



Targeted Adversarial Examples for Black Box Audio Systems

Amog Kamsetty, Rohan Taori, Nikita Vemuri, Brenton Chu

Who Are We?

- Students at UC Berkeley
- Work done at Machine Learning @ Berkeley (ML@B)
 - ml.berkeley.edu
 - Aim to provide AI/ML opportunities at the undergraduate level



A decorative network diagram in the top-left corner, consisting of various sized grey circles (nodes) connected by thin grey lines (edges). Some nodes are solid grey, while others are hollow with a grey outline. The network is sparse and irregular, extending from the top-left towards the center.

What is an Adversarial Example?

A decorative network diagram in the bottom-right corner, similar to the one in the top-left. It features a cluster of grey nodes connected by lines, with some nodes being solid and others hollow. The network is more dense and complex than the one in the top-left, extending from the bottom-right towards the center.

Adversarial Example



x

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

=



$x +$


$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence


Untargeted vs Targeted

A decorative network diagram in the top right corner of the slide. It consists of several interconnected nodes, represented by circles of varying sizes and shades of gray, connected by thin lines. Some nodes are highlighted with a darker gray or blue color.

- Untargeted: Provide input to the model such that it misclassifies the adversarial input
 - Targeted: Provide input to the model so it classifies it as a predetermined target class
- 
- A decorative network diagram in the bottom left corner of the slide. It consists of several interconnected nodes, represented by circles of varying sizes and shades of gray, connected by thin lines. Some nodes are highlighted with a darker gray or blue color.

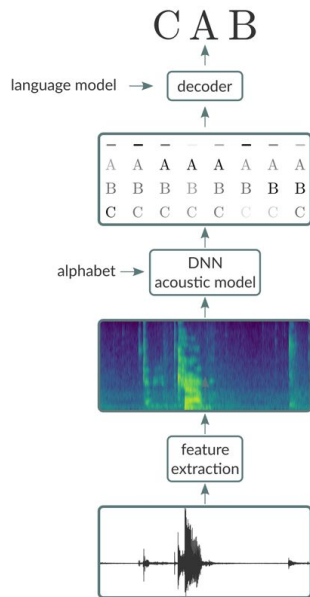
White Box vs Black Box

A decorative network diagram in the top right corner, consisting of interconnected nodes and lines, rendered in a light blue and grey color scheme.

- White box: complete knowledge of model architecture and parameters, allows for gradient computation
 - Black box: no knowledge of model or parameters except for output logits of model
- 
- A decorative network diagram in the bottom left corner, consisting of interconnected nodes and lines, rendered in a light blue and grey color scheme.

Why does this matter?


- Black box attacks can be of particular interest in ASR systems



- If we can create an adversarial audio file, we can trick the model into translating what we want
- If we do this with a black box approach, we can apply this to proprietary systems (ex. Google or IBM APIs)

Classical Adversarial Attacks

A decorative network diagram in the top right corner, consisting of various sized grey circles (nodes) connected by thin grey lines (edges). Some nodes are highlighted with a darker grey or blue color.

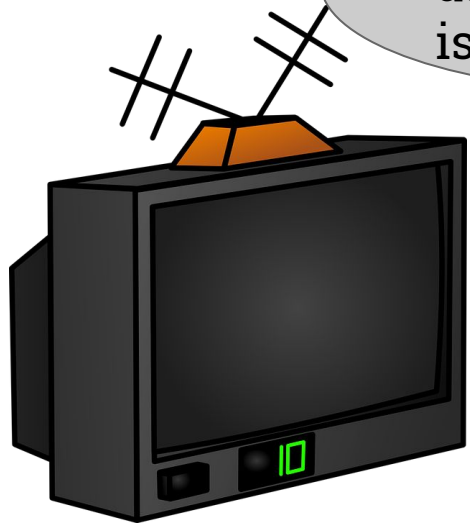
- Taking gradient iteratively
 - FGSM - Fast Gradient Sign Method
 - Houdini
- 
- A decorative network diagram in the bottom left corner, similar to the one in the top right, with grey nodes and edges, some nodes highlighted in darker shades.

Prior Work in Audio

- UCLA - Black box genetic algorithm on single word classes → softmax loss
- Carlini & Wagner: white box attack
 - CTC loss allows for comparison with arbitrary length translations
- Our project: Black box genetic algorithm on sentences using CTC Loss

Problem Statement

without the
dataset the article
is useless



+



**Adversarial
noise**

ok google
browse to
evil.com

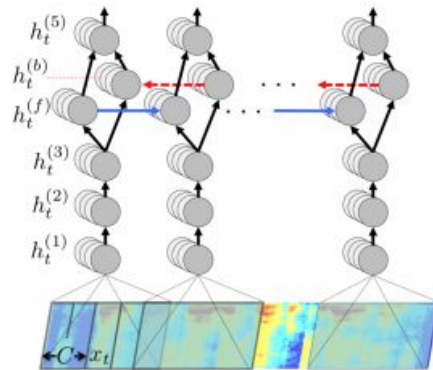


Problem Statement

- Black-box Targeted Attack
 - Given a target \mathbf{t} , a benign input \mathbf{x} , and model \mathbf{M} , perturb \mathbf{x} to form $\mathbf{x}' = \mathbf{x} + \delta$
 - S.t. $\mathbf{M}(\mathbf{x}') = \mathbf{t}$ while maximizing $cross_correlation(\mathbf{x}, \mathbf{x}')$
 - Only have access to logits of \mathbf{M}
 - Not given gradients!

Datasource: DeepSpeech

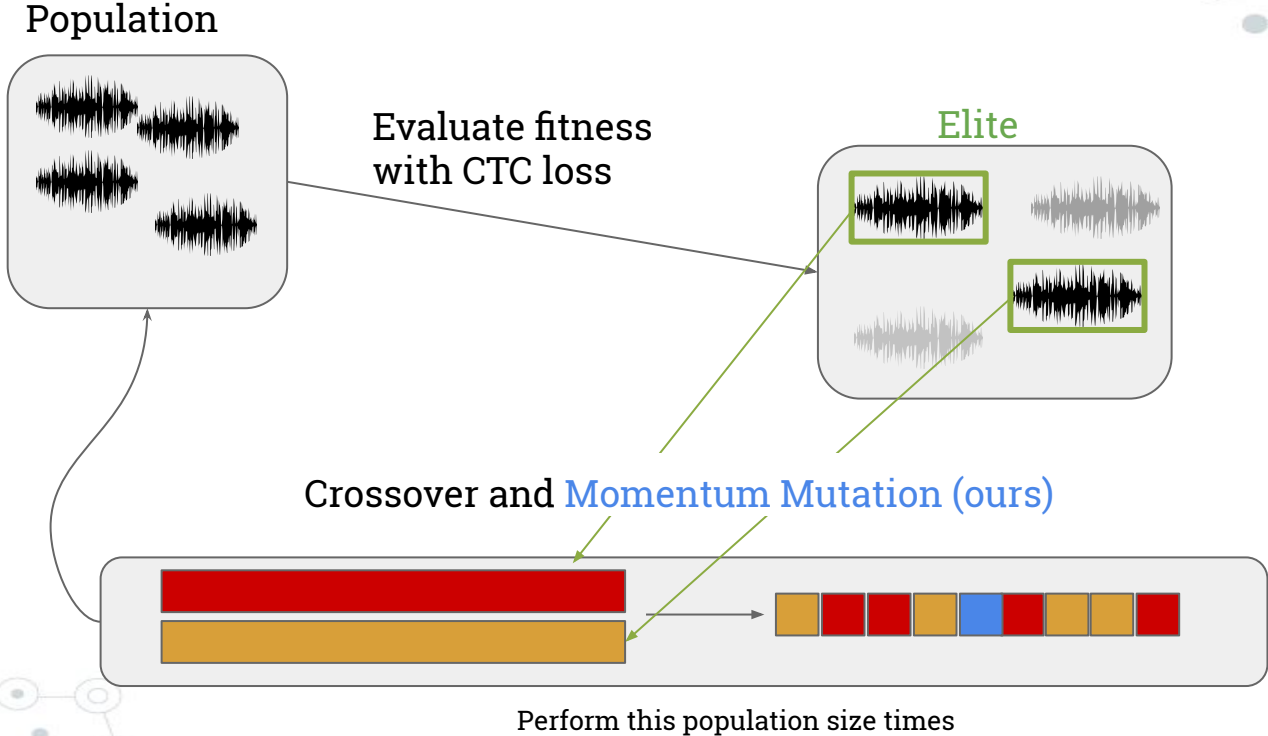
- The model we are targeting is DeepSpeech
 - Architecture created by Baidu
 - Tensorflow implementation by Mozilla; available on Github
- Utilize Common Voice dataset by Mozilla
 - Consists of voice samples
 - Sampling rate of 16 KHz



Final Algorithm: Guided Selection

- Genetic Algorithm approach
- Given the benign input, generate population of size 100
- On each iteration select the best 10 samples using scoring function
- Perform crossover and **momentum mutation**
- Apply high pass filter to added noise

Genetic Algorithm with Momentum



Momentum Mutation

$$p_{new} = (\alpha \times p_{old}) + \frac{\beta}{|currScore - prevScore|}$$

- Probability of mutation is function of difference in scores across iterations
- If little increase in score between iterations, increase “momentum” by increasing probability of mutation
- Encourages decoding to build up to target after making input similar to silence

Decodings while training

and you know it
and he nowit
nd he now
d he now
e now
a eload
elord
heloword
hello world

Gradient Estimation

- Genetic algorithms work best when search space is large
- However, when adversarial sample is near target, only few key perturbations are necessary
- Apply gradient estimation at 100 random indices

$$FD_x(g(x), \delta) = \begin{bmatrix} (g(x_1 + \delta) - g(x_1))/\delta \\ \vdots \\ (g(x_n + \delta) - g(x_n))/\delta \end{bmatrix}$$

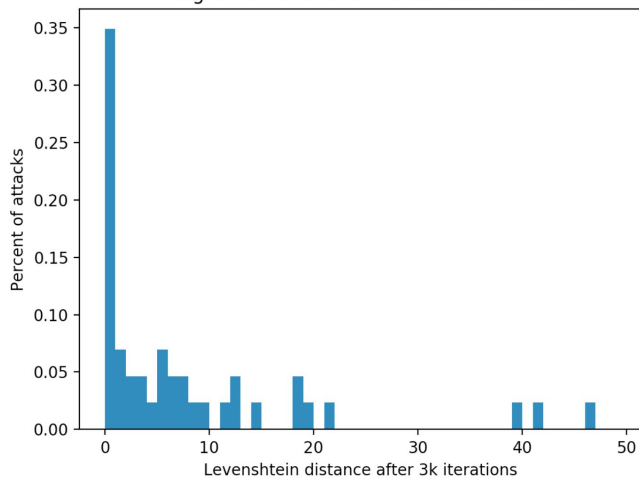
Results

- ⊙ Tested on first 100 samples of CommonVoice dataset
- ⊙ Randomly generated 2 target words
- ⊙ Targeted attack similarity: **89.25%**
 - Algorithm could almost always reach the target
- ⊙ Average similarity score: **94.6%**
 - Computed via wav-file cross-correlation

Results

Metric	White Box Attacks	Our Method	Single Word Black Box
Targeted attack success rate	100%	35%	87%
Average similarity score	99.9%	94.6%	89%
Similarity score method	cross-correlation	cross-correlation	human study
Loss used for attack	CTC	CTC	Softmax
Dataset tested on	Common Voice	Common Voice	Speech Commands
Target phrase generation	Single sentence	Two word phrases	Single word

Histogram of levenshtein distances of attacks



Example

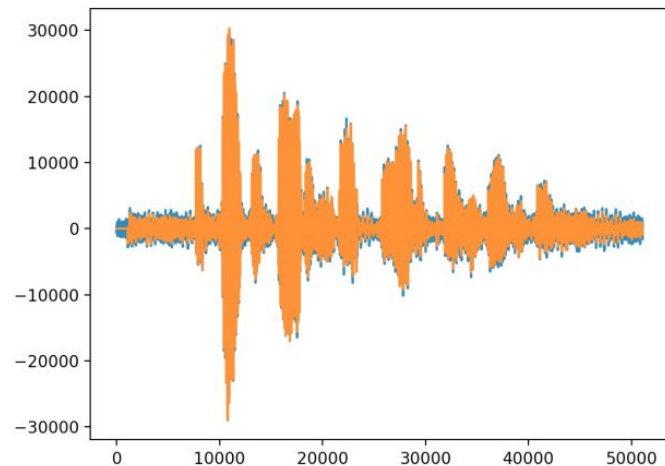
Original file: “and you know it”



Adversarial target: “hello world”



Audio Similarity: 94%
(cross-correlation)



Future Work

- Attack a broader range of models
 - Transferability across models
- Increasing sample efficiency to target
 - API call costs can be prohibitive
- Computational Efficiency

Thank You!

Code and samples:

<https://github.com/rtaori/Black-Box-Audio>