# How CLEVER is your neural network?
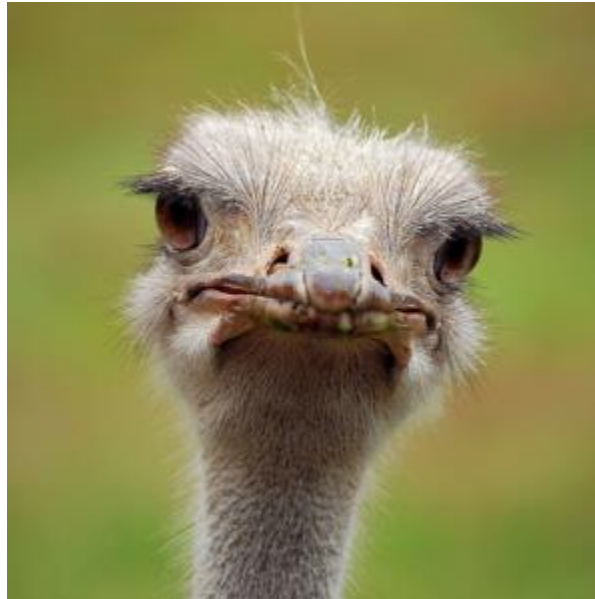## Robustness evaluation against adversarial examples

Pin-Yu Chen

IBM Research AI

O'Reilly AI Conference @ London 2018
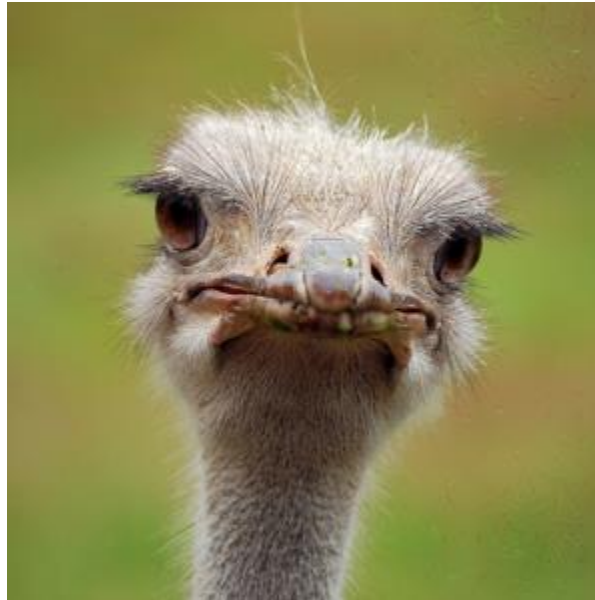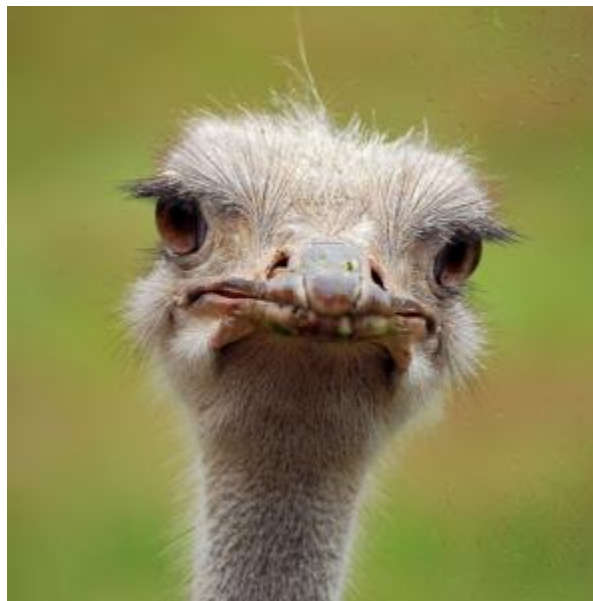
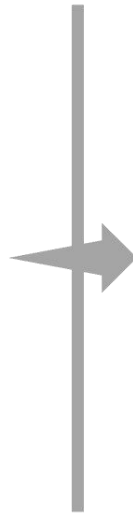# Label it!

# Label it! AI model says:

## ostrich

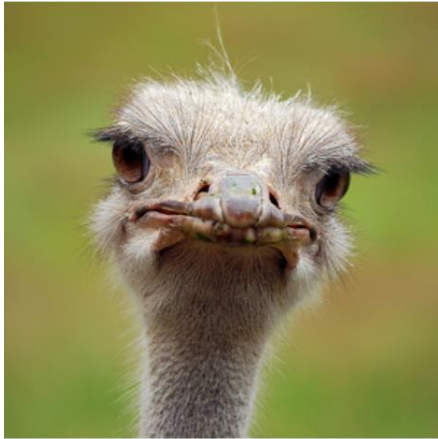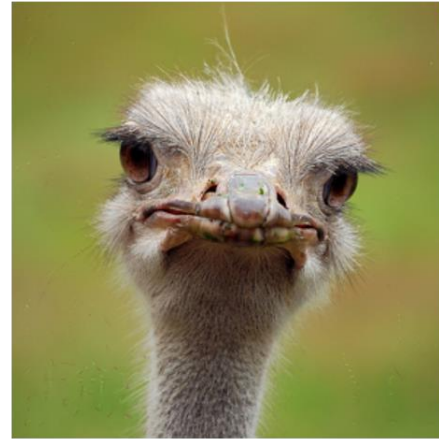# How about this one?
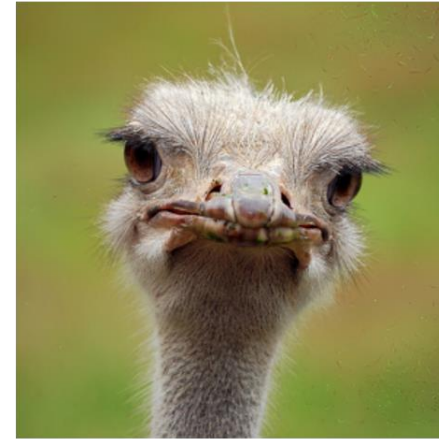
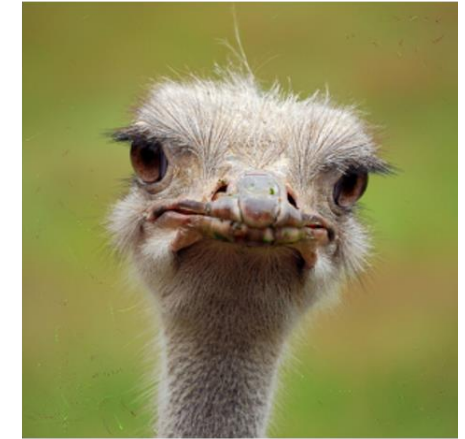# Surprisingly, AI model says:

## shoe shop

ostrich → safe | shoe shop | vacuum

# What is wrong with this AI model?

- This model is one of the BEST image classifier using neural networks

# Adversarial examples: the evil doublegangers



source: Google Images

# Why do adversarial examples matter?

- Adversarial attacks on an AI model deployed at test time (aka evasion attacks)

# Adversarial examples in different domains

- Images
- Videos
- Texts
- Speech/Audio
- Data analysis
- Electronic health records
- Malware
- Online social network
- and many others



**Original Top-3 inferred captions:**

1. A red stop sign sitting on the side of a road.
2. A stop sign on the corner of a street.
3. A red stop sign sitting on the side of a street.

**Adversarial Top-3 captions:**

1. A brown teddy bear laying on top of a bed.
2. A brown teddy bear sitting on top of a bed.
3. A large brown teddy bear laying on top of a bed.

"it was the best of times, it was the worst of times"

× 0.001

**AI model**

"it is a truth universally acknowledged that a single"

# Adversarial examples in image captioning



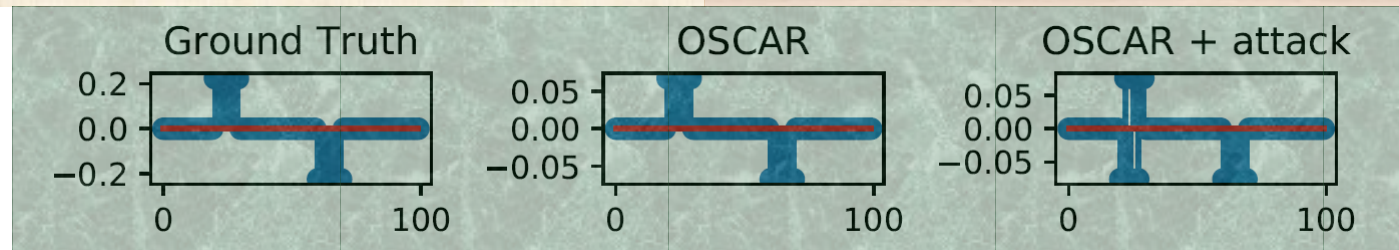Input: image

AI model

Output: caption

**Original Top-3 inferred captions:**
1. A red stop sign sitting on the side of a road.
2. A stop sign on the corner of a street.
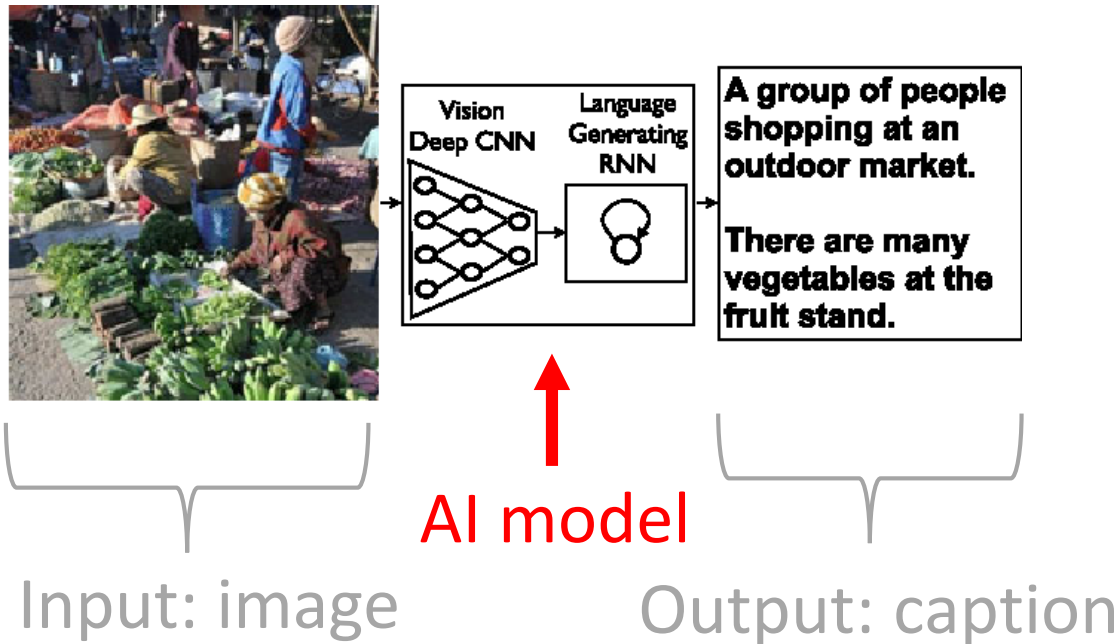3. A red stop sign sitting on the side of a street.

**Adversarial Top-3 captions:**
1. A brown teddy bear laying on top of a bed.
2. A brown teddy bear sitting on top of a bed.
3. A large brown teddy bear laying on top of a bed.

Show and Tell: Lessons Learned from the 2015 MSCOCO Image Captioning Challenge, Oriol Vinyals, AlexanderToshev, Samy Bengio, and Dumitru Erhan, T-PAMI 2017
Attacking Visual Language Grounding with Adversarial Examples: A Case Study on Neural Image Captioning, Hongge Chen*, Huan Zhang*, Pin-Yu Chen, Jinfeng Yi, and Cho-Jui Hsieh, ACL 2018
IBM Research AI

# Adversarial examples in speech recognition



AI model

without the dataset the article is useless

What did your hear?

IBM Research AI

# Adversarial examples in speech recognition



AI model

without the dataset the article is useless

What did your hear?

okay google browse to evil.com

Audio Adversarial Examples: Targeted Attacks on Speech-to-Text, **Nicholas Carlini** and David Wagner, Deep Learning and Security Workshop 2018
IBM Research AI

# Adversarial examples in data regression



Data

Model

Analysis

Factor identification



Ground Truth | OSCAR | OSCAR + attack

# Adversarial examples in physical world

- Real-time traffic sign detector
- 3D-printed adversarial turtle



- Adversarial eye glasses

# Adversarial examples in physical world (1)

- Real-time traffic sign detector



**Robust Physical-World Attacks on Deep Learning Visual Classification**

Kevin Eykholt[*,1], Ivan Evtimov[*,2], Earlence Fernandes[2], Bo Li[3],
Amir Rahmati[4], Chaowei Xiao[1], Atul Prakash[1], Tadayoshi Kohno[2], and Dawn Song[3]

# Adversarial examples in physical world (2)

- **3D-printed adversarial turtle**



classified as turtle   classified as rifle   classified as other

**Synthesizing Robust Adversarial Examples**

Anish Athalye [*1 2]   Logan Engstrom [*1 2]   Andrew Ilyas [*1 2]   Kevin Kwok [2]

# Adversarial examples in physical world (3)

- Adversarial eye glasses that fool face detector



**Accessorize to a Crime: Real and Stealthy Attacks on State-of-the-Art Face Recognition**

Mahmood Sharif
Carnegie Mellon University
Pittsburgh, PA, USA
mahmoods@cmu.edu

Sruti Bhagavatula
Carnegie Mellon University
Pittsburgh, PA, USA
srutib@cmu.edu

Lujo Bauer
Carnegie Mellon University
Pittsburgh, PA, USA
lbauer@cmu.edu

Michael K. Reiter
University of North Carolina
Chapel Hill, NC, USA
reiter@cs.unc.edu

- Adversarial sticker



**Adversarial Patch**

Tom B. Brown, Dandelion Mané, Aurko Roy, Martín Abadi, Justin Gilmer
{tombrown,dandelion,aurkor,abadi,gilmer}@google.com

IBM Research AI

# Adversarial examples in black-box models

- **White-box setting**: adversary knows everything about your model

- **Black-box setting:** craft adversarial examples with limited knowledge about the target model
    - ❖ Unknown training procedure/data/model
    - ❖ Unknown output classes
    - ❖ Unknown model confidence



Black-box attack via iterative model query (ZOO)

Targeted black-box attack on Google Cloud Vision

ZOO: Zeroth Order Optimization based Black-box Attacks to Deep Neural Networks without Training Substitute Models, P.-Y. Chen*, H. Zhang*, Y. Sharma, J. Yi, and C.-J. Hsieh, AI-Security 2017
Black-box Adversarial Attacks with Limited Queries and Information, Andrew Ilyas*, Logan Engstrom*, Anish Athalye*, and Jessy Lin*, ICML 2018
Source: https://www.labsix.org/partial-information-adversarial-examples/
IBM Research AI

# Growing concerns about safety-critical settings with AI

Autonomous cars that deploy AI model for traffic signs recognition

IBM Research AI

# But with adversarial examples...

IBM Research AI

# Where do adversarial examples come from?

- What is the common theme of adversarial examples in different domains?

# Neural Networks: The Engine for Deep Learning

- Applications of neural networks
- ❑ Image processing and understanding
- ❑ Object detection/classification
- ❑ Chatbot, Q&A
- ❑ Machine translation
- ❑ Speech recognition
- ❑ Game playing
- ❑ Robotics
- ❑ Bioinformatics
- ❑ Creativity
- ❑ Drug discovery
- ❑ Reasoning
- ❑ And still a long list…

neural network

outcome (prediction)

2% (traffic light)

**90% (French bulldog)**

3% (basketball)

5% (bagel)

input task

trainable neurons;
usually large and deep

Source: Paishun Ting

IBM Research AI

# The ImageNet Accuracy Revolution and Arms Race



Source: http://image-net.org/challenges/talks_2017/imagenet_ilsvrc2017_v1.0.pdf
Source: https://qz.com/1034972/the-data-that-changed-the-direction-of-ai-research-and-possibly-the-world/

IBM Research AI

# Accuracy ≠ Adversarial Robustness

- Solely pursuing for high-accuracy AI model may get us in trouble…



Tradeoff between Accuracy and $\ell_\infty$ CLEVER Score

Robustness

Accuracy

Our benchmark on 18 ImageNet models reveals a tradeoff in accuracy and robustness

# How can we <u>measure</u> and improve adversarial robustness of my AI/ML model?

An explanation of origins of adversarial examples

The CLRVER score for robustness evaluation

# Learning to classify is all about drawing a line



Classified as 🔵

Classified as ✕

Labeled datasets

——— Decision boundary w/ 100% accuracy

- - - Decision boundary w/ <100% accuracy

IBM Research AI

# Connecting adversarial examples to model robustness

Classified as ⬤

Classified as ✖

Ostrich   shoe shop   vacuum

$x_0$     $x_a$     $x_{a'}$

$\Delta$ = Minimum distortion

$x_a$

adversarial example

$\Delta$

$x_0$

Certified robustness within the grey region

Ostrich

adversarial example

$x_{a'}$

Decision boundary 3

Decision boundary 2

Decision boundary 1

$L_p$ space

bagel   +   =   grand piano

⬤ ✖➜✖ ✖

- Robustness evaluation: how close a refence input is to the (closest) decision boundary

IBM Research AI

# Robustness evaluation is NOT easy

- We still don't fully understand how neural nets learn to predict

❑ calling for interpretable AI

- Training data could be noisy and biased

❑ calling for robust and fair AI

- Neural network architecture could be redundant and leading to vulnerable spots

❑ calling for efficient and secure AI model

- Need for human-like machine perception and understanding

❑ calling for bio-inspired AI model

- Attacks can also benefit and improve upon the progress in AI

❑ calling for attack-independent evaluation

Labeled datasets

**Adversarial Examples Are Not Easily Detected: Bypassing Ten Detection Methods**

Nicholas Carlini     David Wagner

**Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples**

**Anish Athalye** [*1]   **Nicholas Carlini** [*2]   **David Wagner** [2]

# How do we evaluate adversarial robustness?

- ## Game-based approach

❑Specify a set of players (attacks and defenses)

❑Benchmark the performance against each attacker-defender pair

o The metric/rank could be exploited; ⚠ No guarantee on unseen threats and future attacks



Research Prediction Competition

## NIPS 2017: Defense Against Adversarial Attack
Create an image classifier that is robust to adversarial attacks

Google Brain · 107 teams · 3 months ago

IBM Research AI

- ## Verification-based approach

❑Attack-independent: does not use attacks for evaluation

❑Can provide a robustness certificate for safety-critical or reliability-sensitive applications: e.g., no attacks can alter the decision of the AI model if the attack strength is limited

⚠ Optimal verification is provably difficult for large neural nets – computationally impractical

- Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks, Guy Katz, Clark Barrett, David Dill, Kyle Julian, Mykel Kochenderfer, CAV 2017
- Efficient Neural Network Robustness Certification with General Activation Functions, Huan Zhang*, Tsui-Wei Weng*, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel, NIPS 2018

# CLEVER: a tale of two approaches

- An <u>attack-independent</u>, <u>model-agnostic</u> robustness metric that is <u>efficient to compute</u>

- Derived from theoretical robustness analysis for verification of neural networks: <u>C</u>ross <u>L</u>ipschitz <u>E</u>xtreme <u>V</u>alue for n<u>E</u>twork <u>R</u>obustness

- Use of extreme value theory for efficient estimation of minimum distortion

- Scalable to large neural networks

- Open-source codes: https://github.com/IBM/CLEVER-Robustness-Score



Ostrich shoe shop vacuum

$x_0$ $x_a$ $x_{a'}$

adversarial example $x_a$

$\Delta$ = Minimum distortion $\approx$ **CLEVER score**

Certified robustness within the grey region

$x_0$

Ostrich

adversarial example $x_{a'}$

Decision boundary 3

Decision boundary 2

Decision boundary 1

$L_p$ space

$g(x)$

Slope $= -L_q$

$g(x_0)$

Slope $= L_q$

$x_0$ $x_0 + \delta$ $x$

$g(x_0) + L_q\|\delta\|_p$

Bounded by Local Lipschitz constant $L_q$

$g(x_0 + \delta)$
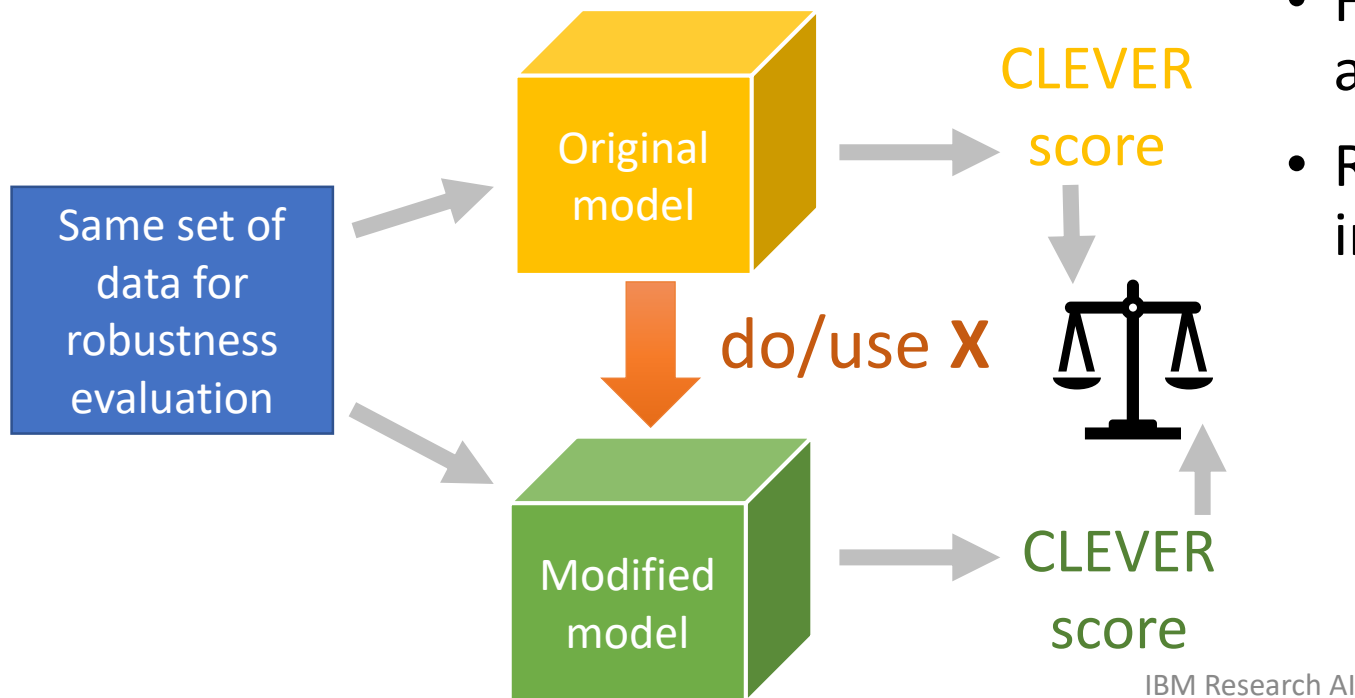
$g(x_0) - L_q\|\delta\|_p$

**input-output perturbation analysis of neural net**

Evaluating the Robustness of Neural Networks: An Extreme Value Theory Approach, Tsui-Wei Weng*, Huan Zhang*, Pin-Yu Chen, Jinfeng Yi, Dong Su, Yupeng Guo, Cho-Jui Hsieh, and Luca Daniel, ICLR 2018
On Extensions of CLEVER: a Neural Network Robustness Evaluation Algorithm, Tsui-Wei Weng*, Huan Zhang*, Pin-Yu Chen, Aurelie Lozano, Cho-Jui Hsieh, and Luca Daniel, GlobalSIP 2018

IBM Research AI
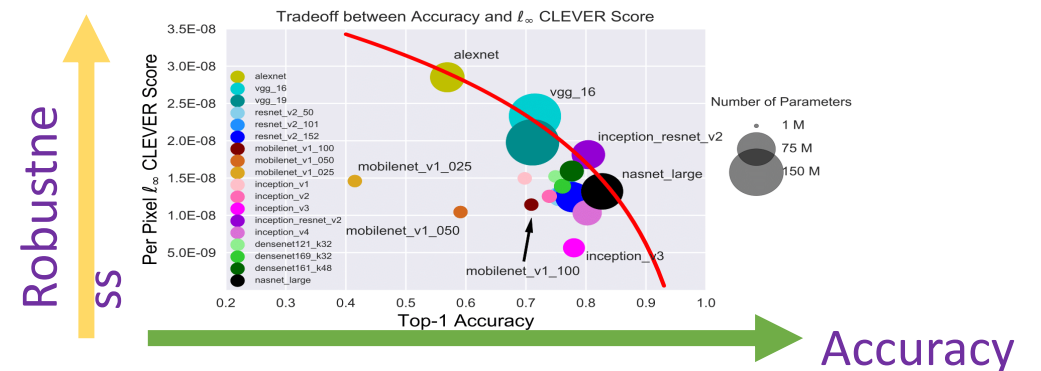
# How do we use CLEVER?

**Before-After robustness comparison**

- Will my model become more robust if I do/use **X**?



**Other use cases**

- Characterize the behaviors and properties of adversarial examples

- Hyperparameter selection for adversarial attacks and defenses
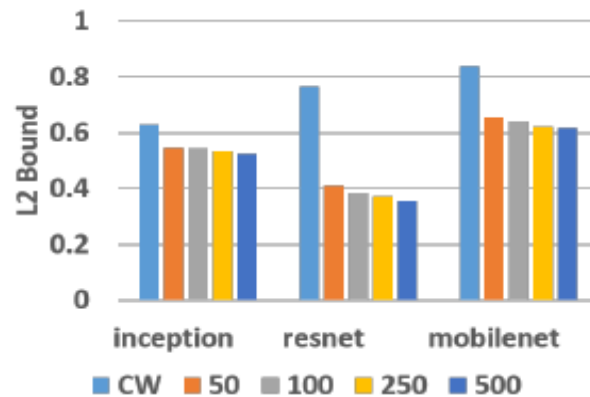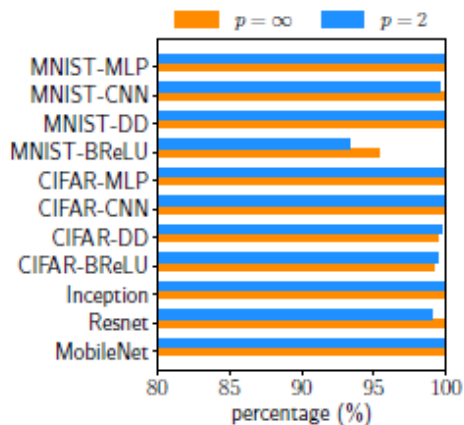
- Reward-driven model robustness improvement



IBM Research AI

# Examples of CLEVER

- CLEVER enables robustness comparison between <u>different</u>

❑ Threat models

❑ Datasets

❑ Neural network architectures

❑ Defense mechanisms

# Where to Find CLEVER? It's ART



## Adversarial Robustness Toolbox (ART)

External: https://github.com/IBM/adversarial-robustness-toolbox

- Python library, 7K lines of code
- State-of-the-art attacks, defences and robustness metrics

```
from keras.datasets import mnist
from keras.models import load_model

from art.attacks import CarliniL2Attack
from art.classifier import KerasClassifier
from art.metrics import loss_sensitivity

# Load data
(_, _), (x_test, y_test) = mnist.load_data()

# Load model and build classifier
model = load_model('my_favorite_keras_model.h5')
classifier = KerasClassifier((0, 1), model)

# Perform attack
attack = CarliniL2Attack(classifier)
adv_x_test = attack.generate(x_test)

# Compute metrics on model robustness
print(loss_sensitivity(classifier, x_test))
```

Load ART modules
Load classifier model (Keras, TF, PyTorch etc)
Perform attack
Evaluate robustness

Open-source release @ RSA 2018:

~ 3.5K+ views of IBM blogs
~ 100+ news outlets covering release
~ 1.3M+ Social Media potential impressions
~ 5K+ views of GitHub repo

**siliconANGLE**
Attackers can fool AI programs. Here's how developers can fight back

**ZDNet**
IBM launches open-source library for securing AI systems
The framework-agnostic software library contains attacks, defenses, and benchmarks for securing artificial intelligence systems.

IBM ENTWICKELT WERKZEUGE GEGEN HACKERANGRIFFE DURCH "BÖSE" KI

IBM、AIシステムを保護するオープンソースライブラリ「Adversarial Robustness Toolbox」

Выпущена Adversarial Robustness Toolbox, открытая библиотека от IBM для защиты ИИ

Adversarial Robustness Toolbox : IBM propose une boite à outils open source pour sécuriser l'intelligence artificielle

IBM Adversarial Robustness Toolbox beschermt tegen kwaadaardige AI

| Evasion attacks | Evasion defenses | Poisoning detection | Robustness metrics |
|---|---|---|---|
| • FGSM | • Feature squeezing | • Detection based on clustering activations | • CLEVER |
| • JSMA | • Spatial smoothing | | • Empirical robustness |

Also available at https://github.com/IBM/CLEVER-Robustness-Score

IBM Research AI

# Take-aways

- Adversarial robustness is a new AI standard

❑Robustness does not come for free: adversarial examples exist in digital space, physical world, and different domains

❑High accuracy ≠ Good robustness

❑Arms race: adversary-aware AI v.s. AI for adversary

- How to evaluate the robustness of my AI model?

❑CLEVER: an attack-independent robustness score

❑Robustness comparison in before-after setting

❑Where to find CLEVER? It's ART!

Human

AI

Data

Robustness

# Beyond Robustness: Trusted AI

**Trusted AI**

IBM Research is building and enabling AI solutions people can trust

As AI advances, and humans and AI systems increasingly work together, it is essential that we trust the output of these systems to inform our decisions. Alongside policy considerations and business efforts, science has a central role to play: developing and applying tools to wire AI systems for trust. IBM Research's comprehensive strategy addresses multiple dimensions of trust to enable AI solutions that inspire confidence.

## Robustness

We are working to ensure the security and reliability of AI systems by exposing and fixing their vulnerabilities: identifying new attacks and defense, designing new adversarial training methods to strengthen against attack, and developing new metric to evaluate robustness.

View publications

## Fairness

To encourage the adoption of AI, we must ensure it does not take on and amplify our biases. We are creating methodologies to detect and mitigate bias through the life cycle of AI applications.

View publications

## Explainability

Knowing how an AI system arrives at an outcome is key to trust, particularly for enterprise AI. To improve transparency, we are researching local and global interpretability of models and their output, training for interpretable models and visualization of information flow within models, and teaching explanations.

View publications

## Lineage

Lineage services can infuse trust in AI systems by ensuring all their components and events are trackable. We are developing services like instrumentation and event generation, scalable event ingestion and management, and efficient lineage query services to manage the complete lifecycle of AI systems.

View publications

IBM Research AI

# Acknowledgement

- Collaborators: Tsui-Wei Weng(MIT), Luca Daniel(MIT), Honnge Chen(MIT) Huan Zhang(UCLA), Cho-Jui Hsieh(UCLA), Jinfeng Yi(JD AI), Yupeng Gao(IBM), Bhanukiran Vinzamuri(IBM), Sijia Liu(IBM), Yash Sharma, Su Dong, Chun-Chen Tu(UMich), Paishun Ting(Umich)

- MIT-IBM Watson AI Lab: David Cox, Lisa Amini

- IBM Research AI – Learning Group: Payel Das, Saska Mojsilovic

- IBM AI-Security Group: Ian Molloy, Mathieu Sinn, and their teams

- IBM Big Check Demo: Casey Dugan and her team

- IBM DLaaS Group: Evelyn Duesterwald and her team

❑Personal Website: www.pinyuchen.com

❑Twitter: pinyuchen.tw