

Research Review 2021

Train, but Verify: Towards practical AI robustness

September 2021

Nathan VanHoudnos

Copyright 2021 Carnegie Mellon University.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

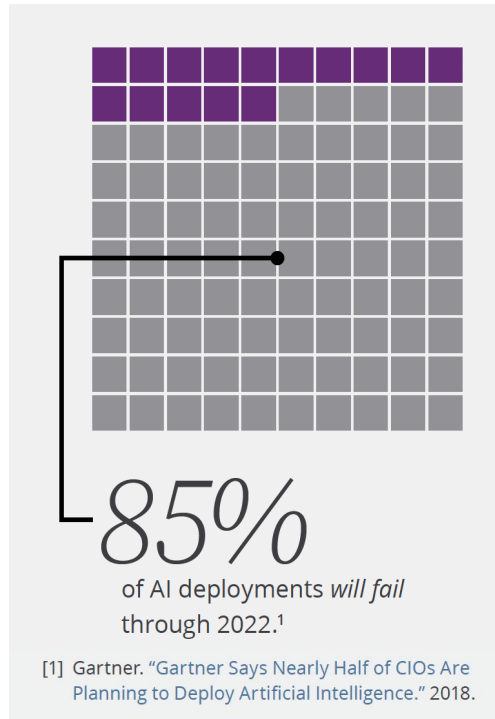
NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other use. Requests for permission should be directed to the Software Engineering Institute at permission@sei.cmu.edu.

DM21-0869

What is AI Engineering?



AI Engineering:

- field of research and practice
- integrates software engineering, systems, CS, and human-centered design
- builds AI responsive to human needs and mission outcomes.

Human-centered



Works with and for people

Scalable



Size, speed, & complexity of mission needs

Robust and Secure

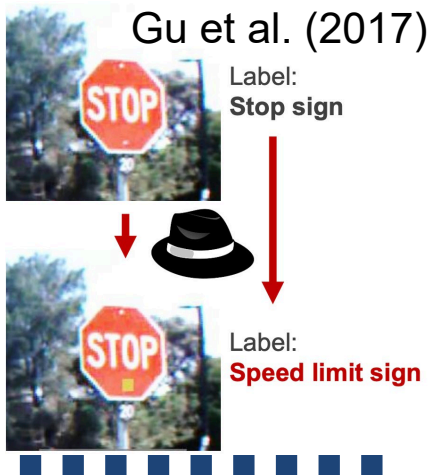


Reliable when under uncertainty or **threat**

<https://www.sei.cmu.edu/our-work/artificial-intelligence-engineering/>

Beieler (2018): An adversary can make an ML component...

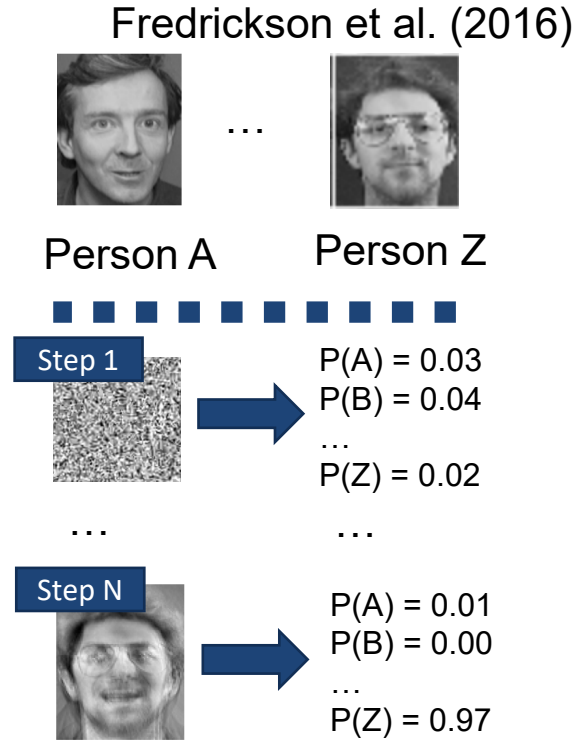
Learn the Wrong Thing



Do the Wrong Thing



Reveal the Wrong Thing



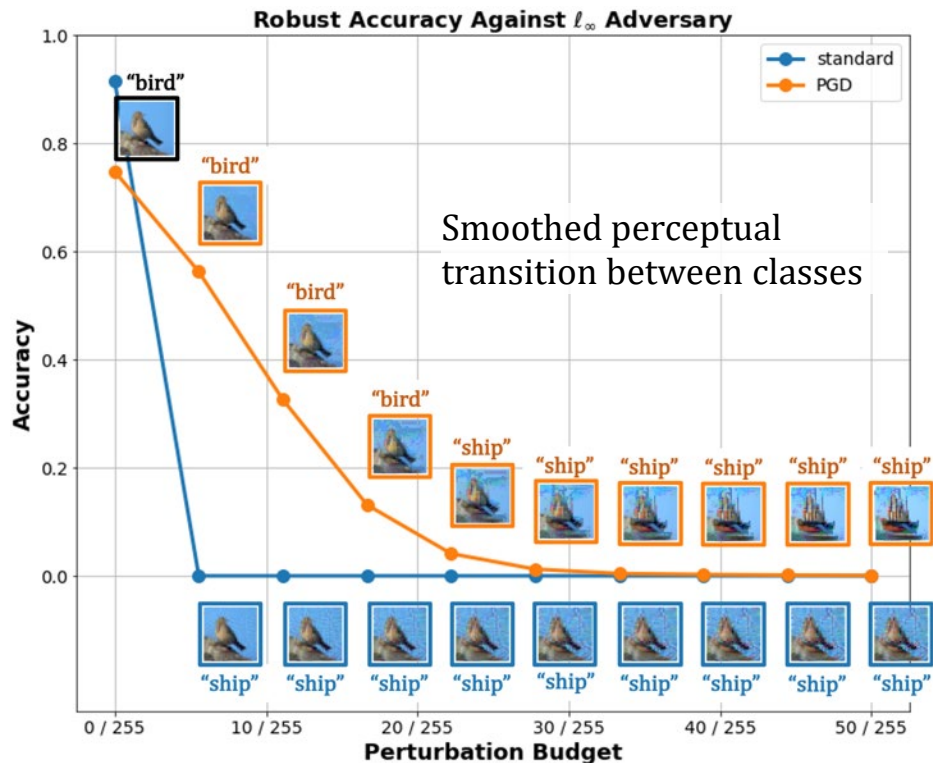
Train, but Verify

| Train \ Verify | Verify “learn” policy | Verify “do” policy | Verify “reveal” policy |
|----------------------------------|----------------------------|--------------------|------------------------|
| Train to enforce “learn” policy | IARPA TrojAI DARPA GARD | | |
| Train to enforce “do” policy | | DARPA GARD | ? |
| Train to enforce “reveal” policy | | | NGA GURU |

Problem:

- AI promises capability for the DoD, but today is untrustworthy.
- Most defensive work focuses on one security policy, but the DoD has wider concerns.
 - What if a system makes high stakes decisions (do policy) and is trained on sensitive data (reveal policy)?

Defenses for do policies reveal information about the data



[Helland et al. 2020 – On The Human Recognizability Phenomenon of Adversarially Trained Deep Image Classifiers]

Model: [Engstrom 2019 – Robustness (Python Library)]

Examples extracted from defended model

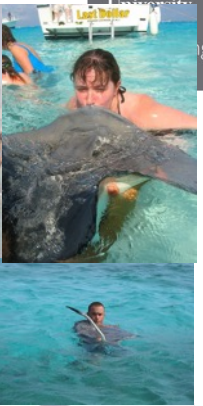
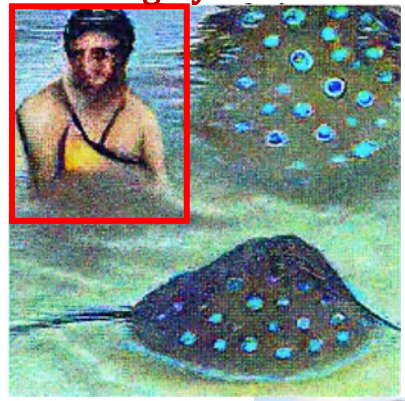


Examples extracted from undefended model



stingray 0.9920

Stingrays and people? Sure thing.



Images from [Deng et al. 2009 - ImageNet A Large-Scale Hierarchical Image Database]

Purple cauliflower? You bet.



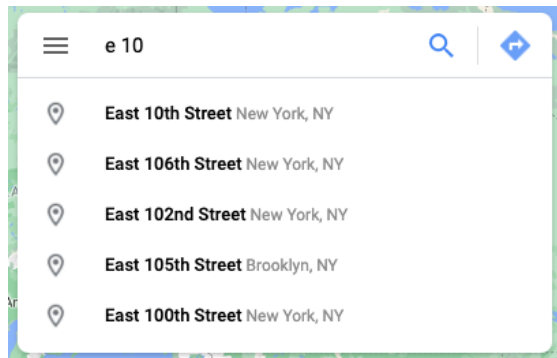
Images from [Deng et al. 2009 - ImageNet A Large-Scale Hierarchical Image Database]

Let's look at the **street sign** class

E 10?



Google Maps



Google Maps Street View



Does my dataset contain photographs in **New York**?

Yes, quite a few.



Images from [*Deng et al. 2009 - ImageNet A Large-Scale Hierarchical Image Database*]

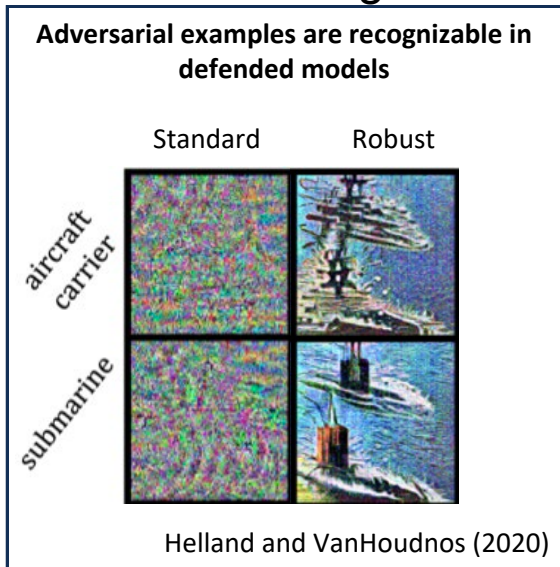
Train, but Verify: Towards practical AI robustness
© 2021 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution

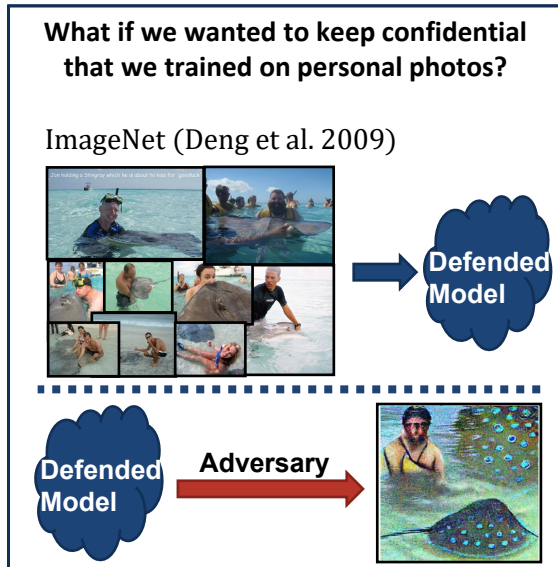
Images from [*Deng et al. 2009 - ImageNet A Large-Scale Hierarchical Image Database*]

Secure AI Engineering for DoD defends from multiple attacks

State-of-the-art methods to enforce “doing the right thing” can leak information about the training data.



A defended model may (unintentionally) reveal critical information about the data.



Organizations with high stakes systems trained on sensitive data need new defense methods.

The Train, but Verify grid

| Verify \ Train | Learned correctly | Did correctly | No Revealed secrets |
|-----------------------|-------------------|---|---------------------|
| To Learn correctly | | | |
| To Do correctly | | <p>Goal Satisfy both Do and Reveal</p> | |
| To not Reveal secrets | | | |

Train, but Verify: FY2021 Goals & Milestones

DoD needs secure AI across multiple policies.

| Verify Train | Learned correctly | Did correctly | No Revealed secrets |
|-----------------------|-------------------|------------------------------------|---------------------|
| To Learn correctly | | | |
| To Do correctly | | Goal Satisfy both Do and Reveal | |
| To not Reveal secrets | | | |

Impact: Allow for use of sensitive data in high stakes environments.

Quantify attacks to reveal policies.

- [Under Review]: Property Inference Attacks in Robust and Private Models

Develop new methods for do defenses and do attacks.

- [Under Review]: Self-Repairing Neural Networks
- [Under Review]: Constrained Gradient Descent: Strong Attacks Against Neural Networks

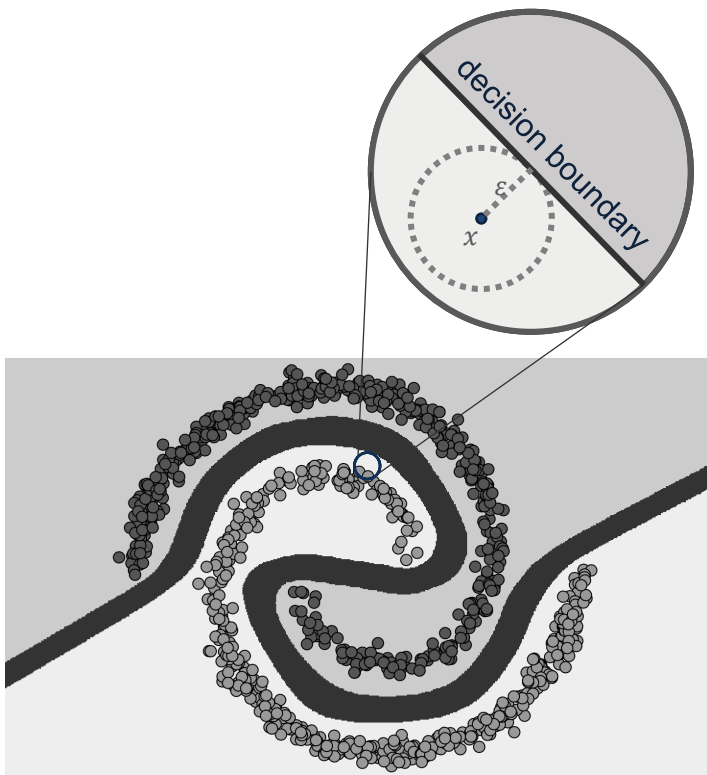
Develop new methods to verify do policies

- **ICML '21: Globally-Robust Neural Networks**
- [Under Review]: Relaxing Local Robustness
- ICLR '21: Fast Geometric Projections for Local Robustness Certification

Release AI Engineering tools

- **Juneberry 0.5 released to GitHub**

ICML '21: Globally-Robust Neural Networks (GloRo Nets)



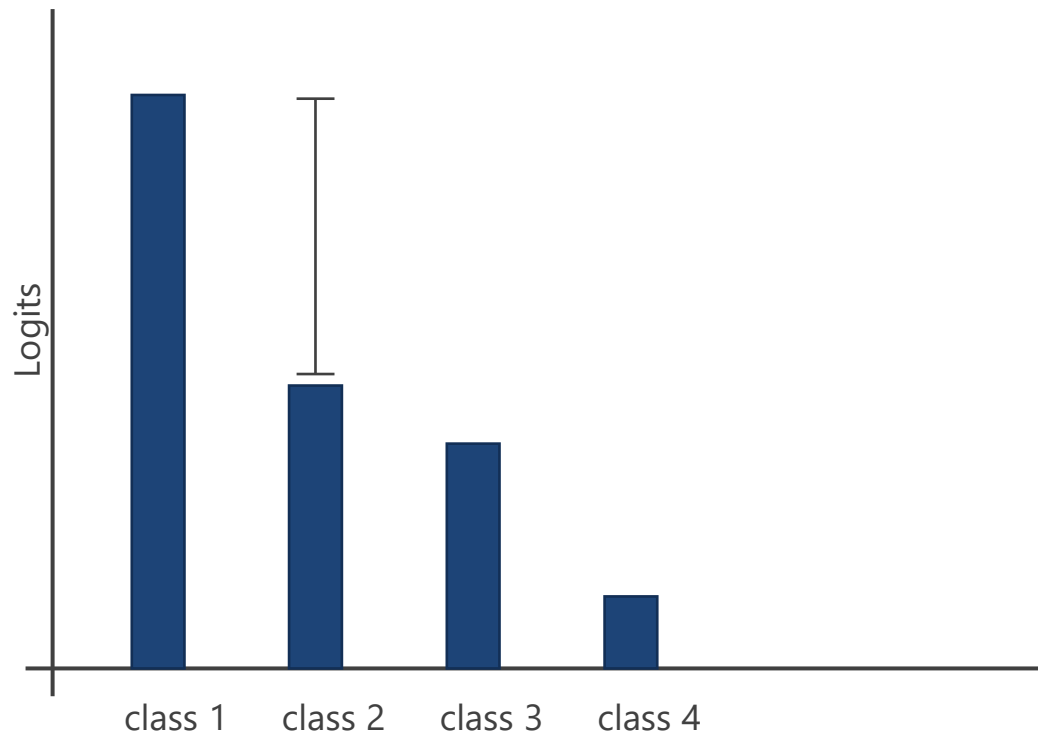
A model F satisfies **local robustness** with robustness radius ϵ on a point x if

$$\forall x': \|x - x'\|_p \leq \epsilon \implies F(x) = F(x')$$

A model F satisfies **global robustness** with robustness radius ϵ if $\forall x$

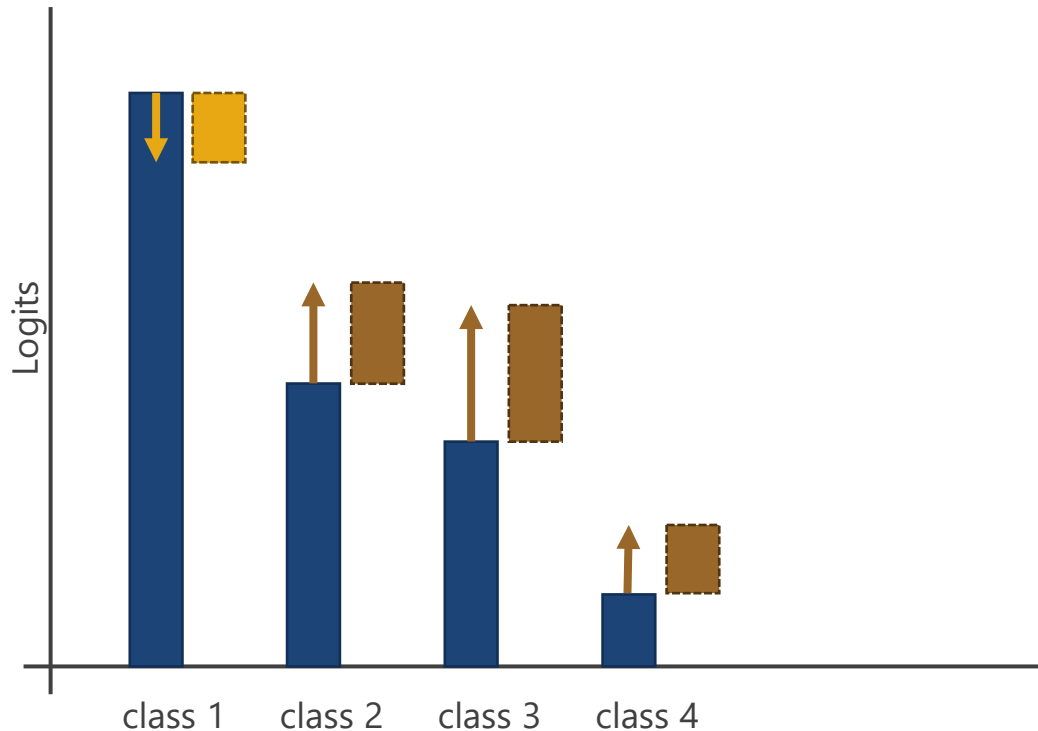
- F is $(\epsilon/2)$ -locally robust at x or
- $F(x) = \perp$ i.e. “No comment”

GloRo Nets: Intuition



If this margin is sufficiently large, a small change to the input will not allow class 2 to surpass class 1

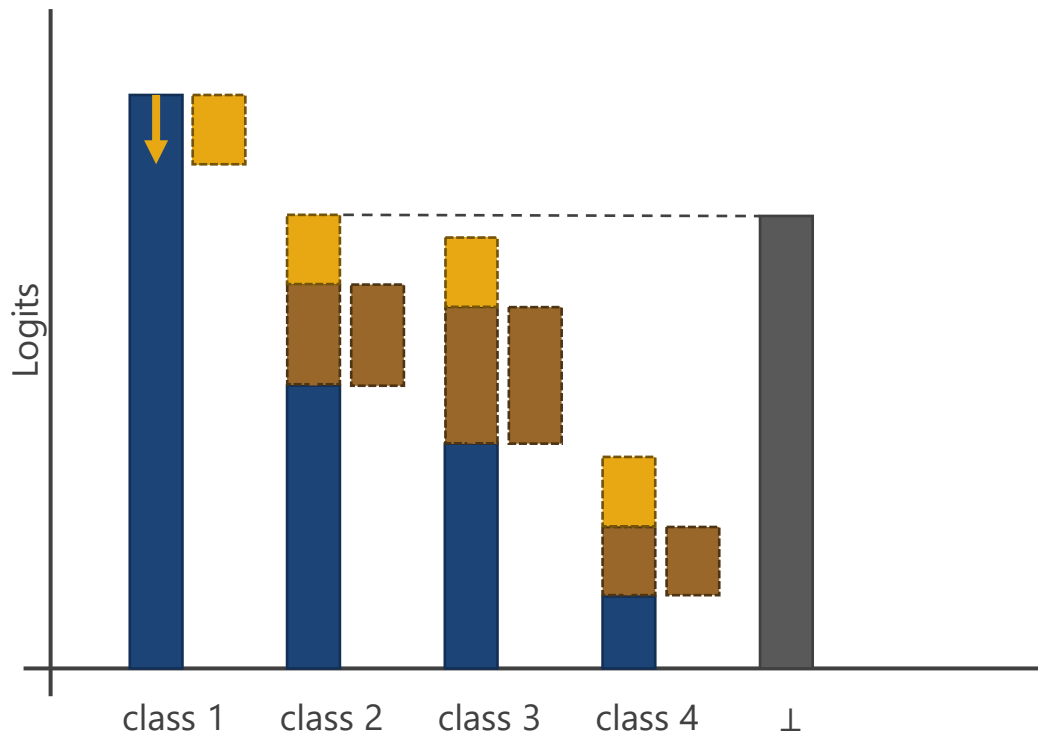
GloRo Nets: Intuition



If this margin is sufficiently large, a small change to the input will not allow class 2 to surpass class 1

The *Lipschitz constant* tells us how much each class can change with a small change to the input in the worst case

GloRo Nets: Intuition



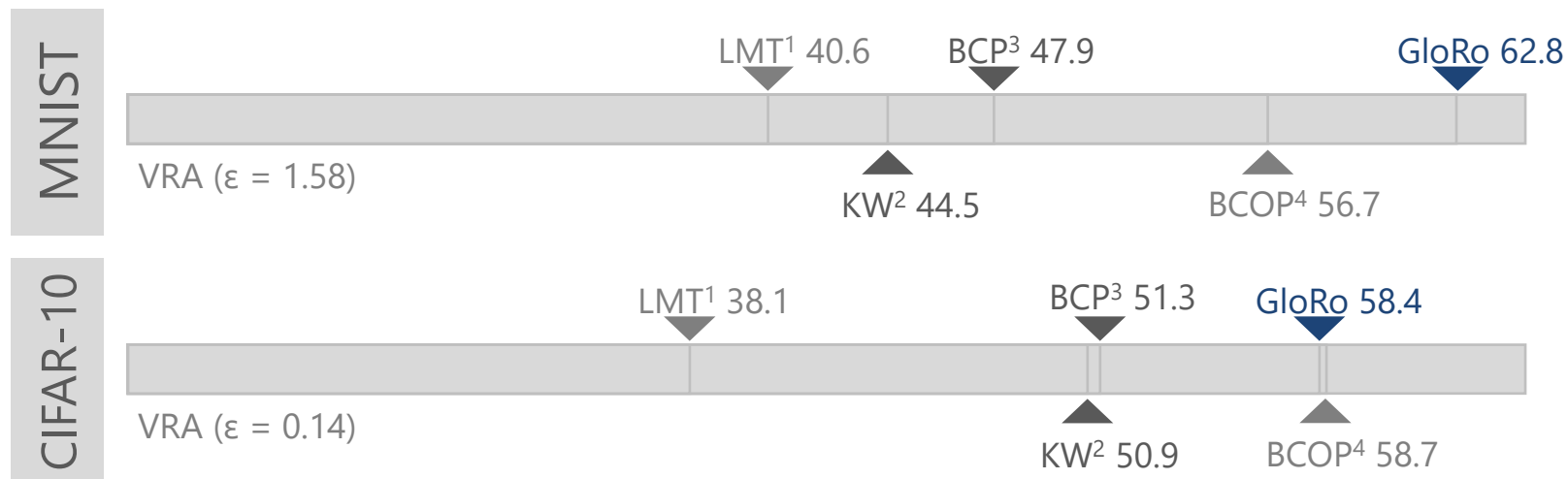
If this margin is sufficiently large, a small change to the input will not allow class 2 to surpass class 1

The *Lipschitz constant* tells us how much each class can change with a small change to the input in the worst case

We add a new class, \perp , which reflects the highest score an adversary can get relative to the top class

GloRo Nets: Results

GloRo Nets match or exceed VRA of previous state-of-the-art deterministic certification methods



¹Tsuzuku et al., 2018; ²Wong & Kolter, 2018; ³Lee et al., 2020; ⁴Li et al., 2019

GloRo Nets: Performance

GloRo Net certification and training is significantly more time and memory efficient than other methods, and more scalable than any other deterministic method

| | method | time to certify test set (s) | memory per instance (MB) |
|----------|------------------|------------------------------|--------------------------|
| CIFAR-10 | GloRo | 0.4 | 1.8 |
| | KW ¹ | 2,500.0 | 1,400.0 |
| | BCP ² | 5.8 | 19.1 |
| | RS ³ | 36,800.0 | 19.8 |

¹Tsuzuku et al., 2018; ²Wong & Kolter, 2018; ³Lee et al., 2020; ⁴Li et al., 2019

Train, but Verify: FY2021 Goals & Milestones

DoD needs secure AI across multiple policies.

| Verify Train | Learned correctly | Did correctly | No Revealed secrets |
|-----------------------|-------------------|------------------------------------|---------------------|
| To Learn correctly | | | |
| To Do correctly | | Goal Satisfy both Do and Reveal | |
| To not Reveal secrets | | | |

Impact: Allow for use of sensitive data in high stakes environments.

Quantify attacks to reveal policies.

- [Under Review]: Property Inference Attacks in Robust and Private Models

Develop new methods for do defenses and do attacks.

- [Under Review]: Self-Repairing Neural Networks
- [Under Review]: Constrained Gradient Descent: Strong Attacks Against Neural Networks

Develop new methods to verify do policies

- **ICML '21: Globally-Robust Neural Networks**
- [Under Review]: Relaxing Local Robustness
- ICLR '21: Fast Geometric Projections for Local Robustness Certification

Release AI Engineering tools

- **Juneberry 0.5 released to GitHub**

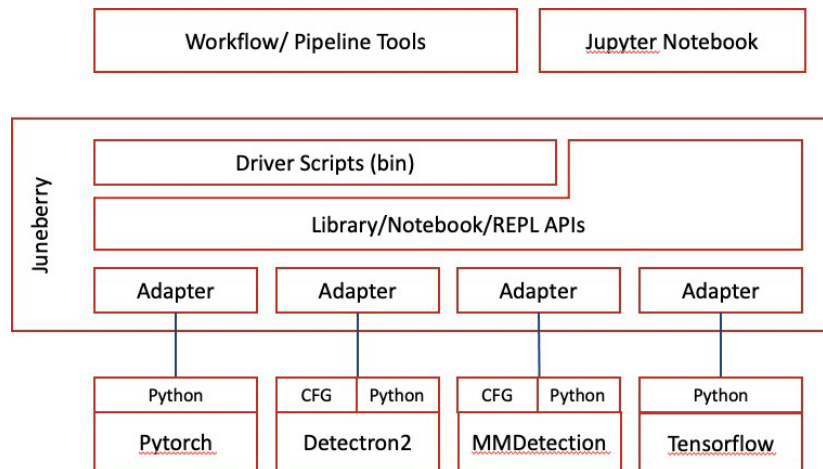
Juneberry: A tool for Robust & Secure AI Engineering



Juneberry

<https://github.com/cmu-sei/Juneberry>

- provides a framework for reproducible ML research
- improves the experience of machine learning experimentation
- is extensible for a variety of ML tasks
 - classification (v. 0.2)
 - object detection (v. 0.4)
 - differential privacy (v. 0.4)
 - certified robustness (v. 0.5)



Juneberry: A tool for Robust & Secure AI Engineering



Juneberry

<https://github.com/cmu-sei/Juneberry>

Vignette: Replicating a Classic Machine Learning Result with Juneberry

- Load data (data config)
- Wrap a model (model factory)
- Replicate a training strategy (model config)
- Train a model (jb_train)
- Evaluate a model (jb_evaluate)
- Replicate a results table (experiment outline)
- Execute an experiment (jb_run_experiment)
- Compare results with the published results

Additional vignettes (certified robustness, object detection, ...) coming soon!

Train, but Verify: Towards Practical AI Robustness

DoD needs secure AI across multiple policies.

| Verify Train | Learned correctly | Did correctly | No Revealed secrets |
|-----------------------|-------------------|------------------------------------|---------------------|
| To Learn correctly | | | |
| To Do correctly | | Goal Satisfy both Do and Reveal | |
| To not Reveal secrets | | | |

Impact: Allow for use of sensitive data in high stakes environments.

FY 2021:

- Quantify attacks to reveal policies.
- Develop new methods for do defenses and do attacks.
- Develop new methods to verify do policies
- Release AI Engineering tools

FY 2022:

- Develop training methods for **do & reveal** that either
 - enforce both
 - trade between them
- Release AI Engineering tools

Train, but Verify: Towards Practical AI Robustness



Nathan VanHoudnos

Senior ML Research Scientist
PI



Andrew Mellinger

Senior Software Developer
Co-PI



Lujo Bauer

Professor
Co-PI



Matt Fredrickson

Associate Professor
Co-PI



Bryan Parno

Associate Professor
Co-PI



Rob Beveridge

Technical Manager



Matthew Churilla

Software Security
Engineer



Bill Shaw

Senior Engineer



Aymeric Fromherz

Post-Doc @ Inria Paris
(formerly PhD Student)



Zifan Wang

PhD Student



Tina Sciuillo-Schade

Research Project
Manager



Violet Turri

Assistant Software
Developer



Annika Horgan

Associate Software
Developer



Weiran Lin

PhD Student



Kevin Li

Masters Student



Jon Helland

Associate ML
Researcher



Nick Winski

Software Developer.



Klas Leino

PhD Student



Clement Fung

PhD Student