

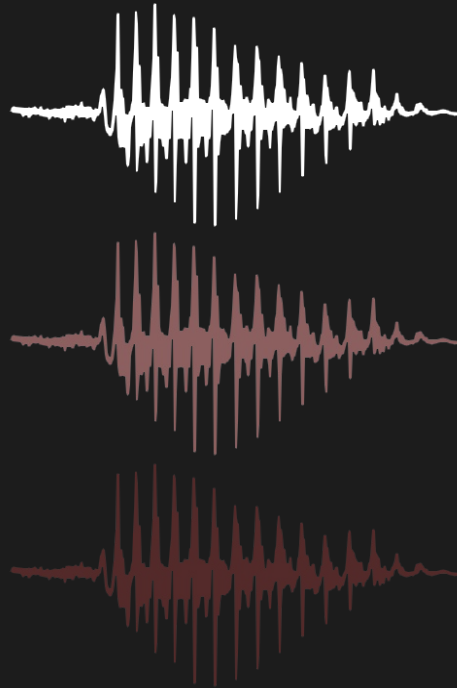
# Adversarial Attacks

on Audio Speech Recognition systems



LINKMEDIA  
TEAM

Student : Yoann Lemesle  
Supervisor : Laurent Amsaleg



# Deep Learning & applications

- **Machine Learning**

ML

Learning **automatically** from a set of data to perform a task without explicit programming.

- **Deep Learning**

DL

Methods of ML that use **neural networks**.



**Input**  
 $x$



**NEURAL  
NETWORK**

**Model**  
 $f$

→ 0.01 | **ASSAULT RIFLE**  
→ 0.05 | **HORSE**  
→ 0.04 | **SHEEP**  
→ 0.11 | **DUCK**  
→ 0.08 | **FROG**  
→ 0.03 | **DOG**  
→ **0.84** | **CAT**

**Predictions**  
 $f(x)$

# Deep Learning & security

- **Adversarial Example** | A seemingly benign input that fools a neural network.



**Input**  
 $x$



**NEURAL  
NETWORK**

**Model**  
 $f$

→ 0.99 | **ASSAULT RIFLE**  
→ 0.05 | **HORSE**  
→ 0.04 | **SHEEP**  
→ 0.11 | **DUCK**  
→ 0.08 | **FROG**  
→ 0.03 | **DOG**  
→ 0.02 | **CAT**

**Predictions**  
 $f(x)$

# Neural Networks

- **Powerful graphs**

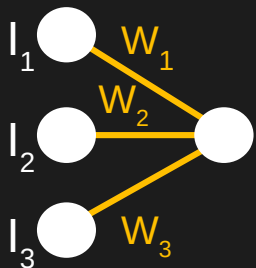
Weighted / Directed

Layers of **neurons** with **weighted connections**.

Acts as an **adjustable** function  $f(\theta, x)$

- **Forward propagation**

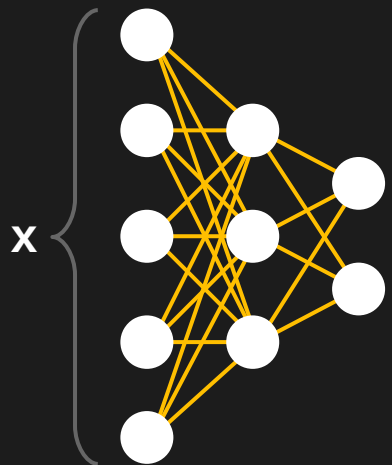
How information is processed.



$$\text{out} = \Phi(I_1 * W_1 + I_2 * W_2 + I_3 * W_3 + \text{Bias})$$

- **Back propagation**

How a neural network is trained.



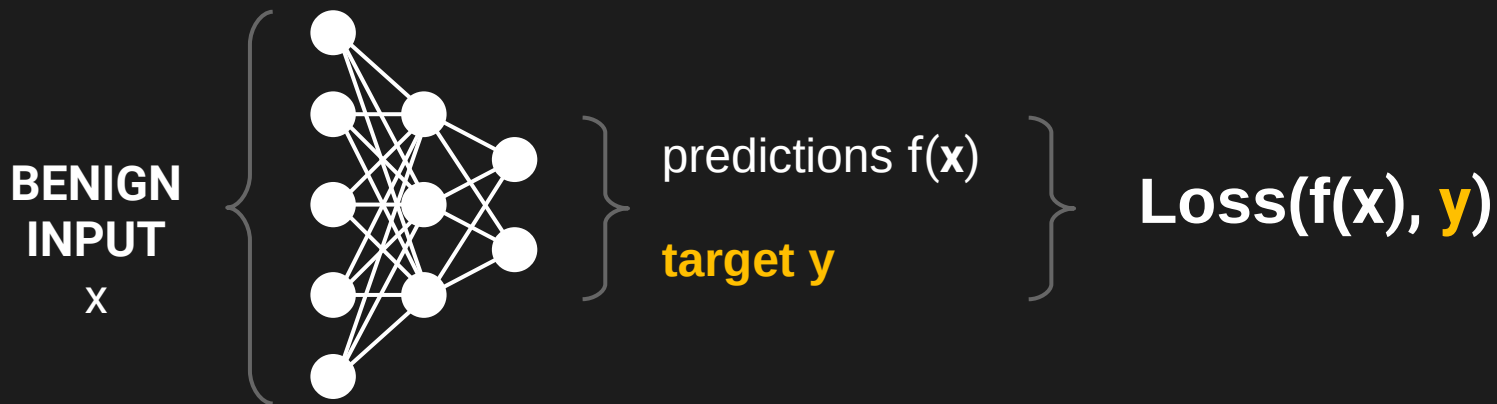
predictions  $f(\theta, x)$

truth  $y$

$$\text{Loss}(\theta, x, y)$$

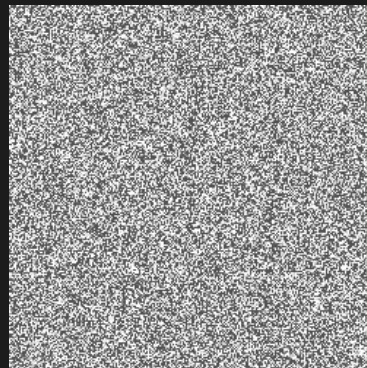
$$\theta' = \theta - \nabla_{\theta} \text{Loss}$$

# Fast Gradient Sign Method



$x$

$- 0.01 *$



$=$

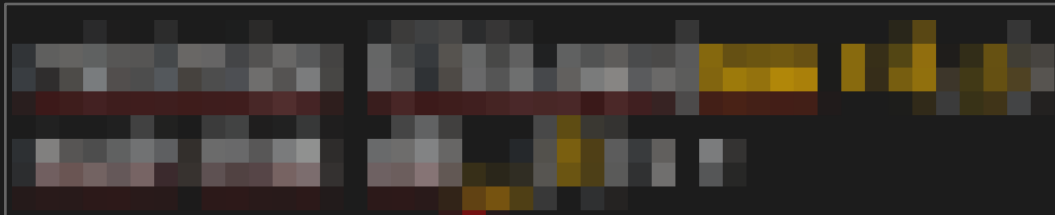


$x'$

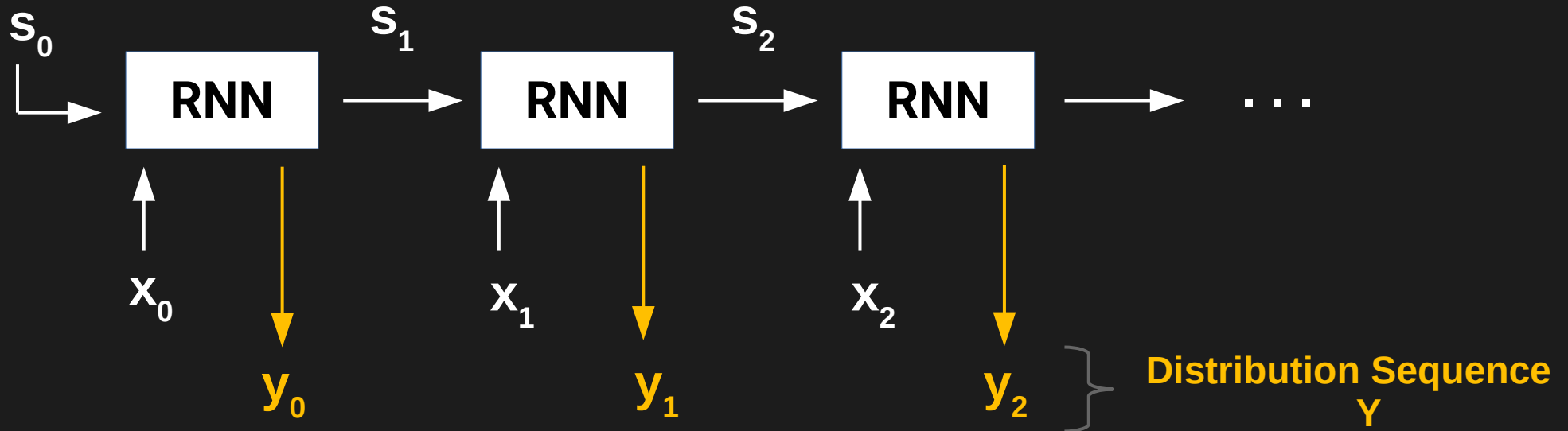
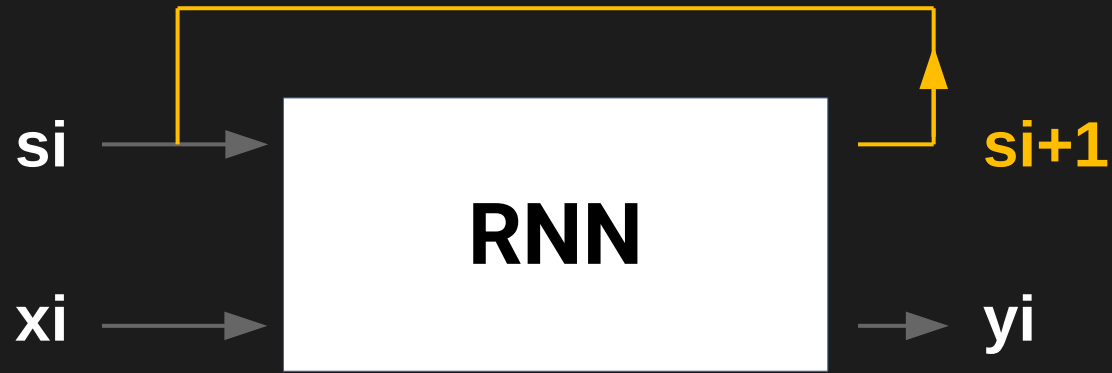
$$\delta = \nabla_x \text{Loss}$$

# Carlini & Wagner's adversarial attack (1)

- **Targeted** | The resulting adversarial example has a **desired** classification.
- **White Box** | Requires **full** knowledge of the model.
- **Minimally Perceptible** | Trying to minimize the **perceptibility** of the **adversarial noise  $\delta$** .



# Recurrent Neural Networks



# Connectionist Temporal Classification

Alignment  $\pi$

Reduces to

Sentence  $p$

BBBεAεεBBεAAA

BABA

---

ε	$\begin{pmatrix} 0.1 \\ 0.2 \\ 0.6 \\ 0.1 \end{pmatrix}$	$\begin{pmatrix} 0.5 \\ 0.1 \\ 0.2 \\ 0.3 \end{pmatrix}$	$\begin{pmatrix} 0.1 \\ 0.7 \\ 0.1 \\ 0.1 \end{pmatrix}$	$\begin{pmatrix} 0.1 \\ 0.8 \\ 0.0 \\ 0.1 \end{pmatrix}$	$\begin{pmatrix} 0.0 \\ 0.1 \\ 0.0 \\ 0.9 \end{pmatrix}$
A					
B					
C					

Distribution Sequence  $\mathbf{Y}$

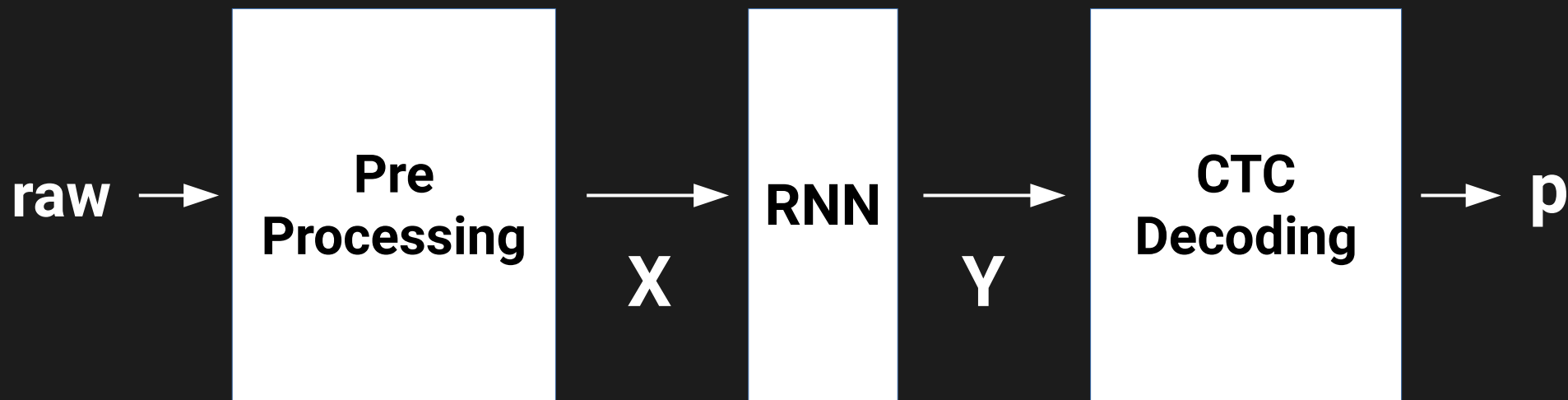
$$Pr(\pi|\mathbf{Y}) = \prod_i \mathbf{Y}_{\pi^i}^i$$

$$Pr(\mathbf{p}|\mathbf{Y}) = \sum_{\pi \in \Pi(\mathbf{p}, \mathbf{Y})} Pr(\pi|\mathbf{Y})$$



# DeepSpeech

---



---

$$\text{CTCLoss}(\text{raw}, \text{p}) = -\log \Pr(\text{p}|\text{Y})$$

## Carlini & Wagner's adversarial attack (2)

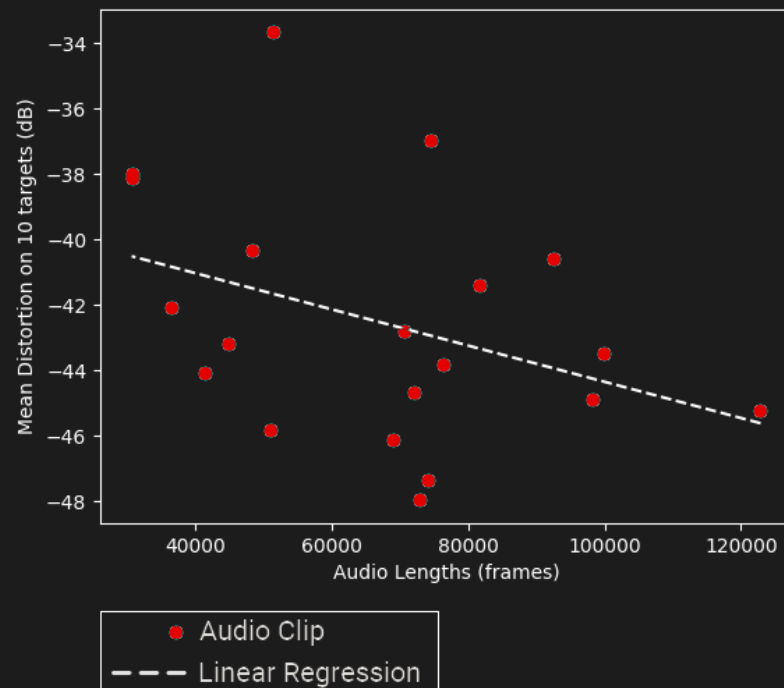
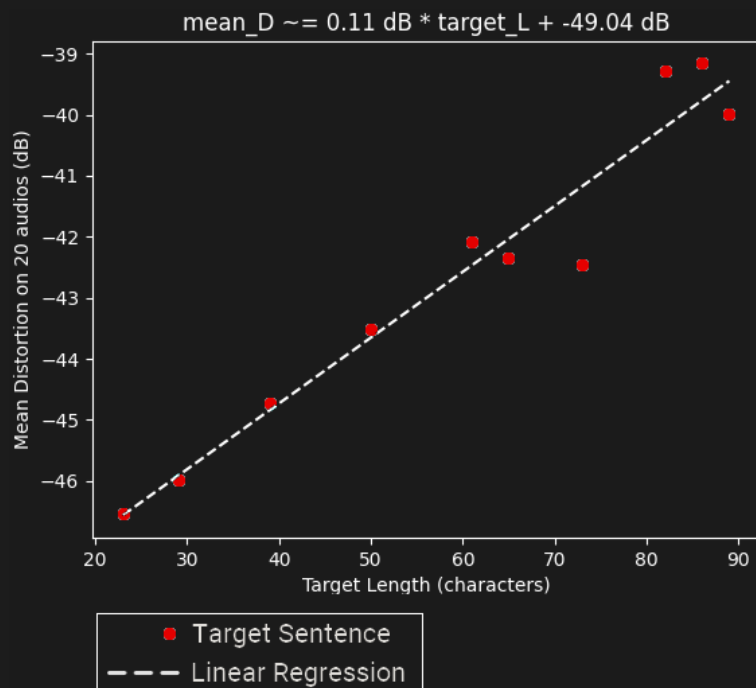
- **Targeted** | The resulting adversarial example has a **desired** classification.
- **White Box** | Requires **full** knowledge of the model.
- **Minimally Perceptible** | Trying to minimize the **perceptibility** of the **adversarial noise  $\delta$** .



$$\begin{aligned} &\text{minimize } \text{CTCLoss}(\text{raw} + \delta, p) \\ &\text{such that } \text{dB}_{\text{raw}}(\delta) \leq \tau \end{aligned}$$

# PyTorch Implementation

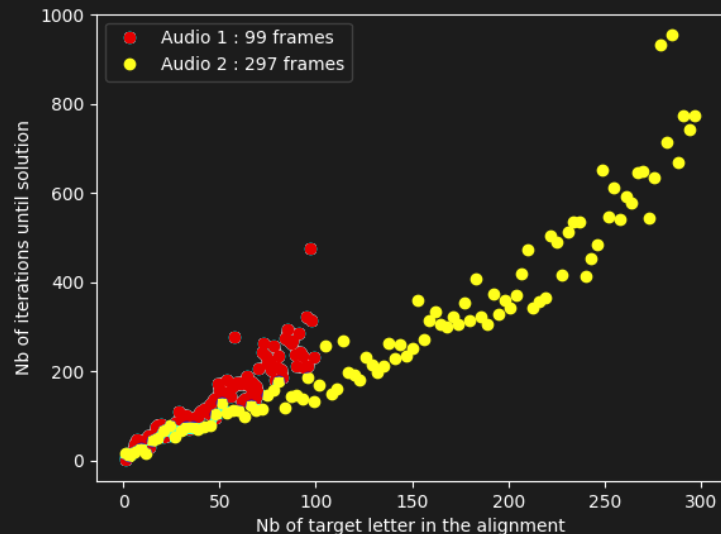
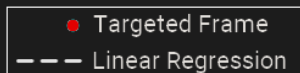
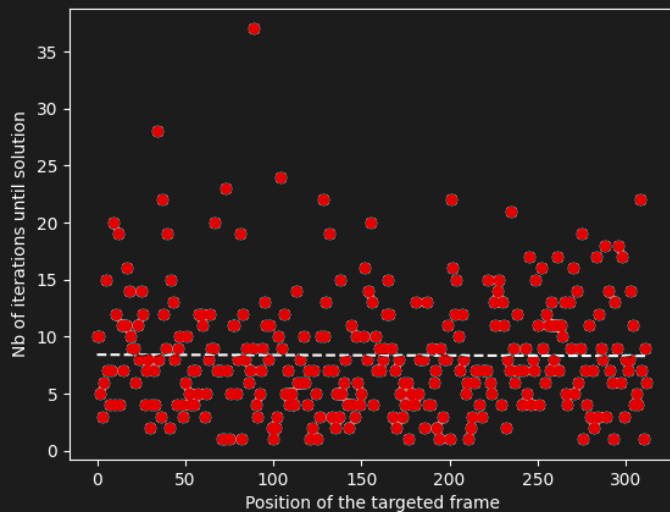
- **Methodology** | Targeting 10 different sentences on 20 audios.



# Specificities of audio adversarial attacks

- **Distortion Metrics** |  $L_\infty$  &  $L_2$  norms work well for images, not for audio !
- **Degrees of non linearity** | Differentiating through the **pre-processing** and **CTC decoding** step is not easy.

- **RNNs ?**



# Conclusion

---

- **Adversarial Attacks are a security and scientific challenge**
- **Neural Networks are not well understood**
- **Adversarial Attacks are harder on sequence input**

# Main References

---

## **Intriguing properties of neural networks (2013)**

C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow and R. Fergus.

## **Explaining and harnessing adversarial examples (2014)**

I. J. Goodfellow, J. Shlens, and C. Szegedy

## **Audio adversarial examples: Targeted attacks on speech-to-text (2018)**

Carlini, N. and Wagner