# MaskMark: Robust Neural Watermarking for Real and Synthetic Speech

Patrick O'Reilly[1], Zeyu Jin[2], Jiaqi Su[2], Bryan Pardo[1]

*IEEE ICASSP 2024*

1. Northwestern University
2. Adobe Research

# MaskMark: Robust Neural Watermarking for Real and Synthetic Speech

(Listening examples)

# TL;DR:

In this work, we show how to **hide a binary vector in audio** that can be recovered even when the audio has been altered significantly.

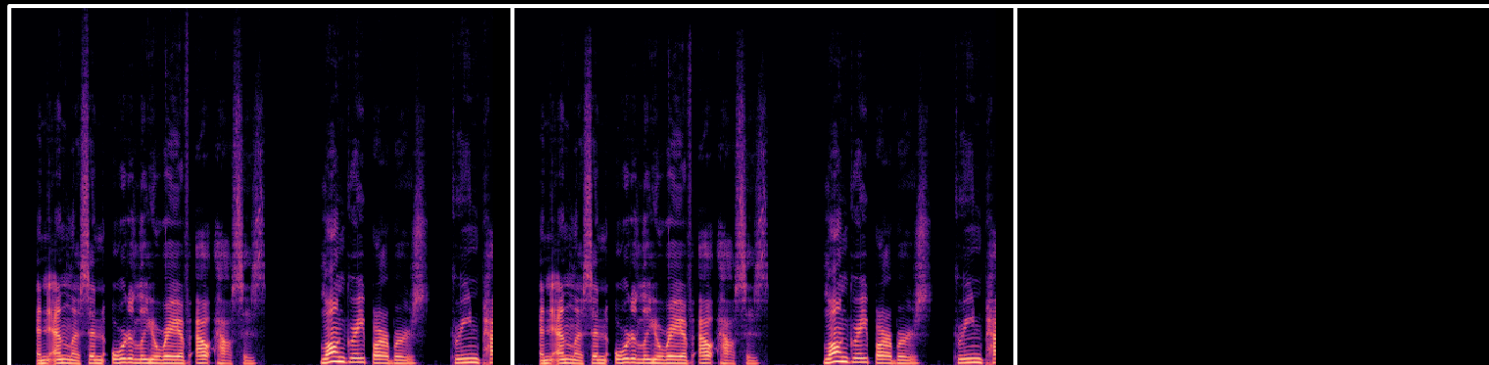Let's look at some examples.

This audio has no hidden vector

Clean

This audio has no hidden vector
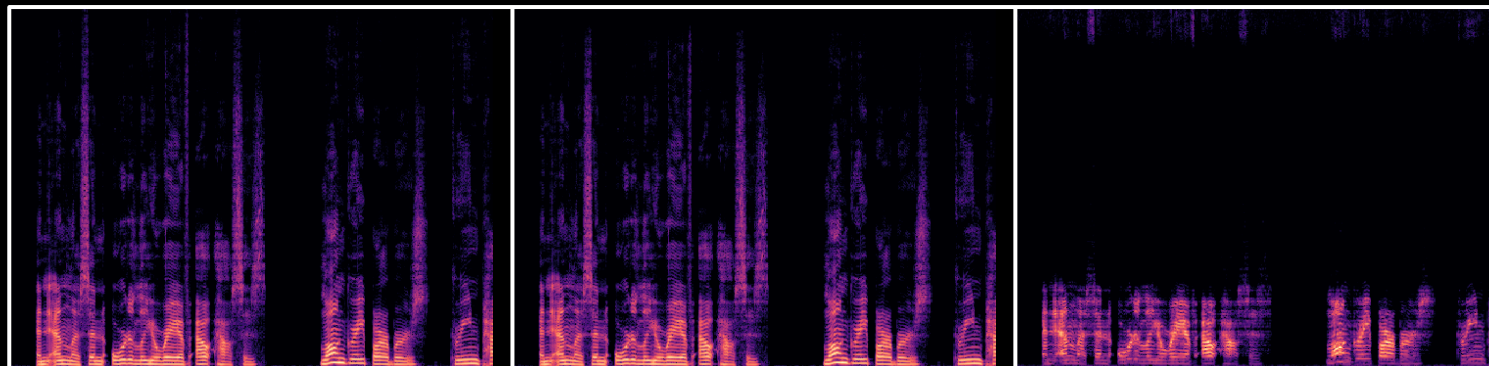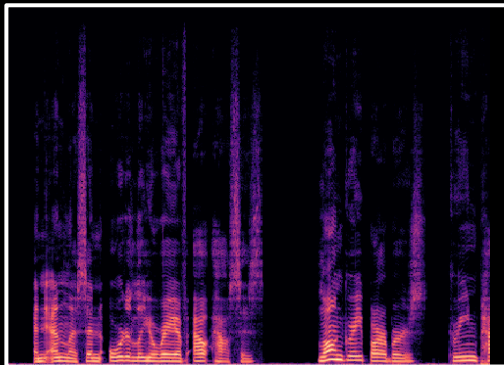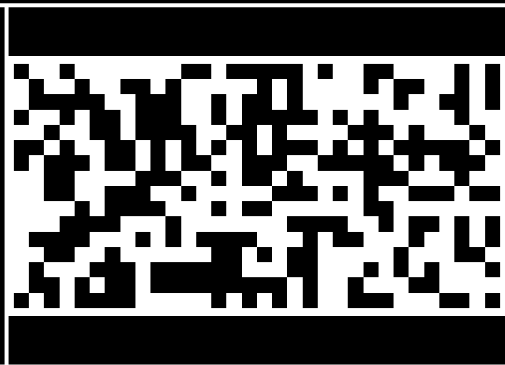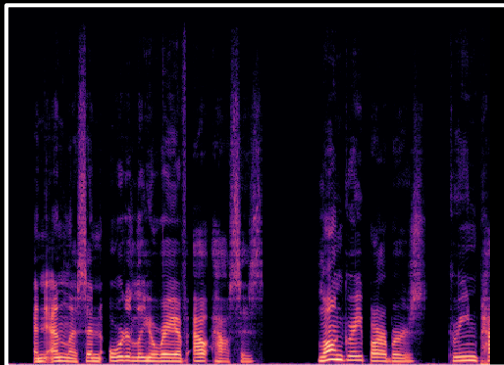
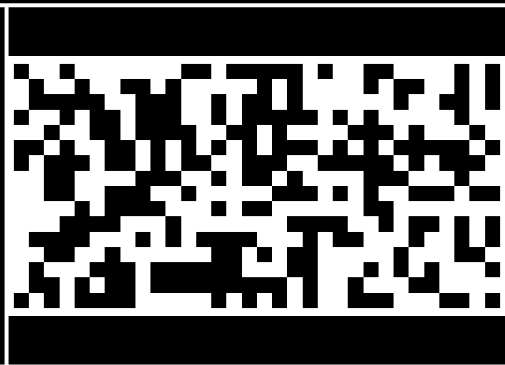This audio has a hidden vector

Clean
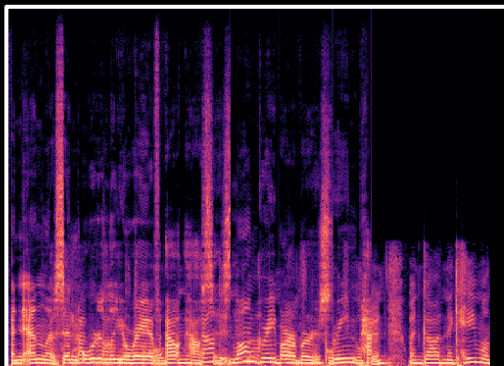
Watermarked

Normalized Difference
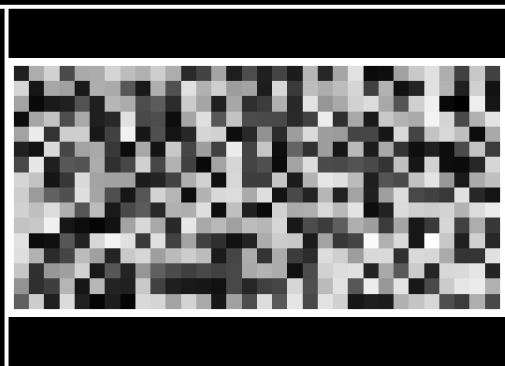
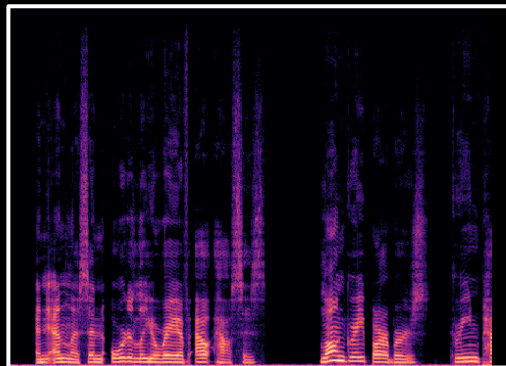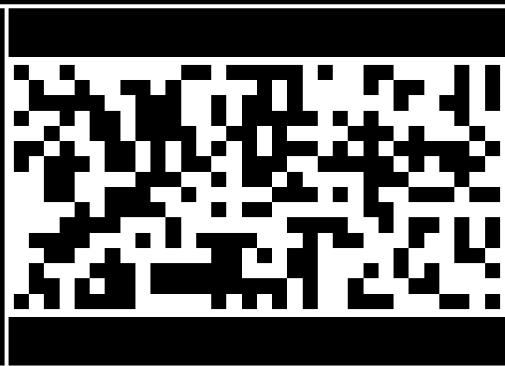**Watermarked**          **Embedded key vector**

**Watermarked**

**Embedded key vector**

**Simulated editing**
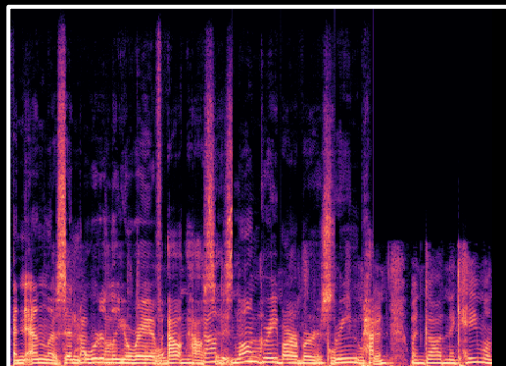
**Recovered key vector
(logits)**

**Watermarked**      **Embedded key vector**

(random chance is 50%)

**100% match**

**Simulated editing**      **Recovered key vector (logits)**      **Recovered key vector (quantized)**

**Watermarked**                    **Embedded key vector**

**Watermarked**

**Embedded key vector**

**HiFiGAN resynthesis**

**Recovered key vector (logits)**

**Watermarked**  **Embedded key vector**

99% match

**HiFiGAN resynthesis**  **Recovered key vector (logits)**  **Recovered key vector (quantized)**

**Watermarked**                    **Embedded key vector**

**Watermarked**

**Embedded key vector**

**Simulated over-the-air**

**Recovered key vector
(logits)**

**Watermarked**      **Embedded key vector**

**83% match**

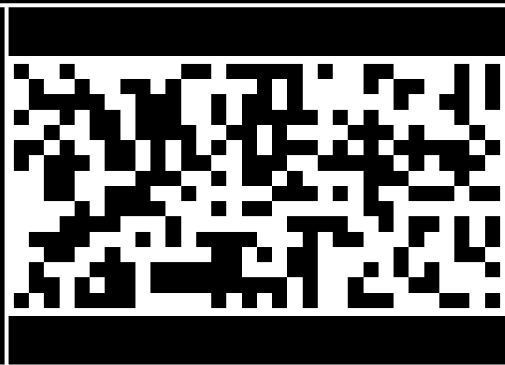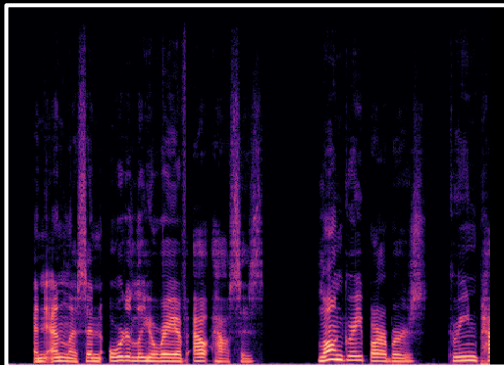**Simulated over-the-air**      **Recovered key vector (logits)**      **Recovered key vector (quantized)**

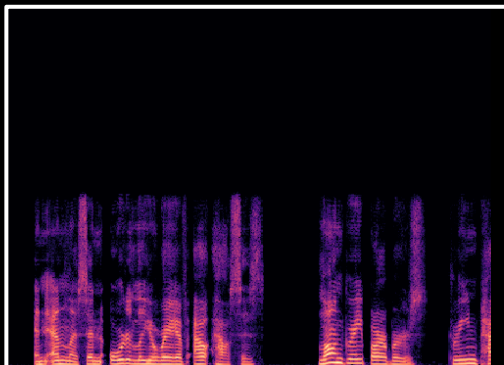Why should we care about hiding binary vectors in audio clips?

# 2016: WaveNet

**Expertise** + **compute** + **a large single-speaker dataset** + **lots of time**

=

# 2016: WaveNet

**Expertise** + **compute** + **a large single-speaker dataset** + **lots of time**

**=**

# 2023: Suno Bark

$0 + 1-10 min. audio + 5 min. editing

=

# 2023: Suno Bark

**$0** + **1-10 min. audio** + **5 min. editing**

**=**

# 4chan users embrace AI voice clone tool to generate celebrity hatespeech



Illustration by Alex Castro / The Verge

/ Free AI voice cloning technology from startup ElevenLabs has been used by trolls to imitate the voices of celebrities. The generated audio ranges in content from memes and erotica to virulent hatespeech.

By James Vincent, a senior reporter who has covered AI, robotics, and more for eight years at The Verge.

Jan 31, 2023, 5:00 AM PST | 💬 7 Comments / 7 New

# AI-Enabled Voice Cloning Anchors Deepfaked Kidnapping

Virtual kidnapping is just one of many new artificial intelligence attack types that threat actors have begun deploying, as voice cloning emerges as a potent new imposter tool.

# What can speech synthesis providers do?

We can **hide** a **message** in all the audio we generate

We can **hide** a **message** in all the audio we generate

We can **check** any audio for the message

We can **hide** a **message** in all the audio we generate

We can **check** any audio for the message

If we find the message, the audio was generated by our system

**"embed"**

We can **hide** a **message** in all the audio we generate

We can **check** any audio for the message

If we find the message, the audio was generated by our system

"embed"

We can **hide** a **message** in all the audio we generate

"watermark key"

We can **check** any audio for the message

If we find the message, the audio was generated by our system

"embed"

We can **hide** a **message** in all the audio we generate

"watermark key"     `[0, 1, 1, 0, 0, 1, … ]`

$n$ bits

We can **check** any audio for the message

If we find the message, the audio was generated by our system

"embed"

We can **hide** a **message** in all the audio we generate

"watermark key"     [0, 1, 1, 0, 0, 1, … ]

$n$ bits

We can **check** any audio for the message

If we find the message, the audio was generated by our system

"embed"

We can **hide** a **message** in all the audio we generate

"watermark key"

`[0, 1, 1, 0, 0, 1, … ]`

$n$ bits

We can **check** any audio for the message

If we find the message, the audio was generated by our system

**"embed"**

**We can hide a message in all the audio we generate**

**"watermark key"**

`[0, 1, 1, 0, 0, 1, … ]`

*n* bits

**We can check any audio for the message**

**"detect"**

**If we find the message, the audio was generated by our system**

# Watermarking

# Watermarking

**Embed** the watermark

# Watermarking

**Embed** the watermark

watermark key    [0, 1, 1, 0, 0, 1, … ]

E

# Watermarking

**Embed** the watermark



watermark key    [0, 1, 1, 0, 0, 1, … ]

E

# Watermarking



**Embed** the watermark

**Detect** the watermark

watermark key  [0, 1, 1, 0, 0, 1, … ]

E

# Watermarking

**Embed** the watermark

**Detect** the watermark

watermark key  [0, 1, 1, 0, 0, 1, … ]

E

D

# Watermarking



**Embed** the watermark

watermark key    [0, 1, 1, 0, 0, 1, … ]

E

**Detect** the watermark

D   → score (big)

# Watermarking

**Embed** the watermark

watermark key  [0, 1, 1, 0, 0, 1, … ]

E

**Detect** the watermark

D

score **(small)**

# Watermarking

**Embed** the watermark

watermark key    `[0, 1, 1, 0, 0, 1, … ]`

E

**Detect** the watermark

≠

D

score **(small)**

# Desiderata

# Desiderata

1.     ≈ 

# Desiderata

1.  ≈ 

2. $\left| \left\{ \phantom{x} \right\} \right|$ is big

# Desiderata

1.  ≈ 

2. $\left|\left\{\text{}\right\}\right|$ **is big**

3.  **is hard to remove from** 

# Desiderata

1.  ≈   **Perceptual transparency**

   *watermark doesn't ruin user experience*

2. $\left| \{ \text{🔴} \} \right|$ **is big**

3. 🔴 **is hard to remove from** 

# Desiderata

**1.**  ≈ 

                                        **Perceptual transparency**

                               *watermark doesn't ruin user experience*

**2.** $\left| \left\{ \bullet \right\} \right|$ **is big**

                                              **Capacity**

                              *can hide info like user IDs in the watermark*

**3.**  **is hard to remove from** 

# Desiderata

**1.**  ≈ 

**Perceptual transparency**

*watermark doesn't ruin user experience*

**2.**  **is big**

**Capacity**

*can hide info like user IDs in the watermark*

**3.**  **is hard to remove from** 

**Robustness**

*watermark works under realistic conditions*

# Desiderata

# Desiderata



Perceptual transparency

Capacity

Robustness

Hard to remove, but low information capacity

# Desiderata

# Desiderata

# Balancing these is hard!

# Robustness

# Robustness

# Robustness



channel

**C**

noise
reverb
compression
pitch shift
…

**D**

score **(big)**

We only need 1 bit to answer "fake or not?"

# Desiderata

How can we **robustly** and **transparently** hide a little information in audio?

# "EigenWatermark" (Tai & Mansour 2019)



**Clean**

# "EigenWatermark" (Tai & Mansour 2019)



Clean

Watermarked

# "EigenWatermark" (Tai & Mansour 2019)



Clean         Watermarked         Normalized Difference

# "EigenWatermark" (Tai & Mansour 2019)



**Watermarked**

**Completely breaks the watermark!**

# "EigenWatermark" (Tai & Mansour 2019)



Watermarked        Speed up 2%

Completely breaks the watermark!

How can we **robustly** and **transparently** hide a little information in audio?

**Let's make** $\boxed{E}$ **and** $\boxed{D}$ **neural networks**



[0, 1, 1, 0, 0, 1, … ]

# Let's make  E  and  D  neural networks

[0, 1, 1, 0, 0, 1, … ]

E

# Let's make [E] and [D] neural networks

# Let's make E and D neural networks

# Let's make E and D neural networks

Let's make **E** and **D** neural networks

cross-entropy loss

[0, 1, 1, 0, 0, 1, … ]

E

D

[-2.7, 10.1, 3.5, -0.1, -20.2, 3.6, … ]

> 0.5

[0,   1,   1,   0,   0,   1,   … ]

STFT loss

What are we missing?

# Let's make E and D neural networks



cross-entropy loss

STFT loss

channel

C

E

D

[0, 1, 1, 0, 0, 1, … ]

[-2.7, 10.1, 3.5, -0.1, -20.2, 3.6, … ]

> 0.5

[0,   1,   1,   0,   0,   1,   … ]

# People have tried this!

| | Less robust ← → More robust | DNN-A (Pavlovic et al. 2022) | EigenWatermark (Tai & Mansour 2019) |
|---|---|---|---|
| Sample rate | | 16kHz | 44.1kHz |
| Required audio length | | 2s | 1s |
| Robustness | | | |
| Signal-processing | | | |
| Neural audio codec | | | |
| Neural vocoder | | | |
| Neural denoiser | | | |

| | Less robust ⟶ More robust | DNN-A (Pavlovic et al. 2022) | EigenWatermark (Tai & Mansour 2019) |
|---|---|---|---|
| Sample rate | | 16kHz | 44.1kHz |
| Required audio length | | 2s | 1s |
| Robustness | | | |
| Signal-processing | | TPR @ 1% FPR = 0.00 | |
| Neural audio codec | | | |
| Neural vocoder | | | |
| Neural denoiser | | | |

| | Less robust     More robust | DNN-A (Pavlovic et al. 2022) | EigenWatermark (Tai & Mansour 2019) |
|---|---|---|---|
| Sample rate | | 16kHz | 44.1kHz |
| Required audio length | | 2s | 1s |
| **Robustness** | | | |
| Signal-processing | | TPR @ 1% FPR = 0.00 | |
| Neural audio codec | | 0.00 | |
| Neural vocoder | | 0.00 | |
| Neural denoiser | | 0.00 | |

| | Less robust → More robust | DNN-A (Pavlovic et al. 2022) | EigenWatermark (Tai & Mansour 2019) |
|---|---|---|---|
| Sample rate | | 16kHz | 44.1kHz |
| Required audio length | | 2s | 1s |
| Robustness | | | |
| Signal-processing | | TPR @ 1% FPR = 0.00 | 0.73 |
| Neural audio codec | | 0.00 | 0.39 |
| Neural vocoder | | 0.00 | 0.01 |
| Neural denoiser | | 0.00 | 1.00 |

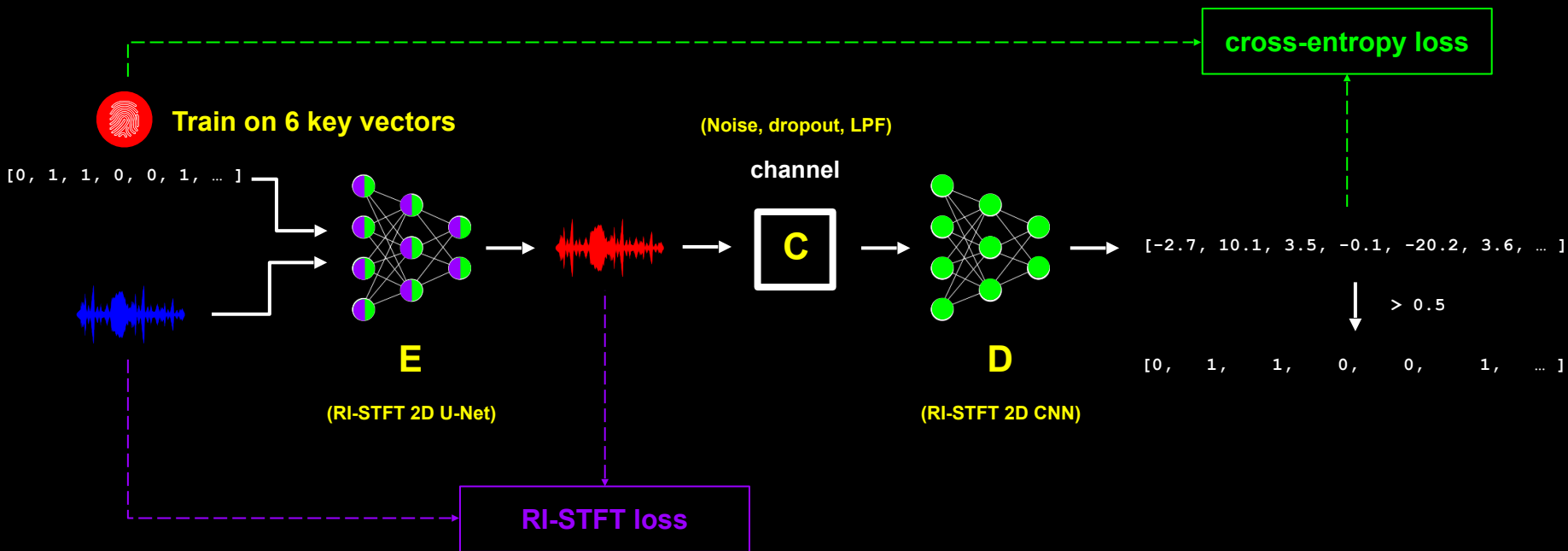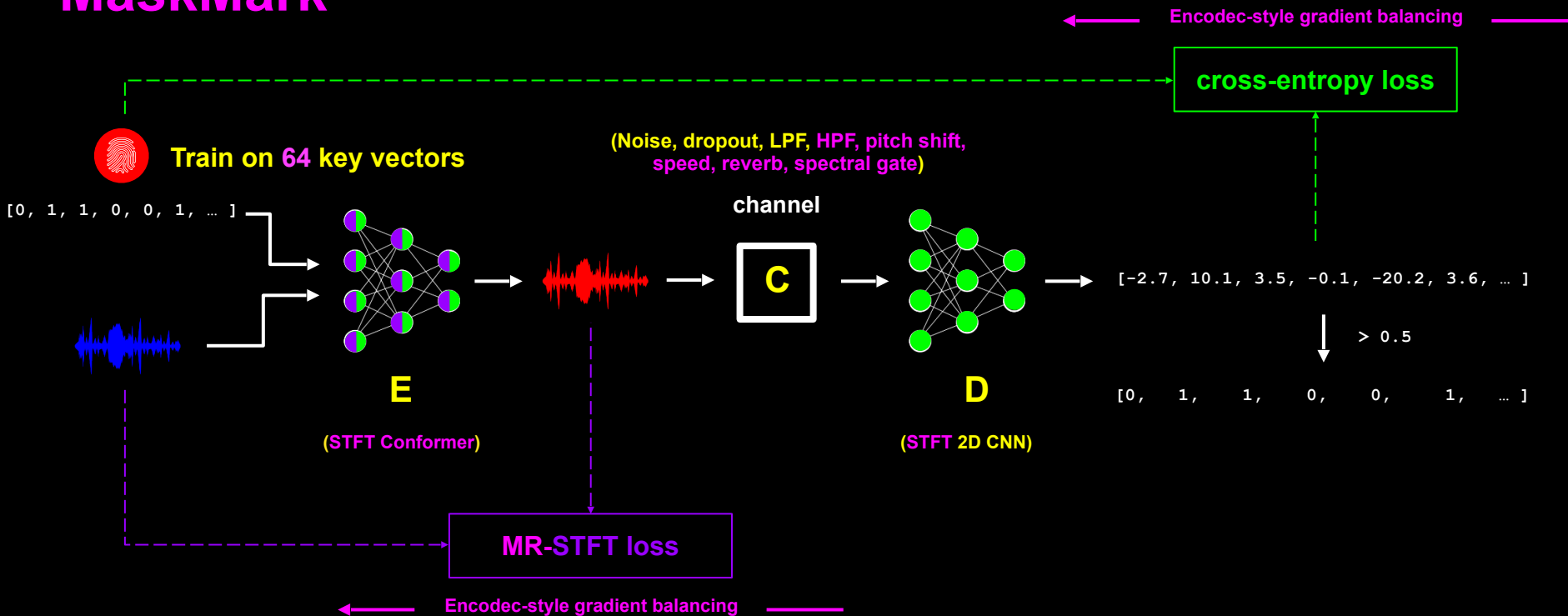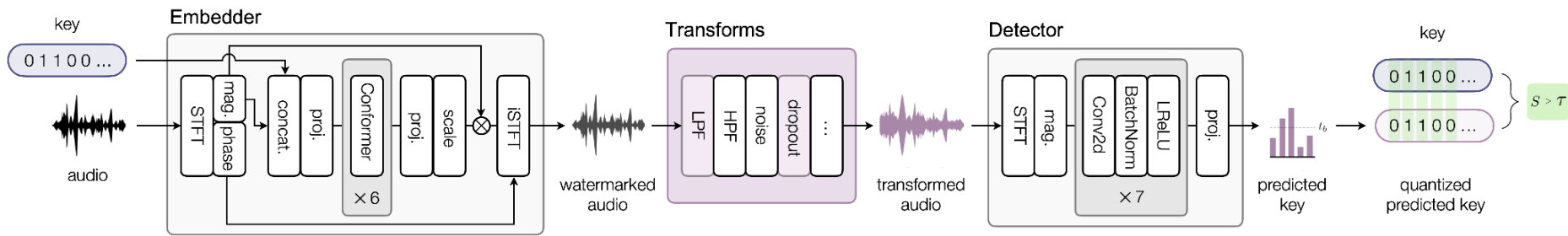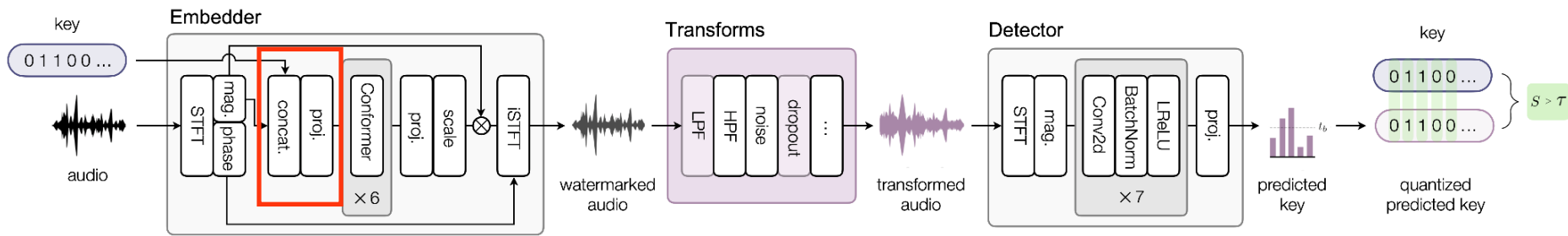| | DNN-A | Eigen | MaskMark |
|---|---|---|---|
| **Less robust** ➞ **More robust** | | | |
| Sample rate | 16kHz | 44.1kHz | 48kHz |
| Required audio length | 2s | 1s | 1s |
| **Robustness** | | | |
| Signal-processing | TPR @ 1% FPR = 0.00 | 0.73 | 0.97 |
| Neural audio codec | 0.00 | 0.39 | 0.45 |
| Neural vocoder | 0.00 | 0.01 | 0.82 |
| Neural denoiser | 0.00 | 1.00 | 0.99 |

DNN-A, Pavlovic et al. (2022)
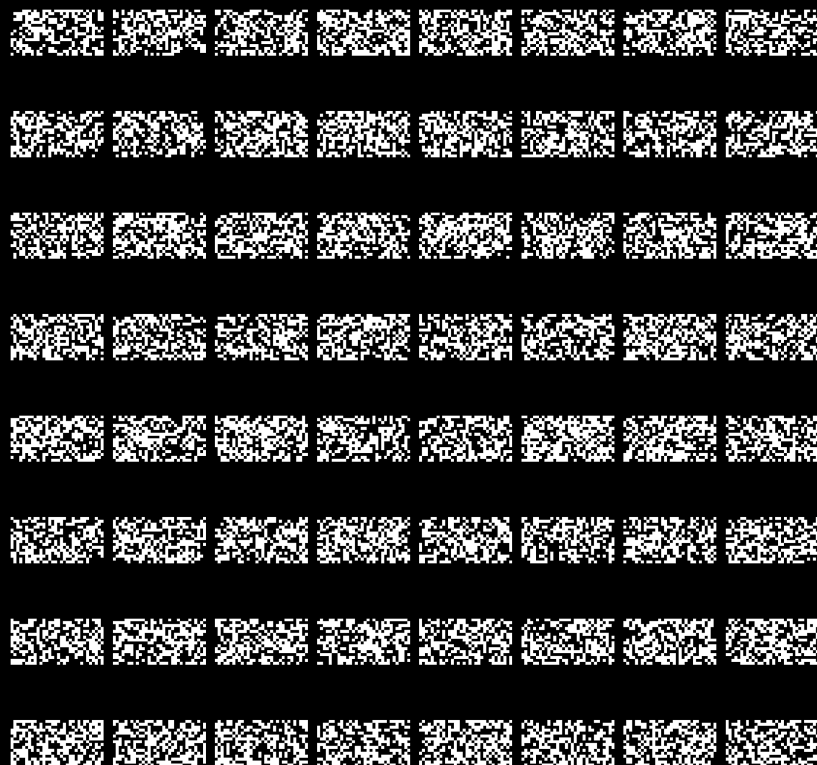
# Architecture details

# Architecture details



We hide the same watermark
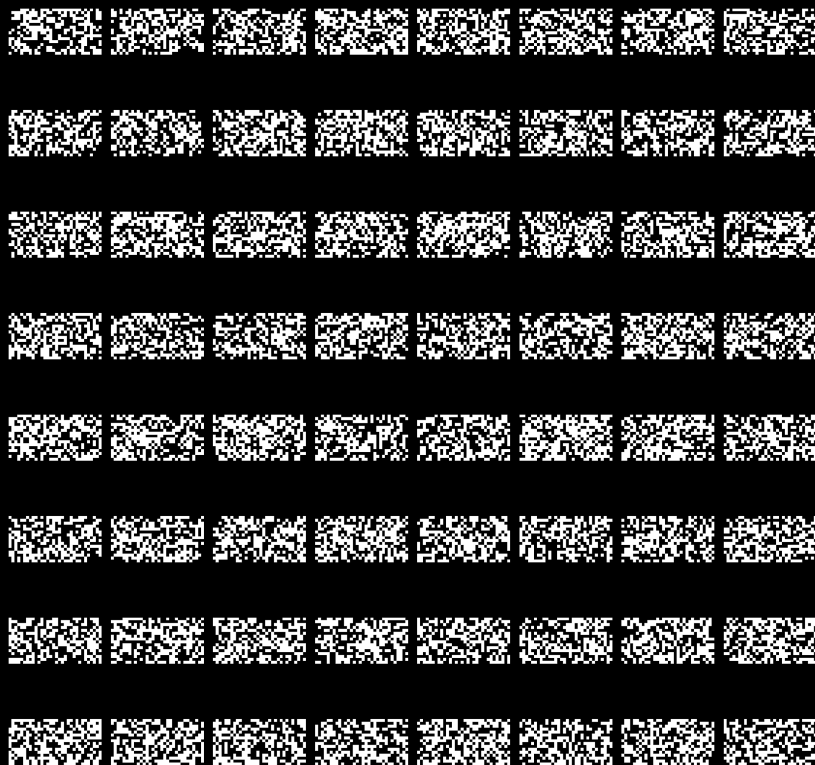information in every frame

No explicit learning of an "un-watermarked" class!

**64 learned watermarks**

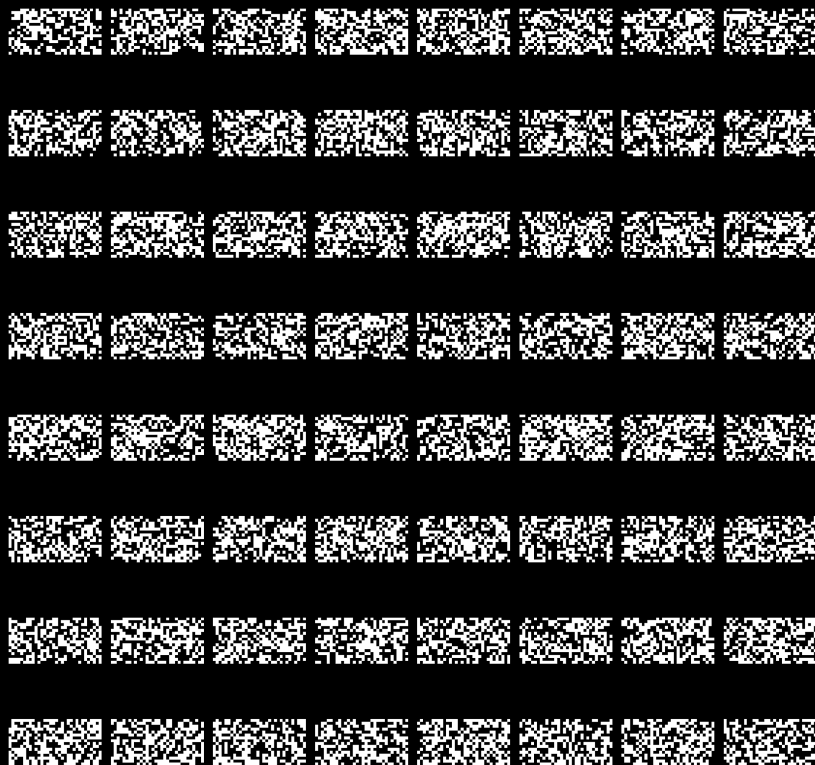Detector prediction for **unwatermarked** audio (essentially random)

64 learned watermarks

~50% bit accuracy

Detector prediction for
**unwatermarked** audio
(essentially random)

64 learned watermarks

~100% bit accuracy
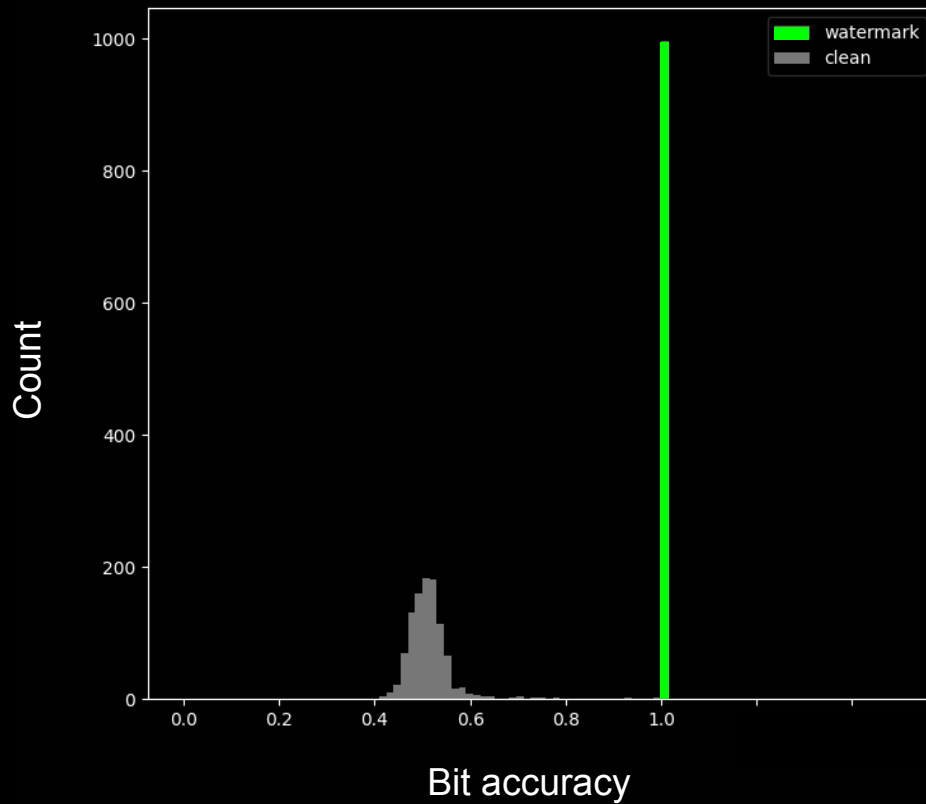
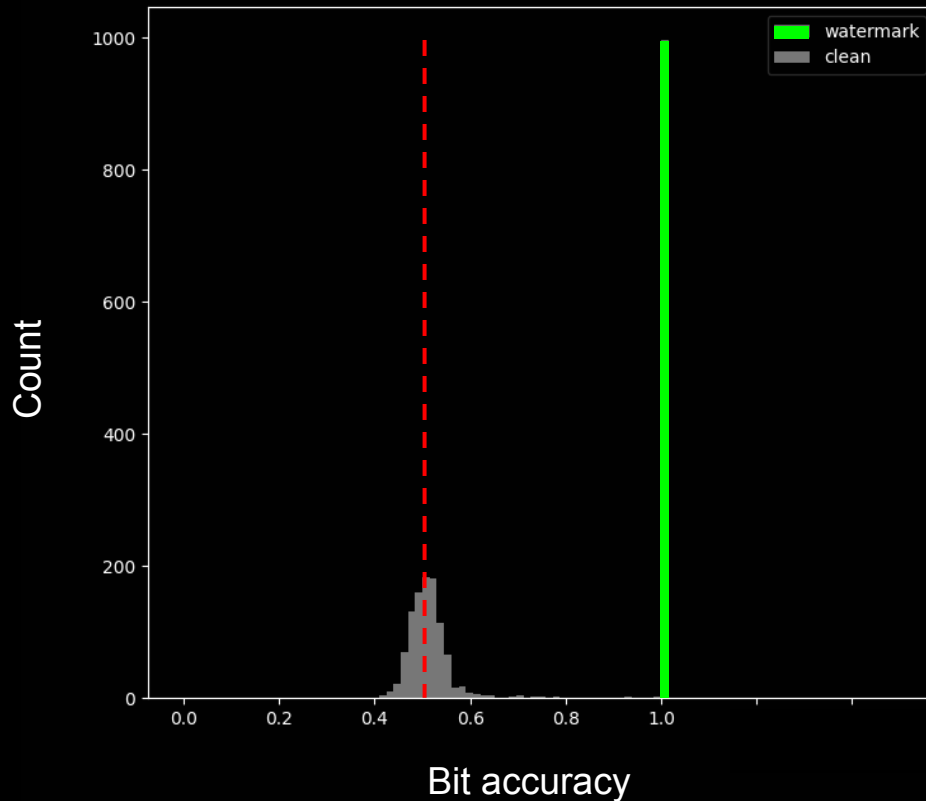Detector prediction for
audio with watermark 0

64 learned watermarks

Bit accuracy vs. known key vector, watermarked & un-watermarked audio

Bit accuracy vs. known key vector, watermarked &
un-watermarked audio

~50% bit accuracy for unwatermarked

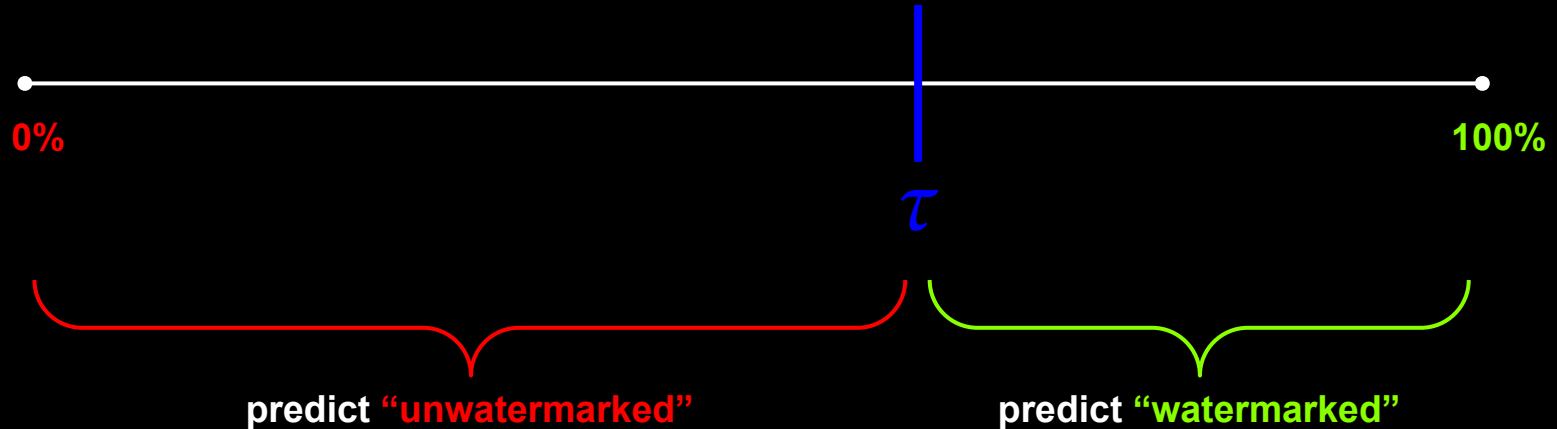~100% bit accuracy for watermarked

This lets us **distinguish between watermarked and unwatermarked audio** using bit accuracy

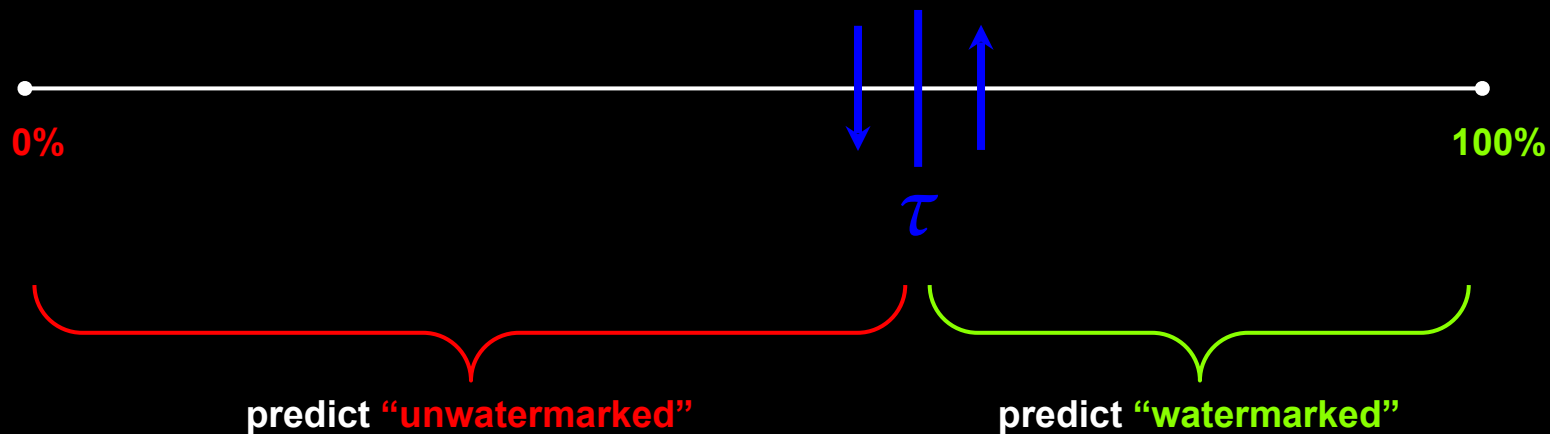This lets us **distinguish between watermarked and unwatermarked audio** using bit accuracy

0%                                                                    100%

When targeting a low (1%) FPR, our approach outperforms recent signal-processing and neural-network watermarks!

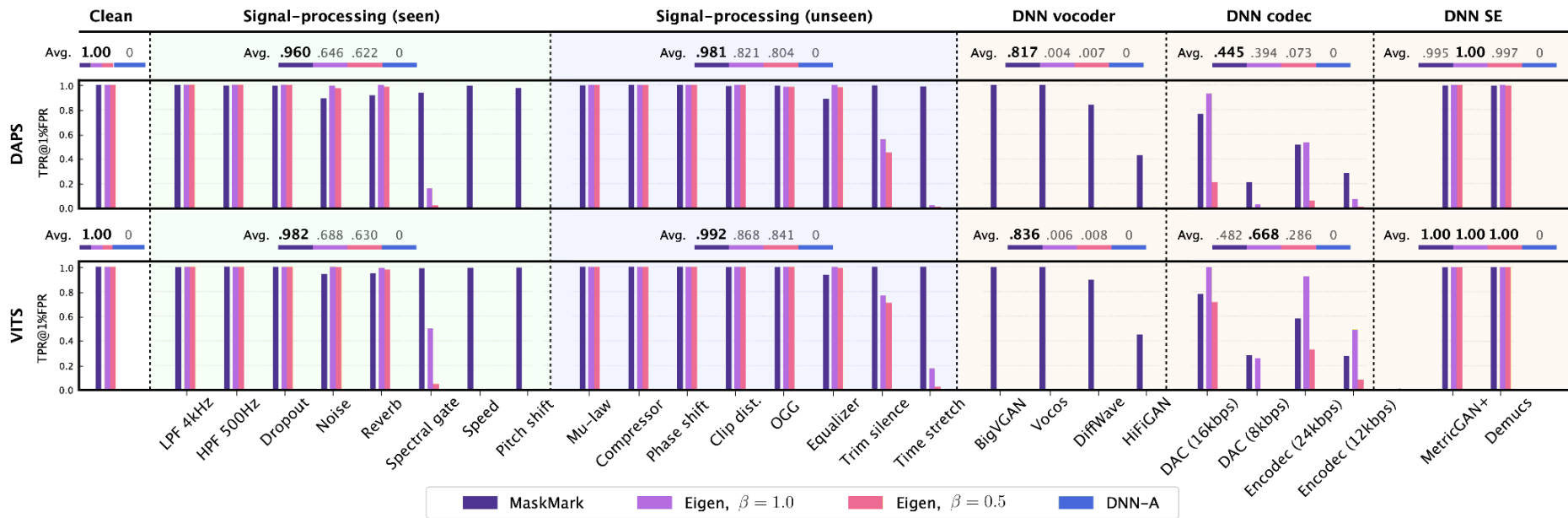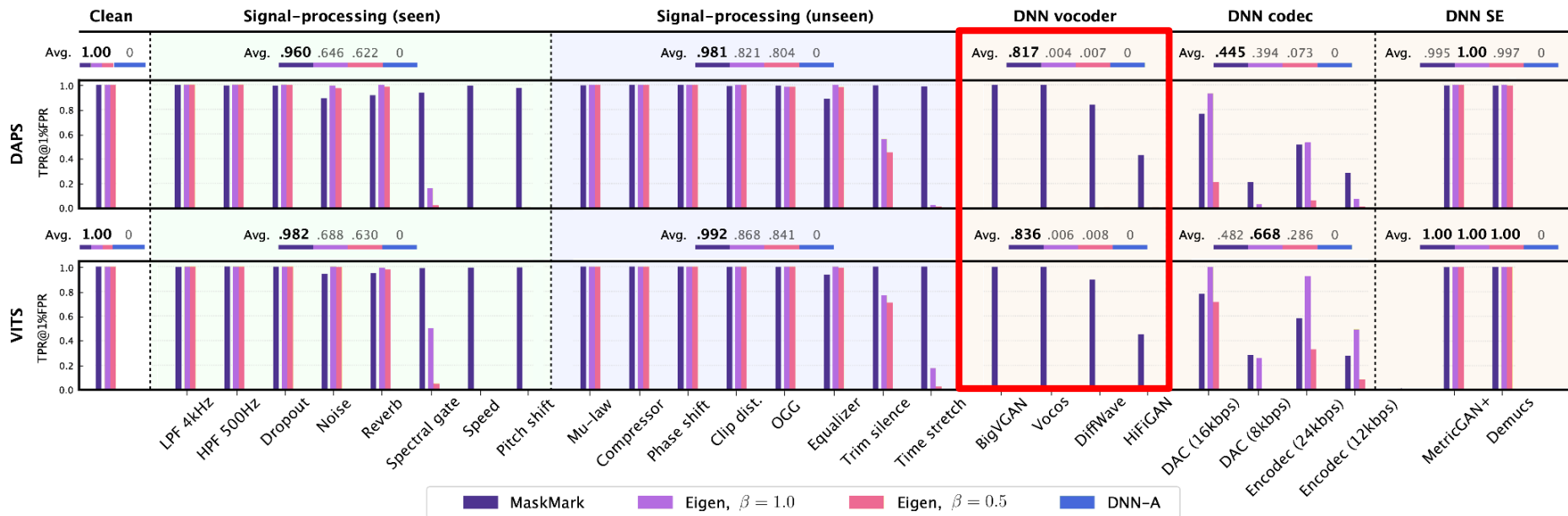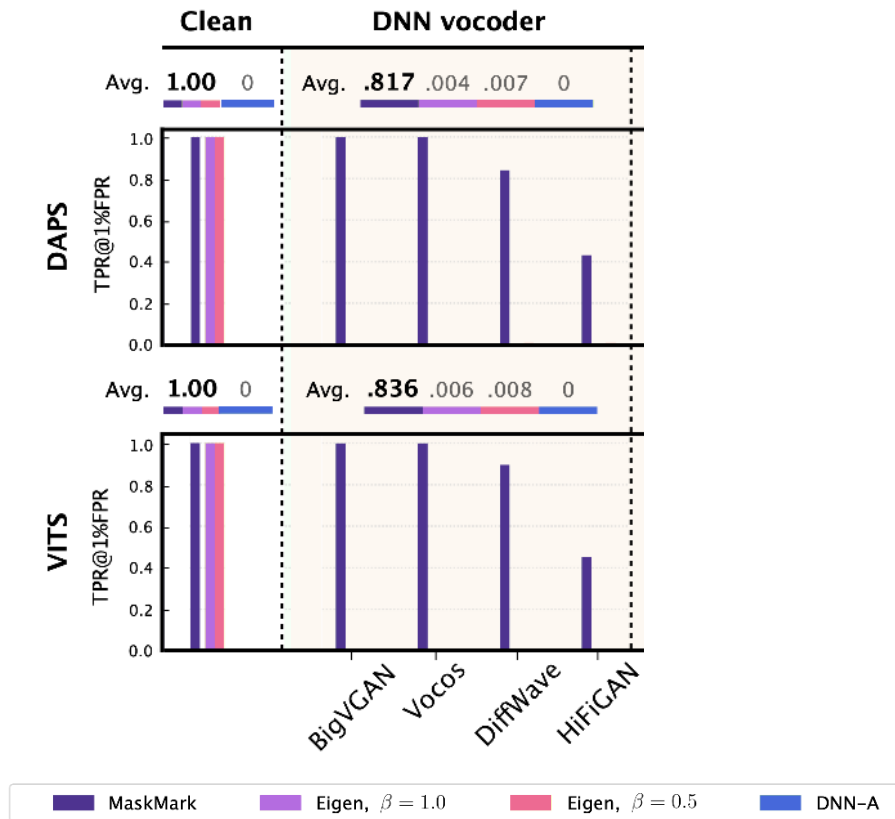|  | DNN-A | Eigen | Proposed |
|---|---|---|---|
| Less robust ← → More robust | | | |
| Sample rate | 16kHz | 44.1kHz | 48kHz |
| Required audio length | 2s | 1s | 1s |
| Robustness | | | |
| Signal-processing | TPR @ 1% FPR = 0.00 | 0.73 | 0.97 |
| Neural audio codec | 0.00 | 0.39 | 0.45 |
| Neural vocoder | 0.00 | 0.01 | 0.82 |
| Neural denoiser | 0.00 | 1.00 | 0.99 |

Neural vocoders can wipe out other watermarks while maintaining high audio quality!

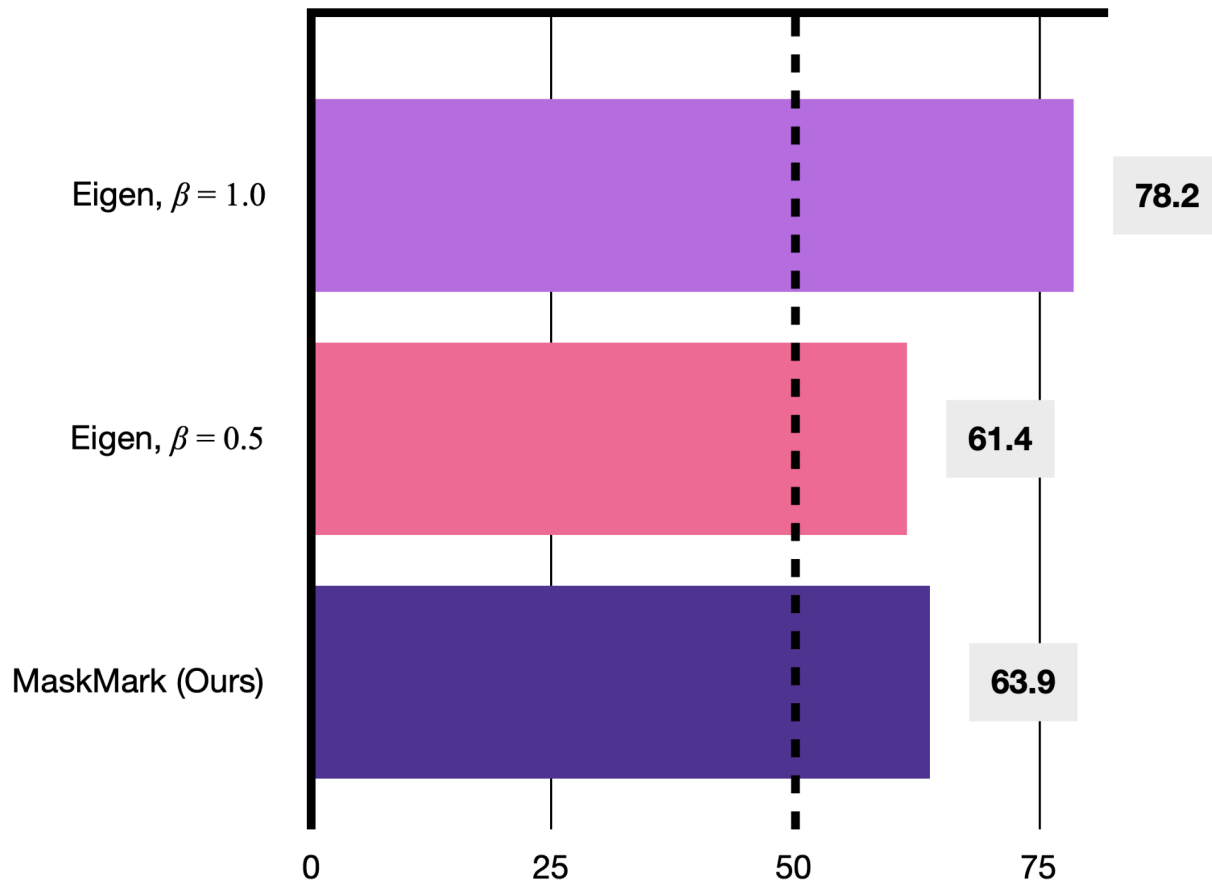**Neural vocoders can wipe out other watermarks while maintaining high audio quality!**

**Redundant frame-level embedding helps against pitch- and time-scale modification**

Our approach preserves audio quality as rated by human listeners.

% Recordings Correctly Identified

Eigen, $\beta = 1.0$: 78.2

Eigen, $\beta = 0.5$: 61.4

MaskMark (Ours): 63.9

# Concurrent works:

- **Timbre Watermark** (Liu et al. 2024) uses a similar network design and also demonstrates robustness against neural network-based transformations

- **WavMark** (Chen et al. 2023) uses invertible neural networks to achieve a higher watermark capacity, but considers a narrower and "gentler" set of transformations

- **AudioSeal** (Roman et al. 2024) embeds a residual signal in the time domain and likewise considers a narrower set of transformations

# Future directions:

- Improved robustness to neural network-based transformations

- Robustness to adversarial (optimization-based) attacks

- Increased information capacity

# MaskMark: Robust Neural Watermarking for Real and Synthetic Speech



(Listening examples)