TEMPEST talks 2022

Original

Cat (99.6)

Dog (0.1)

Panda (0.1)

Horse (0.1)

Bird (0.1)

TEMPEST talks

2022

Original + Perturbation

**Input**

Pre-recorded image

**Attack Building Blocks**

**Laser Attack Modeling**

Camera Modeling
Direct Interpolation | Bilinear Interpolation

Laser Interference Modeling
Photon Filtering | Color Strip Addition

Laser Effect Emulation
Incidence Direction | Lens' Imperfectness

**Parameter Search**
Wavelength
Minimum Strength
Maximum Strength
Incidence Direction
Incidence Function

**Laser Generation**
Wavelength
Power
Pulse Width
Pulse Period
Incidence Angle

**Output**

Real image

Laser diode

Green light

Camera

**Target System**

Captured image

Traffic light detection & recognition

**Result**

Red light

---

**Legitimate Sample**

Stop Sign

STOP

**+**

**Adversarial Perturbation**

**=**

**Adversarial Sample**

Yield Sign

STOP

Source Image | Guide Image | Adversarial Sample | Adversarial Noise

TEMPEST talks
2022

"panda" $\quad +\ .007\ \times\quad$ noise $\quad =\quad$ "gibbon"

TEMPEST talks
2022

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

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

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

TEMPEST talks
2022

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

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

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

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

TEMPEST talks
2022

- Generative Adversarial Networks (GANs)
  - Treina simultaneamente duas redes neurais que competem entre si

- 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

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

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

TEMPEST talks
2022

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

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

$$L_G = -D(x + G(z))$$
$$L_D = D\big(x + G(z)\big) - D(x)$$



Latent Space
$z \sim N(0,1)$

Clean
Samples

GAN Generator

GAN Discriminator

Modulation
Classifier

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

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



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

TEMPEST talks
2022

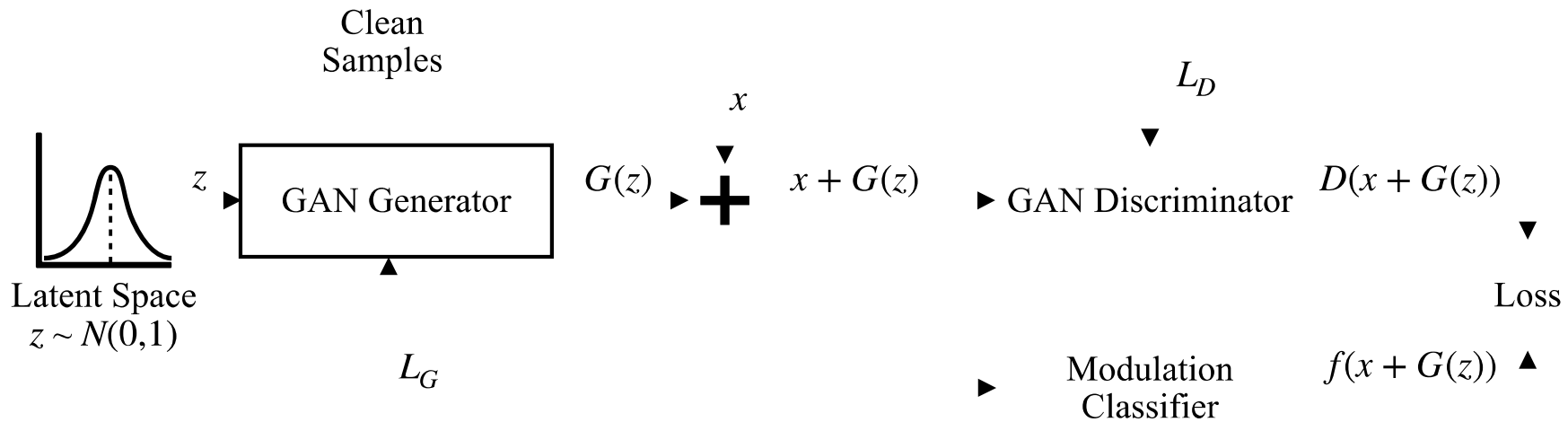- Modificamos a estrutura da GAN para que o gerador produza perturbações adversariais

$$L_{G1} = -D(x + G(z))$$

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

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

$$L_{G1} = -D(x + G(z))$$

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

$$L_G = \alpha L_{G1} + \beta L_{G2}$$

- Multi-Task Loss

$$p(\boldsymbol{y}|\boldsymbol{f^W}(\boldsymbol{x})) = Softmax(\boldsymbol{f^W}(\boldsymbol{x}))$$

- Multi-Task Loss

$$p(\boldsymbol{y}|\boldsymbol{f^W}(\boldsymbol{x})) = Softmax(\boldsymbol{f^W}(\boldsymbol{x}))$$

$$\log p(\boldsymbol{y}|\boldsymbol{f^W}(\boldsymbol{x})) \propto -\frac{1}{2\sigma^2}||\boldsymbol{y} - \boldsymbol{f^W}(\mathbf{x})||^2 - \log\sigma$$

- Multi-Task Loss

$$p(\boldsymbol{y}|\boldsymbol{f}^W(\boldsymbol{x})) = Softmax(\boldsymbol{f}^W(\boldsymbol{x}))$$

$$\log p(\boldsymbol{y}|\boldsymbol{f}^W(\boldsymbol{x})) \propto -\frac{1}{2\sigma^2}||\boldsymbol{y} - \boldsymbol{f}^W(\mathbf{x})||^2 - \log\sigma$$

$$\log p(\boldsymbol{y_1}, \boldsymbol{y_2}|\boldsymbol{f}^W(\boldsymbol{x})) = p(\boldsymbol{y_1}|\boldsymbol{f}^W(\boldsymbol{x})) \cdot p(\boldsymbol{y_2}|\boldsymbol{f}^W(\boldsymbol{x}))$$

- Multi-Task Loss

$$p(\boldsymbol{y}|\boldsymbol{f^W}(\boldsymbol{x})) = Softmax(\boldsymbol{f^W}(\boldsymbol{x}))$$

$$\log p(\boldsymbol{y}|\boldsymbol{f^W}(\boldsymbol{x})) \propto -\frac{1}{2\sigma^2}||\boldsymbol{y} - \boldsymbol{f^W}(\mathbf{x})||^2 - \log\sigma$$

$$\log p(\boldsymbol{y_1}, \boldsymbol{y_2}|\boldsymbol{f^W}(\boldsymbol{x})) = p(\boldsymbol{y_1}|\boldsymbol{f^W}(\boldsymbol{x})) \cdot p(\boldsymbol{y_2}|\boldsymbol{f^W}(\boldsymbol{x}))$$

$$L(\mathbf{W}, \sigma_1, \sigma_2) = -\log p(\boldsymbol{y_1}, \boldsymbol{y_2}|\boldsymbol{f^W}(\boldsymbol{x}))$$
$$\propto \frac{1}{2\sigma_1^2}||\boldsymbol{y_1} - \boldsymbol{f^W}(\mathbf{x})||^2 + \frac{1}{2\sigma_2^2}||\boldsymbol{y_2} - \boldsymbol{f^W}(\mathbf{x})||^2 + \log\sigma_1\sigma_2$$
$$= \frac{1}{2\sigma_1^2}L_1(\boldsymbol{W}) + \frac{1}{2\sigma_2^2}L_2(\boldsymbol{W}) + \log\sigma_1\sigma_2$$

- Multi-Task Loss

$$L_{G1} = -D(x + G(z))$$

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

$$L_G = \frac{1}{2\sigma_1^2} L_{G1} + \frac{1}{2\sigma_2^2} L_{G2} + \log(\sigma_1 \sigma_2)$$
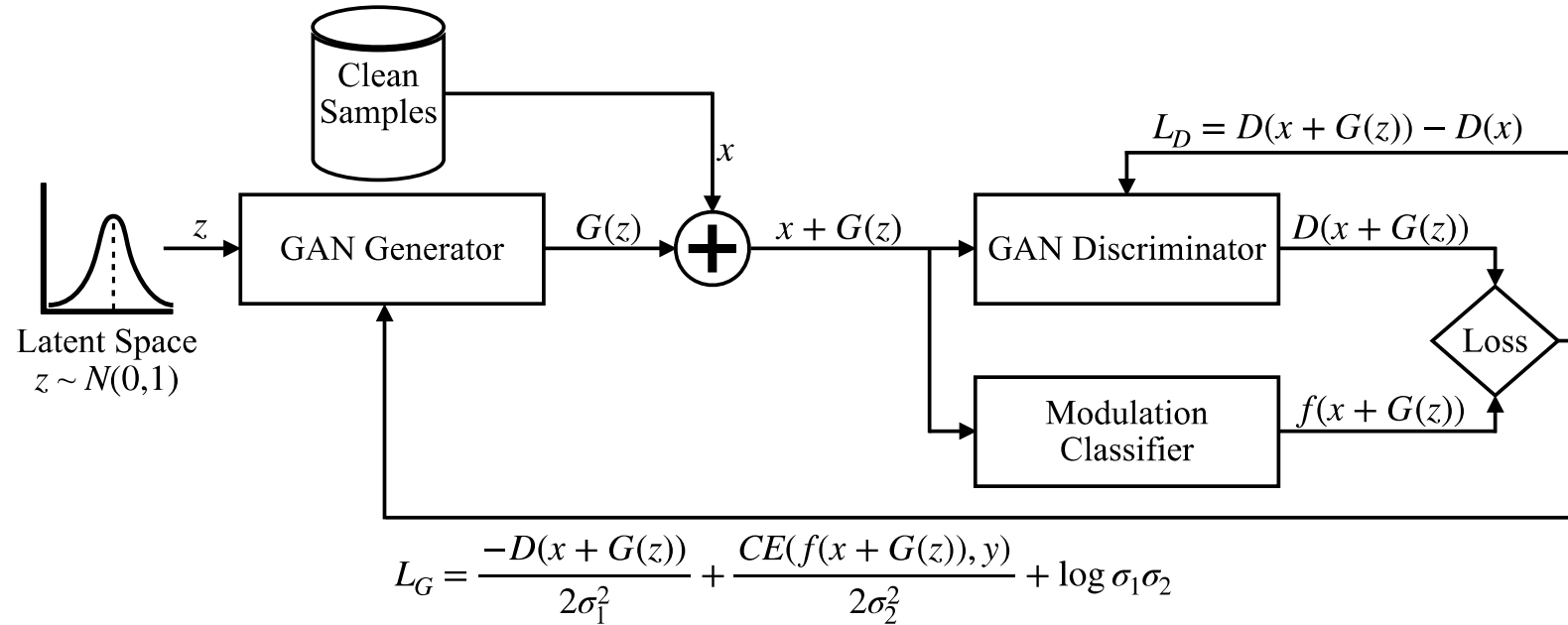
- Multi-Task Loss

$$L_{G1} = -D(x + G(z))$$

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

$$L_G = \frac{1}{2\sigma_1^2} L_{G1} + \frac{1}{2\sigma_2^2} L_{G2} + \log(\sigma_1 \sigma_2)$$

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

TEMPEST talks
2022

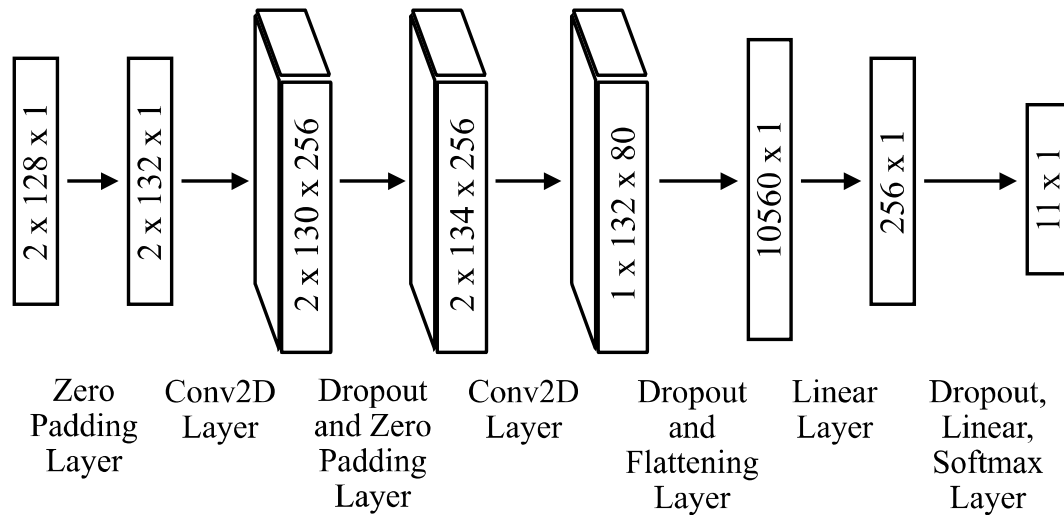- Modificamos a estrutura da GAN e a combinamos com a Multi-Task Loss
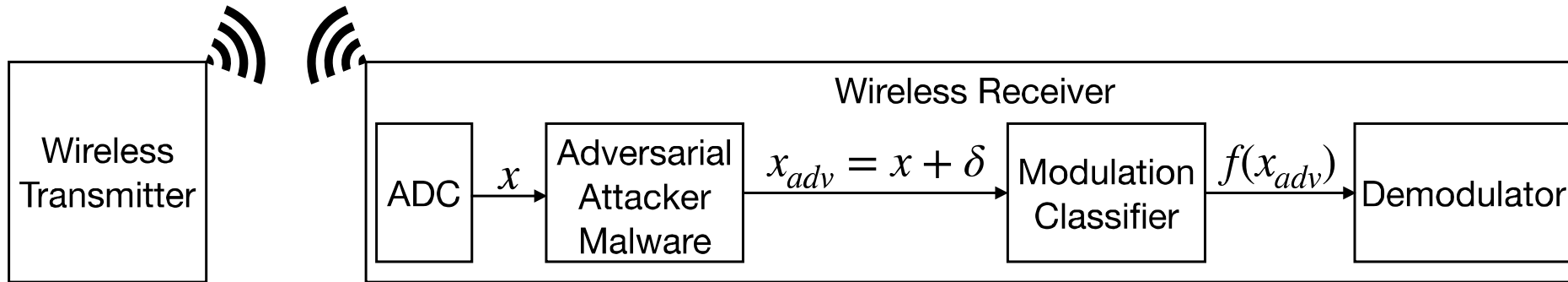


$$L_D = D\big(x + G(z)\big) - D(x)$$

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

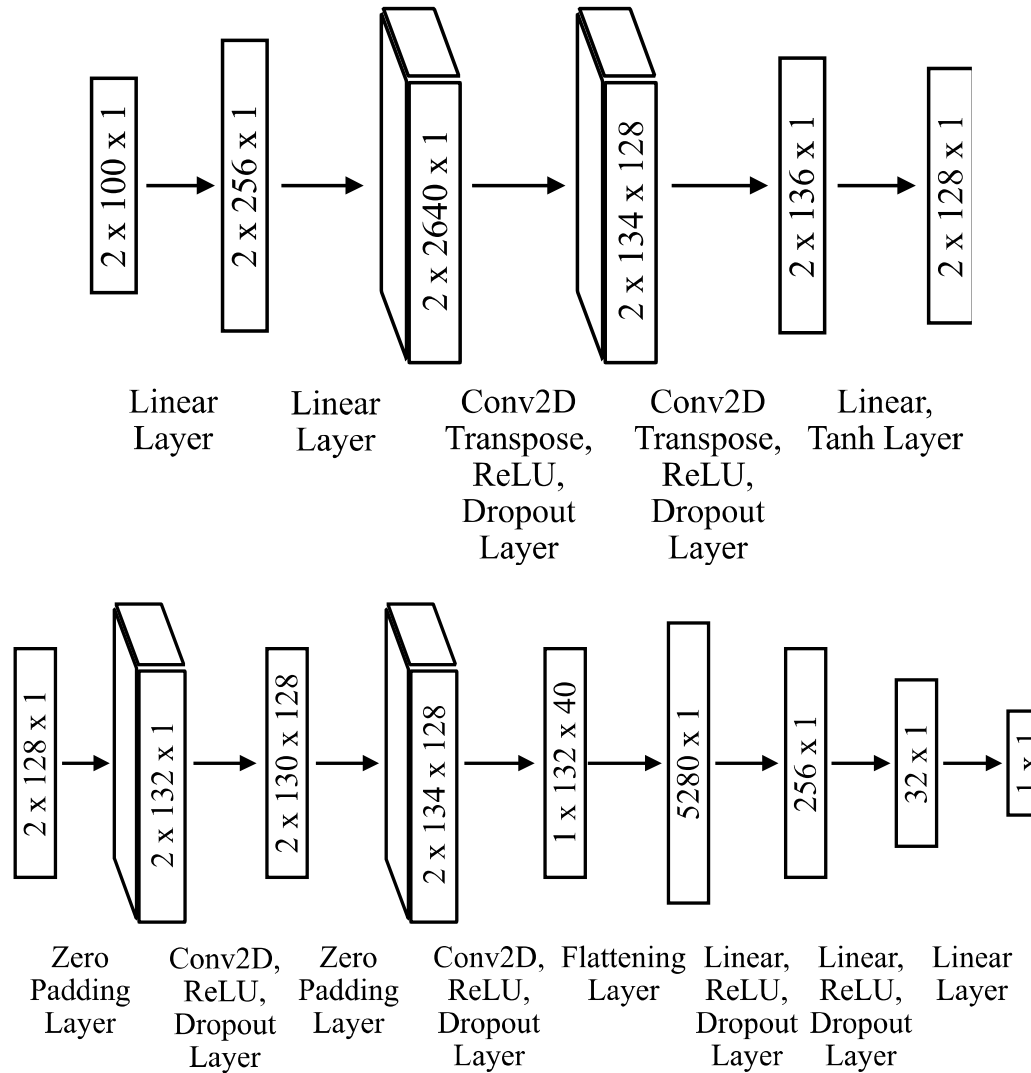- Multi-Objective GAN-Based Adversarial Attack

---

**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 $G(z)$

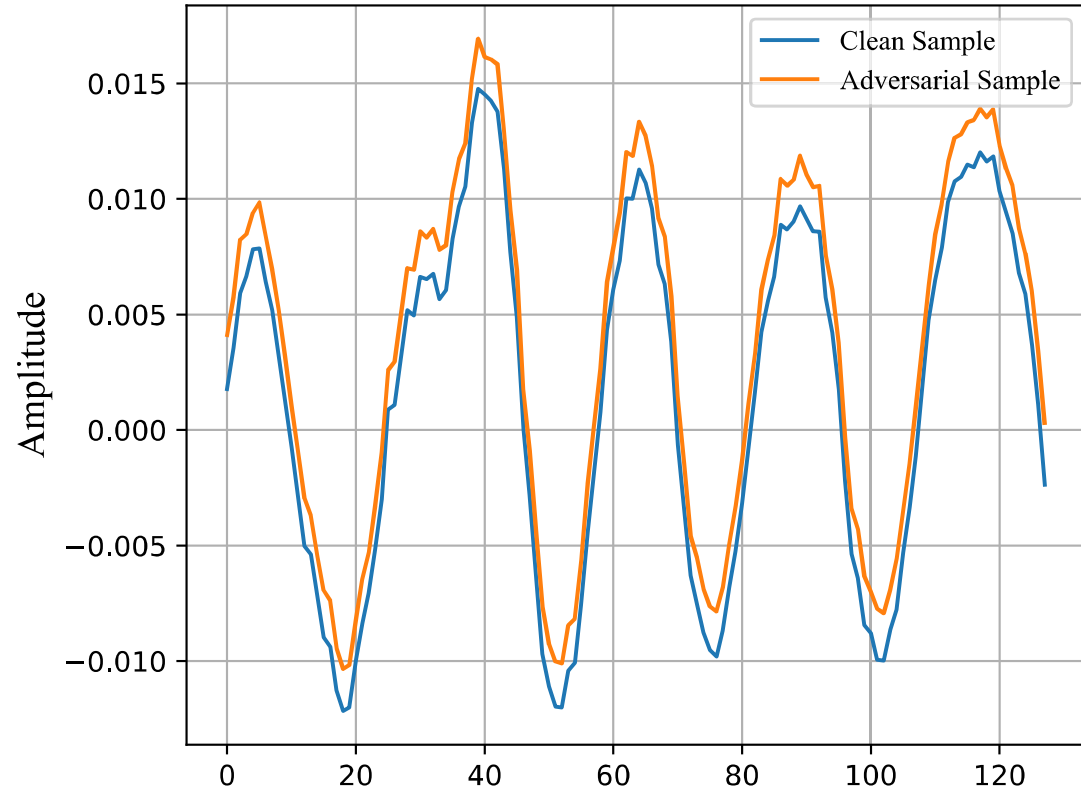4:     Construct the adversarial sample $x_{adv} = x + G(z)$

5: **end for**

---

TEMPEST talks
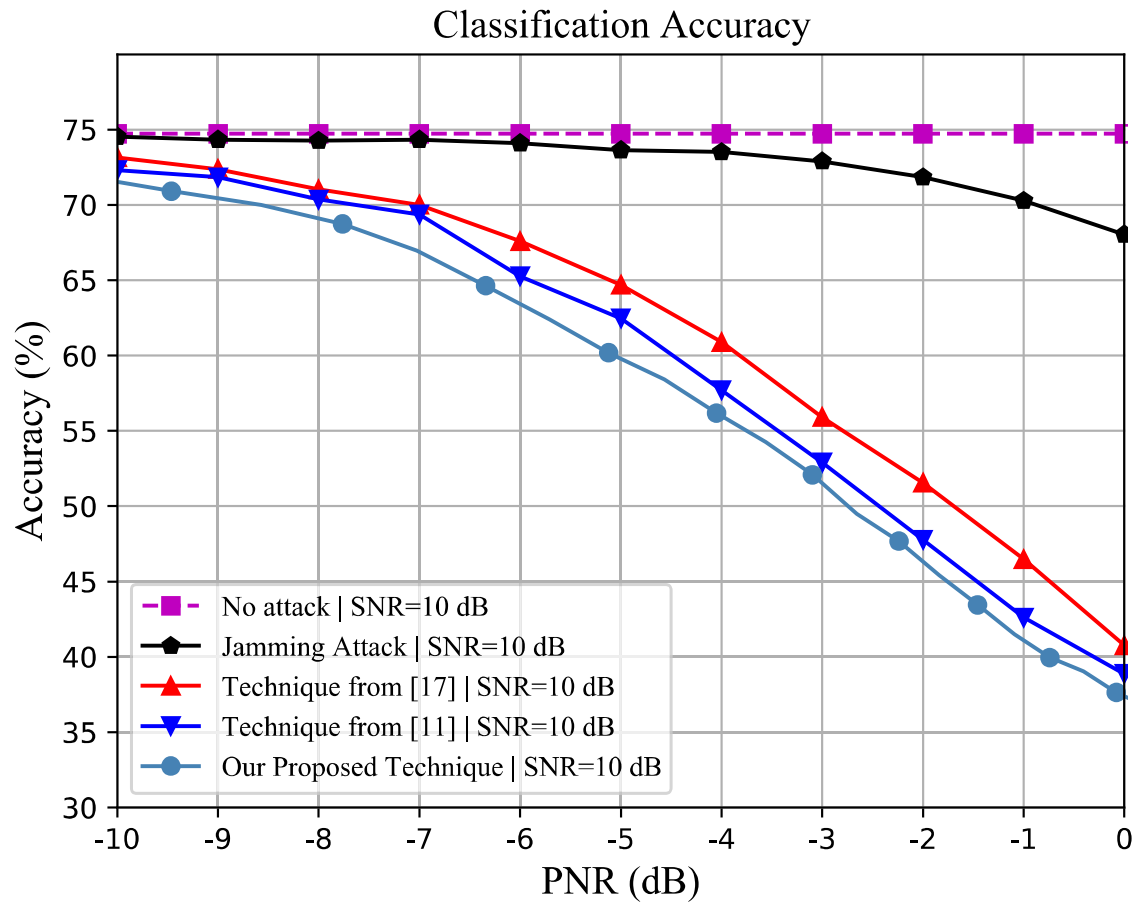
2022

- Multi-Objective GAN-Based Adversarial Attack

- Multi-Objective GAN-Based Adversarial Attack

- Multi-Objective GAN-Based Adversarial Attack
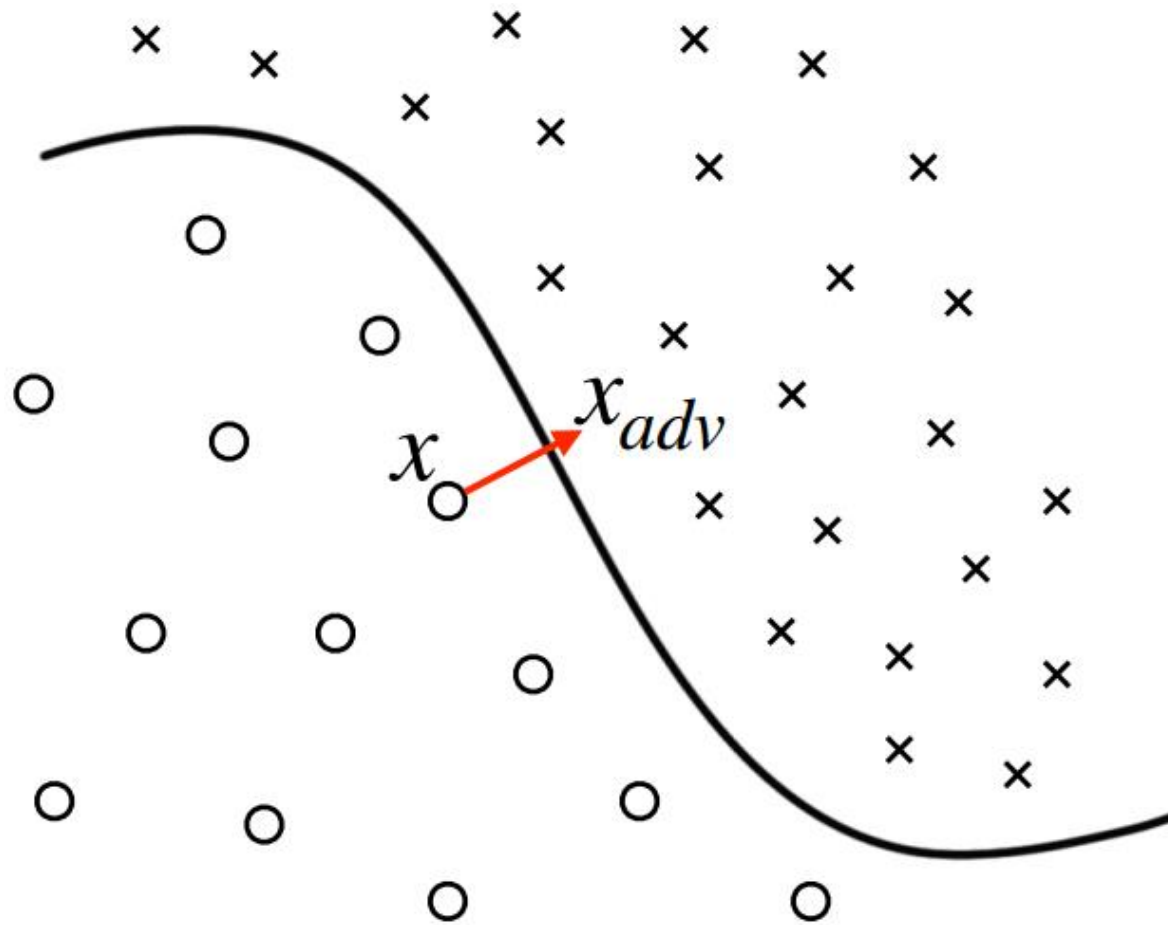
- Multi-Objective GAN-Based Adversarial Attack



Classification Accuracy

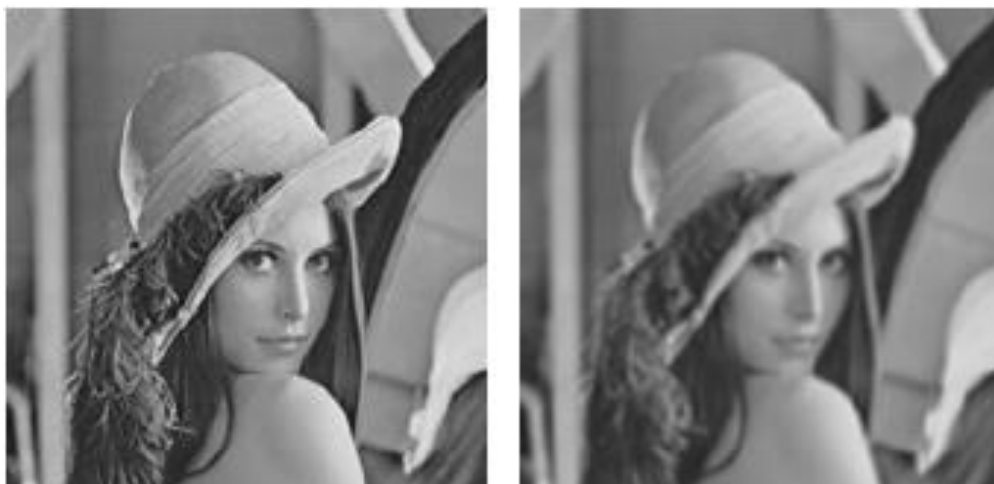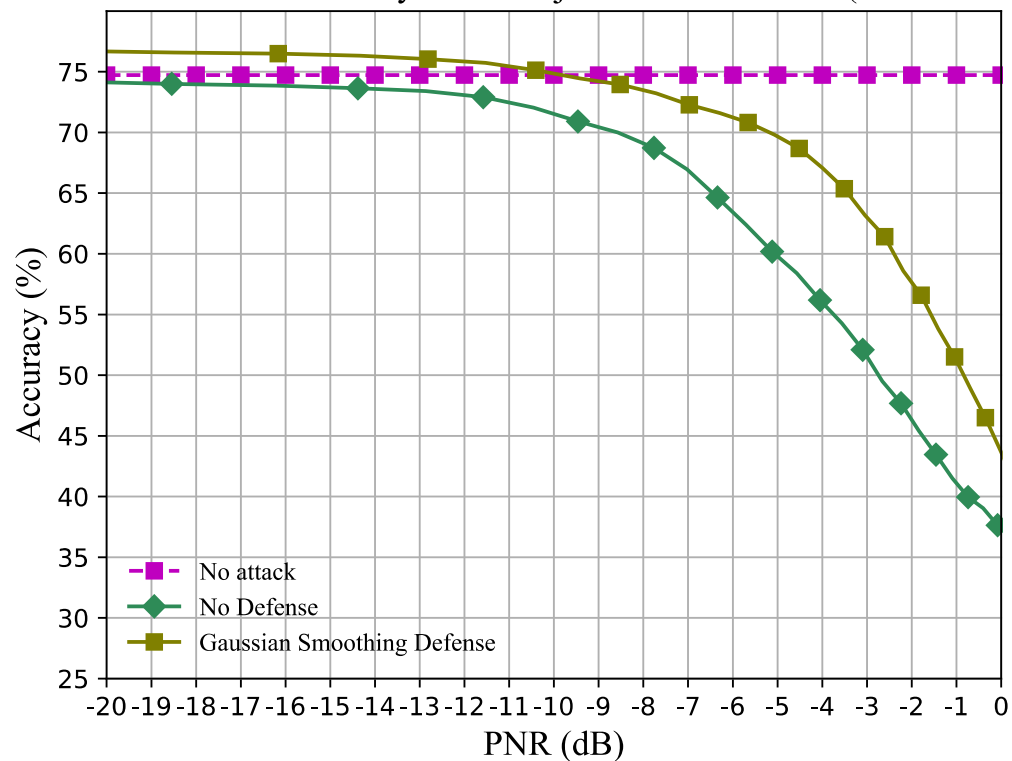| Adversarial Attack Technique | Mean Execution Time per Sample |
|---|---|
| Technique from [17] | 20189 $ms$ |
| Technique from [11] | 234 $ms$ |
| **Our Proposed Technique** | 0.6980 $ms$ |

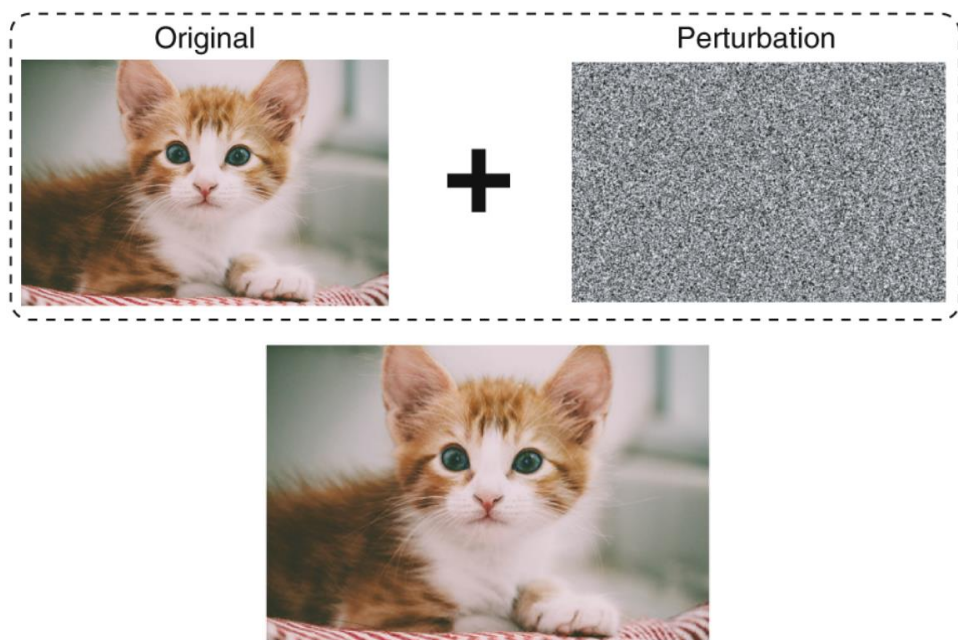- Diminuir a sensibilidade das fronteiras de decisão

- Diminuir a sensibilidade das fronteiras de decisão
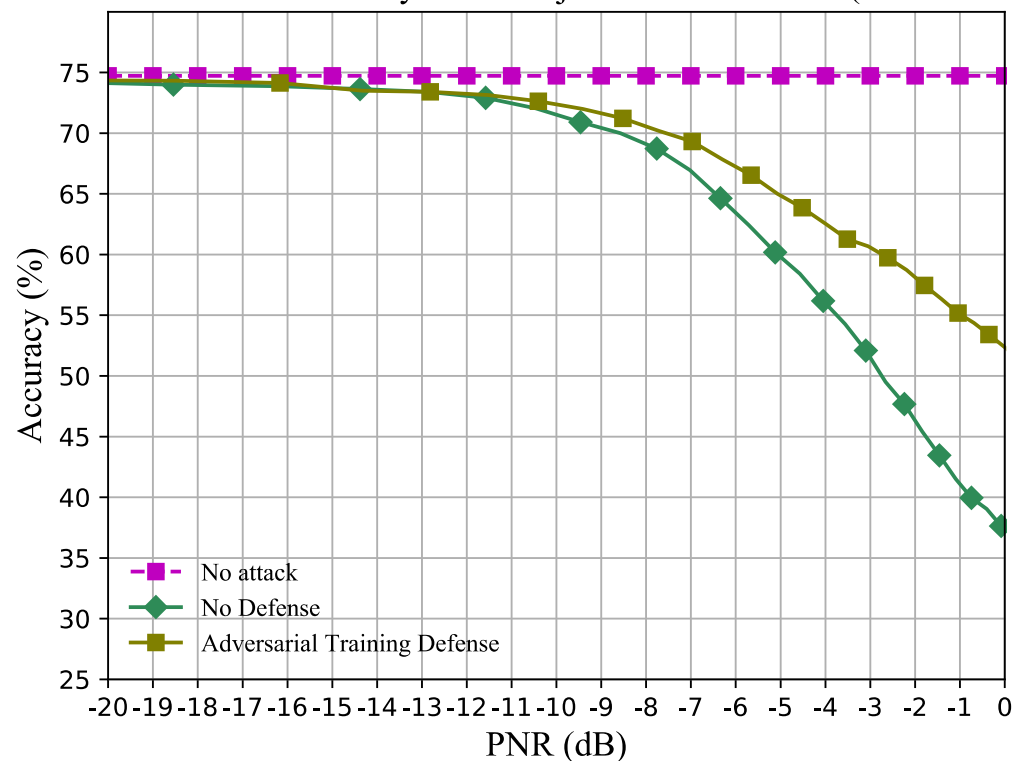


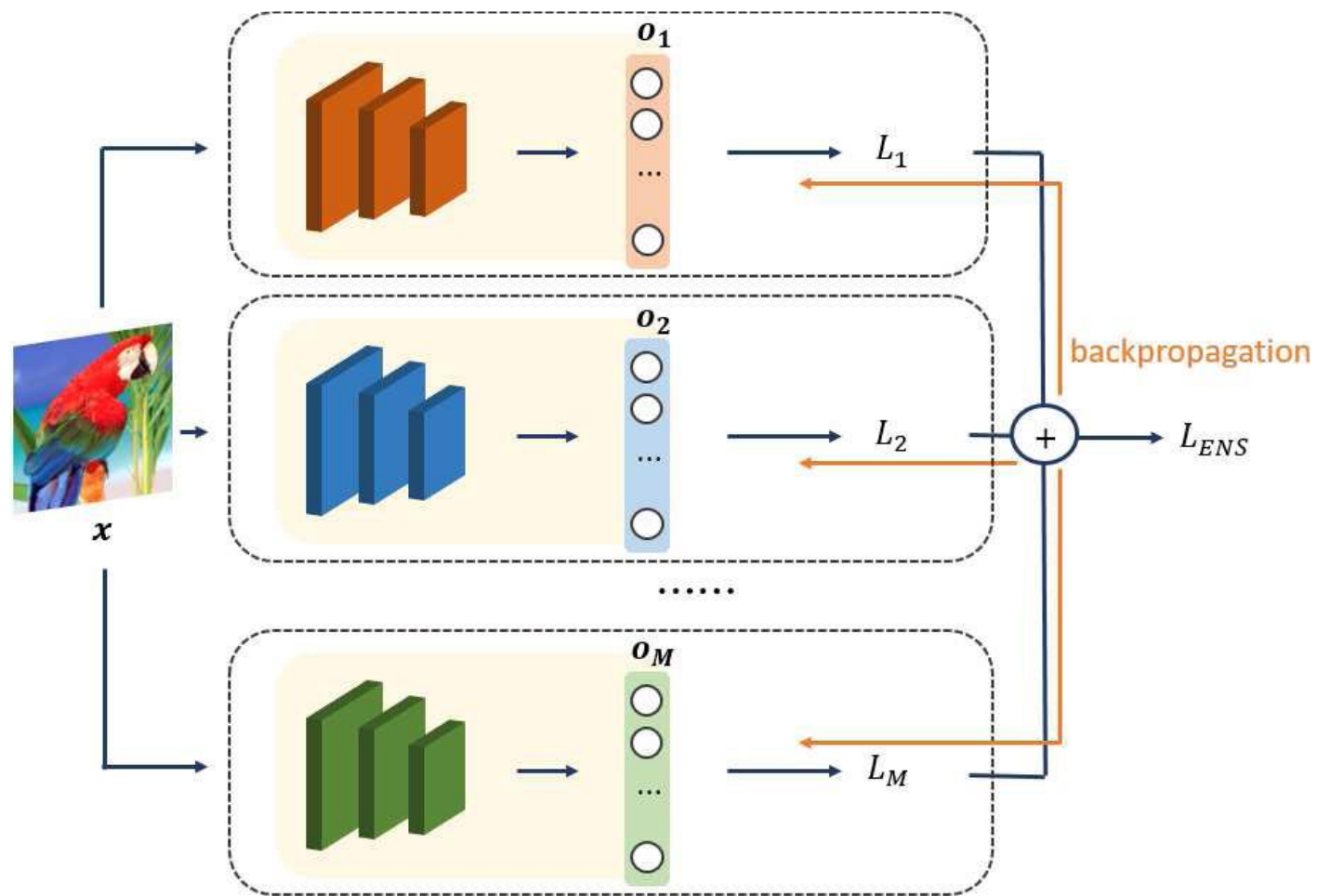Classification Accuracy Multi-Objective GAN Attack (SNR=10dB)

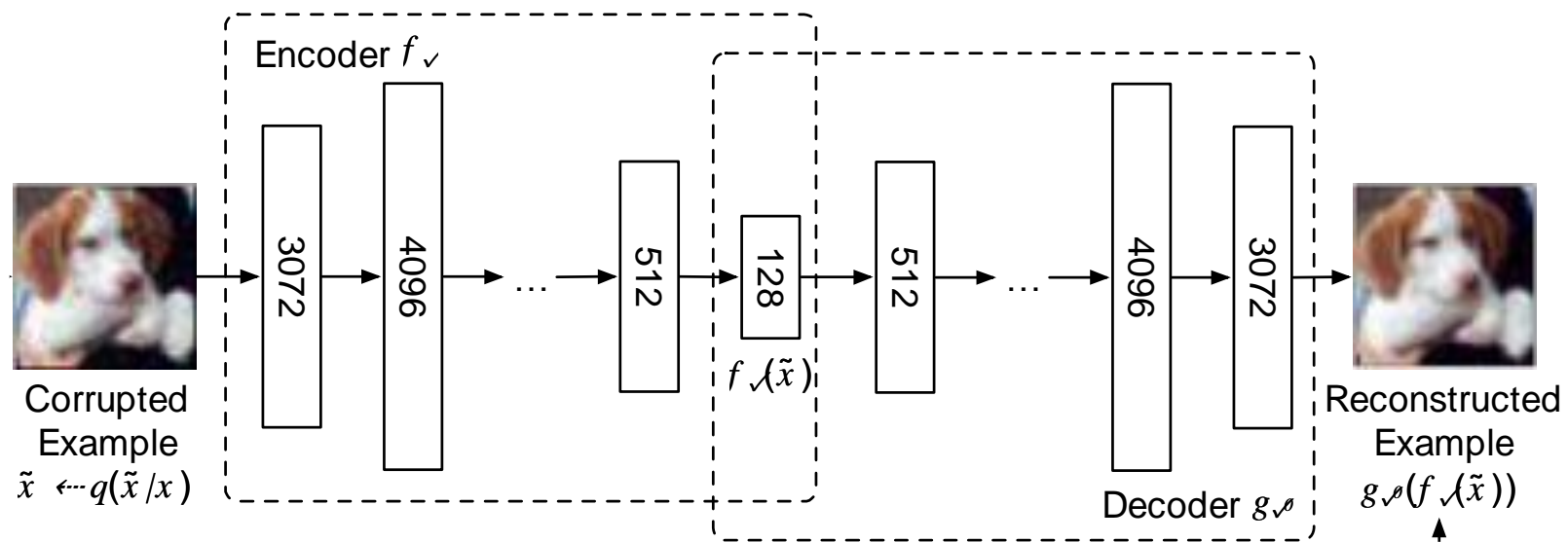- Diminuir a sensibilidade das fronteiras de decisão



Original

Perturbation

+

Classification Accuracy Multi-Objective GAN Attack (SNR=10dB)



- - No attack
- No Defense
- Adversarial Training Defense

Accuracy (%)

PNR (dB)

TEMPEST talks

2022

- Combinação de modelos

- Remoção de ruído e perturbações adversariais

# Obrigado!

**TEMPEST** talks

2022

P. Freitas de Araujo-Filho, G. Kaddoum, M. Naili, E. T. Fapi and Z. Zhu, "Multi-Objective GAN-Based Adversarial Attack Technique for Modulation Classifiers," in IEEE Communications Letters, vol. 26, no. 7, pp. 1583-1587, July 2022, doi: 10.1109/LCOMM.2022.3167368.

Paulo Freitas de Araujo Filho
paulo.freitas@tempest.com.br