iterative optimization**

Repeat until no solution  $\delta$  can be found by reducing  $\tau$ .

$\delta$  : perturbation,  $c$ : regularization coefficient,  $x$  : input,

$T$  : target phrase,  $\|\cdot\|_2$ : 2D p-norm  $(\sum_{i,j} (|\delta_{(i,j)}|^p)^{1/p})$ ,  $\tau$ : hyper parameter



# Evaluation

- **Speech to another speech**
- **Non-speech to speech**
- **Hide speech**

# Evaluation (Speech to Speech)

DeepSpeech transcribes it as

**“later we simply let life proceed  
in its own direction toward its own fate”**



**“the boy looked out at the horizon”**

# Evaluation (Information Density)

**Input waveform is converted into 50 frame/sec**

**→ theoretical maximum density of audio is 50 character/sec**

**→ can generate adversarial examples at this maximum rate**

**→ make short audio clips recognized as longer phrases**

# Evaluation (Non-Speech to Speech)

**DeepSpeech transcribes it as**

**“speech can be embedded in music”**



# Evaluation (Speech to Silence)

- **You can “Hide” speech by adding adversarial noise.**
  - **DeepSpeech transcribes nothing.**
- **Targeting silence is easier than targeting a specific phrase.**
  - **Easier to construct adversarial examples when starting with longer audio waveforms than shorter ones**

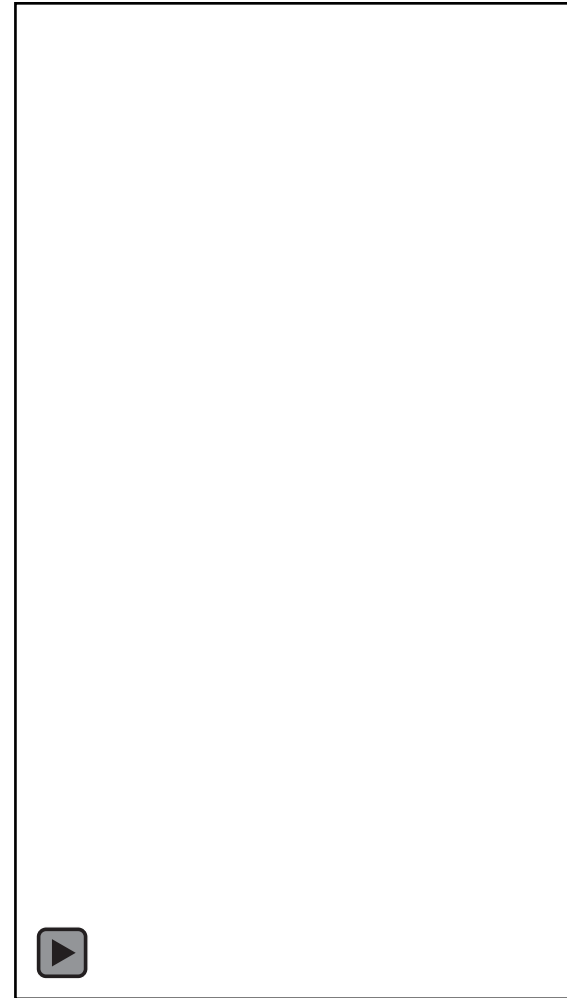
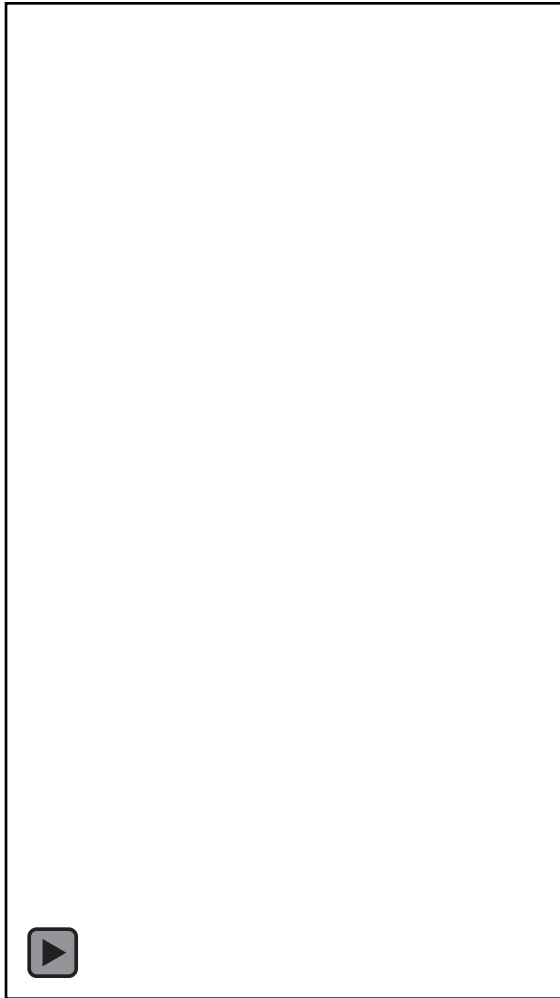
# Conclusion

- **Demonstrate targeted audio adversarial examples are effective on automatic speech recognition**
- **With optimization-based attacks applied end-to-end, turn any audio waveform into any target transcription with 100% success by only adding a slight distortion.**
- **Open Questions**
  - **Can these attacks be played over-the-air?**
  - **Are audio adversarial examples transferable?**
  - **Which existing defenses can be applied audio?**

# Follow-up works

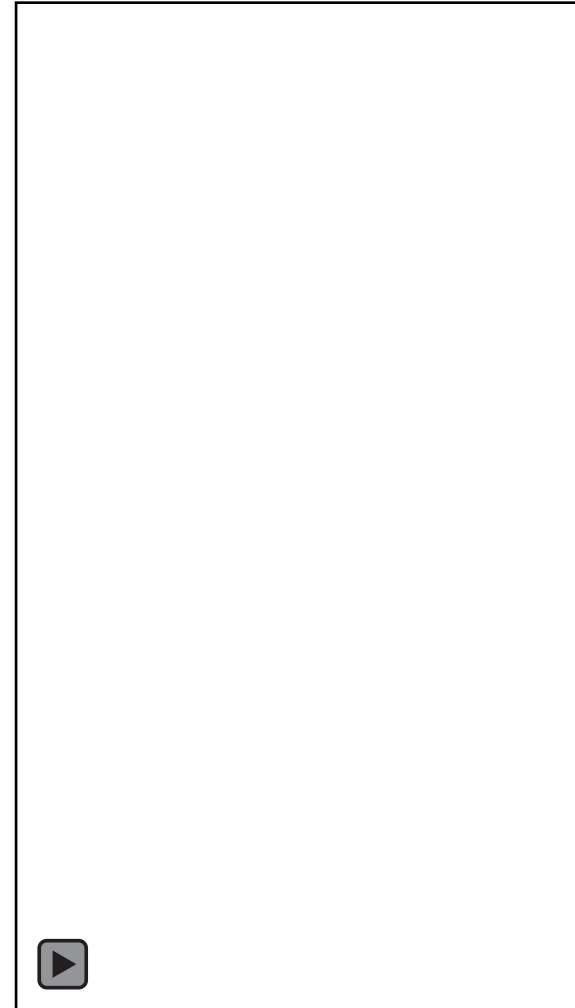
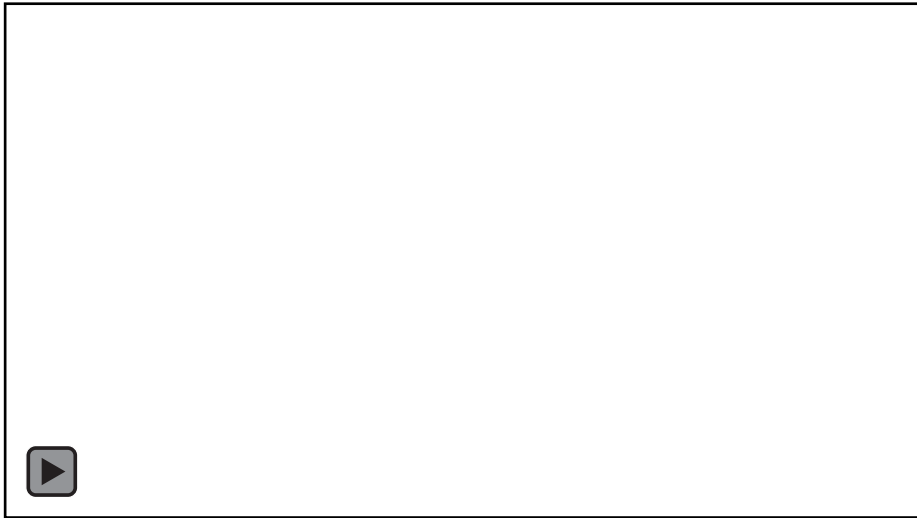
- **Yuan, Xuejing, et al. "CommanderSong: A Systematic Approach for Practical Adversarial Voice Recognition." (USENIX 2018)**
  - Integrate the commands into a song in a way that can be effectively recognized by ASR through the air, in the presence of background noise, while not being detected by a human listener.
- **Chen, Yuxuan, et al. "Devil's whisper: A general approach for physical adversarial attacks against commercial black-box speech recognition devices." (USENIX 2020)**
  - Enhance a simple local model roughly approximating the target black-box platform with a white-box model that is more advanced yet unrelated to the target

# Follow-up works





# Follow-up works



# Questions

## Question 1) Black-box Attack

From 고우영

**Q: This attack is also a white box attack. The DeepSpeech and Kaldi is an open-source model, but if the model weight is not opened, this attack can not be used. Is there any powerful black box ASR adversarial attack, and how much is the performance?**

- **As I mentioned in the follow-up works, there are some works.**
- **CommanderSong attack achieved 100% success rate.**
- **However, none of them accepted at the conference have shown successful demonstrations of commercial speech recognition systems.**

# Questions

## Question 2) Defense Method

From 한상구

**Q: As author said, all defenses to adversarial examples are focused on image domains. Are there any effective defense method published after this?**

- **Kwon, Hyun, Hyunsoo Yoon, and Ki-Woong Park. “ POSTER: Detecting Audio Adversarial Example through Audio Modification.” (CCS 2019)**
- **Yang, Zhuolin, et al. “Characterizing Audio Adversarial Examples Using Temporal Dependency.”**

# Questions

## Question 3) Real-world application

From Tuan Hoang Dinh (Best question)

**Q: The adversarial examples they constructed can not be played over the air, which leads to impractical attacks in the real world. Why people still pay attention on this topic? Does industrial care about this?**

- **Yuan, Xuejing, et al. " CommanderSong: A Systematic Approach for Practical Adversarial Voice Recognition." (USENIX 2018)**
- **CommanderSong attack was originally successful for Apple Siri, but Apple has since updated and blocked it.**

Thank you for listening