# The Space of Adversarial Strategies

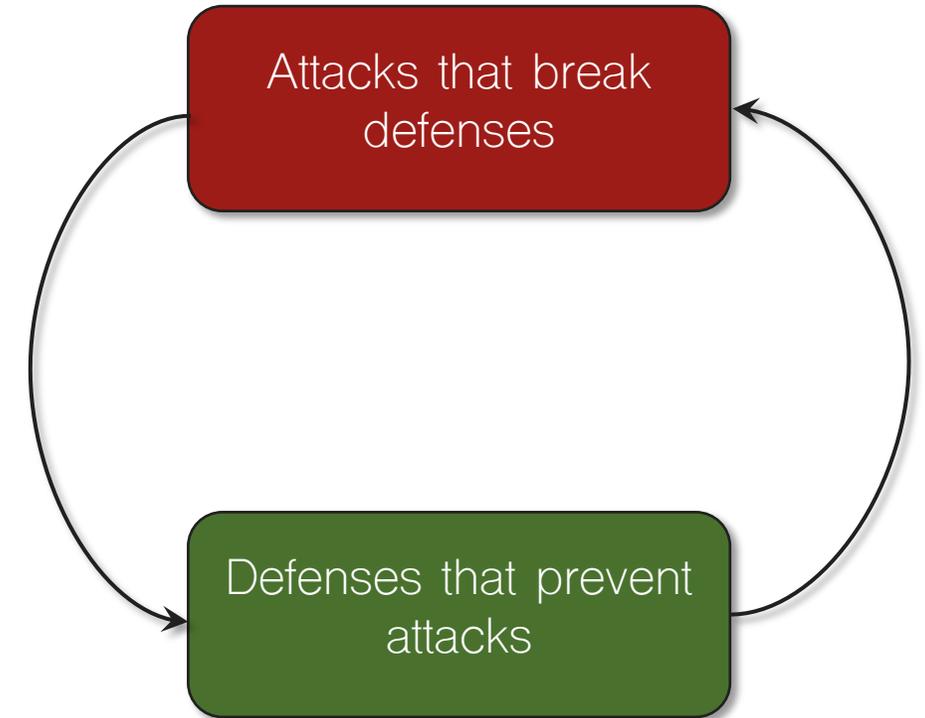Ryan Sheatsley*, **Blaine Hoak*,** Eric Pauley, Patrick McDaniel

*Equal Contribution. Thursday, August 10, 2023

MADS&P

IEEE Conference on Secure and Trustworthy Machine Learning

February 8-10, 2023

Hilton Raleigh North Hills

Raleigh, North Carolina

Sponsored by the IEEE Computer Society Technical Committee of Security and Privacy

Attacks that break defenses

Defenses that prevent attacks

LaserAttack

LowProFool

Basic Iterative Method

Adversarial Texture

Wasserstein Attack

Adversarial Patch

Brendel and Bethge Attack

NewtonFool

Projected Gradient Descent

Targeted Universal Perturbation Attack

Frame Saliency Attack

Auto Conjugate Attack

Square Attack

Geometric Decision Based Attack

Carlini and Wagner $l_0$

Carlini and Wagner $l_2$

Malware Gradient Descent

DPatch

Elastic Net Attack

DeepFool

Spatial Transforms Attack

Fast Gradient Method

RobustDPatch

Sign-OPT

ThresholdAttack

Over The Air Flickering Attack

Auto Attack

PixelAttack

Virtual Adversarial Method

GRAPHITE

Carlini and Wagner $l_\infty$

High Confidence Low Uncertainty Attack

Universal Perturbation Attack

Boundary Attack

Carlini and Wagner ASR

Decision Tree Attack

Jacobian Saliency Map Approach

HopSkipJump

Zeroth-Order Optimization

Feature Adversaries

Imperceptible ASR

Shadow Attack

Simple Black-box Adversarial Attack

LaserAttack

LowProFool

Basic Iterative Method

Adversarial Texture

Wasserstein Attack

Adversarial Patch

Brendel and Bethge Attack

NewtonFool

Projected Gradient Descent

Targeted Universal Perturbation Attack

Frame Saliency Attack

Auto Conjugate Attack

Square Attack

Geometric Decision Based Attack

Carlini and Wagner $l_0$

Malware Grad

- How can we systematically represent and evaluate attacks?

Over T

Sign-OPT

versarial Method

GRAPHITE

Carlini and Wagner $l_\infty$

High Confidence Low Uncertainty Attack

Boundary Attack

Universal Perturbation Attack

Carlini and Wagner ASR

HopSkipJump

Decision Tree Attack

Jacobian Saliency Map Approach

Zeroth-Order Optimization

Feature Adversaries

Imperceptible ASR

Shadow Attack

Simple Black-box Adversarial Attack

*Surfaces* compute the perturbation

*Travelers* apply perturbations and other techniques that update the input

**BIM**
$$X_0^{adv} = X$$
$$X_{N+1}^{adv} = Clip_{X,\epsilon}\left\{X_N^{adv} + \alpha\,\text{sign}(\nabla_X J(X_N^{adv}, y_{true}))\right\}$$

**FAB**
$$s \leftarrow \arg\min_{l \neq c} \frac{|f_l(x^{(i)}) - f_c(x^{(i)})|}{\|\nabla f_l(x^{(i)}) - \nabla f_c(x^{(i)})\|_q}$$
$$\delta^{(i)} \leftarrow \text{proj}_p(x^{(i)}, \pi_s, C)$$
$$\delta_{\text{orig}}^{(i)} \leftarrow \text{proj}_p(x_{\text{orig}}, \pi_s, C)$$
compute $\alpha$ as in Equation (9)
$$x^{(i+1)} \leftarrow \text{proj}_C\left((1-\alpha)\left(x^{(i)} + \eta\delta^{(i)}\right) + \alpha(x_{\text{orig}} + \eta\delta_{\text{orig}}^{(i)})\right)$$

**JSMA**
Compute forward derivative $\nabla\mathbf{F}(\mathbf{X}^*)$
$p_1, p_2 = \texttt{saliency\_map}(\nabla\mathbf{F}(\mathbf{X}^*), \Gamma, \mathbf{Y}^*)$
Modify $p_1$ and $p_2$ in $\mathbf{X}^*$ by $\theta$

**APGD**
$$z^{(k+1)} \leftarrow P_{\mathcal{S}}\left(x^{(k)} + \eta\nabla f(x^{(k)})\right)$$
$$x^{(k+1)} \leftarrow P_{\mathcal{S}}\left(x^{(k)} + \alpha(z^{(k+1)} - x^{(k)})\right.$$
$$\left.+ (1-\alpha)(x^{(k)} - x^{(k-1)})\right)$$

**DeepFool**
$$\hat{l} \leftarrow \arg\min_{k \neq \hat{k}(\boldsymbol{x}_0)} \frac{|f_k'|}{\|\boldsymbol{w}_k'\|_2}$$
$$\boldsymbol{r}_i \leftarrow \frac{|f_{\hat{l}}'|}{\|\boldsymbol{w}_{\hat{l}}'\|_2^2}\boldsymbol{w}_{\hat{l}}'$$
$$\boldsymbol{x}_{i+1} \leftarrow \boldsymbol{x}_i + \boldsymbol{r}_i$$
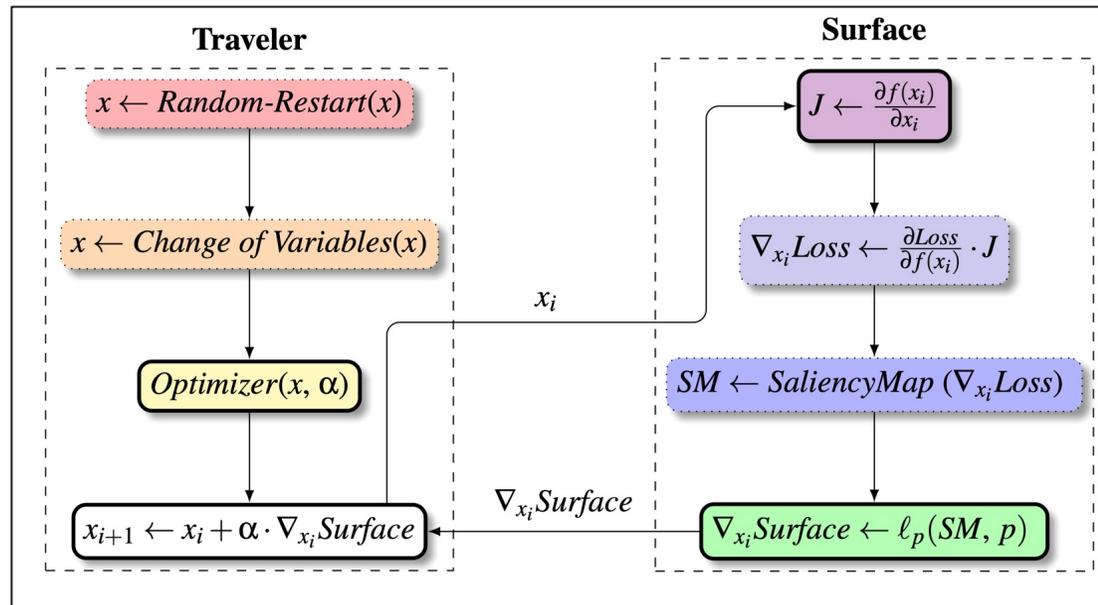
**PGD**
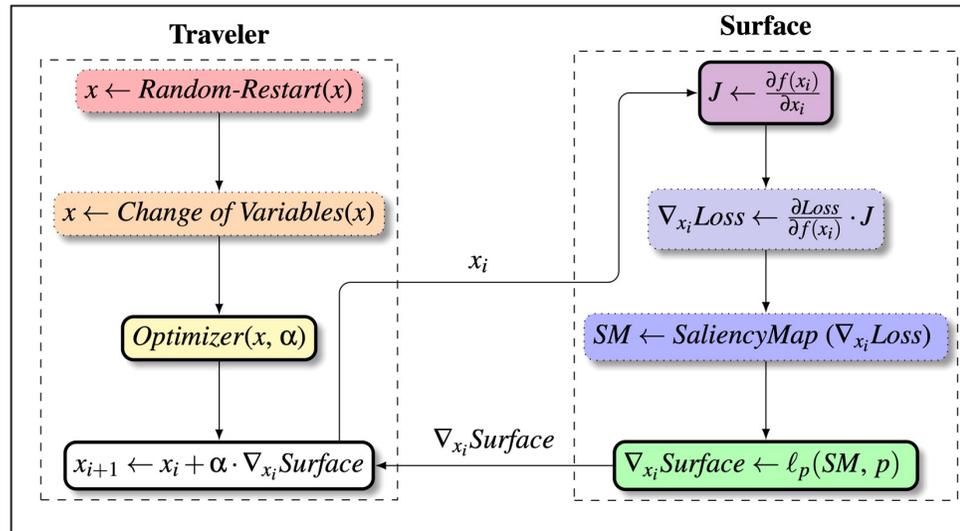$$x^{t+1} = \Pi_{x+\mathcal{S}}\left(x^t + \alpha\,\text{sgn}(\nabla_x L(\theta, x, y))\right)$$

*Components of the traveler and surface. Arrows represent the progression of an an attack through components at each iteration.*

**Traveler**
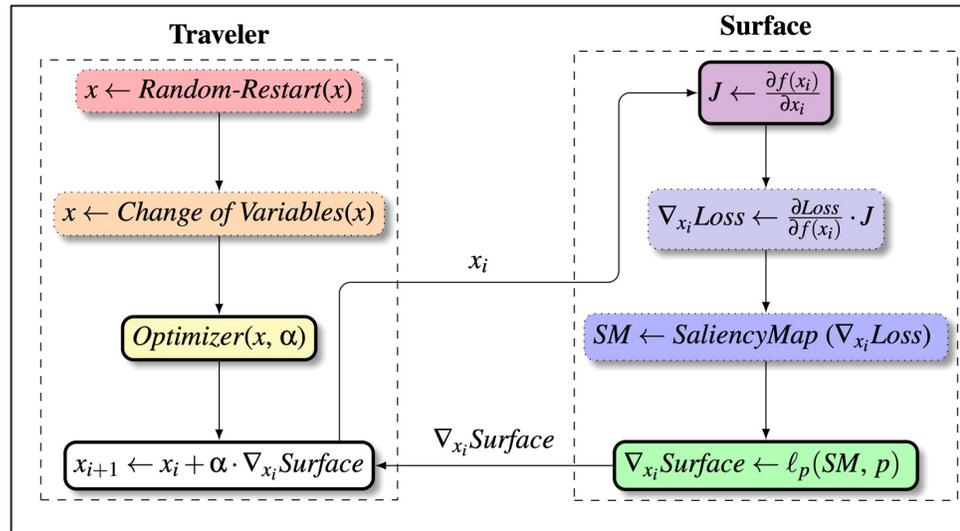
$x \leftarrow Random\text{-}Restart(x)$

$x \leftarrow Change\ of\ Variables(x)$

$Optimizer(x, \alpha)$

$x_{i+1} \leftarrow x_i + \alpha \cdot \nabla_{x_i}Surface$

$x_i$

$\nabla_{x_i}Surface$

**Surface**

$J \leftarrow \dfrac{\partial f(x_i)}{\partial x_i}$

$\nabla_{x_i}Loss \leftarrow \dfrac{\partial Loss}{\partial f(x_i)} \cdot J$

$SM \leftarrow SaliencyMap\,(\nabla_{x_i}Loss)$

$\nabla_{x_i}Surface \leftarrow \ell_p(SM, p)$

**Attack Algorithms**

| | Surface Components | | Traveler Components | |
|---|---|---|---|---|
| *Losses:* | Cross-Entropy | | *Random-Restart*: | Enabled, Disabled |
| | Carlini-Wagner Loss | | | |
| | Identity Loss | | | |
| | Difference of Logits Ratio Loss | | | |
| *Saliency Maps:* | $SM_J, SM_D, SM_I$ | | *Change of Variables*: | Enabled, Disabled |
| $\ell_p$-*norm* | $\ell_0, \ell_2, \ell_\infty$ | | *Optimizer*: | SGD, Adam, MBS, BWSGD |

| | CE | CWL | IL | DLR | $SM_J$ | $SM_D$ | $SM_I$ | $\ell_0$ | $\ell_2$ | $\ell_\infty$ | RR | CoV | SGD | Adam | MBS | BWSGD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BIM | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ |
| PGD | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ● | ○ | ● | ○ | ○ | ○ |
| JSMA | ○ | ○ | ● | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ● | ○ | ○ | ○ |
| DF | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ● | ○ | ○ | ○ |
| CW | ○ | ● | ○ | ○ | ○ | ○ | ● | ○ | ● | ○ | ○ | ● | ○ | ● | ○ | ○ |
| APGD-CE | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ● | ○ | ○ | ○ | ● | ○ |
| APGD-DLR | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ● | ○ | ○ | ○ | ● | ○ |
| FAB | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ○ | ○ | ○ | ● |

PGD: $\nabla \mathrm{CE}(x, y)$

DeepFool: $\nabla f_y(x)$

**Traveler**

$x \leftarrow \textit{Random-Restart}(x)$

$x \leftarrow \textit{Change of Variables}(x)$

$\textit{Optimizer}(x, \alpha)$

$x_{i+1} \leftarrow x_i + \alpha \cdot \nabla_{x_i} Surface$

**Surface**

$J \leftarrow \frac{\partial f(x_i)}{\partial x_i}$

$\nabla_{x_i} Loss \leftarrow \frac{\partial Loss}{\partial f(x_i)} \cdot J$

$SM \leftarrow SaliencyMap(\nabla_{x_i} Loss)$

$\nabla_{x_i} Surface \leftarrow \ell_p(SM, p)$

$x_i$

$\nabla_{x_i} Surface$

**Attack Algorithms**

| | Surface Components | | Traveler Components | |
|---|---|---|---|---|
| *Losses:* | Cross-Entropy | | *Random-Restart:* | Enabled, Disabled |
| | Carlini-Wagner Loss | | | |
| | Identity Loss | | | |
| | Difference of Logits Ratio Loss | | | |
| *Saliency Maps:* | $SM_J$, $SM_D$, $SM_I$ | | *Change of Variables:* | Enabled, Disabled |
| $\ell_p\text{-norm}$ | $\ell_0$, $\ell_2$, $\ell_\infty$ | | *Optimizer:* | SGD, Adam, MBS, BWSGD |

| | CE | CWL | IL | DLR | $SM_J$ | $SM_D$ | $SM_I$ | $\ell_0$ | $\ell_2$ | $\ell_\infty$ | RR | CoV | SGD | Adam | MBS | BWSGD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BIM | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ |
| PGD | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ● | ○ | ● | ○ | ○ | ○ |
| JSMA | ○ | ○ | ● | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ○ | ● | ○ | ○ | ○ |
| DF | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ● | ○ | ○ | ○ |
| CW | ○ | ● | ○ | ○ | ○ | ○ | ● | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ |
| APGD-CE | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ● | ○ | ○ |
| APGD-DLR | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ● | ○ | ○ |
| FAB | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ○ | ○ | ○ | ● |

PGD: $\qquad \nabla \text{CE}(x, y)$

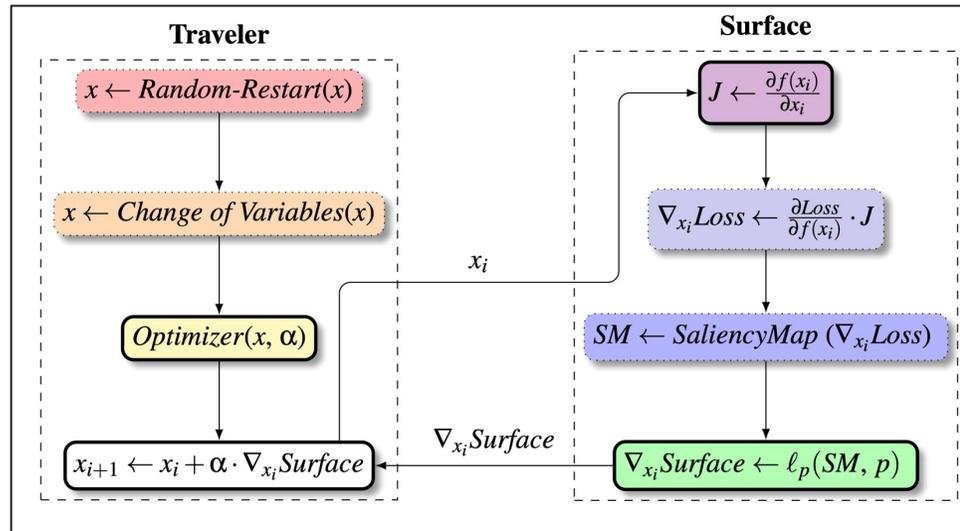DeepFool: $\quad \dfrac{|f_y(x) - f_k(x)|}{\nabla f_y(x) - \nabla f_k(x)} \quad (\nabla f_y(x) - \nabla f_k(x))$

**Traveler**

$x \leftarrow$ *Random-Restart*$(x)$

$x \leftarrow$ *Change of Variables*$(x)$

*Optimizer*$(x, \alpha)$

$x_{i+1} \leftarrow x_i + \alpha \cdot \nabla_{x_i} Surface$

**Surface**

$J \leftarrow \frac{\partial f(x_i)}{\partial x_i}$

$\nabla_{x_i} Loss \leftarrow \frac{\partial Loss}{\partial f(x_i)} \cdot J$

$SM \leftarrow SaliencyMap\,(\nabla_{x_i} Loss)$

$\nabla_{x_i} Surface \leftarrow \ell_p(SM, p)$

$x_i$

$\nabla_{x_i} Surface$

**Attack Algorithms**

| | Surface Components | | Traveler Components | |
|---|---|---|---|---|
| *Losses:* | Cross-Entropy | | *Random-Restart*: | Enabled, Disabled |
| | Carlini-Wagner Loss | | | |
| | Identity Loss | | | |
| | Difference of Logits Ratio Loss | | | |
| *Saliency Maps:* | $SM_J$, $SM_D$, $SM_I$ | | *Change of Variables*: | Enabled, Disabled |
| $\ell_p$-*norm* | $\ell_0, \ell_2, \ell_\infty$ | | *Optimizer*: | SGD, Adam, MBS, BWSGD |

| | CE | CWL | IL | DLR | $SM_J$ | $SM_D$ | $SM_I$ | $\ell_0$ | $\ell_2$ | $\ell_\infty$ | RR | CoV | SGD | Adam | MBS | BWSGD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BIM | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ |
| PGD | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ● | ○ | ● | ○ | ○ | ○ |
| JSMA | ○ | ○ | ● | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ○ | ● | ○ | ○ | ○ |
| DF | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ● | ○ | ○ | ○ |
| CW | ○ | ● | ○ | ○ | ○ | ○ | ● | ○ | ● | ○ | ○ | ● | ○ | ● | ○ | ○ |
| APGD-CE | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ● | ● | ● | ○ | ○ | ○ | ● | ○ |
| APGD-DLR | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ● | ● | ● | ○ | ○ | ○ | ● | ○ |
| FAB | ○ | ○ | ● | ○ | ○ | ○ | ● | ○ | ● | ○ | ○ | ○ | ○ | ○ | ○ | ● |

PGD: $\delta = \alpha \cdot \mathrm{sgn}(\nabla \mathrm{CE}(x, y))$
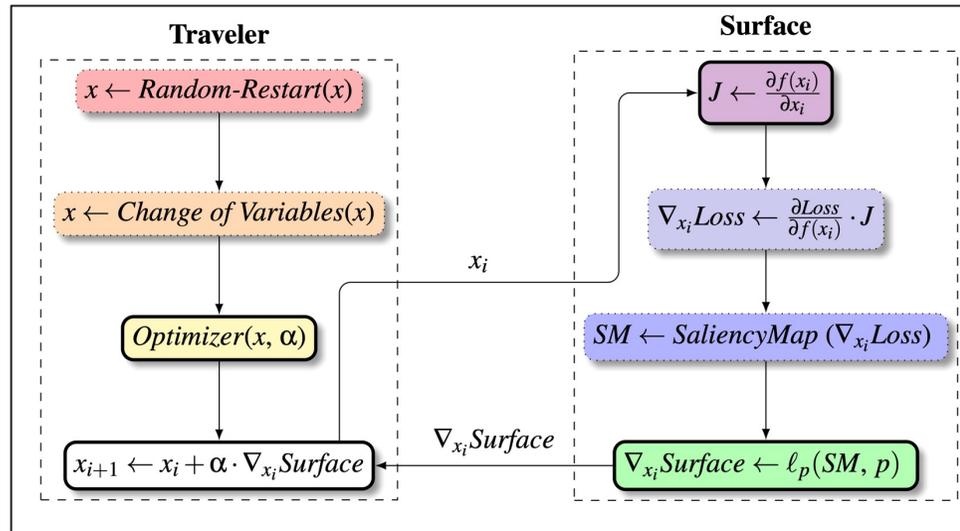
DeepFool: $\delta = \dfrac{|f_y(x) - f_k(x)|}{\|\nabla f_y(x) - \nabla f_k(x)\|_2^2} \cdot (\nabla f_y(x) - \nabla f_k(x))$
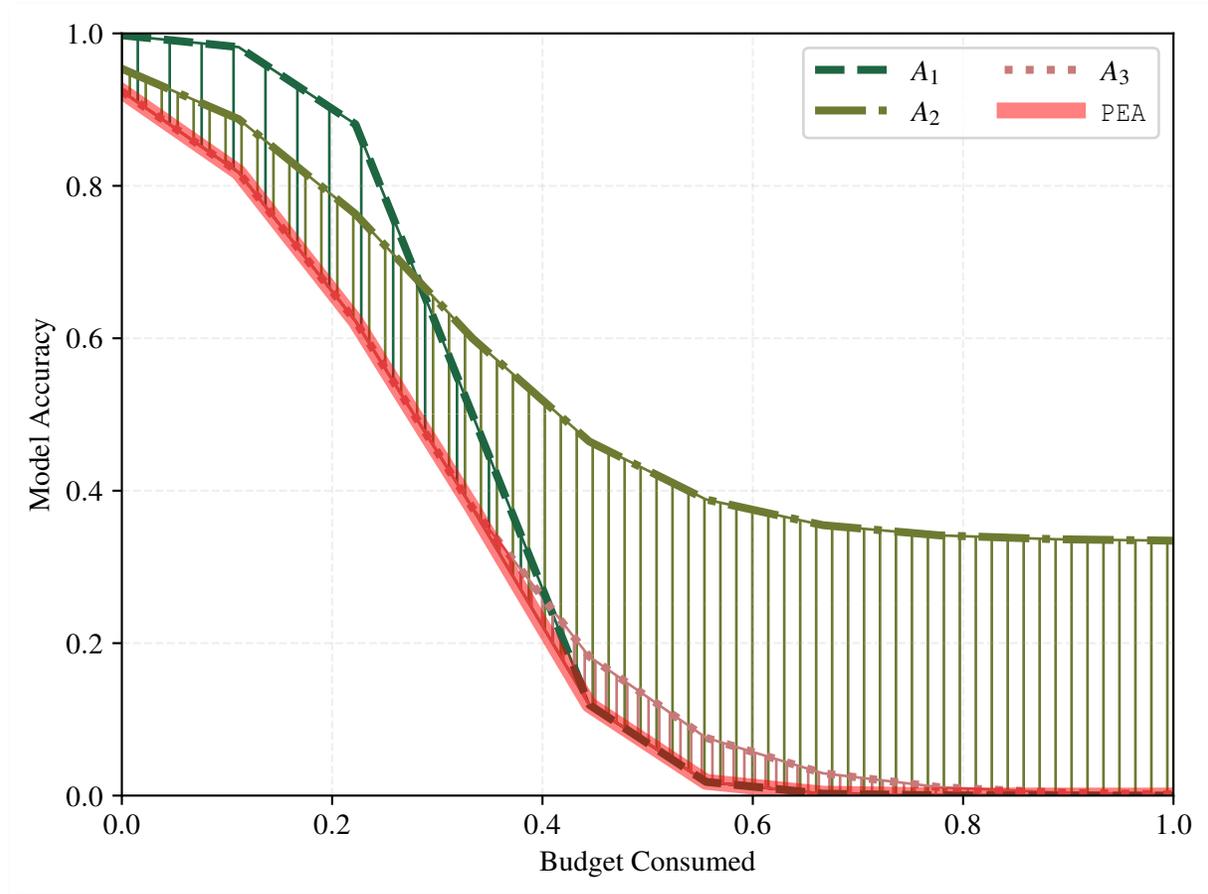
Traveler

$x \leftarrow Random\text{-}Restart(x)$

$x \leftarrow Change\ of\ Variables(x)$

$Optimizer(x, \alpha)$

$x_{i+1} \leftarrow x_i + \alpha \cdot \nabla_{x_i} Surface$

$x_i$

$\nabla_{x_i} Surface$

Surface

$J \leftarrow \frac{\partial f(x_i)}{\partial x_i}$

$\nabla_{x_i} Loss \leftarrow \frac{\partial Loss}{\partial f(x_i)} \cdot J$

$SM \leftarrow SaliencyMap\ (\nabla_{x_i} Loss)$

$\nabla_{x_i} Surface \leftarrow \ell_p(SM, p)$

**Attack Algorithms**

| | Surface Components | | Traveler Components | |
|---|---|---|---|---|
| *Losses:* | Cross-Entropy | | *Random-Restart:* | Enabled, Disabled |
| | Carlini-Wagner Loss | | | |
| | Identity Loss | | | |
| | Difference of Logits Ratio Loss | | | |
| *Saliency Maps:* | $SM_J, SM_D, SM_I$ | | *Change of Variables:* | Enabled, Disabled |
| $\ell_p$-*norm* | $\ell_0, \ell_2, \ell_\infty$ | | *Optimizer:* | SGD, Adam, MBS, BWSGD |

| | CE | CWL | IL | DLR | SM$_J$ | SM$_D$ | SM$_I$ | $\ell_0$ | $\ell_2$ | $\ell_\infty$ | RR | CoV | SGD | Adam | MBS | BWSGD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BIM | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ |
| PGD | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ● | ○ | ● | ○ | ○ | ○ |
| JSMA | ○ | ○ | ● | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ○ | ● | ○ | ○ | ○ |
| DF | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ● | ○ | ○ | ○ |
| CW | ○ | ● | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ● | ○ | ○ |
| APGD-CE | ● | ○ | ○ | ○ | ○ | ○ | ● | ○ | ○ | ● | ● | ○ | ○ | ○ | ● | ○ |
| APGD-DLR | ○ | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ● | ○ | ○ | ○ | ● | ○ |
| FAB | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ● | ○ | ○ | ○ | ○ | ○ | ○ | ● |

$$\text{PGD: } \delta = \alpha \cdot \text{sgn}(\nabla \text{CE}(x, y)), x_0 = x + \mathcal{U}(-\epsilon, \epsilon), x_{i+1} = x_i + \delta$$

$$\text{DeepFool: } \delta = \frac{|f_y(x) - f_k(x)|}{\|\nabla f_y(x) - \nabla f_k(x)\|_2^2} \cdot (\nabla f_y(x) - \nabla f_k(x)), x_0 = x, x_{i+1} = x_i + \delta$$

Our ***extensible*** decomposition of ***mutually compatible and independent*** components allows us to build a vast attack space containing 576 attacks.

# The Pareto Ensemble Attack

# Evaluation

- Questions
  - When and why are attacks performant?
- Setup
  - Adversary has access to model parameters
  - CIC-MalMem2022, Malware Detection, 58k total (k-Fold), 4 classes
  - CIFAR-10, Object Classification, 50k train, 10k test, 10 classes
  - Fashion-MNIST, Clothing Classification, 60k train, 10k test, 10 classes
  - MNIST, Digit Recognition, 60k train, 10k test, 10 classes
  - NSL-KDD, Network Intrusion Detection, 125k train, 22k test, 5 classes
  - Phishing Websites, Phishing Detection, 10k total (k-Fold), 2 classes
  - UNSW-NB15, Network Intrusion Detection, 101k train, 53k test, 10 classes

# Hypothesis Testing Illuminates Effective Strategies

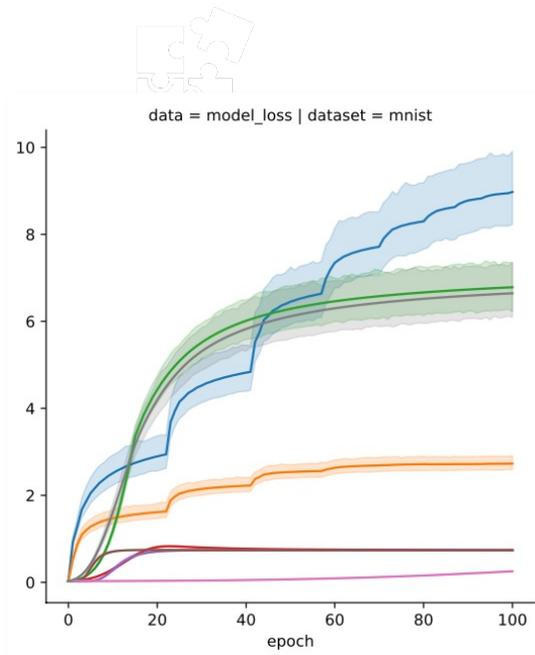| | Component H₁ | | Component H₂ | | Condition | p-value | Effect Size |
|---|---|---|---|---|---|---|---|
| 1. | SGD | is better than | BWSGD | when | $\texttt{Dataset} = \texttt{MNIST}$ | $<2.2 \times 10^{-308}$ | $99\%$ |
| 2. | Adam | is better than | BWSGD | when | $\texttt{Dataset} = \texttt{MNIST}$ | $<2.2 \times 10^{-308}$ | $99\%$ |
| | | | ⋮ | | ⋮ | | |
| 84. | Identity Loss | is better than | Difference of Logits Ratio Loss | when | $\texttt{Dataset} = \texttt{NSL-KDD}$ | $<2.2 \times 10^{-308}$ | $93\%$ |
| 85. | SGD | is better than | BWSGD | when | $\texttt{SaliencyMap} = \texttt{Jacobian Saliency Map}$ | $<2.2 \times 10^{-308}$ | $92\%$ |
| | | | ⋮ | | ⋮ | | |
| 393. | DeepFool Saliency Map | is better than | Jacobian Saliency Map | when | $\texttt{Dataset} = \texttt{FMNIST}$ | $<5 \times 10^{-6}$ | $66\%$ |
| 394. | Cross-Entropy | is better than | Carlini-Wagner Loss | when | *Change of Variables = Disabled* | $<5 \times 10^{-6}$ | $61\%$ |
| | | | ⋮ | | ⋮ | | |
| 1689. | $\ell_0$ | is better than | $\ell_2$ | when | $\texttt{Threat Model} = \ell_2 + 1.0$ | $9.8 \times 10^{-1}$ | $50\%$ |
| 1690. | Identity Saliency Map | is better than | DeepFool Saliency Map | when | $\texttt{Threat Model} = \ell_\infty + 0.4$ | $1.0$ | $49\%$ |

## *Hypothesis Testing Illuminates Effective Strategies*

- Change of Variables → 100% disabled, 0% enabled
- Optimizers → 50% Adam, 33% SGD, 16% MBS, 1% BWSGD
- Random Restart → 61% enabled, 39% disabled
- Saliency Maps → 70% no Saliency Map, 30% either DeepFool or JSMA Saliency Map
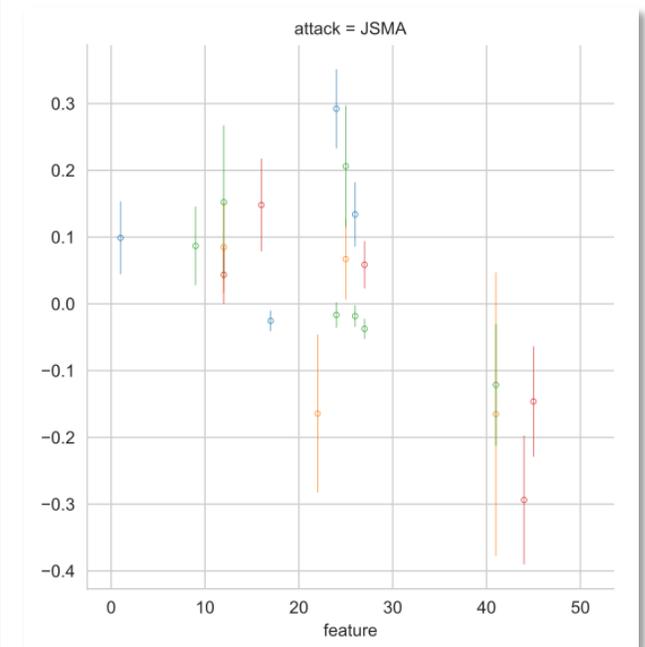- Loss → 47% Identity Loss, 34% Cross Entropy, 18% Carlini Wagner Loss, 1% DLR Loss

## *Hypothesis Testing Illuminates Effective Strategies*

- Change of Variables → 100% disabled, 0% enabled
- Optimizers → 50% Adam, 33% SGD, 16% MBS, 1% BWSGD
- **Random Restart → 61% enabled, 39% disabled**
- Saliency Maps → 70% no Saliency Map, 30% either DeepFool or JSMA Saliency Map
- Loss → 47% Identity Loss, 34% Cross Entropy, 18% Carlini Wagner Loss, 1% DLR Loss

## *Hypothesis Testing Illuminates Effective Strategies*

- Change of Variables → 100% disabled, 0% enabled
- Optimizers → 50% Adam, 33% SGD, 16% MBS, 1% BWSGD
- Random Restart → 61% enabled, 39% disabled
- Saliency Maps → 70% no Saliency Map, 30% either DeepFool or JSMA Saliency Map
- **Loss → 47% Identity Loss, 34% Cross Entropy, 18% Carlini Wagner Loss, 1% DLR Loss**

# GitHub Repo: https://github.com/sheatsley/attacks

# Thank you

🌐 https://hoak.me

✉️ bhoak@cs.wisc.edu

🐦 @blaine_hoak

🐙 blainehoak

WISCONSIN
UNIVERSITY OF WISCONSIN-MADISON

MADS&P