# Dompteur: Taming Audio Adversarial Examples
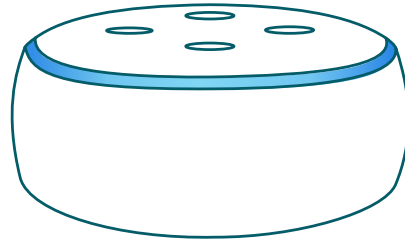
Thorsten Eisenhofer, Lea Schönherr, Joel Frank,
Lars Speckemeier, Dorothea Kolossa, Thorsten Holz

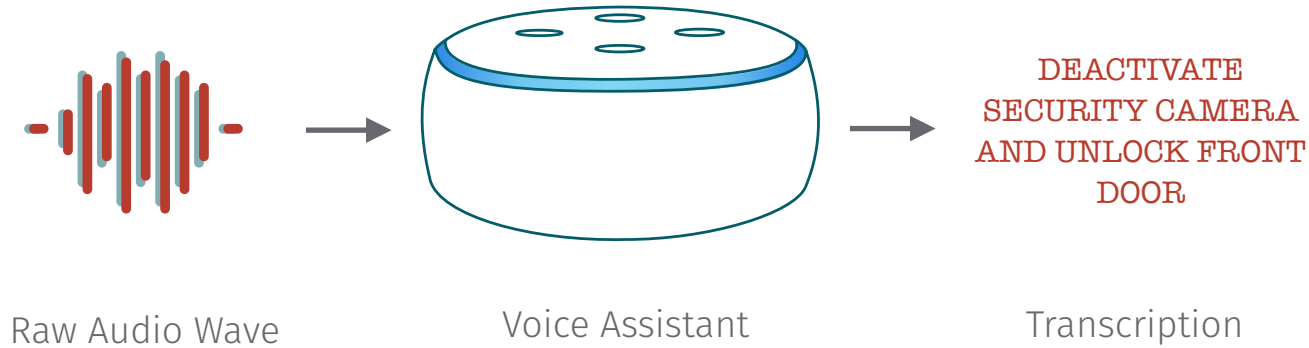USENIX Security Symposium 2021

CASA
Cyber Security in the Age
of Large-Scale Adversaries

RUHR
UNIVERSITÄT
BOCHUM

RUB

Raw Audio Wave

Voice Assistant

BIDS TOTALING SIX
HUNDRED FIFTY ONE
MILLION DOLLARS
WERE SUBMITTED

Transcription

Raw Audio Wave

Voice Assistant

DEACTIVATE SECURITY CAMERA AND UNLOCK FRONT DOOR

Transcription

Raw Audio Wave

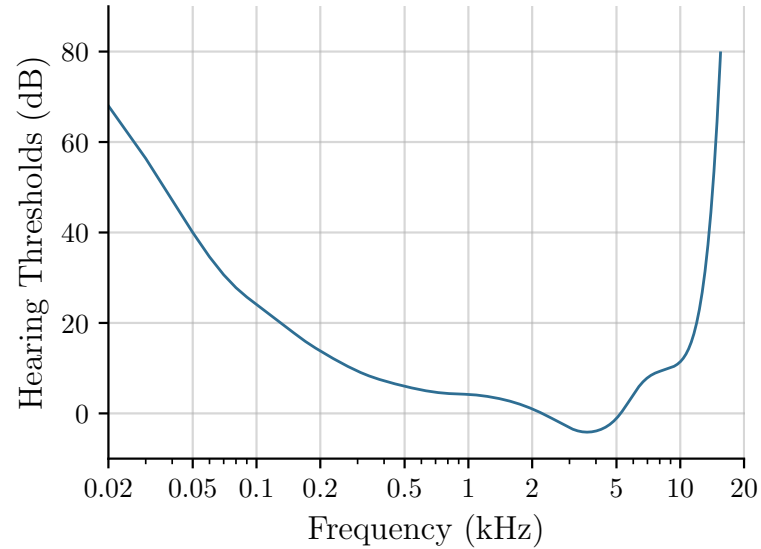Voice Assistant

DEACTIVATE SECURITY CAMERA AND UNLOCK FRONT DOOR

Transcription

*When we **accept** that **adversarial examples** exist, what **else** can we do?*
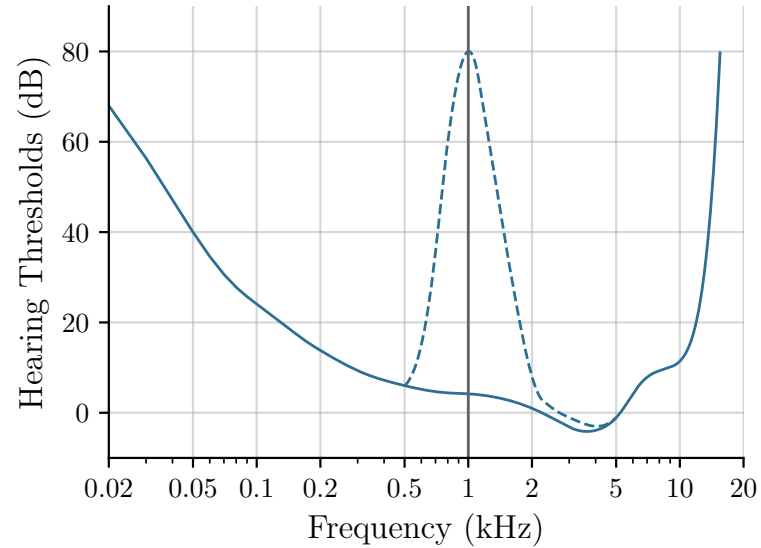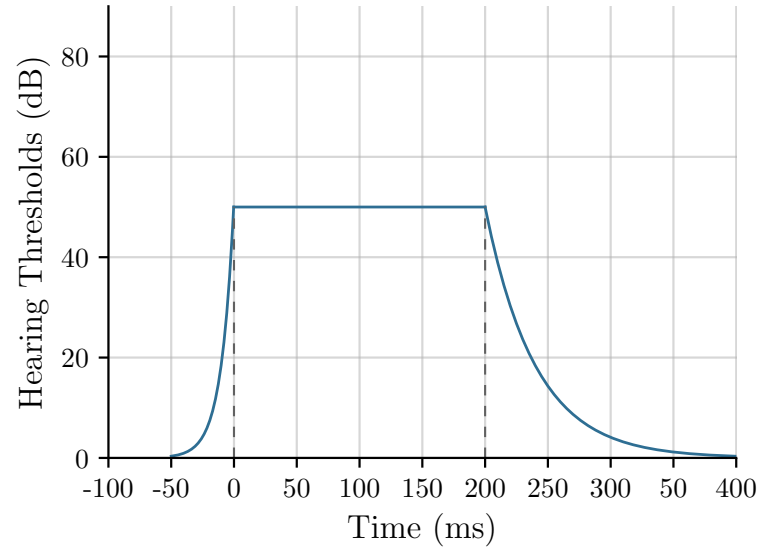
**Absolute Hearing Thresholds**
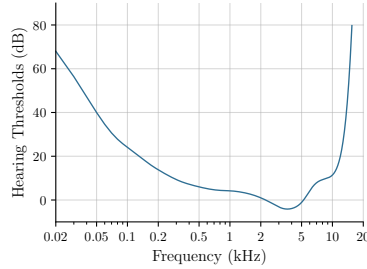
Gustav Fechner
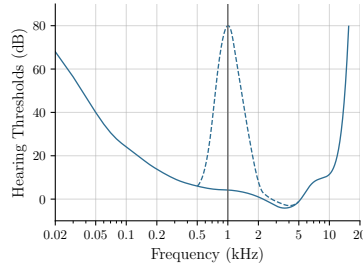1801-1887

**Frequency Masking**

Gustav Fechner
1801-1887

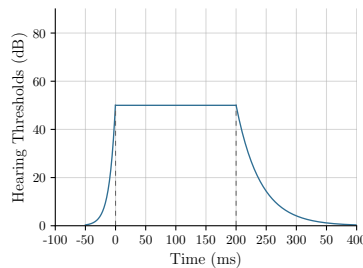**Temporal Masking**

Gustav Fechner
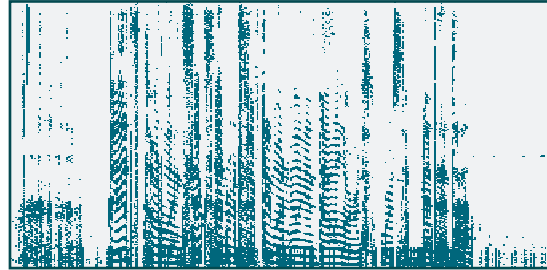1801-1887

Absolute Hearing Thresholds
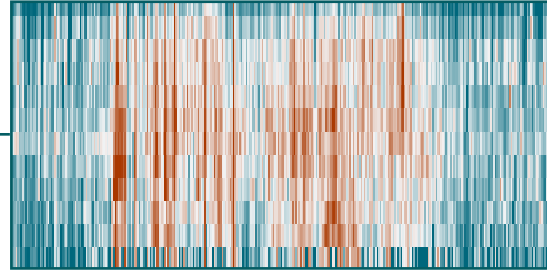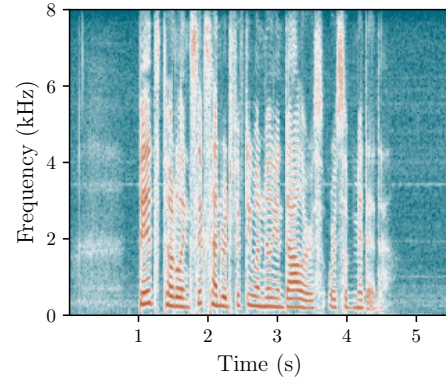
Frequency Masking
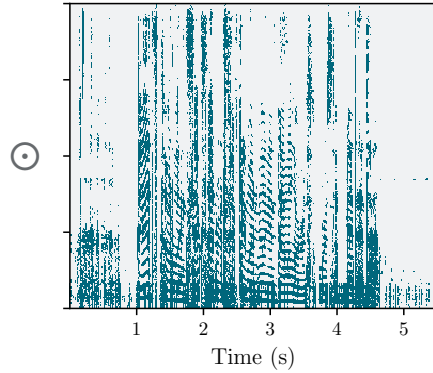
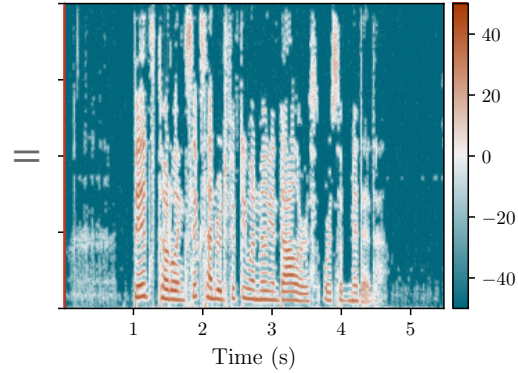Temporal Masking

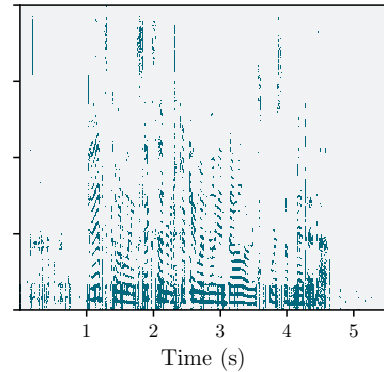Psychoacoustic Hearing Thresholds

Mask  M

Raw Audio Signal $\mathbf{S}$

Mask $\mathbf{M}$ $(\Phi = 0)$

Filtered $\mathbf{T} = \mathbf{S} \odot \mathbf{M}$

Mask $\mathbf{M}$ $(\Phi = 12)$

Filtered $\mathbf{T} = \mathbf{S} \odot \mathbf{M}$

Raw Audio Signal $\mathbf{S}$    Mask $\mathbf{M}$ $(\Phi = 0)$    Filtered $\mathbf{T} = \mathbf{S} \odot \mathbf{M}$

**Band-Pass Filter**

Raw Audio Wave → **Psychoacoustic Filtering** → **Band-Pass Filter** → Voice Assistant → Transcription: I SOLEMNLY SWEAR I AM UP TO NO GOOD

Implemented DOMPTEUR for **Kaldi** toolkit

Standard Input → Kaldi → 5.90% / 8.74%

Processed Input → Dompteur → 6.33% / 6.10%

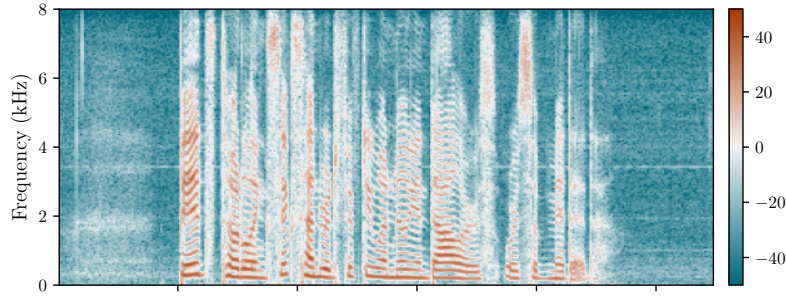Word Error Rate (WER)

**Adversarial Robustness**
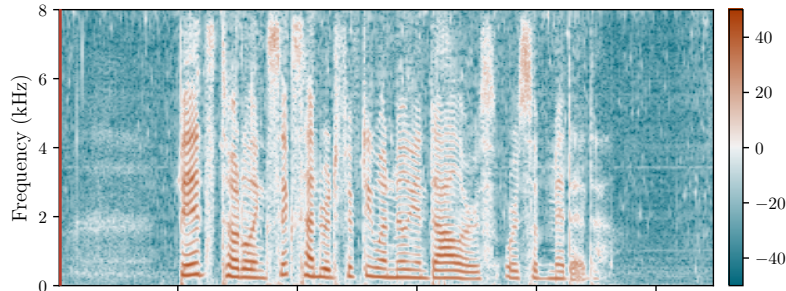
Strong adaptive, **white-box** attacker

**Successful** at computing adversarial
examples against **DOMPTEUR**

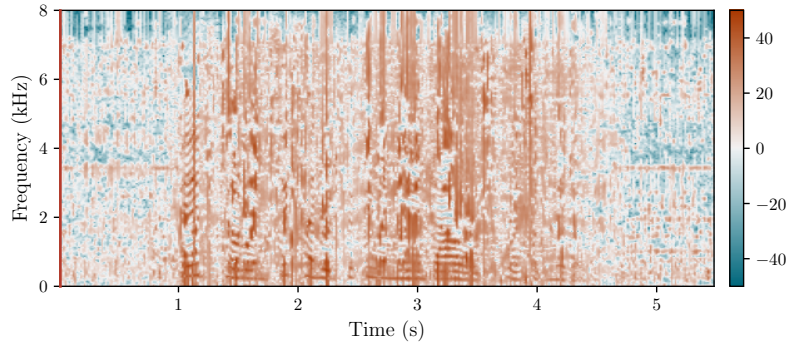**But** attack forced into audible
ranges and **clearly perceivable**

**Unmodified Signal**

BIDS TOTALING SIX
HUNDRED FIFTY ONE
MILLION DOLLARS WERE
SUBMITTED

**KALDI**

SEND SECRET FINANCIAL
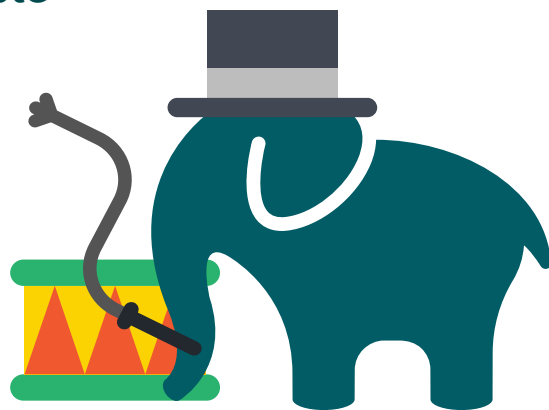REPORT

**DOMPTEUR** $\Phi = 12$

SEND SECRET FINANCIAL
REPORT

# Takeaways

**Adversarial examples** seem to be **inevitable**

New Perspective: Make attack **noticeable**

**Psychoacoustics** effective to force attack into **audible ranges**

Code, Examples and Models available at
github.com/rub-syssec/dompteur

# Thank you!