



Tempest

ACADEMY

Conference
2023

Attacking and Defending Machine Learning-based Systems

Prof. Paulo Freitas
pfreitas@cin.ufpe.br



- **Graduação**
 - Engenharia Eletrônica @UFPE
 - University of Massachusetts
- **Mestrado**
 - Ciência da Computação @UFPE
- **Doutorado**
 - Computer Science @UFPE
 - Eng Elétrica @ÉTS



- Intrepid Control Systems
- Ericsson
- Microsoft Research



Professor
@CIn UFPE

Security
Researcher





Tempest

ACADEMY

Conference

01 The AI Spring

02 The Risk

03 The Aftermath

04 The Future



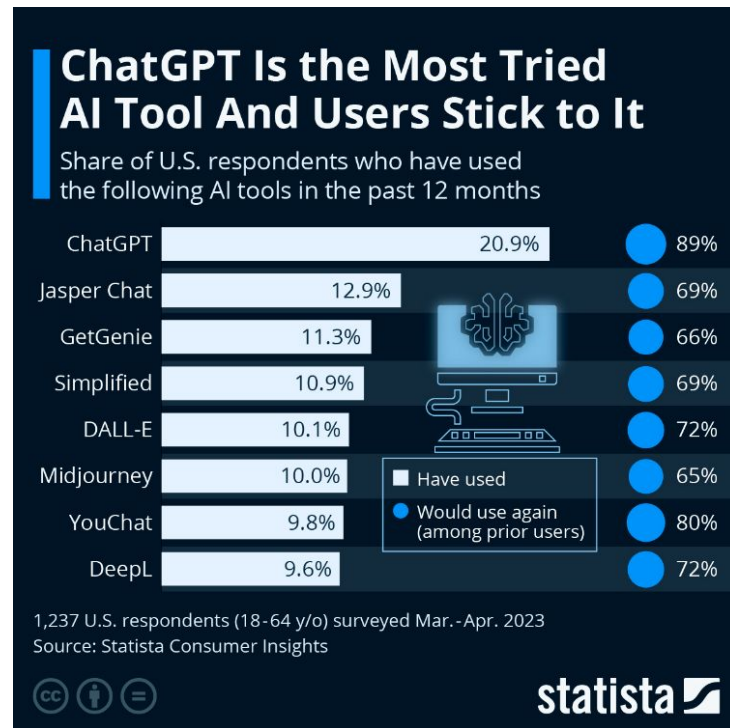
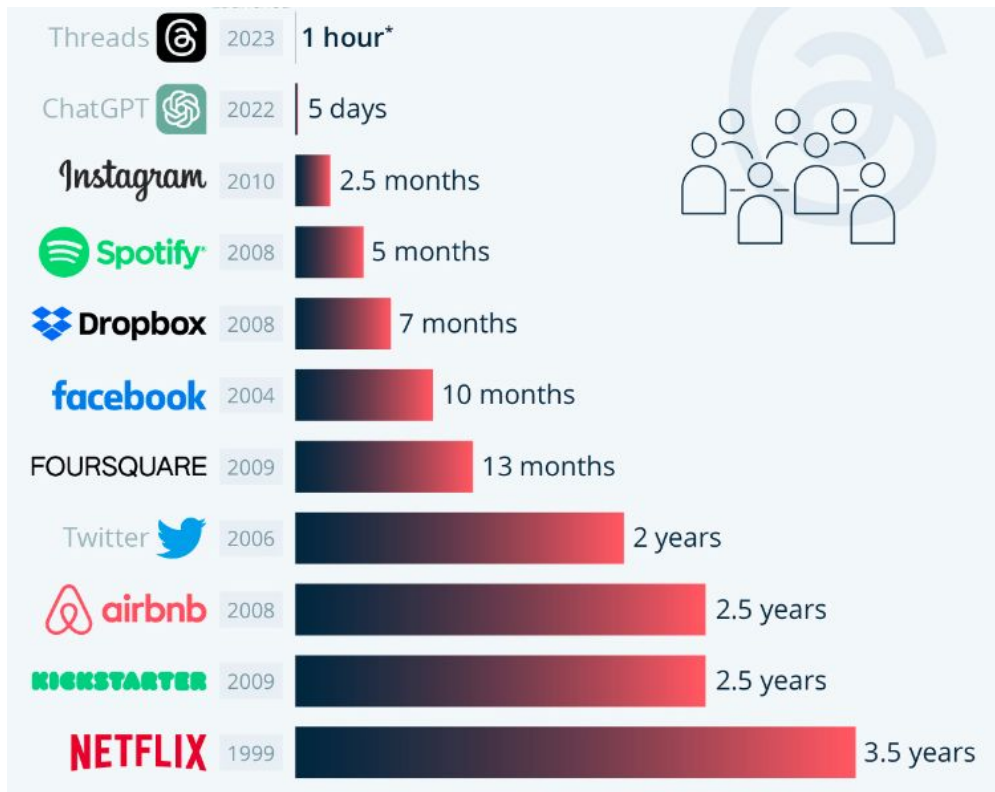
Tempest

ACADEMY

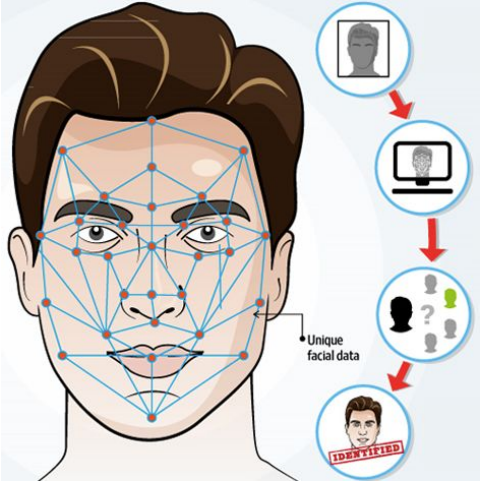
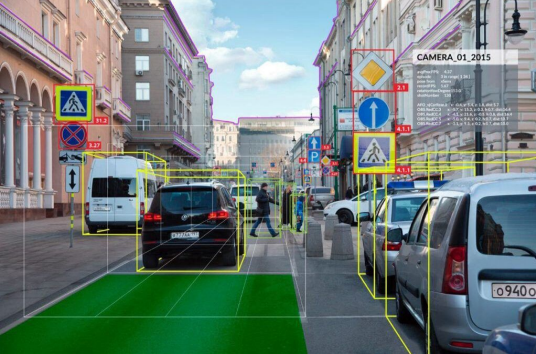
Conference

The AI Spring

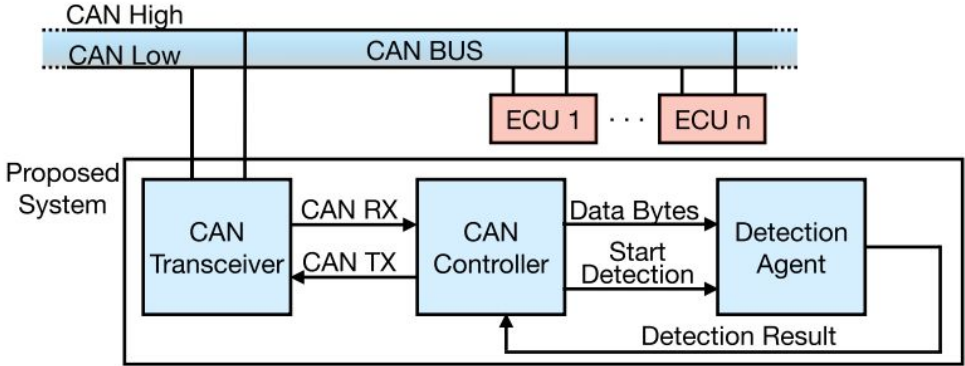
The AI Spring



The AI Spring



The AI Spring



The AI Spring

100 STARTUPS USING ARTIFICIAL INTELLIGENCE TO TRANSFORM INDUSTRIES

CONVERSATIONAL AI/ BOTS



VISION



AUTO



ROBOTICS



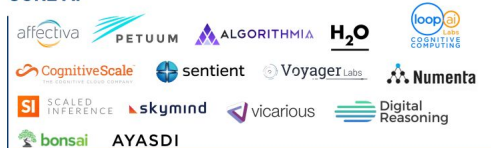
CYBERSECURITY



BUSINESS INTELLIGENCE & ANALYTICS



CORE AI



AD, SALES, CRM



HEALTHCARE



TEXT ANALYSIS/ GENERATION



IOT/IIoT



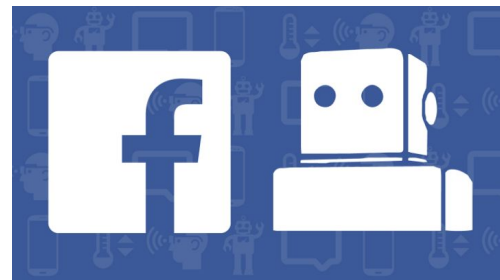
COMMERCE



FINTECH & INSURANCE



OTHER



But why do we use AI at all?

Why not stick to the “classic” algorithms?

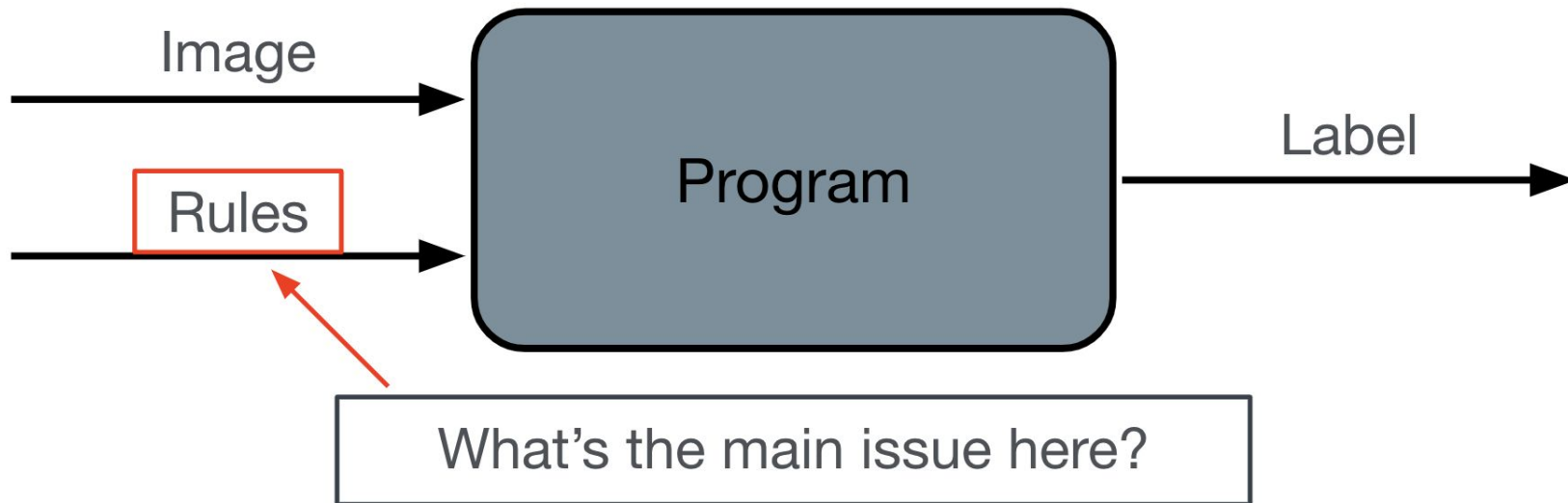
Let's dive into an example...

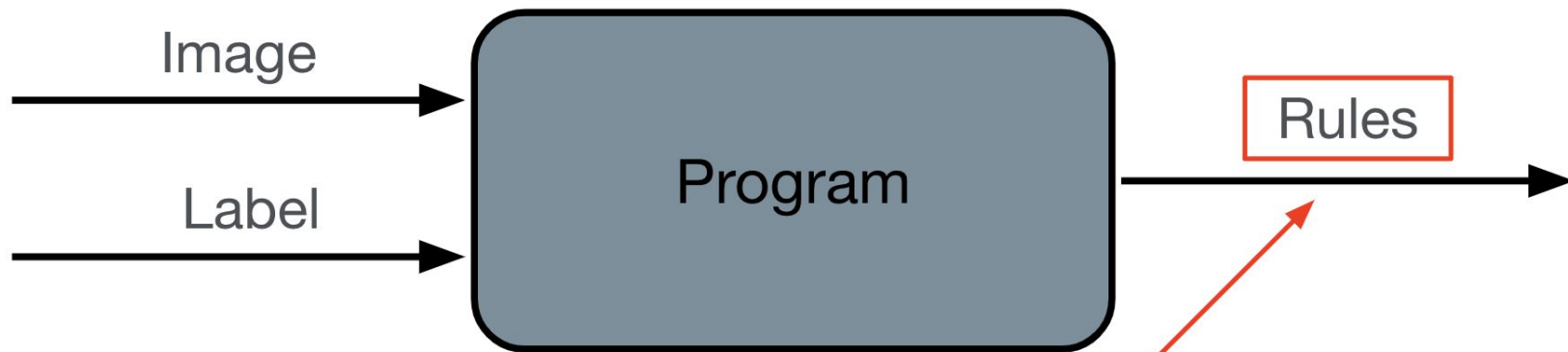
Let's assume we have the following task: given an arbitrary image, decide if it is a panda or not.

Is this a Panda?



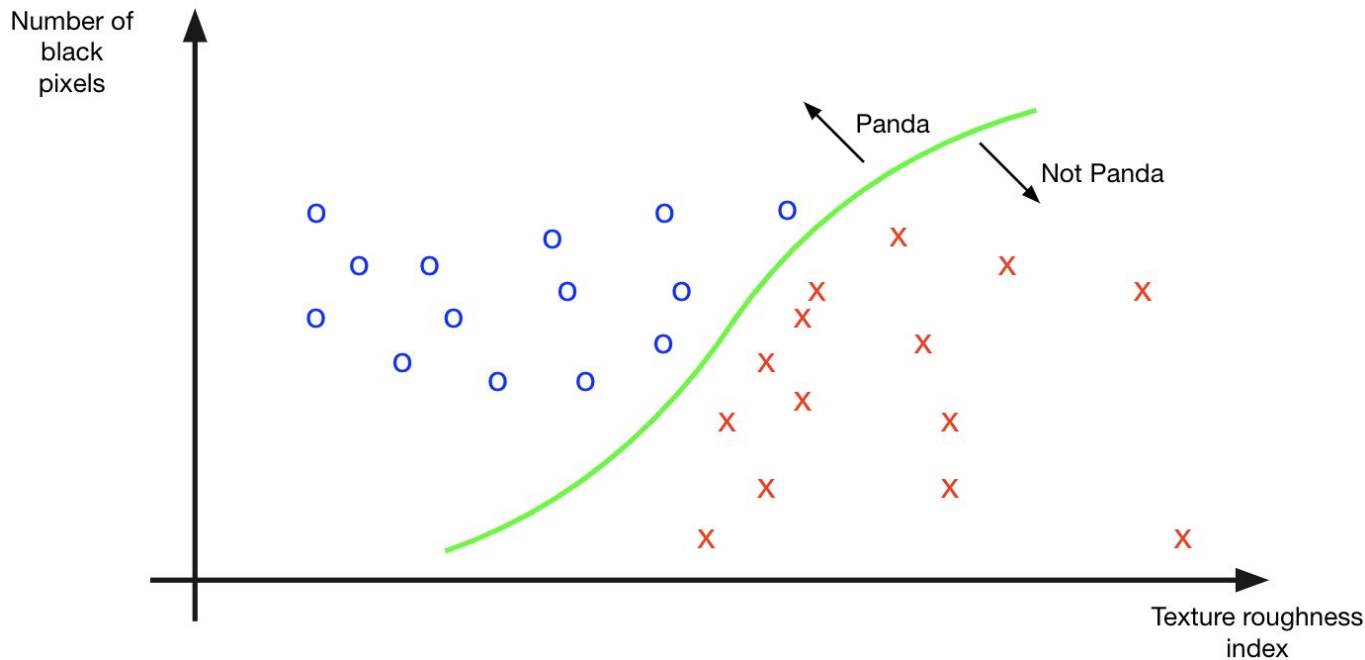
- Read each pixel;
- Analyze the pixels for the colors;
- Look for round black shapes;
- Find black heart-shaped figure just below and between the eyes;
- Define a texture that matches a panda skin;
- Remove other objects;
- Etc.





What if we change the order of the program?

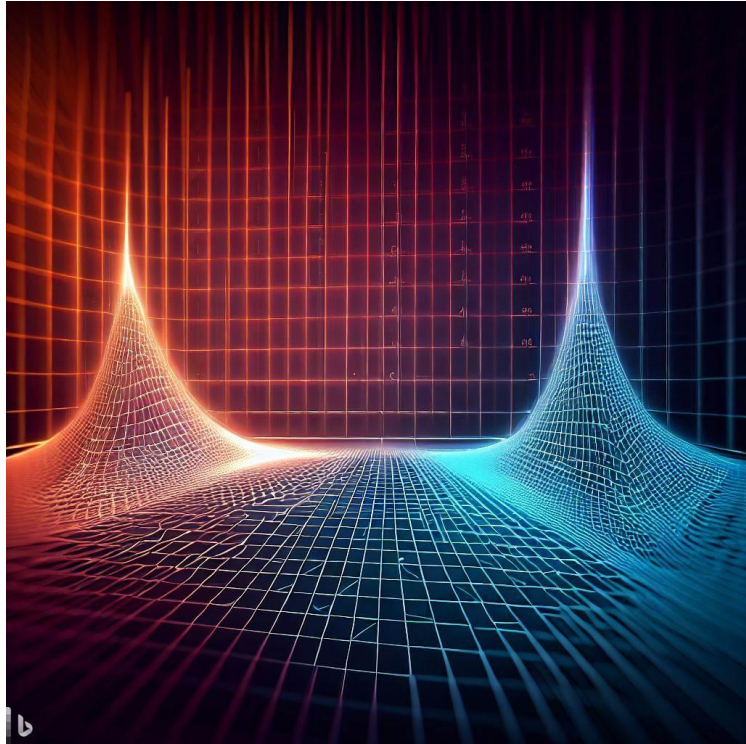
Fronteira de Decisão



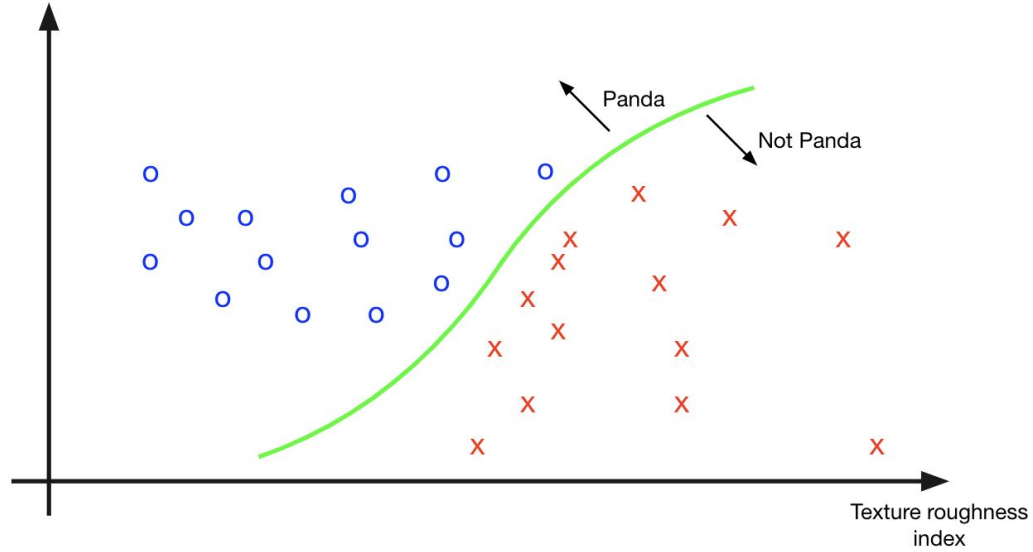
Fronteiras de decisão



Fronteira de Decisão



Number of
black
pixels





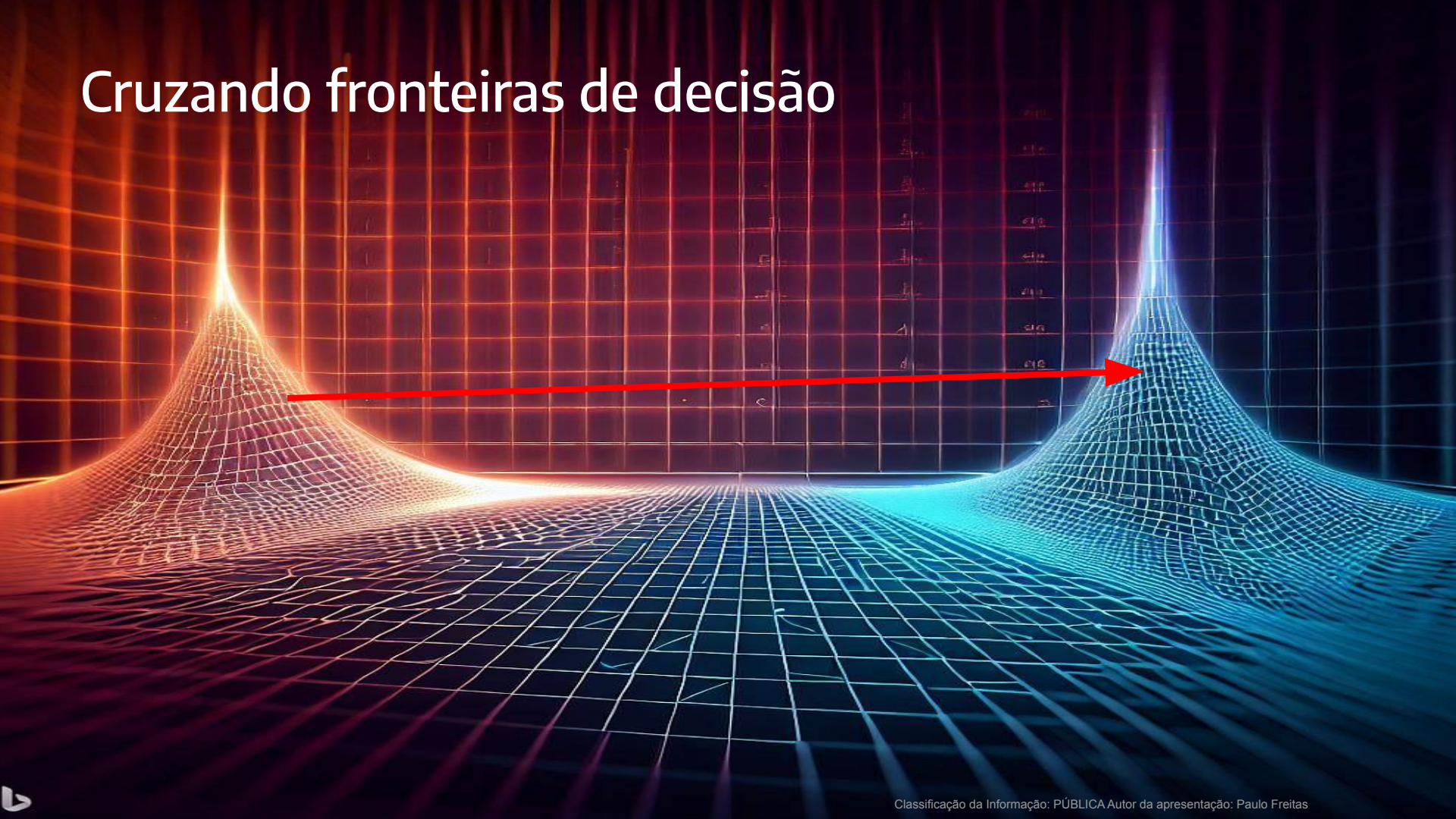
Tempest

ACADEMY

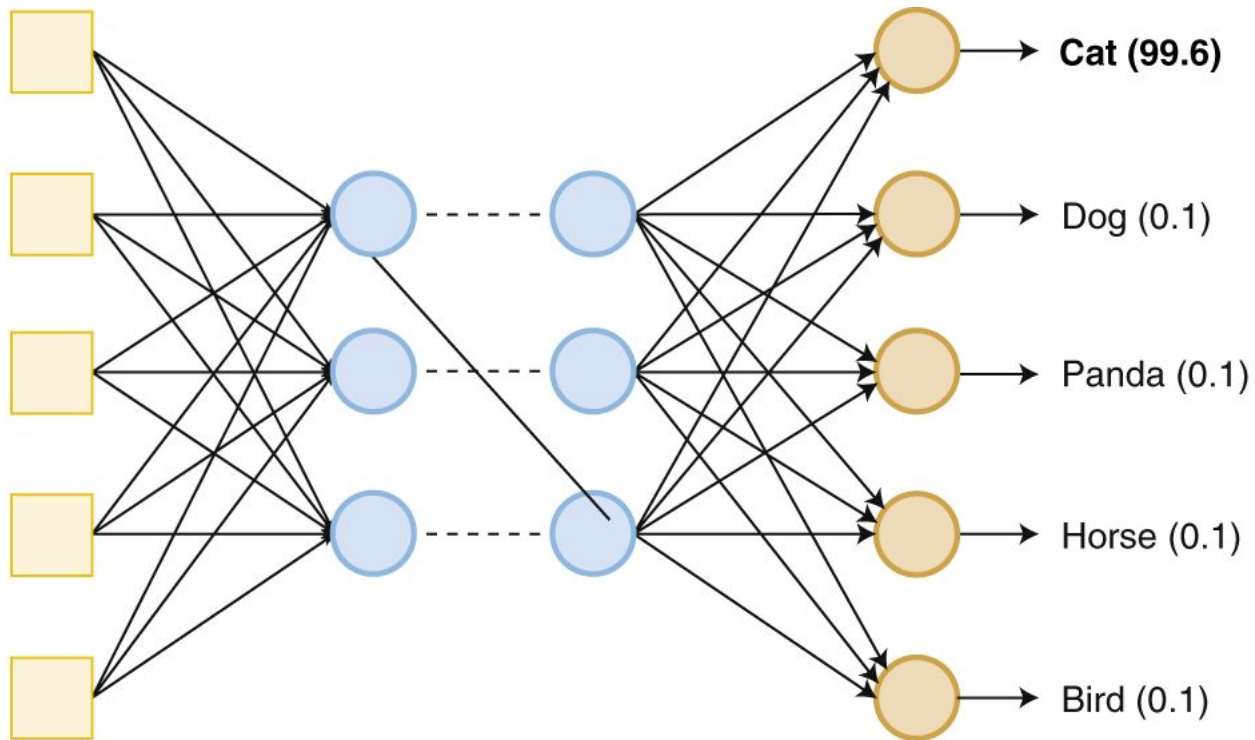
Conference

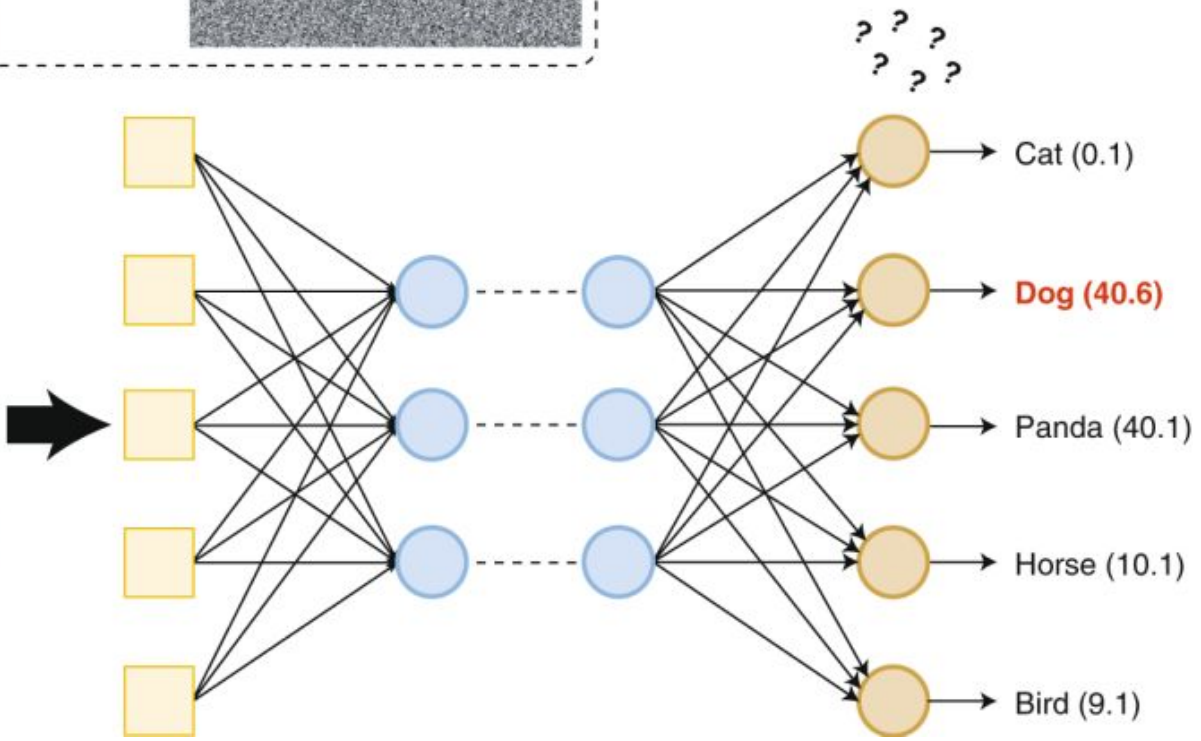
The Risk

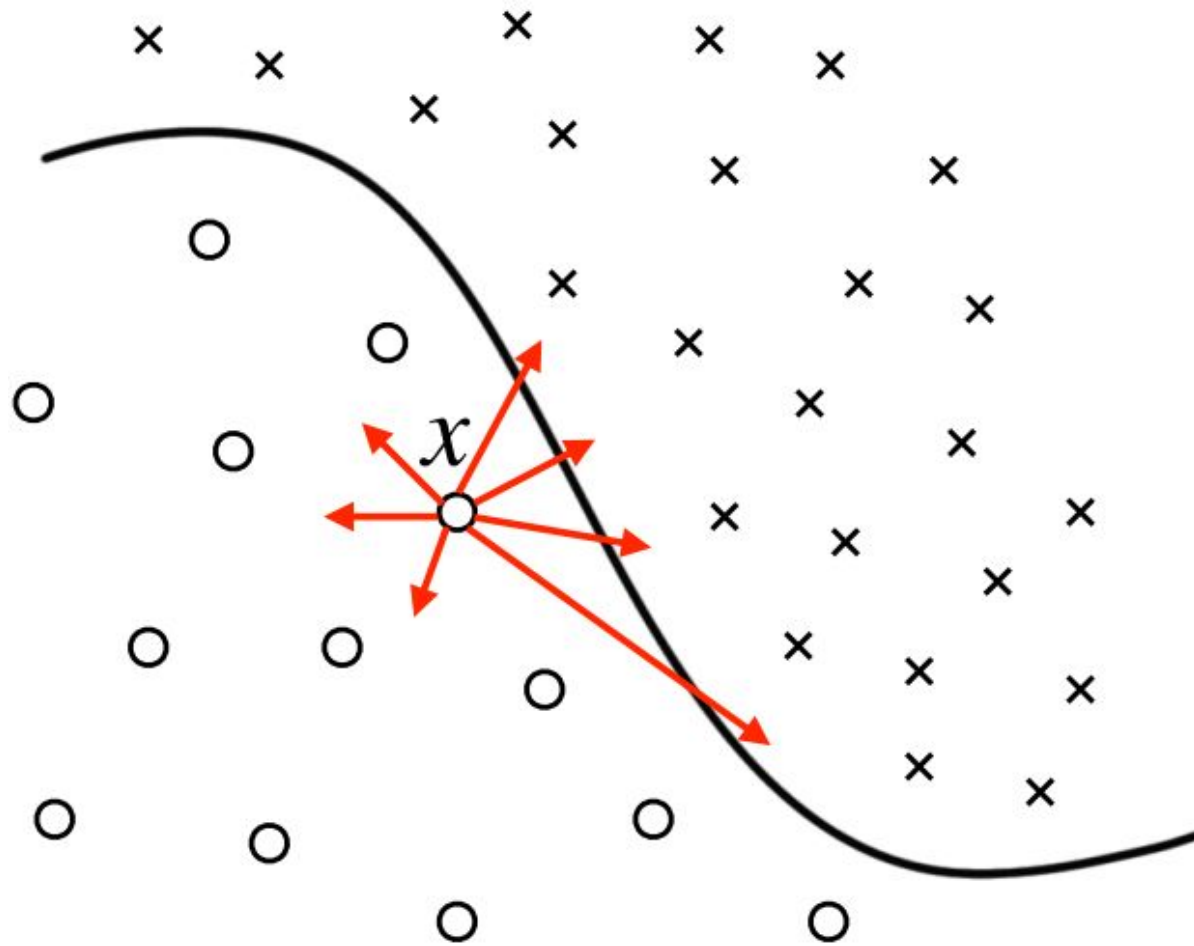
Cruzando fronteiras de decisão

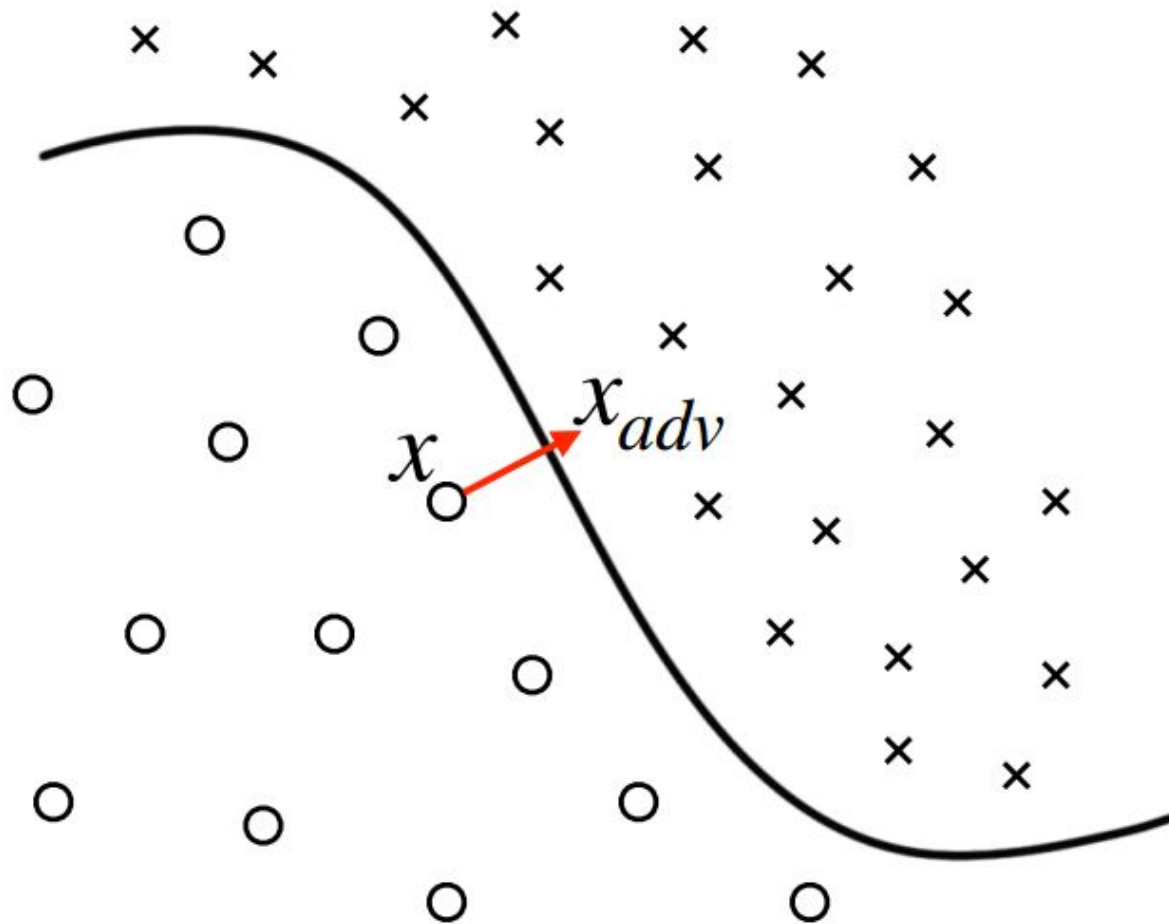


Original









$$x_{adv} = x + \delta$$

$$\min \|x_{adv} - x\| < \rho$$

$$f(x_{adv}) \neq f(x)$$

Forging adversarial perturbations

Multi-Objective GAN-Based Adversarial Attack Technique for Modulation Classifiers

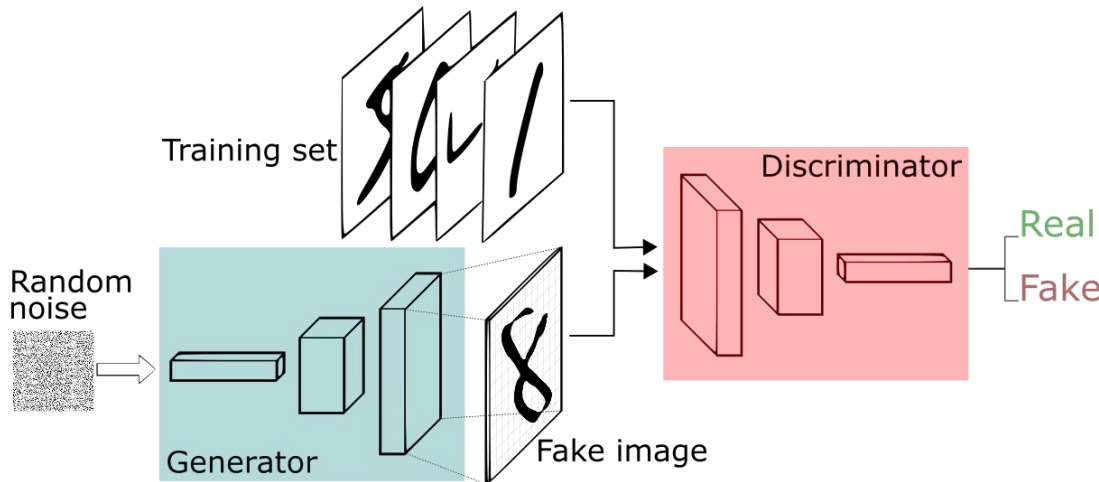
Paulo Freitas de Araujo-Filho^{id}, Georges Kaddoum^{id}, *Senior Member, IEEE*, Mohamed Naili,
Emmanuel Thepie Fapi^{id}, and Zhongwen Zhu^{id}, *Senior Member, IEEE*

$$f(x_{adv}) \neq f(x)$$

$$\min ||x_{adv} - x|| < \rho$$

Generative Adversarial Networks (GANs)

- Gerador G
 - Treinado para produzir amostras sintéticas de dados que sejam reconhecidos para reais
 - Aprende a distribuição de probabilidade de dos dados reais
 - Implicitamente modela o sistema
- Discriminador D
 - Treinado para distinguir as amostras reais daquelas produzidas pelo gerador



$$L_G = -D(G(z))$$

$$L_D = D(G(z)) - D(x)$$

Multi-Objective GAN-Based Adv Attack

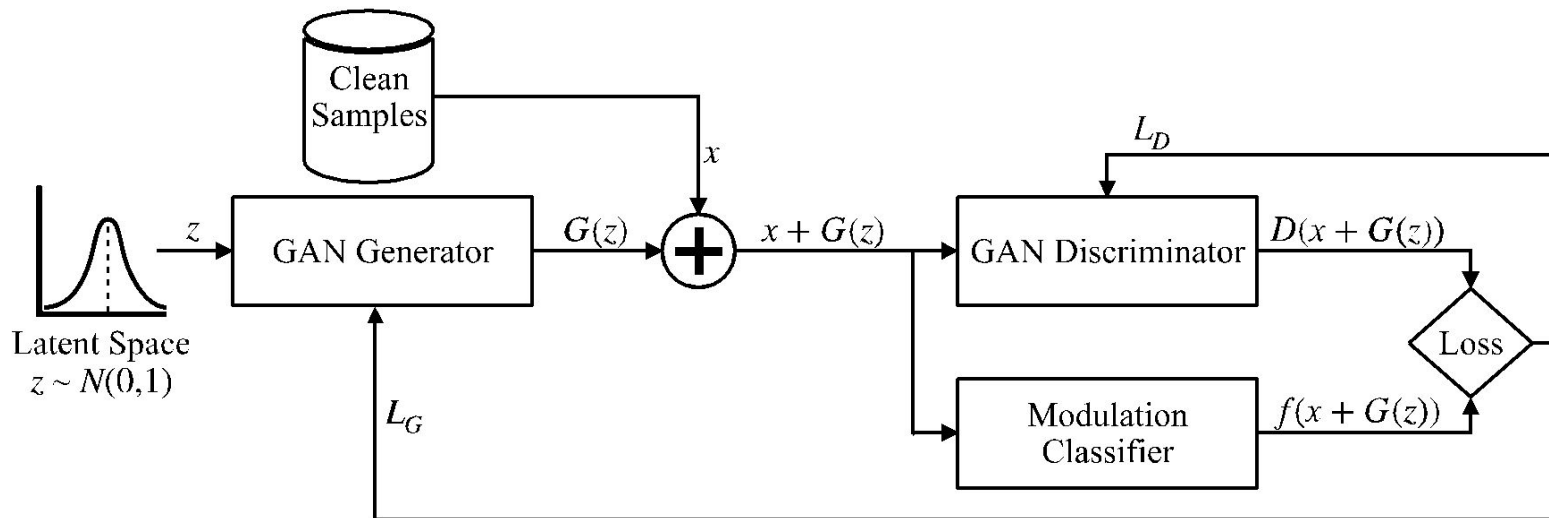
- Modificamos a estrutura da GAN para que o gerador produza perturbações adversariais

$$\delta = G(z)$$

$$x_{adv} = x + G(z)$$

$$L_G = -D(x + G(z))$$

$$L_D = D(x + G(z)) - D(x)$$

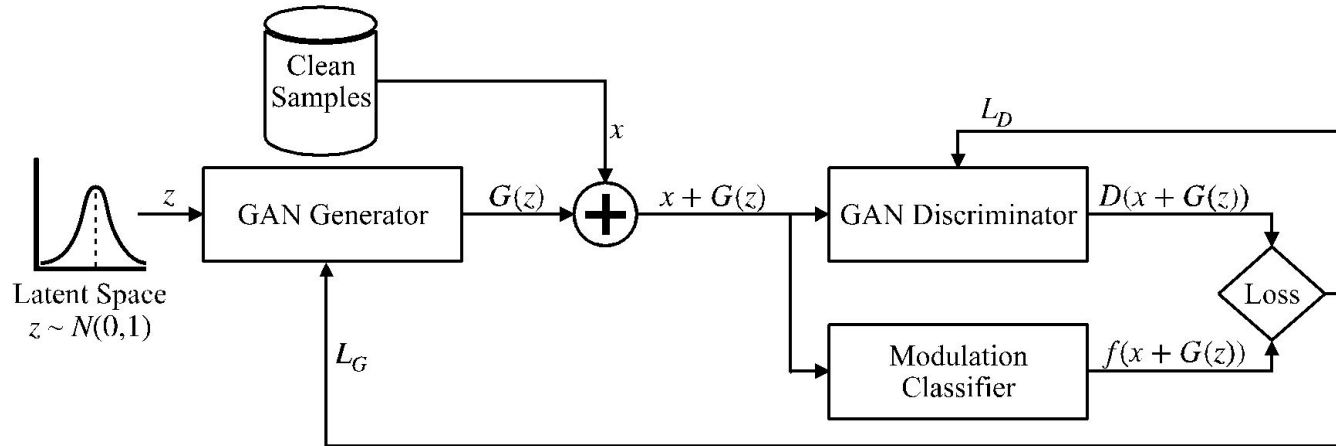


Multi-Objective GAN-Based Adv Attack

- Modificamos a estrutura da GAN para que o gerador produza perturbações adversariais

$$L_D = D(x + G(z)) - D(x)$$

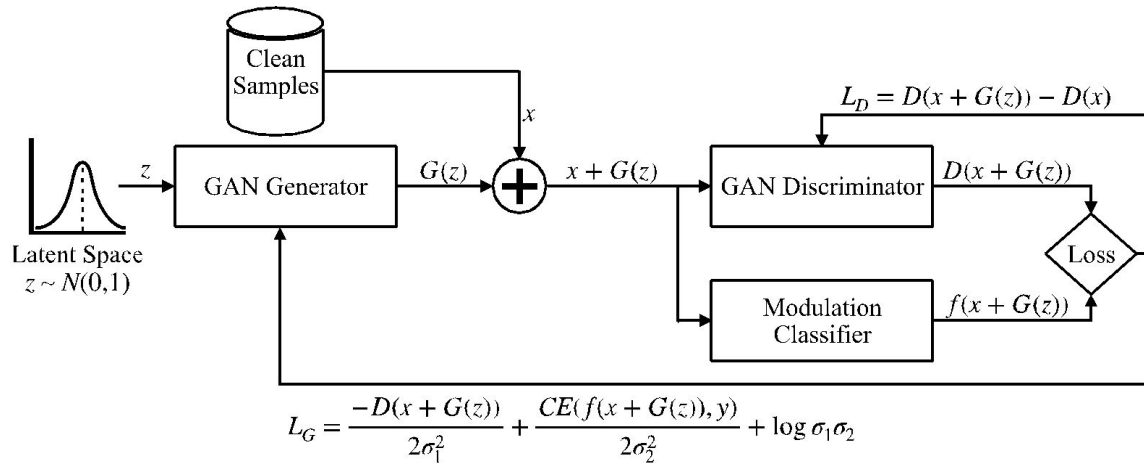
$$L_G = -D(x + G(z))$$



$$L_{G2} = CE(f(x + G(z)), y) = - \sum_{i=1}^n y_i \log(f_i(x + G(z)))$$

Multi-Objective GAN-Based Adv Attack

- Modificamos a estrutura da GAN e a combinamos com a Multi-Task Loss



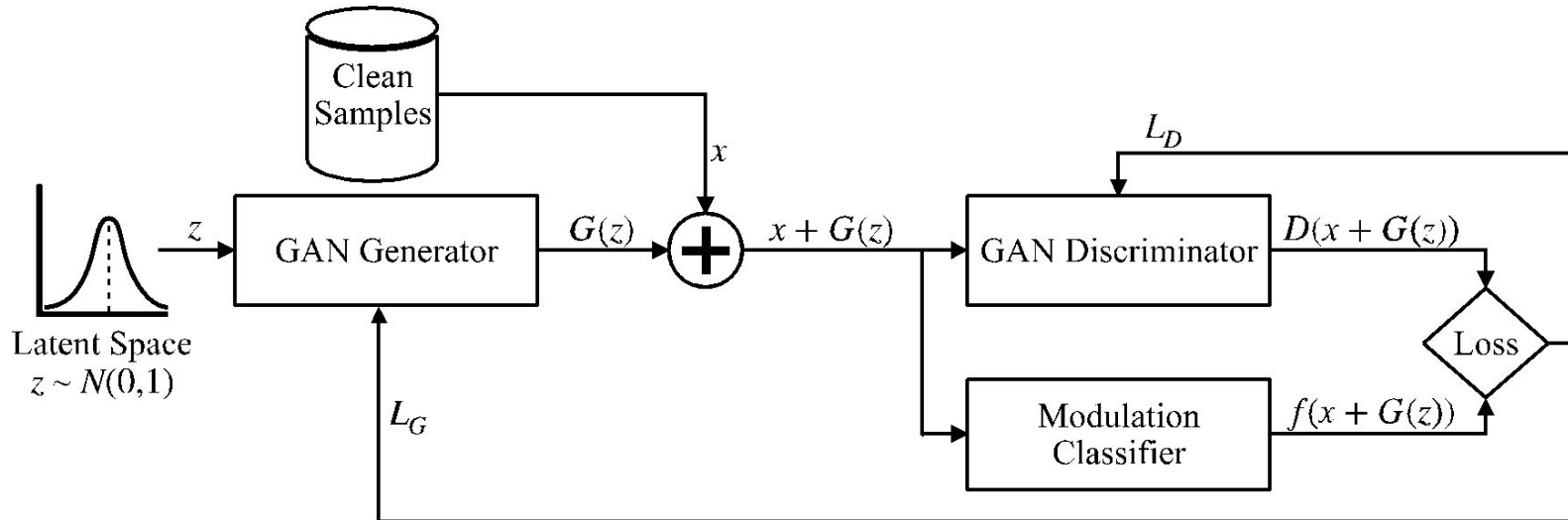
$$L_D = D(x + G(z)) - D(x)$$

$$L_G = -\frac{D(x + G(z))}{2\sigma_1^2} + \frac{CE(f(x + G(z)), y)}{2\sigma_2^2} + \log(\sigma_1 \sigma_2)$$

Multi-Objective GAN-Based Adv Attack

- Modificamos a estrutura da GAN para que o gerador produza perturbações adversariais

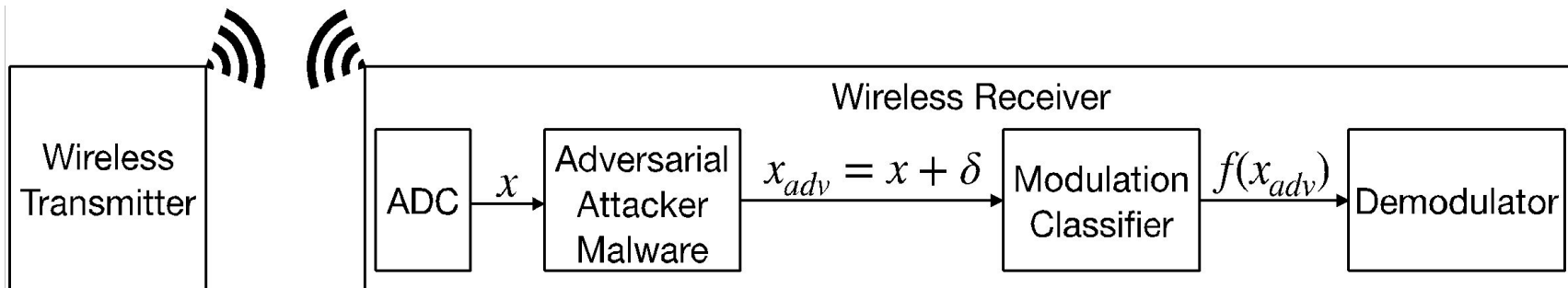
$$L_D = D(x + G(z)) - D(x) \quad L_G = -\frac{D(x+G(z))}{2\sigma_1^2} + \frac{CE(f(x+G(z)),y)}{2\sigma_2^2} + \log(\sigma_1\sigma_2)$$



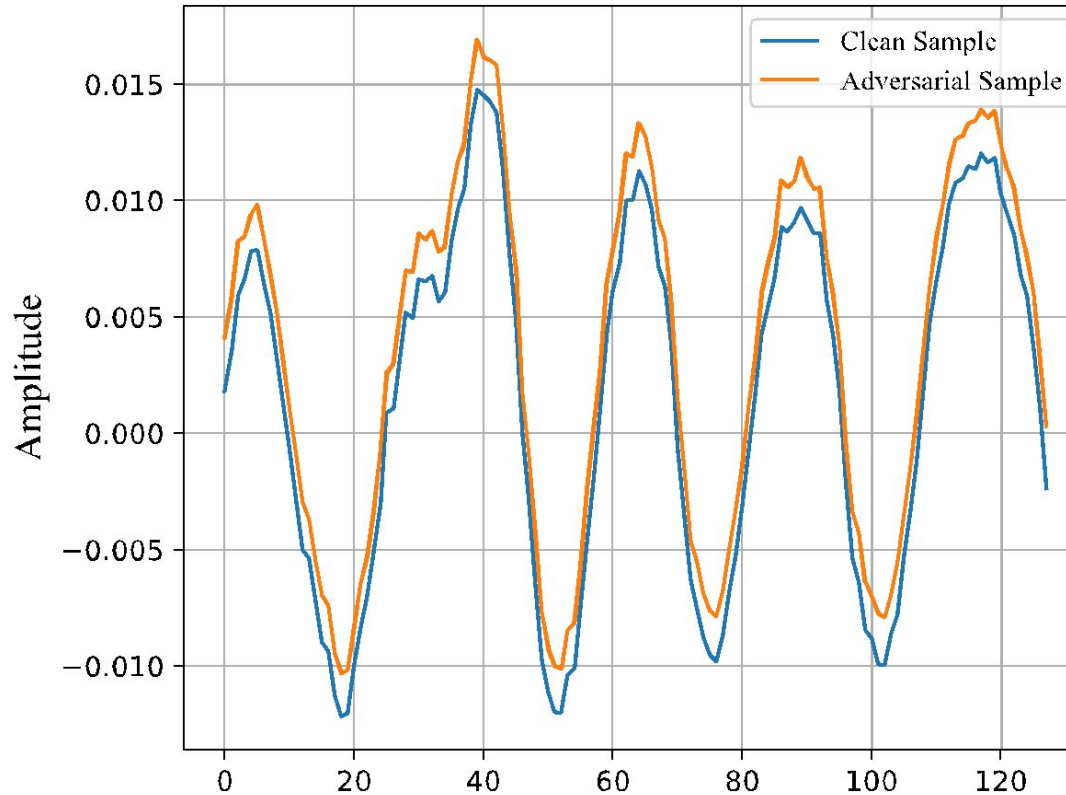
Multi-Objective GAN-Based Adv Attack

Algorithm 1 Proposed Adversarial Attack Technique

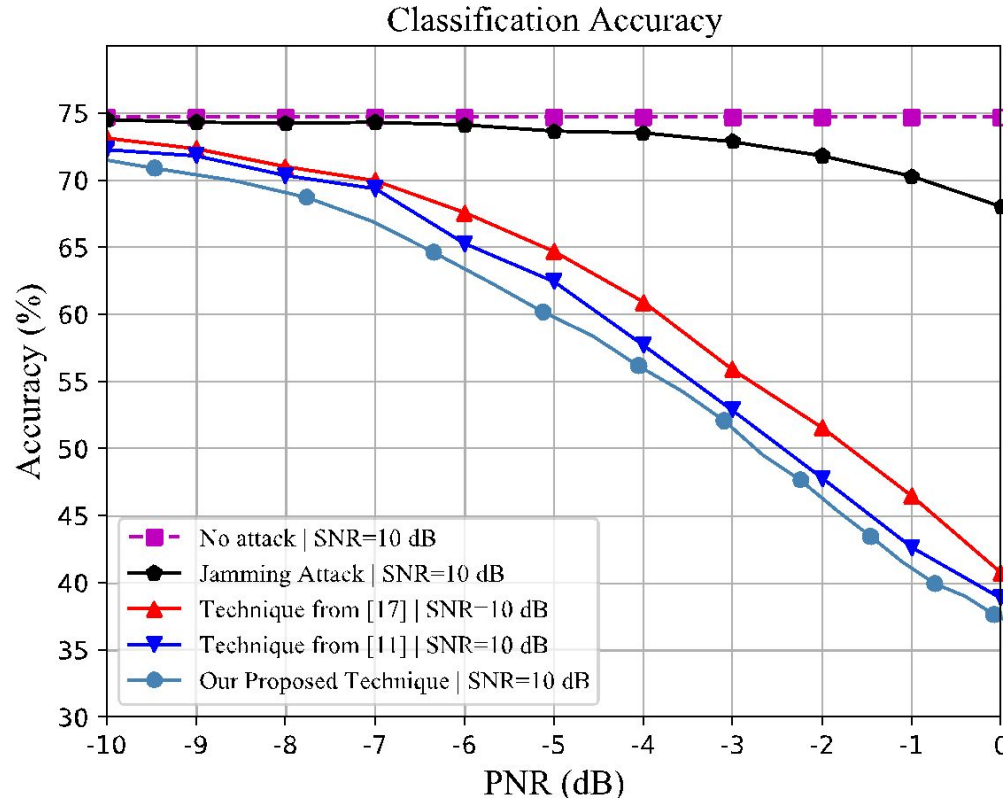
- 1: Train a GAN according to equations (4) and (5)
 - 2: **for** Each incoming sample x **do**
 - 3: Compute $G(z)$
 - 4: Construct the adversarial sample $x_{adv} = x + G(z)$
 - 5: **end for**
-



Multi-Objective GAN-Based Adv Attack



Multi-Objective GAN-Based Adv Attack



Adversarial Attack Technique	Mean Execution Time per Sample
Technique from [17]	20189 <i>ms</i>
Technique from [11]	234 <i>ms</i>
Our Proposed Technique	0.6980 <i>ms</i>



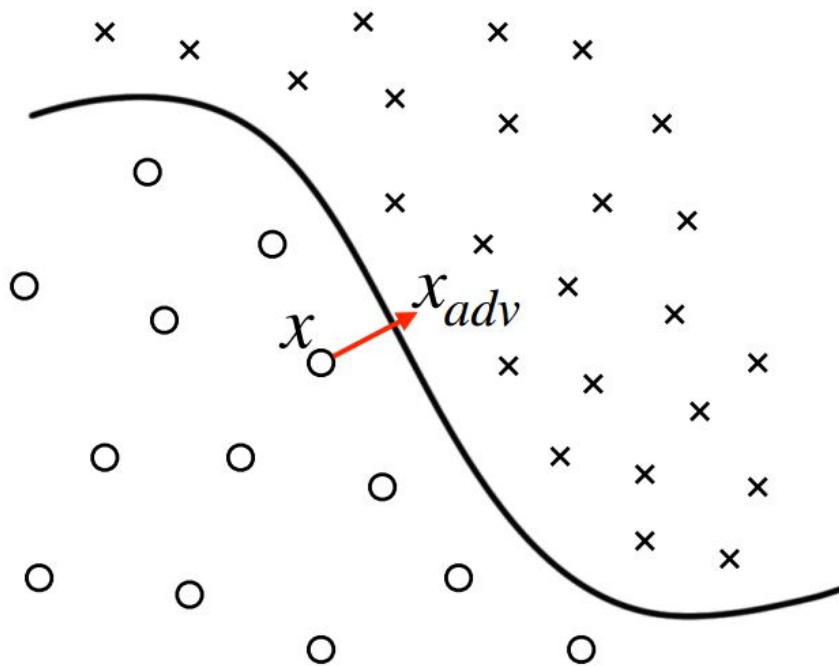
Tempest

ACADEMY

Conference

The Aftermath

Como mitigar?

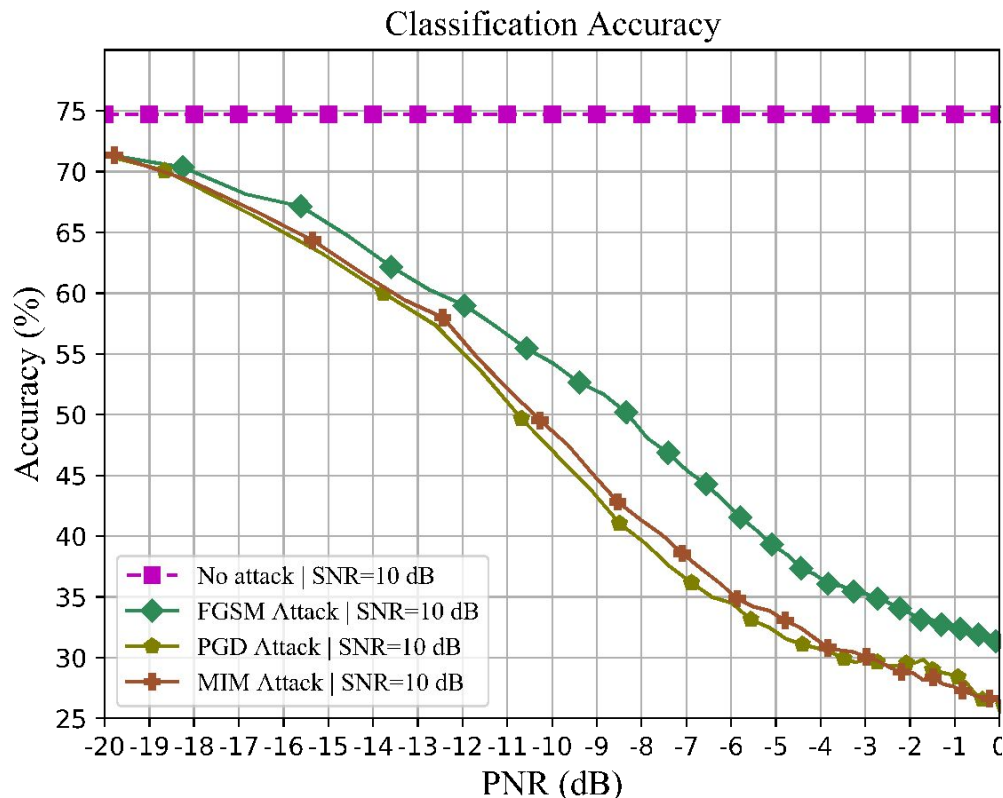


- Dificultar acesso
- Esconder informações
- Detectar perturbações
- Remover perturbações
- Suavizar fronteiras de decisão
- Combinar modelos

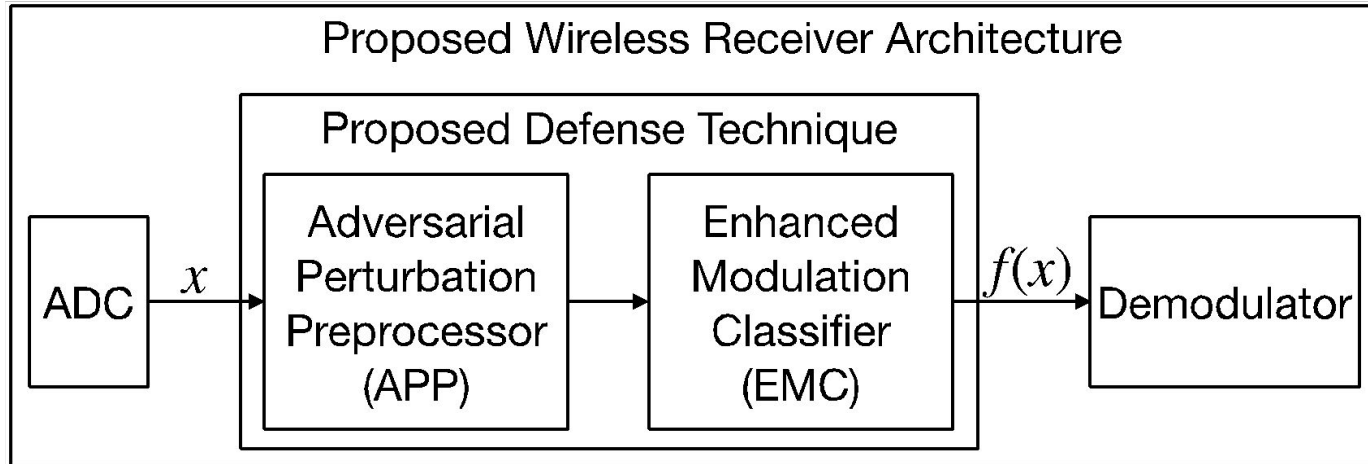
Defending Wireless Receivers Against Adversarial Attacks on Modulation Classifiers

Paulo Freitas de Araujo-Filho^{ID}, Georges Kaddoum^{ID}, *Senior Member, IEEE*, Mohamed Chiheb Ben Nasr^{ID}, Henrique F. Arcoverde, and Divanilson R. Campelo^{ID}, *Member, IEEE*

The Defense



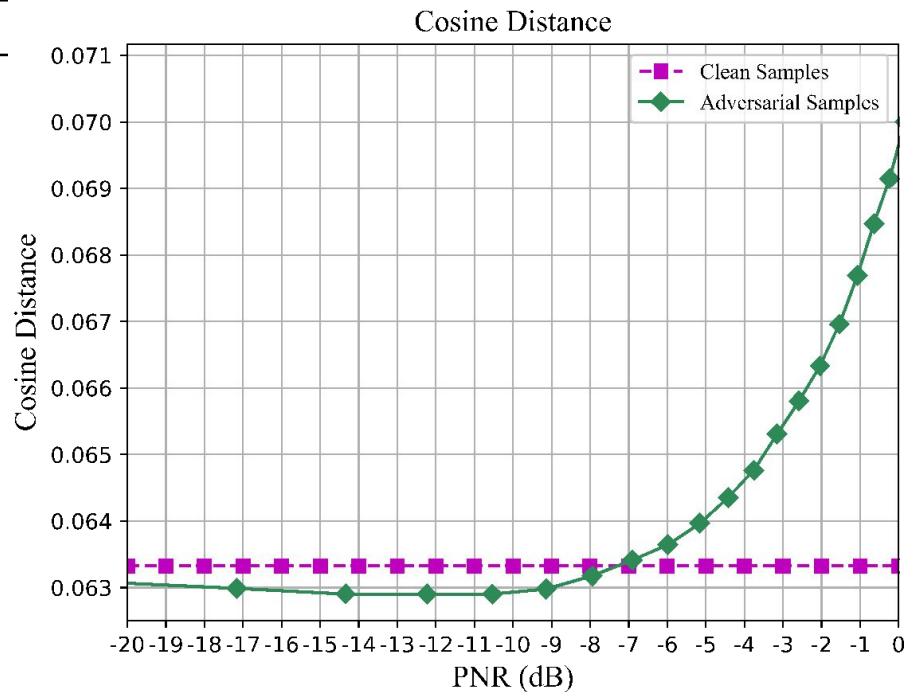
The Defense



The Defense

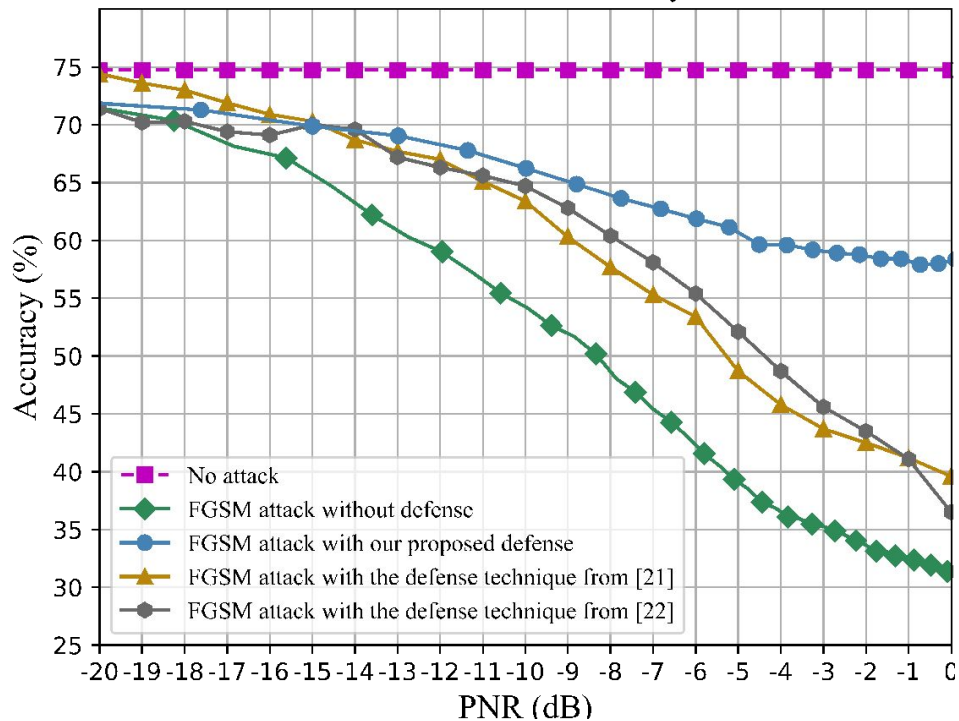
Algorithm 1: Proposed Defense Technique

- 1: Train a DAE with samples tampered with Gaussian noise and adversarial perturbations
 - 2: Train a EMC with samples tampered with Gaussian noise and adversarial perturbations
 - 3: **for** Each incoming sample x_i **do**
 - 4: Compute $x_o = DAE(x_i)$
 - 5: Compute $\beta = CD(x_i, x_o)$
 - 6: **if** $\beta \geq t$ **then**
 - 7: Preprocess data sample $x = x_o$
 - 8: **else**
 - 9: Do not preprocess data sample $x = x_i$
 - 10: **end if**
 - 11: Classify data sample $y = f(x)$
 - 12: **end for**
-

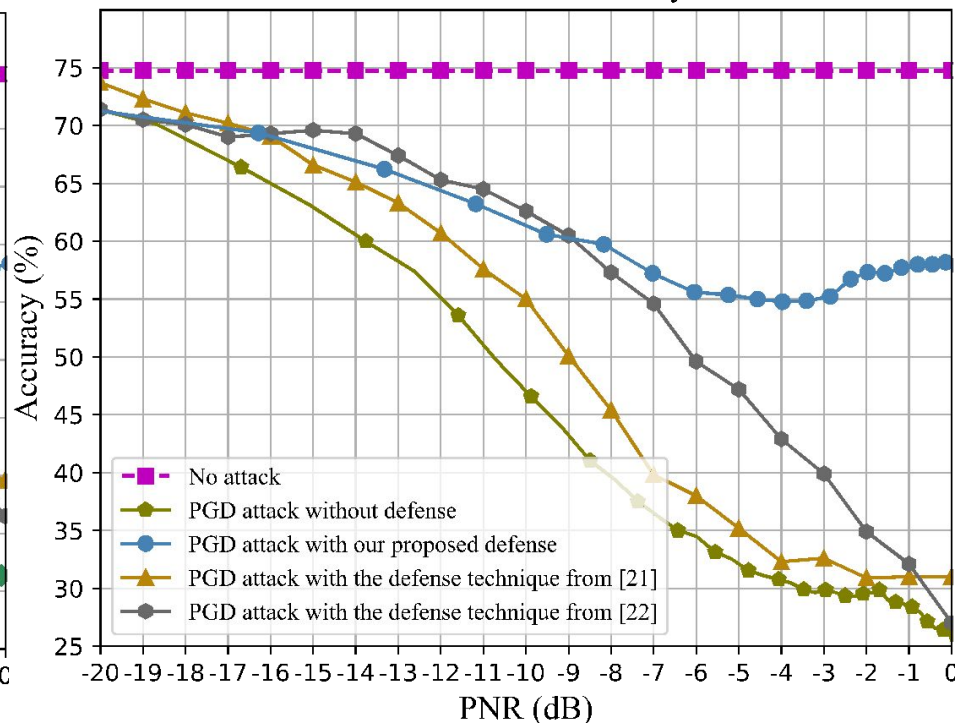


The Defense

Classification Accuracy



Classification Accuracy





Tempest

ACADEMY

Conference

The Future

- AI brings great changes
- Adversarial attacks are a real problem that can cause devastating consequences to many different applications
- Defense mechanisms exist but cannot yet fully protect ML-based systems.

**New and better
defense mechanisms are needed!**



Tempest

ACADEMY

Conference

2023

