



truera

Tutorial | SIGKDD 2021

# Machine Learning *Explainability* and *Robustness* → Connected at the Hip

[Anupam Datta](#)  
[Matt Fredrikson](#)  
[Klas Leino](#)  
[Kaiji Lu](#)  
[Shayak Sen](#)  
[Zifan Wang](#)



*Anupam Datta*



*Matt Fredrikson*



*Caleb Lu*



*Klas Leino*



*Shayak Sen*



*Zifan Wang*

# Machine Learning Systems are Ubiquitous



# Google

April 3, 2013, Vol 309, No. 13 >

< Previous Article    Next Article >

Viewpoint | April 3, 2013

## The Inevitable Application of Big Data to Health Care

Travis B. Murdoch, MD, MSc; Allan S. Detsky, MD, PhD

[+] Author Affiliations



Big Data in Government, Defense and Homeland Security 2015 - 2020

[Share](#) [G+](#) [Twitter](#) [LinkedIn](#) [Pinterest](#) [Email](#)

NEW YORK, May 12, 2015 /PRNewswire/

## How Big Data Could Replace Your Credit Score

Credit scores are useful in determining who gets loans, but they're far from perfect. AvantCredit determines loan-worthiness based on all sorts of factors, including your use of social media and prepaid cell phones.

## Big Data in Education

Learn how and when to use key methods for educational data mining and learning analytics on large-scale educational data.

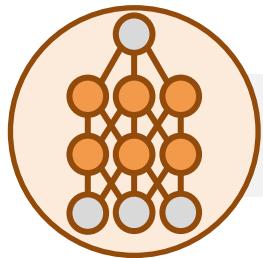
TEACHERS COLLEGE  
COLUMBIA UNIVERSITY

amazon

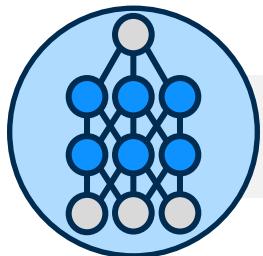
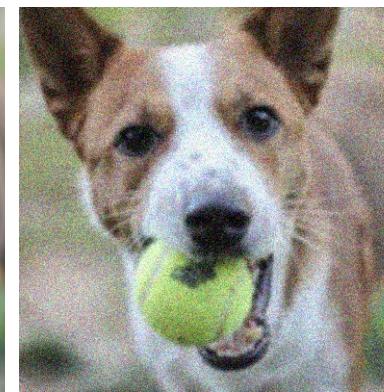
facebook

bing

# Motivating example



model 1

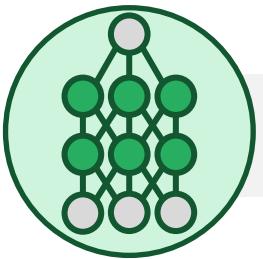


model 2

dog  
dog  
dog

sports car  
dog  
dog

cat  
dog  
cat



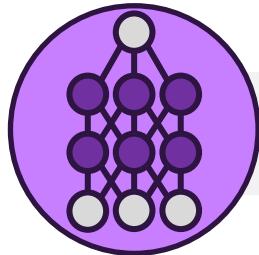
model 3

dog

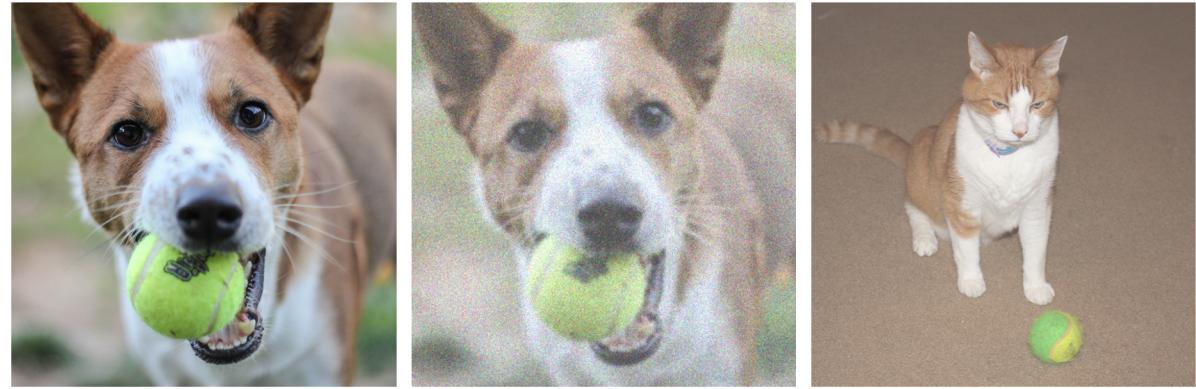
dog

cat

# Can we trust a new model based on its predictions?



model 4



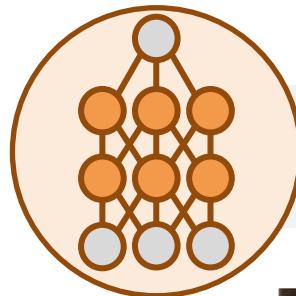
prediction agrees  
with all previous models  
on our validation point;  
how do we know which  
model's behavior model 4  
will best resemble?

dog
dog
dog
dog

sports car
dog
dog
?

cat
dog
cat
?

# Problem: lack of robustness

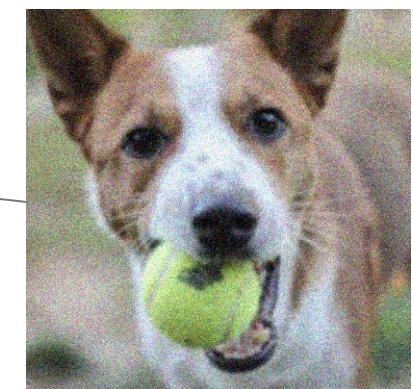
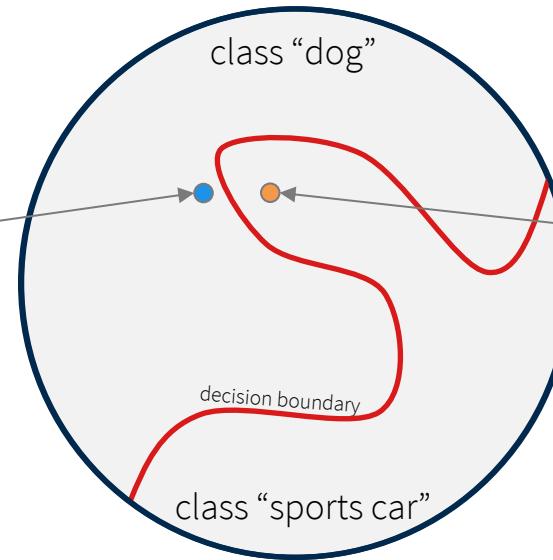


model 1

model 1 is not *robust*

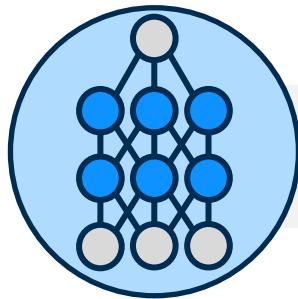


dog



sports car

# Problem: lack of conceptual soundness

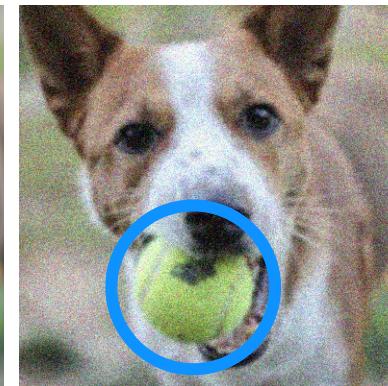


model 2

model 2 is not *conceptually sound*



dog

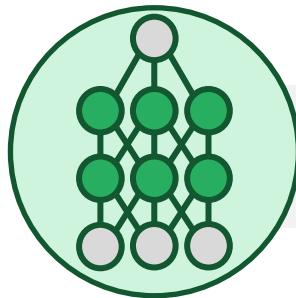


dog



dog

# Ideal: robustness and conceptual soundness



model 3

model 3 is *robust* and *conceptually sound*



dog

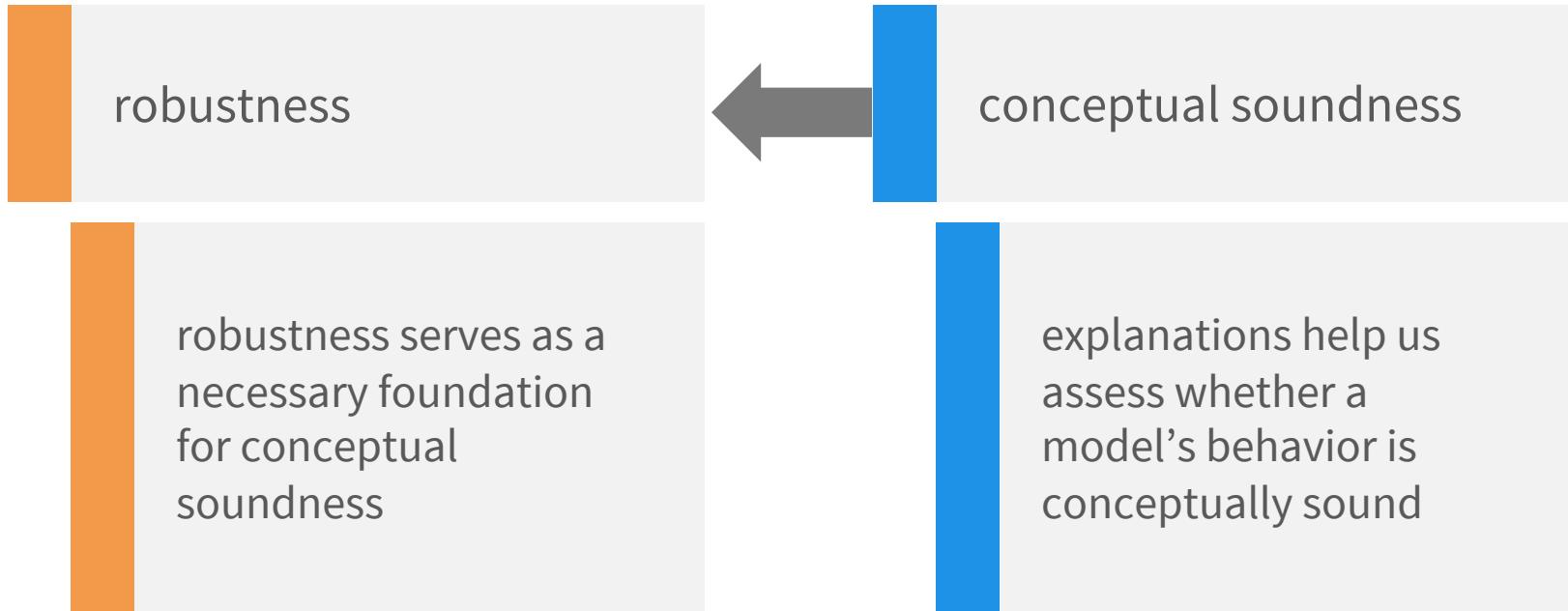


dog

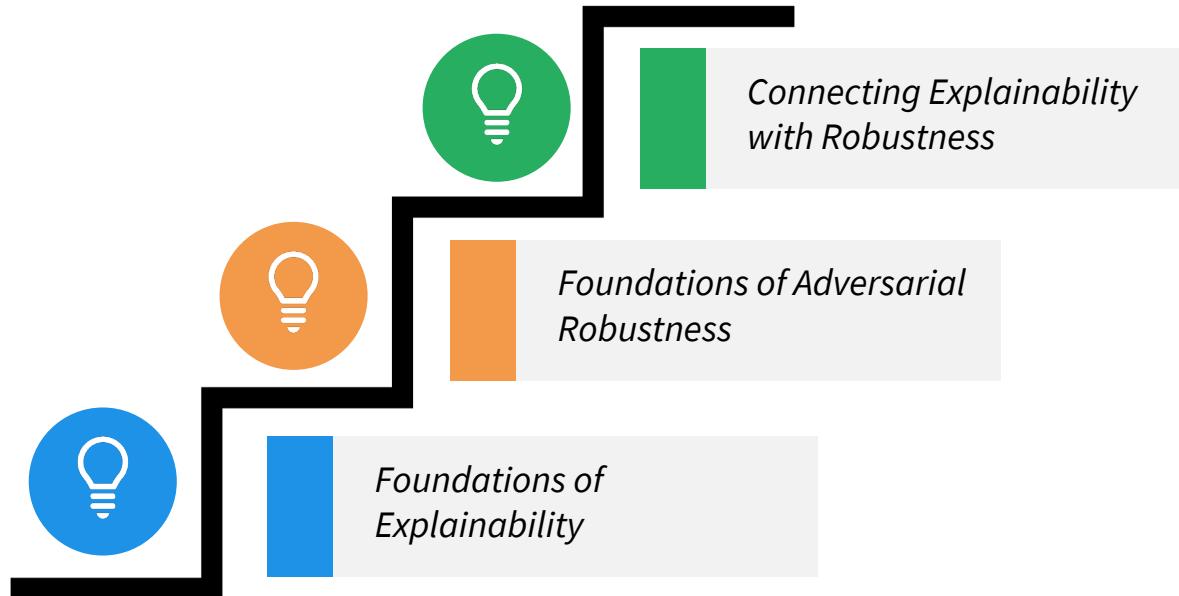


cat

# Takeaways



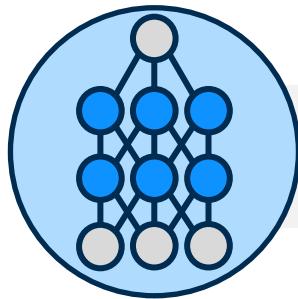
# Tutorial Roadmap



# Section I

## Foundations of XAI

# Explainability helps assess conceptual soundness



model 2

model 2 is not *conceptually sound*



dog

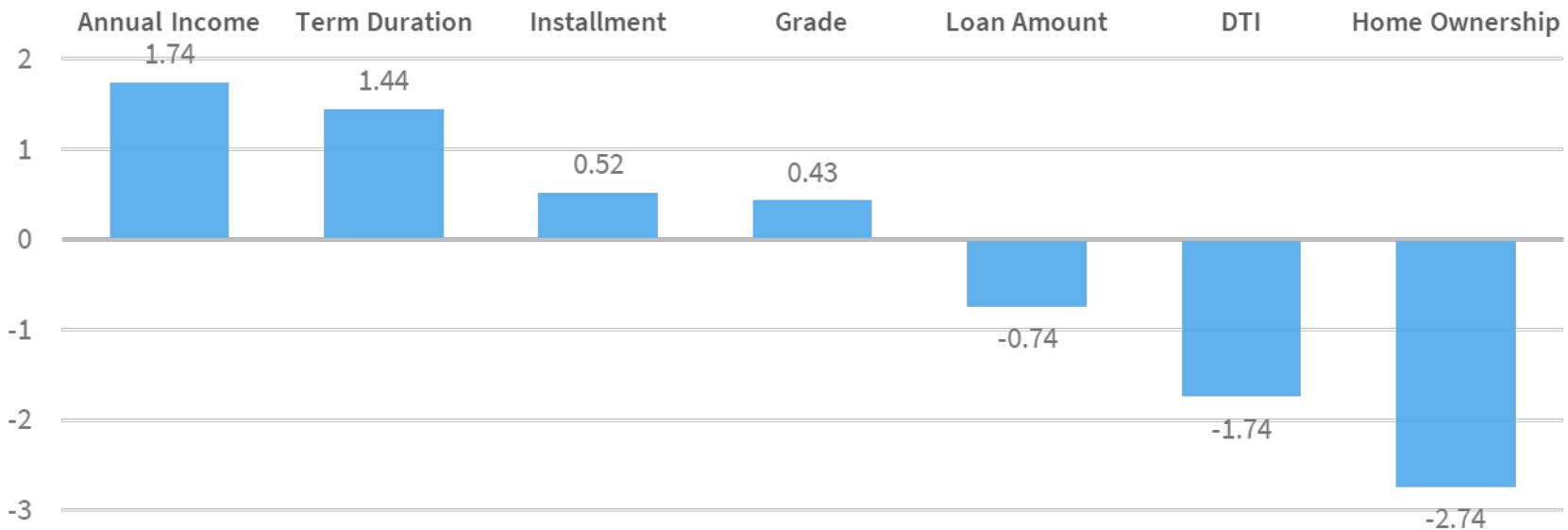


dog



dog

# Input feature importance for a tree model



# Internal explanation for a deep network

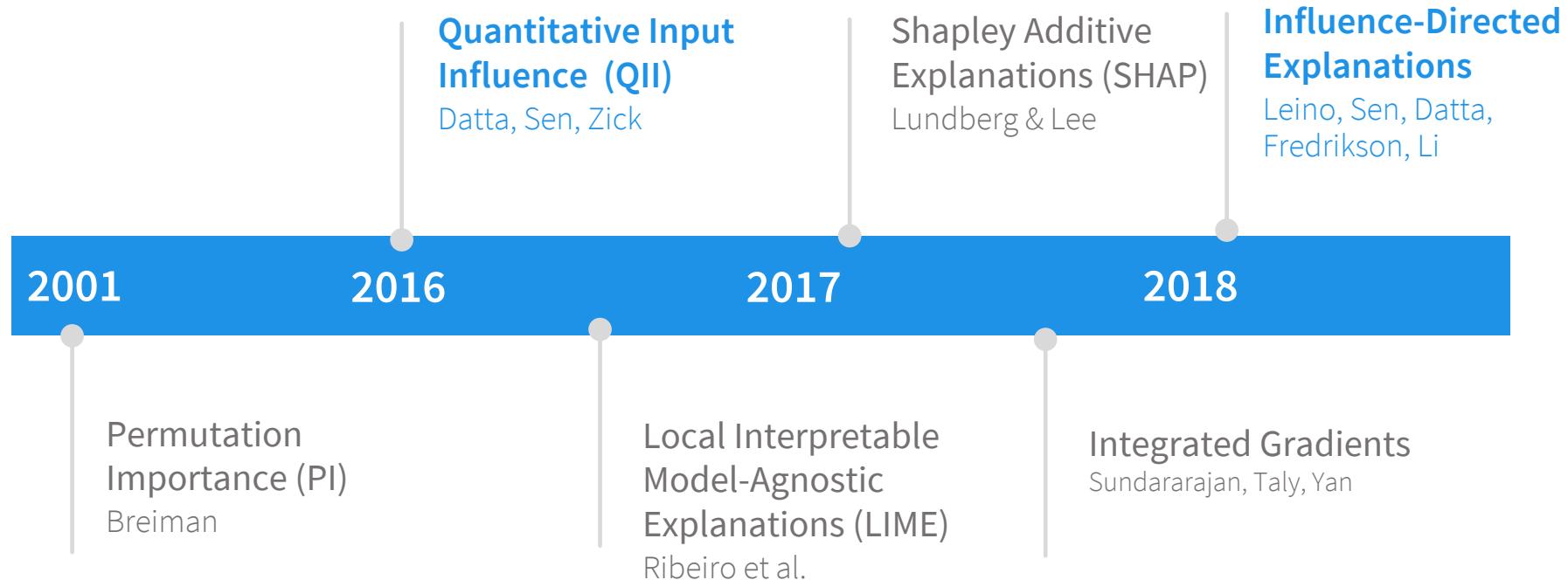


What makes Orlando Bloom Orlando Bloom?

**Influence-Directed  
Explanations**

Leino, Sen, Fredrikson, Datta, Li, ITC '18

# Methods for Computing Feature Importance

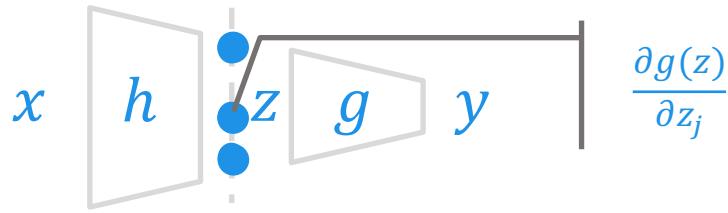


# Requirements for Good Explanations

- Explanation Accuracy
  - Accurately reflect reasons for model behavior
- Generality
  - Answer rich set of questions
- Interpretation devices
  - To make explanations meaningful to human users

Satisfied by a class of gradient-based explanation methods for deep neural networks

# Accuracy of Gradient Explanations



$$y = F(x) = g(z), z = h(x)$$

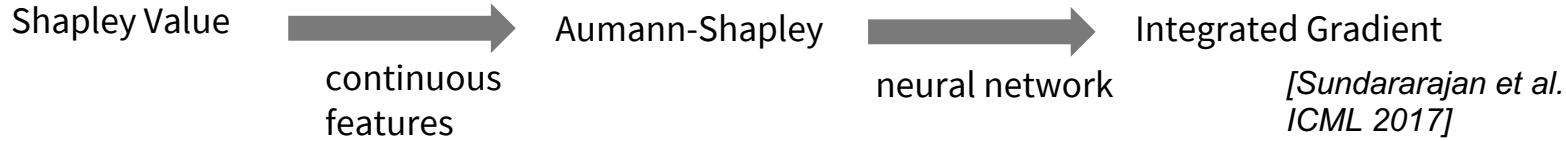
Feature Importance = Average gradient over a relevant distribution

- **Causality.** Gradients capture causal impact of features on model behavior
- **Distributional Marginality.** If gradients of output wrt a feature are equal for two different networks on a relevant distribution, then that feature is equally important for both networks

Influence-Directed  
Explanations

Leino, Sen, Fredrikson, Datta, Li, ITC '18

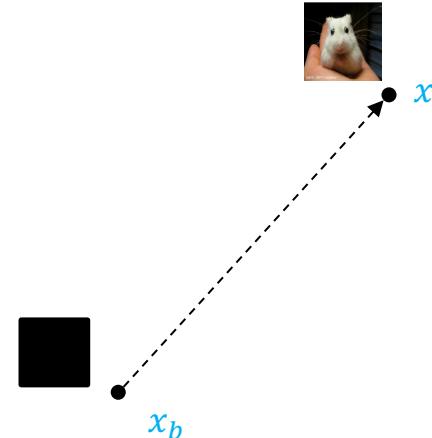
# Integrated Gradient



$$IG(x; x_b, F) = (x - x_b) \int_0^1 \frac{\partial F(\gamma(\alpha; x, x_b))}{\partial \gamma} d\alpha$$

where  $\gamma(\alpha; x, x_b) = x_b + \alpha(x - x_b)$

Average gradient of points on a linear path  
from a baseline to the target input



# Integrated Gradient

Shapley Value



continues  
features

Aumann Shapley



differentiable  
output

Integrated Gradient

[Sundararajan et al.  
ICML 2017]

Integrated Gradient is the **only** path method  
that satisfies

- Symmetry
- Dummy
- Efficiency(Completeness)
- Additivity



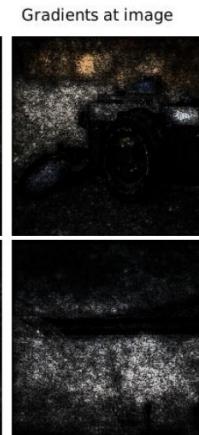
Original image

Top label and score  
Top label: reflex camera  
Score: 0.993755

Top label: fireboat  
Score: 0.999961

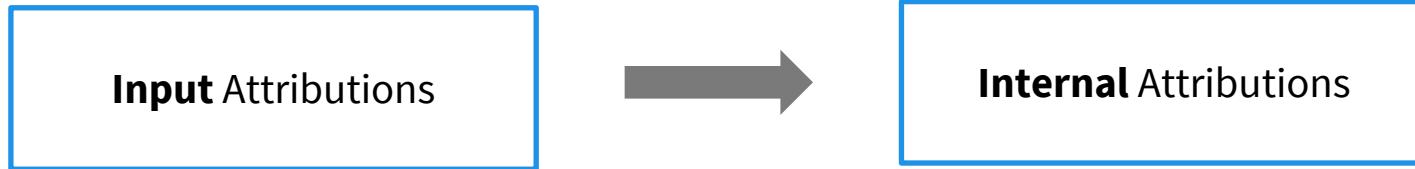


Integrated gradients



Gradients at image

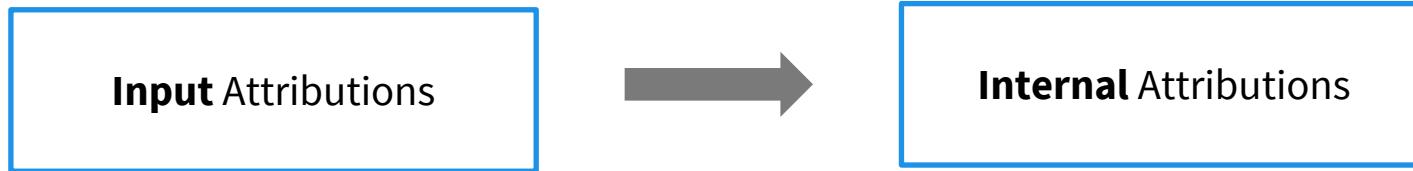
# Now It's Time to Dive Deeper...



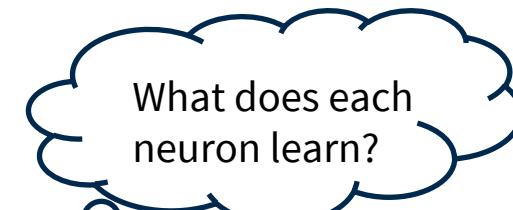
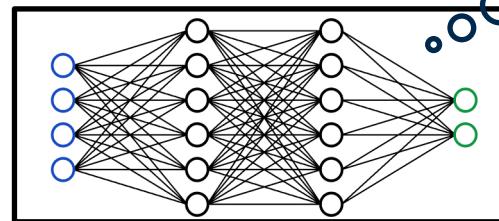
Why are we interested in internal representations?



# Now It's Time to Dive Deeper...

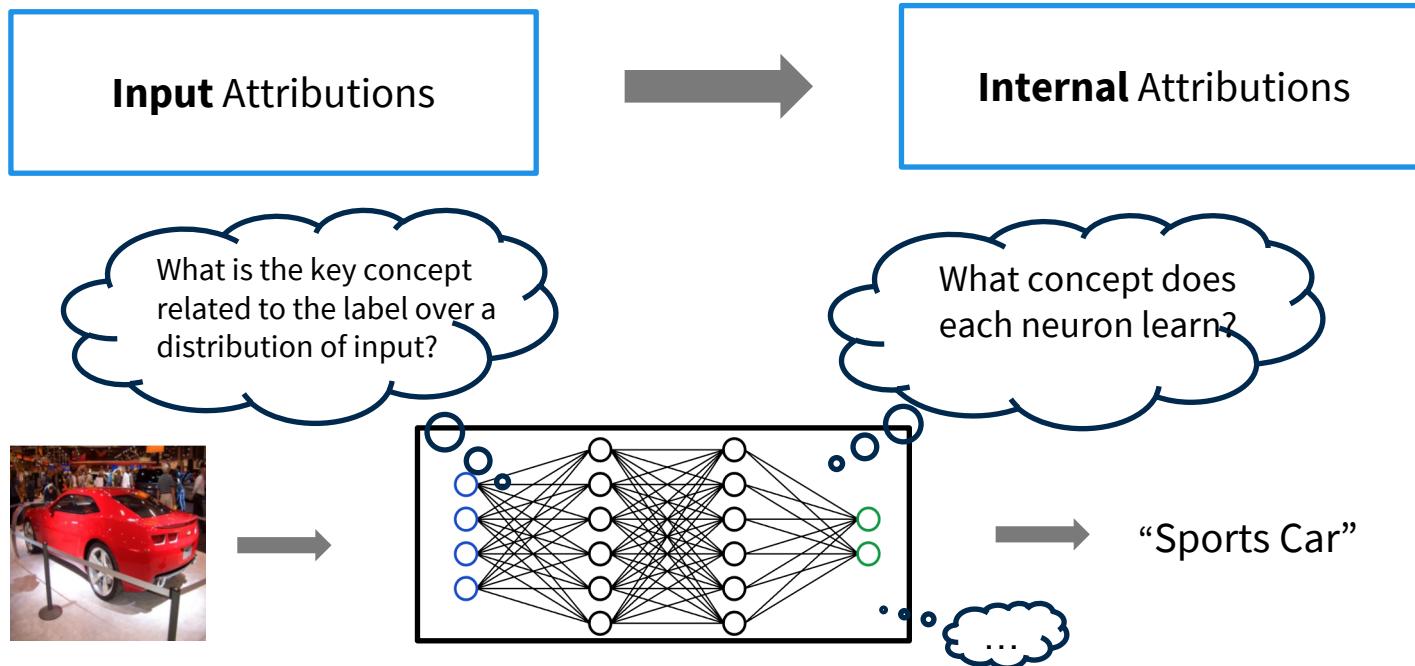


Why are we interested in internal representations?



→ “Sports Car”

# Now It's Time to Dive Deeper...



# What Makes Orlando Bloom Orlando Bloom?



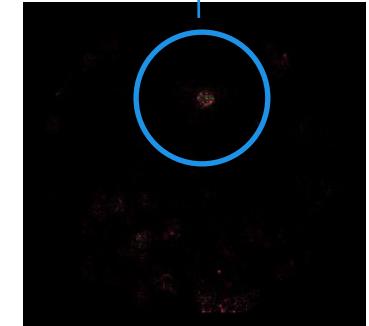
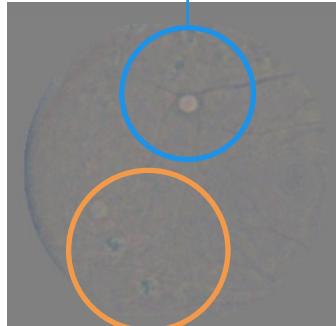
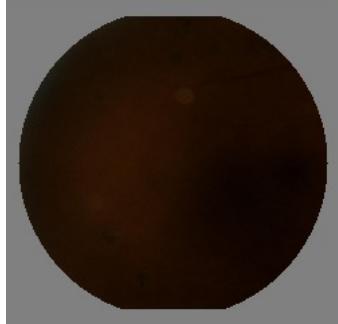
Internal explanation for a deep network

**Influence-Directed  
Explanations**

Leino, Sen, Fredrikson, Datta, Li, ITC '18

# Detecting Diabetic Retinopathy Stage 5

Optical Disk



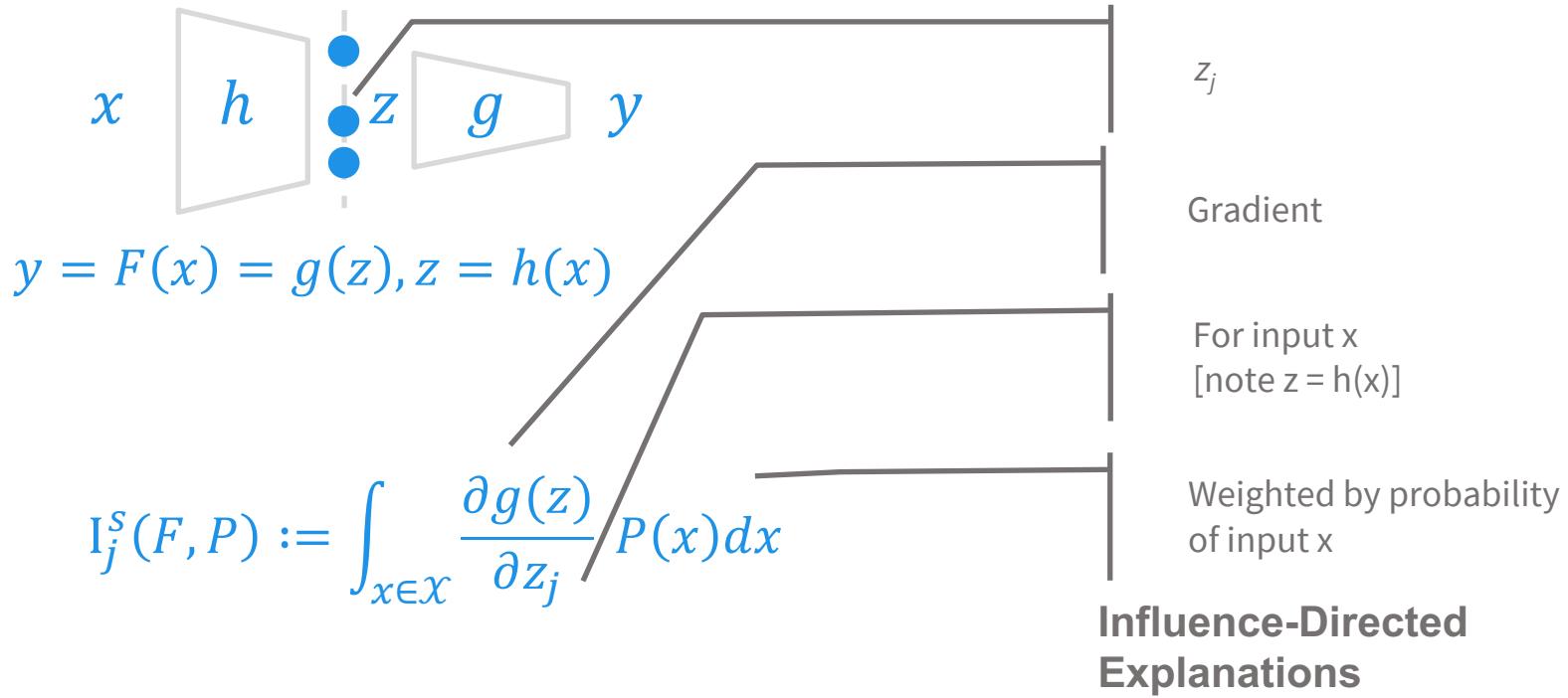
Lesions

Influence-Directed  
Explanations

Leino, Sen, Fredrikson, Datta, Li 2018

# Distributional Influence

Influence = average gradient over distribution of interest

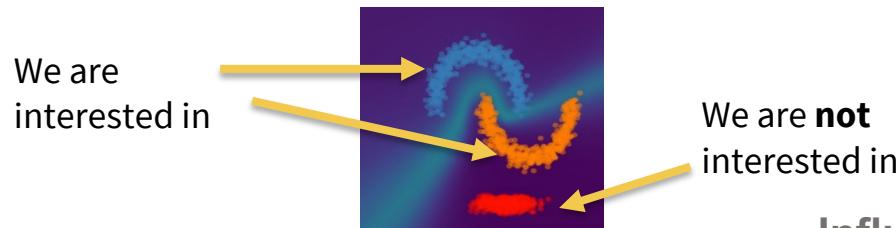


# Axiomatic Foundation for Distributional Influence

$$I_j^s(F, P) := \int_{x \in \mathcal{X}} \frac{\partial g(z)}{\partial z_j} P(x) dx$$

When  $s$  is the input slice ( $h(x) = x$ ), Distributional Influence satisfies:

- **Distributional Marginality**: If the partial derivatives w.r.t. an input feature are identical for  $F_1, F_2$  over the distribution of interest, then  $I_j^s(F_1, P) = I_j^s(F_2, P)$
- **Linear Agreement**
- **Completeness**
- ...



**Influence-Directed Explanations**

Leino, Sen, Fredrikson, Datta, Li ITC '18

# Distributional Influence Generalizes Existing Methods

$$I_j^s(F, P) := \int_{x \in \mathcal{X}} \frac{\partial g(z)}{\partial z_j} P(x) dx$$

When  $s$  is the input slice( $h(x) = x$ )

- and  $\mathcal{X}$  is a set of points (uniformly) distributed on a linear path from a baseline input to the target input



multiplying  $I_j^s(F, P)$   
with  $(x - x_b)$

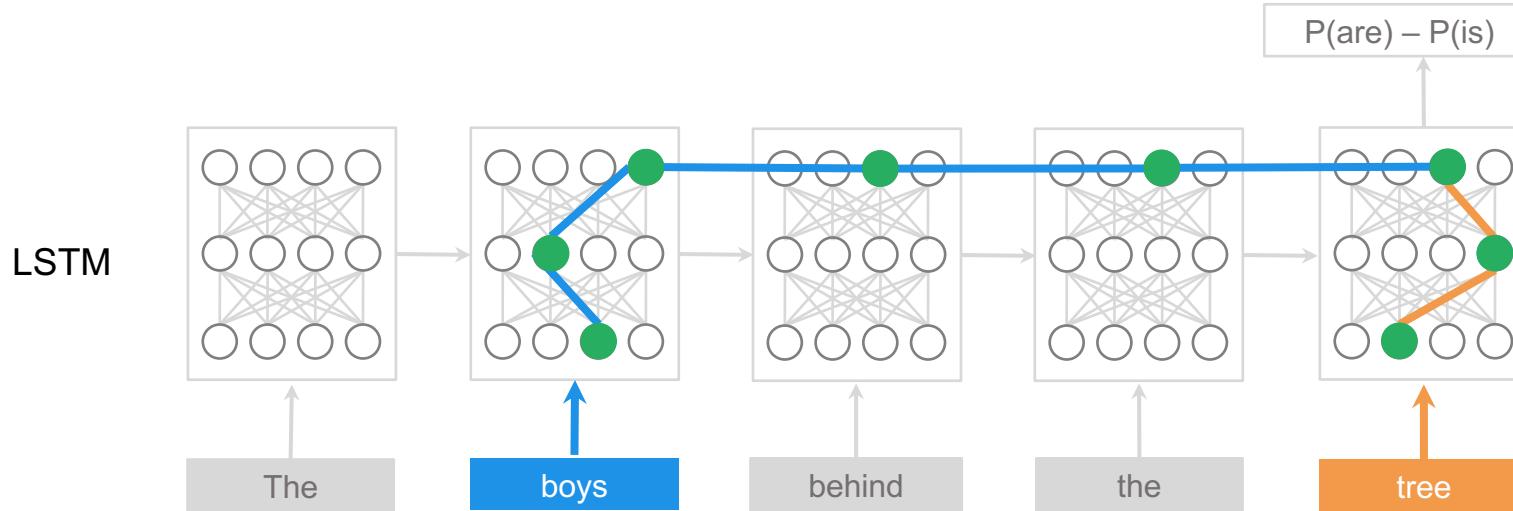
Integrated Gradient  
[Sundararajan et al. 2017]



Smooth Gradient  
[Smilkov et al. 2017]

⋮

# Internal Explanations via Influence Paths



- Influence paths provide insights into misclassifications
- Model can be compressed down the influential paths without changing the utility of the model

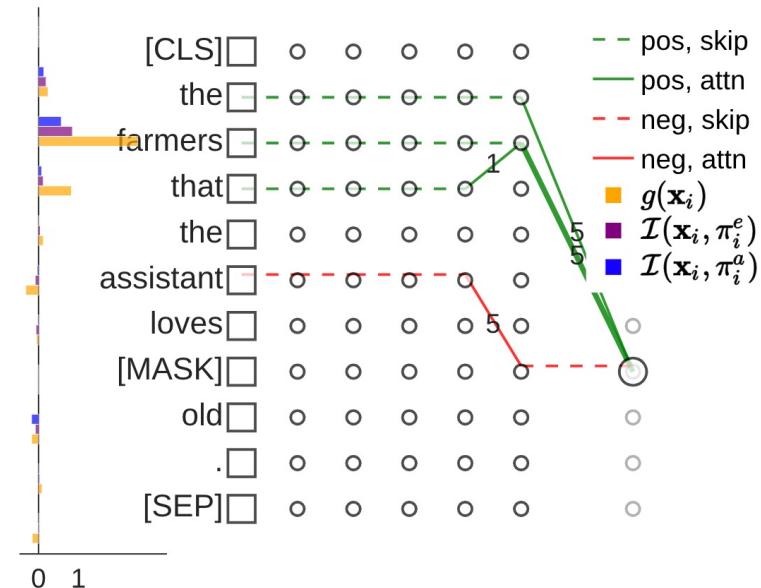
## Influence Paths

Lu, Mardziel, Leino, Fedrikson,  
Datta, ACL '20

# Influence Graphs for BERT

## BERT v.s. LSTM

- Scaling up method to identify influential paths
- Prevalence of “copy” and “transfer” operations to carry context



# Related Work

	Explanation Framework Properties			Influence Properties	
	Quantity	Distribution	Internal	Marginality	Sensitivity
Influence-Directed Explanation [Leino et al. ITC '18]	✓	✓	✓	✓	✓*
Conductance [Dhamdhere et al. ICLR '19]		✓-	✓	✓	✓
Integrated Gradient [Sundararajan et al. ICML '17]		✓-		✓	✓
Smooth Gradient [Smilkov et al. 2017]		✓-		✓	✓
Simple Taylor [Bach et al. 2015 PLOS ONE]		✓-		✓	
Deconvolution [Zeiler et al. ECCV '14]				✓†	
Guided Backpropagation [Springenberg et al. 2015 ICLR Workshop]				✓†	✓
Layer-wise Relevance Propagation [Bach et al. 2015 PLOS ONE]		✓-	✓†	✓*	✓*

✓ Supports

✓- Limited flexibility

✓\* Supports under some parameterizations

✓† Internal influence as an intermediate step

# Takeaways: Requirements for Good Explanations

- Explanation Accuracy
  - Accurately reflect reasons for model behavior
- Generality
  - Answer rich set of questions
- Interpretation devices
  - To make explanations meaningful to human users

Satisfied by a class of gradient-based explanation methods for deep neural networks

# Demo TruLens

Library containing attribution and interpretation methods for deep nets.

`pip install trulens`

TruLens is supporting models built with



# Q & A [15 min]

# Break I

# Section II

## Foundations of Adversarial Robustness

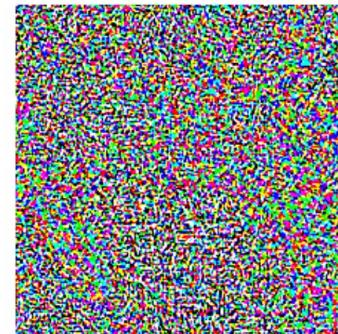
# Intriguing Properties of Neural Networks [Szegedy et al. 2014]

*“We find that applying an imperceptible non-random perturbation to a test image, it is possible to arbitrarily change the network’s prediction”*



“panda”

$+ .007 \times$

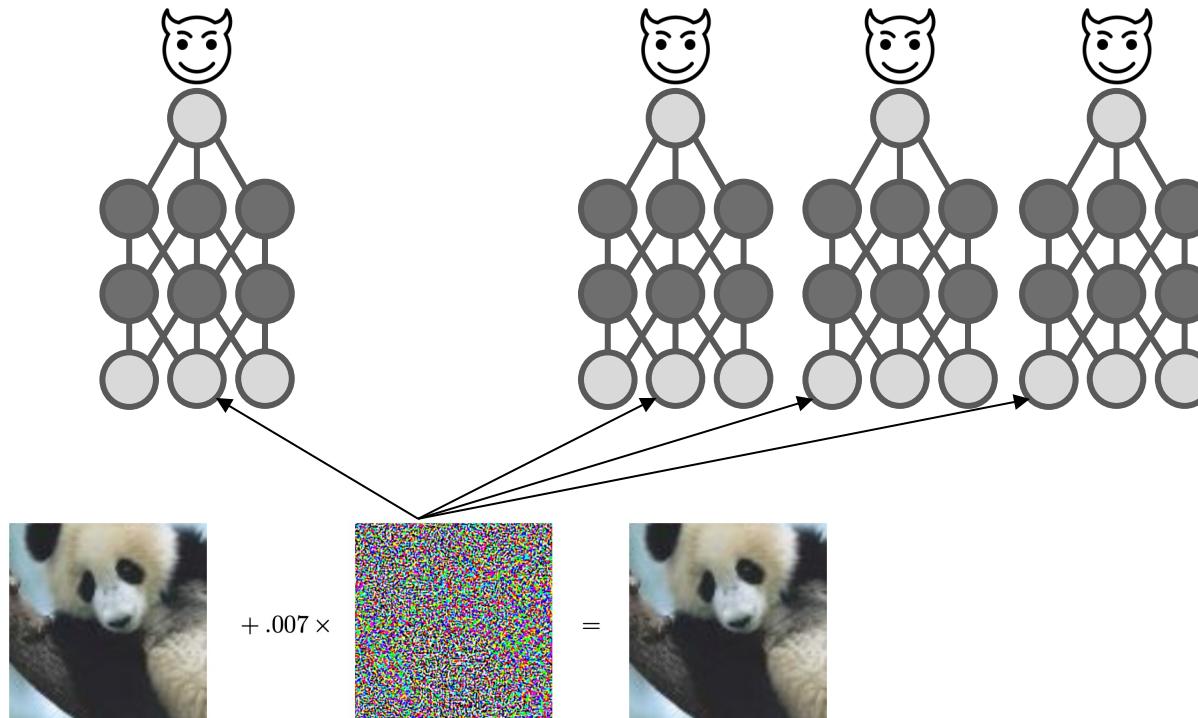


adversarial perturbation



“gibbon”

# Adversarial examples generalize



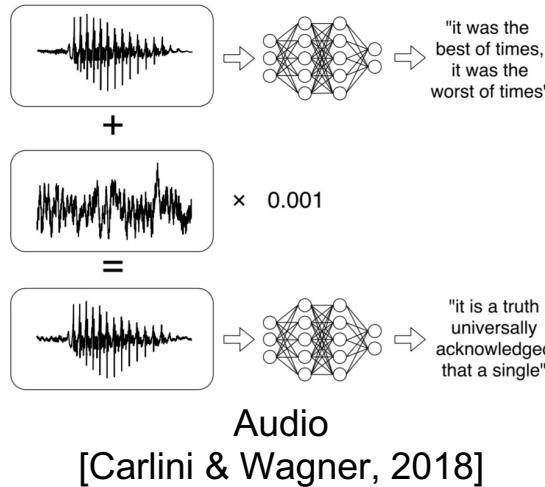
# It's not just images!

```

0x4587: ... jmp 0x4800    ...
0x458c: mov cx, cx      (6689c9)
0x458f: cmp ax, bx      (6639d8)
...
0x4800: add ax, 0x10    ...
0x4804: sub bx, 0x10    ...
0x4808: nop             (90)
0x4805: pushfd          (9c)
0x4806: push ebx         (53)
0x4807: add ebx, 0x1a    (83c31a)
0x480a: pop ebx          (5b)
0x480b: popfd            (9d)
0x480d: jmp 0x458c      (e97afdf000)
...

```

Malware Detection  
[Lucas et al., 2021]



Audio  
[Carlini & Wagner, 2018]

---

South Africa's historic Soweto township marks its 100th birthday on Tuesday in a mood of optimism.  
**57% World**

South Africa's historic Soweto township marks its 100th birthday on Tuesday in a mood of optimism.  
**95% Sci/Tech**

---

Chancellor Gordon Brown has sought to quell speculation over who should run the Labour Party and turned the attack on the opposition Conservatives.  
**75% World**

Chancellor Gordon Brown has sought to quell speculation over who should run the Labour Party and turned the attack on the opposition Conservatives.  
**94% Business**

---

Text  
[Ebrahimi et al., 2018]

# Is this really a problem?

Perhaps models just learn different concepts than humans

Some applications need more assurance

Secure biometrics [Sharif et al., 2016]



Safe autonomy

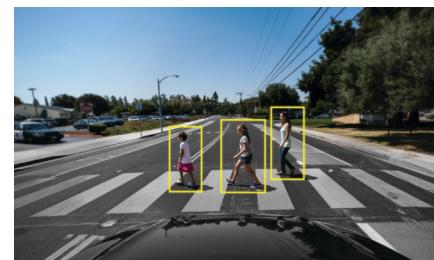
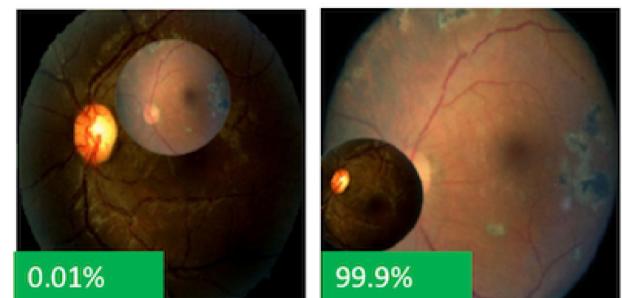


Image source: nvidia.com

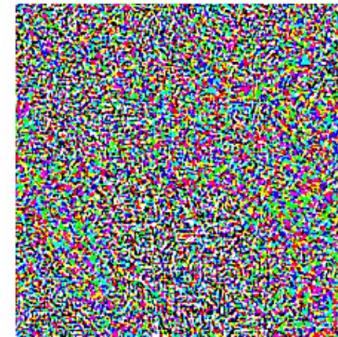
Trustworthy Diagnostics  
[Finlayson et al., 2019]



# Can we explain vulnerable models?



+ .007 ×



=



“panda”

Important features?

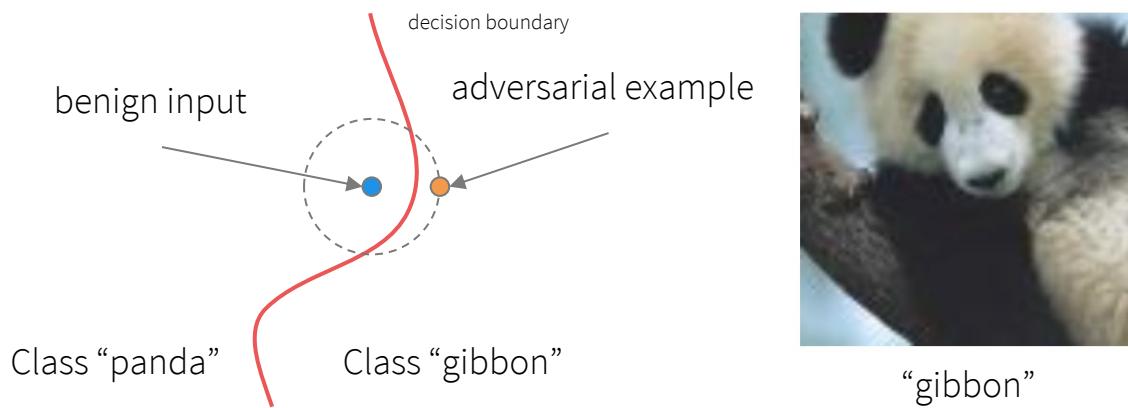
“gibbon”

*Adversarial examples are incompatible with conceptually-sound models*

# Adversarial examples are a violation of *local robustness*



“panda”



“gibbon”

# Local robustness

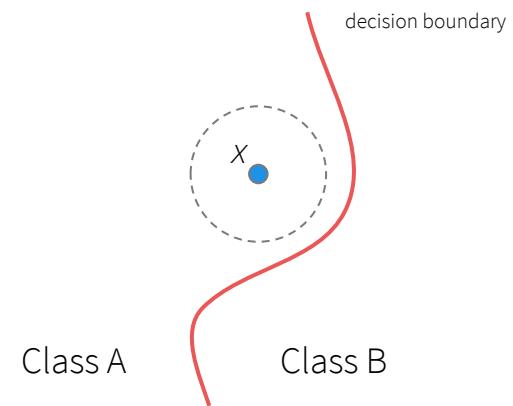
## Definition

A classifier,  $F$ , is  $\epsilon$ -locally-robust at  $x$  if  $\forall x'$ ,

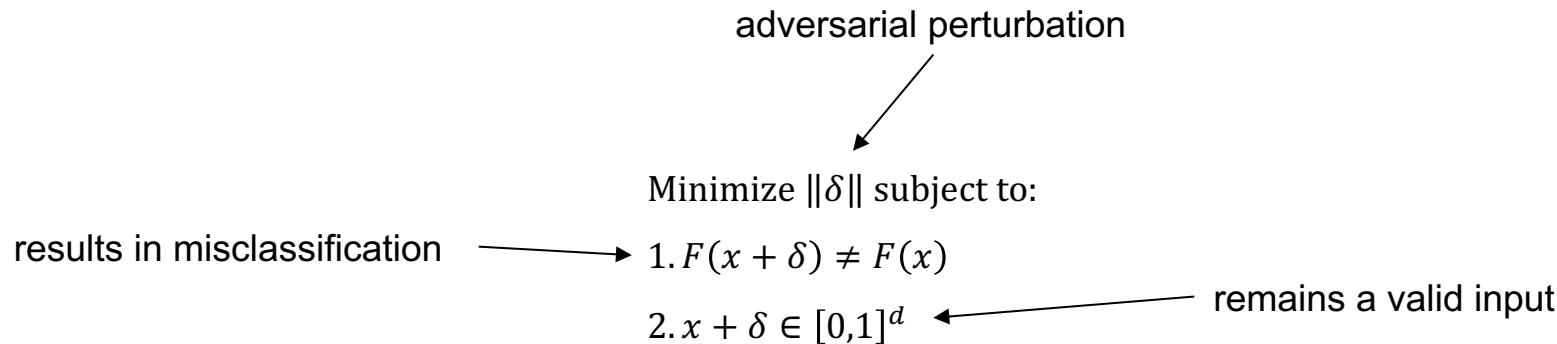
$$\|x - x'\| \leq \epsilon \Rightarrow F(x) = F(x')$$

$L_p$  norm, typically  $p \in \{1, 2, \infty\}$

i.e., the model makes the same prediction on all points in the  $\epsilon$ -ball centered at  $x$

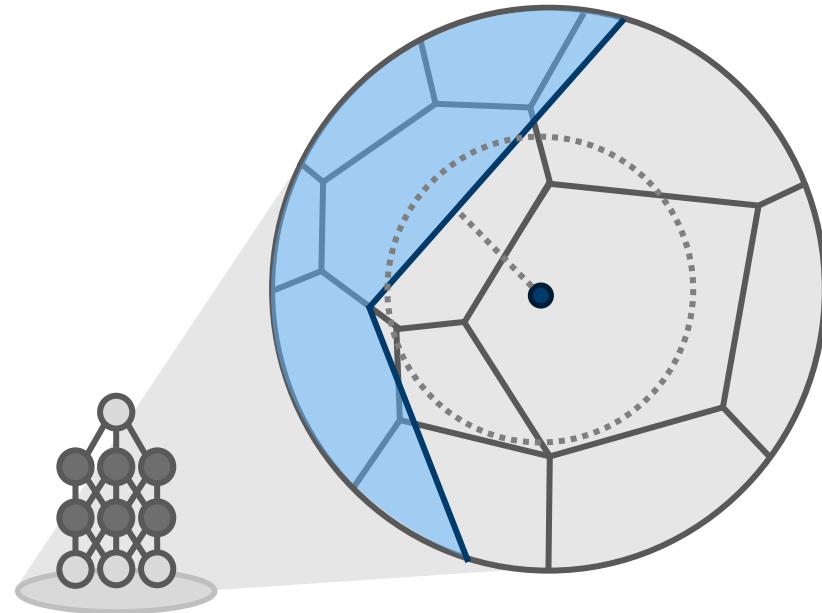


# Finding violations



# Exhaustive search

- ReLU networks are *piecewise-linear*
- Components comprise a *polyhedral complex*
- Each *activation pattern* of neurons gives a polyhedral region with linear boundaries
- Use constraint-solving to find distance to boundary for a point [Jordan et al. 2019]



**This is computationally infeasible for most networks**

# Finding violations, efficiently

Minimize  $\|\delta\|$  subject to:

1.  $F(x + \delta) \neq F(x)$
2.  $x + \delta \in [0,1]^d$



Maximize  $L(x + \delta, F(x))$  subject to  $\|\delta\| < \epsilon, x + \delta \in [0,1]^d$

- Loss function  $L$  is cross-entropy
- New formulation solved with *projected gradient descent*
- Despite non-convexity, this works well in practice
- Tune by varying number of PGD steps, step size

# Building robust models

- Goal: push boundaries away from data manifold
- Not of interest: detecting or filtering potential attacks, modifying inputs, stochastic smoothing

[Goodfellow et al., 2015]  
 [Madry et al., 2016]

**Standard Training**

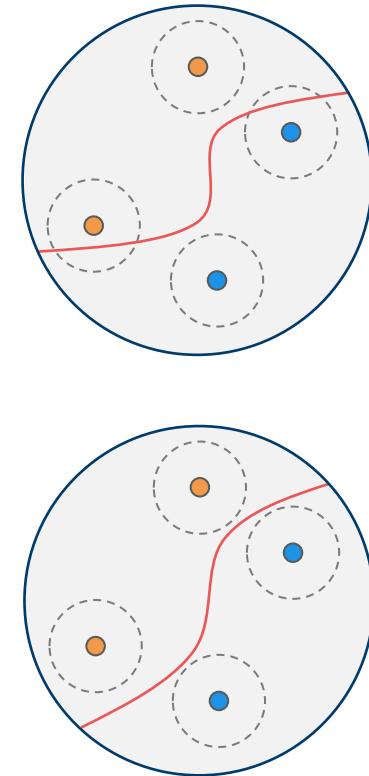
$$\min_{\theta} \mathbb{E}_{x,y \sim D} [L(\theta, x, y)]$$

- $D$  is the training data
- $\theta$  are the model parameters

**Adversarial Training**

$$\min_{\theta} \mathbb{E}_{x,y \sim D} [\max_{\|\delta\| < \epsilon} L(\theta, x + \delta, y)]$$

- $D$  is the training data
- $\theta$  are the model parameters
- $\epsilon$  is the local robustness radius



# Are adversarially-trained models robust?

## Adversarial Training

$$\min_{\theta} \mathbb{E}_{x,y \sim D} [\max_{\|\delta\| < \epsilon} L(\theta, x + \delta, y)]$$

*Small loss guarantees that no attack will succeed*  
[Madry et al., 2017]

Practical hurdles remain:

- Unclear whether we can sufficiently minimize loss
- This statement applies to the **training data**
- What can we say about test data?

In practice, adversarial training does not guarantee a specific degree of robustness

# Uncertain results

Reported degree of protection depends on the attack used to evaluate it!

Takeaways for AT:

- Difficult to train a robust model
- More difficult to evaluate its robustness on new data

<b>CIFAR-10 - <math>l_{\infty}</math> - <math>\epsilon = 8/255</math></b>		<b>originally reported</b>	<b>new attack</b>
1	(Carmon et al., 2019)	89.69	62.5 <b>-2.97</b>
2	(Alayrac et al., 2019)	86.46	56.30 <b>-0.27</b>
3	(Hendrycks et al., 2019)	87.11	57.4 <b>-2.48</b>
4	(Rice et al., 2020)	85.34	58 <b>-4.58</b>
5	(Qin et al., 2019)	86.28	52.81 <b>0.03</b>
6	(Engstrom et al., 2019)	87.03	53.29 <b>-4.04</b>
7	(Kumari et al., 2019)	87.80	53.04 <b>-3.92</b>
8	(Mao et al., 2019)	86.21	50.03 <b>-2.62</b>
9	(Zhang et al., 2019a)	87.20	47.98 <b>-3.15</b>
10	(Madry et al., 2018)	87.14	47.04 <b>-3.00</b>
11	(Pang et al., 2020)	80.89	55.0 <b>-11.52</b>
12	(Wong et al., 2020)	83.34	46.06 <b>-2.85</b>
13	(Shafahi et al., 2019)	86.11	46.19 <b>-4.72</b>
14	(Ding et al., 2020)	84.36	47.18 <b>-5.74</b>
15	(Moosavi-Dezfooli et al., 2019)	83.11	41.4 <b>-2.90</b>
16	(Zhang & Wang, 2019)	89.98	60.6 <b>-23.96</b>
17	(Zhang & Xu, 2020)	90.25	68.7 <b>-32.25</b>
18	(Jang et al., 2019)	78.91	37.40 <b>-2.45</b>
19	(Kim & Wang, 2020)	91.51	57.23 <b>-23.01</b>
20	(Moosavi-Dezfooli et al., 2019)	80.41	36.3 <b>-2.60</b>
21	(Wang & Zhang, 2019)	92.80	58.6 <b>-29.25</b>
22	(Wang & Zhang, 2019)	92.82	66.9 <b>-39.97</b>
23	(Mustafa et al., 2019)	89.16	32.32 <b>-32.04</b>
24	(Chan et al., 2020)	93.79	15.5 <b>-15.24</b>

# Provable robustness

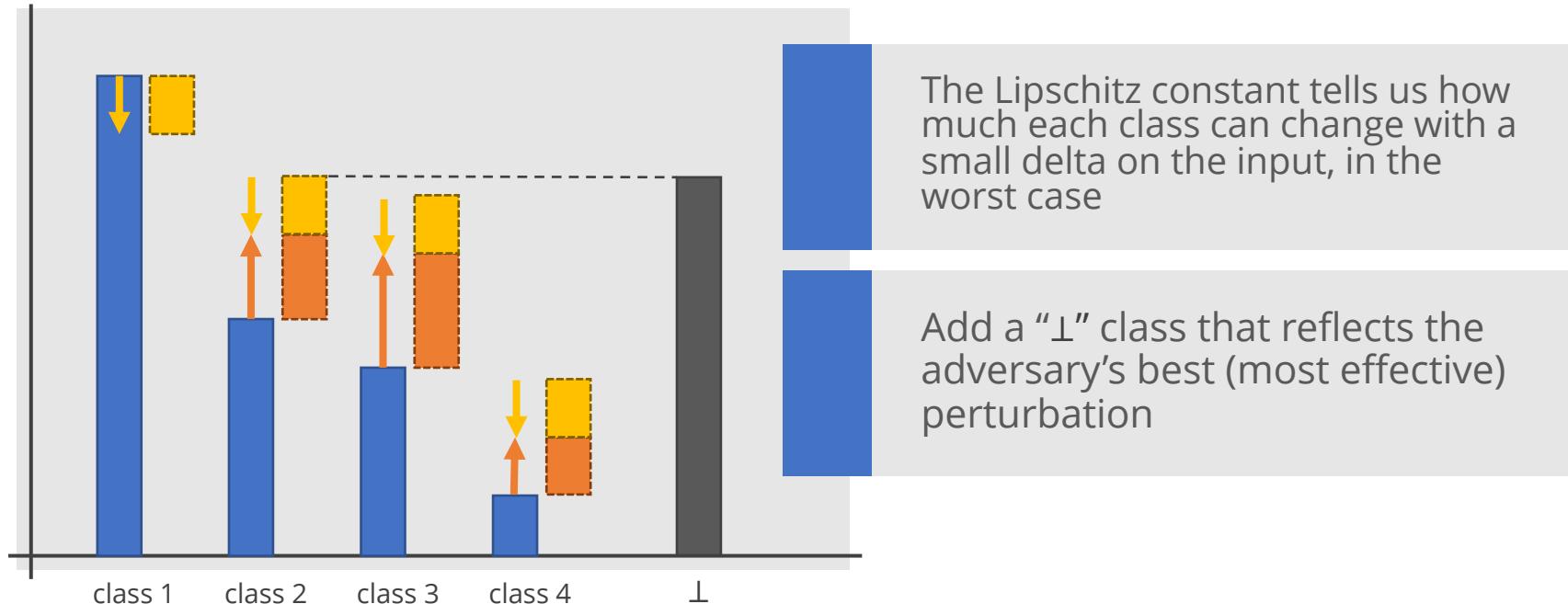
Goal: train the model so that it can be proven robust on new points

A function  $f$  is  $(K, \delta)$  locally-Lipschitz continuous iff for all  $x$ , there exists a neighborhood  $U_\delta$  around  $x$  such that  $f|_{U_\delta}$  is  $K$ -Lipschitz continuous:

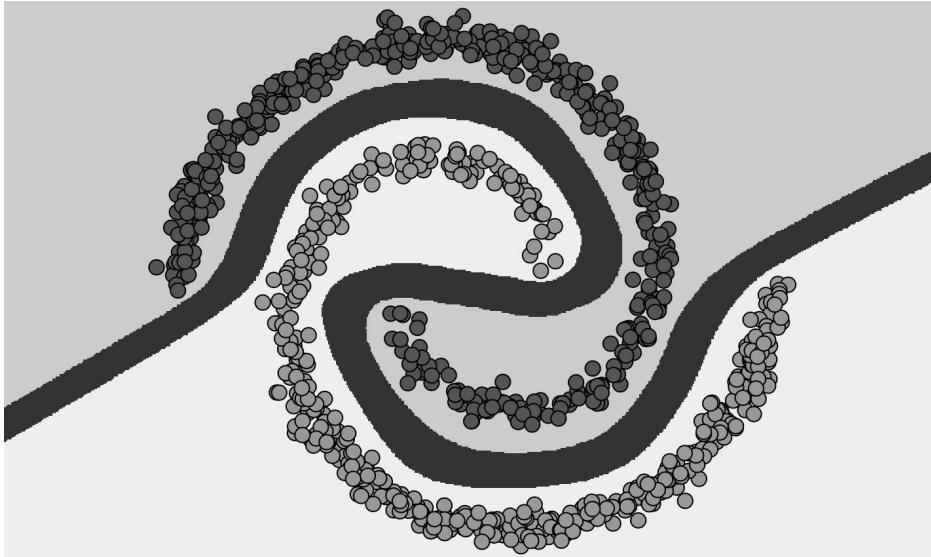
$$\forall x, x'. \|x - x'\| < \delta \Rightarrow \|f(x) - f(x')\| < K\|x - x'\|$$

**Lipschitz continuity is key to robust predictions!**

# Building a provably-robust model



# Training with $\perp$ [Leino et al. 2021]



The “ $\perp$ ” class lets us train the model using standard loss

Points classified as  $\perp$  incrementally corrected with each update

This leads to networks that:

- are sufficiently Lipschitz
- have boundaries with margins away from the manifold

# Calculating the Lipschitz bound



The **global** Lipschitz constant can be efficiently bounded by taking a product of the spectral norm of each layer



Any layer with a bounded Lipschitz constant can be used

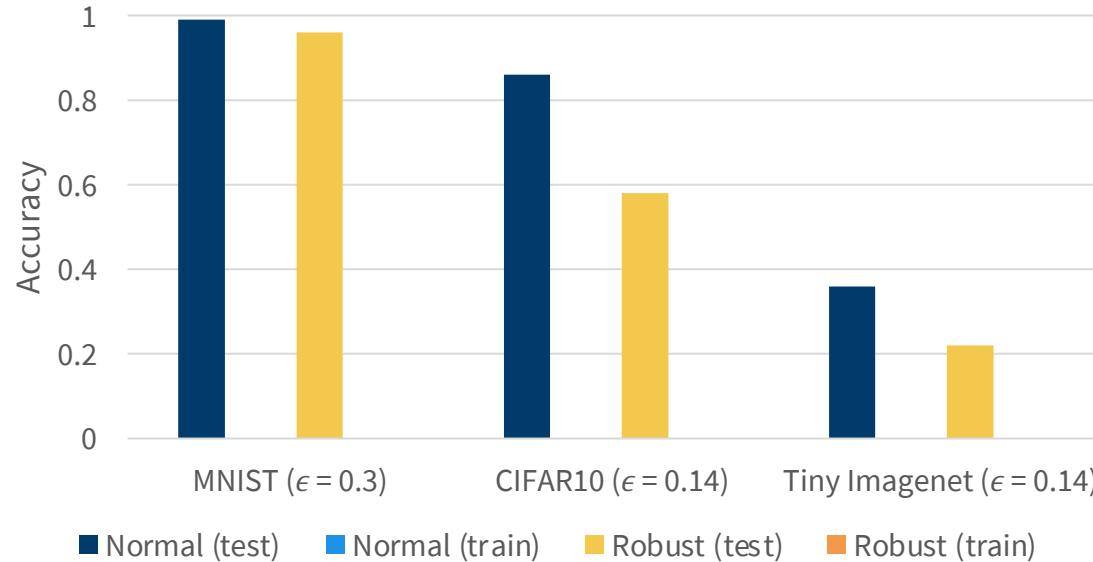


Little training overhead, compatible with normal optimization



Final Lipschitz constant computed **once** after training—no cost at test time!

# Robust training results



**Key challenge: improve robust generalization**

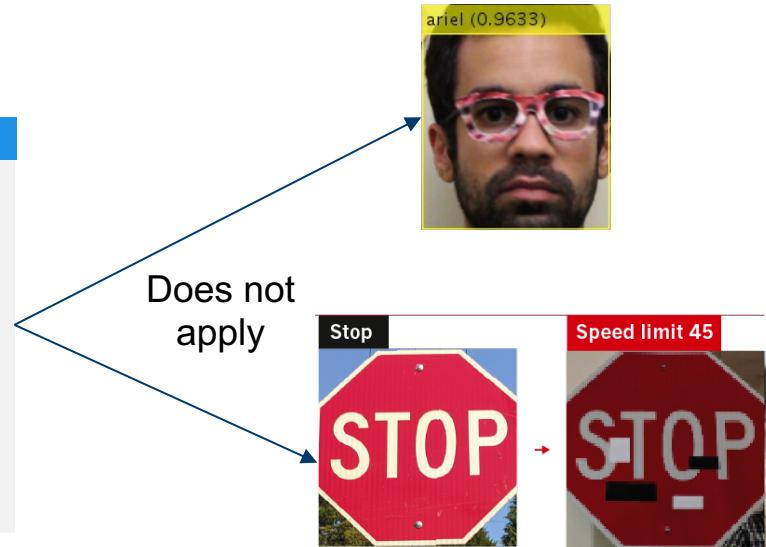
# Beyond $L_p$ -robustness

## Definition

A classifier,  $F$ , is  $\epsilon$ -locally-robust at  $x$  if  $\forall x'$ ,

$$\|x - x'\| \leq \epsilon \Rightarrow F(x) = F(x')$$

Does not apply



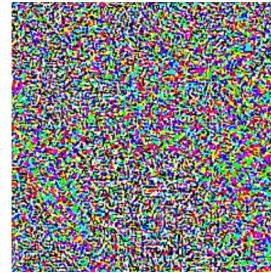
**Key challenge: “semantic” robustness**

# Takeaways

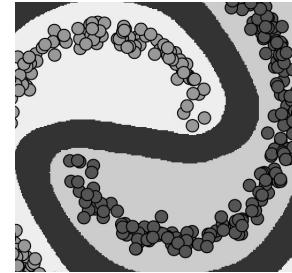
“gibbon”



*Adversarial examples* are ubiquitous in deep learning



They raise important questions about models' *conceptual soundness* and *explainability*



*Local robustness* addresses this problem, and robust training is becoming more effective and scalable



Far from being solved, extending robust learning to “semantic” attacks is an outstanding challenge

Up next: how does robustness bear on explainability?

# Q & A [15 min]

# Break II

# Section III

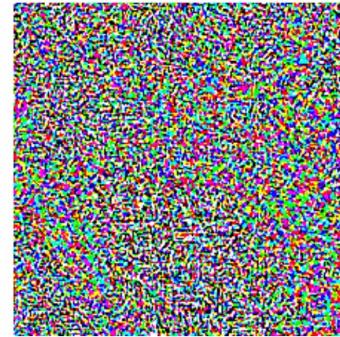
## Connecting Explainability with Robustness

# Recall: deep networks are easy to fool



“panda”

$+ .007 \times$

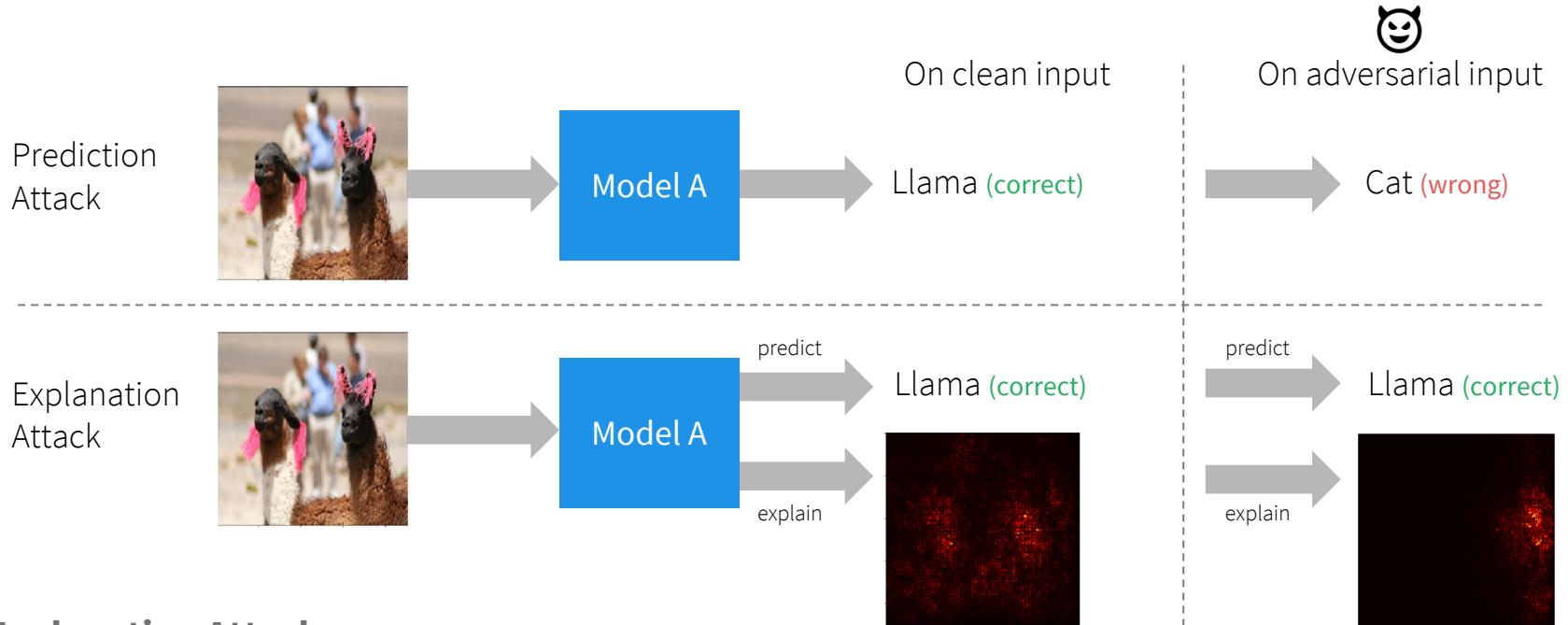


adversarial perturbation



“gibbon”

# Explanations can also be manipulated adversarially



## Explanation Attacks

Ghorbani et al. AAAI '19\*

Dombrowski et al. NeurIPS '19

Wang et al. NeurIPS '20



attribution map changes  
significantly

# Can we trust explanations?

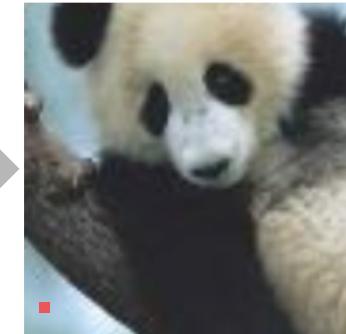
- If explanations can be manipulated, can we trust them?
- Is there something wrong with the explanation method that produces these anomalies?

# Can we trust explanations?

suppose that  
changing just one  
pixel in this region  
prevents the model  
from predicting  
“panda”



“panda”



not “panda”



possible explanation



Is it really wrong to assign influence to the pixel that can be modified to change the model’s prediction?

*If it weren’t for this pixel, this point would not be classified as “panda”*

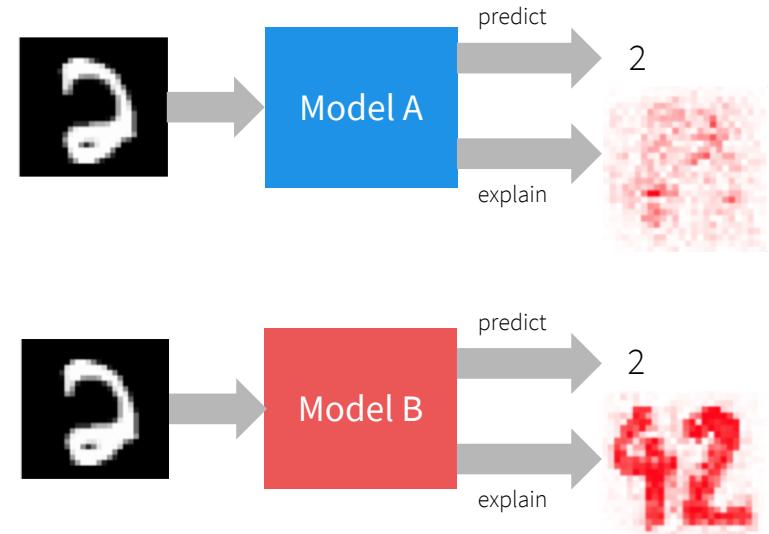
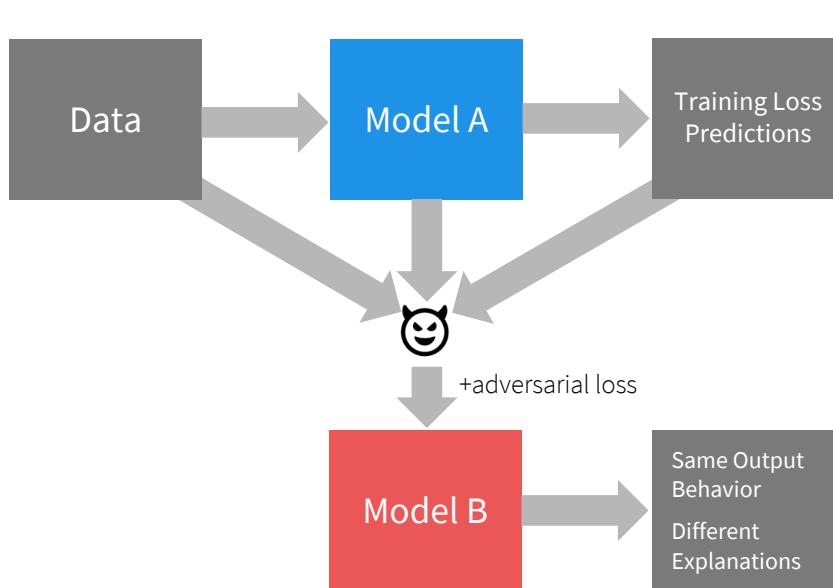
# Proposition



## Key Idea

“bugs” in *accurate* explanations are evidence of model quality issues

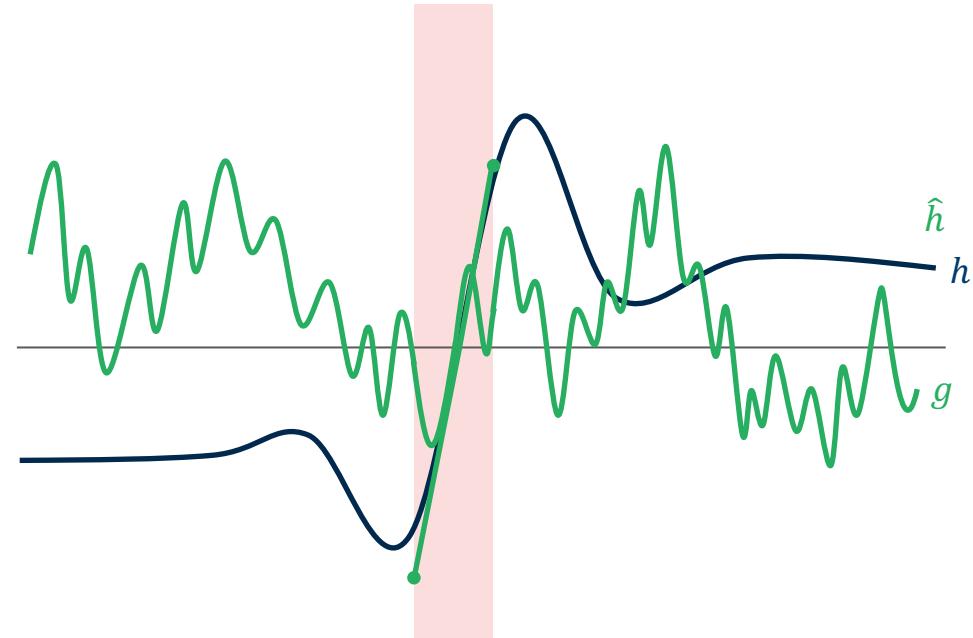
# Model-based attacks on explanations



# Same predictions, different gradients

## Theorem (Black et al. 2021)

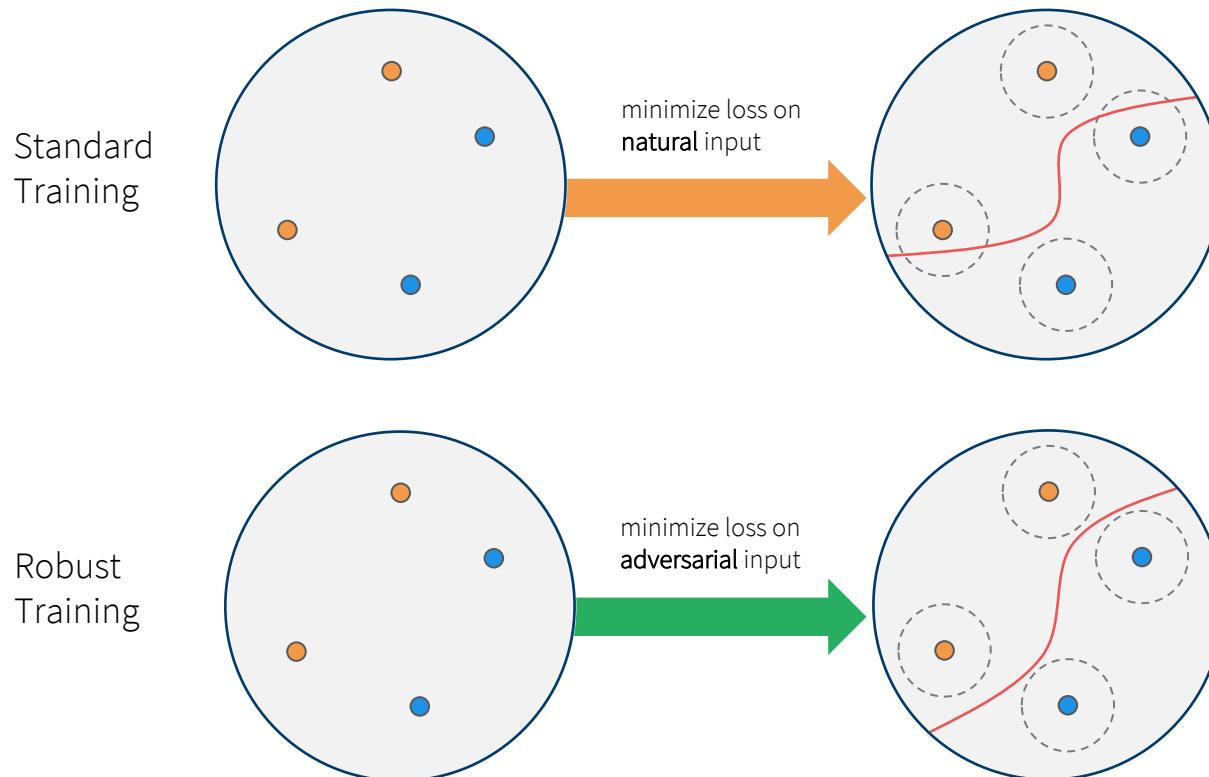
Let  $h$  be a binary classifier and let  $g$  be an unrelated bounded function with arbitrary gradients. Then there exists some classifier  $\hat{h}$  such that (1)  $\hat{h}$  makes the same predictions as  $h$ , and (2)  $\hat{h}$  has the same gradients as  $g$  almost everywhere.



# Now what?

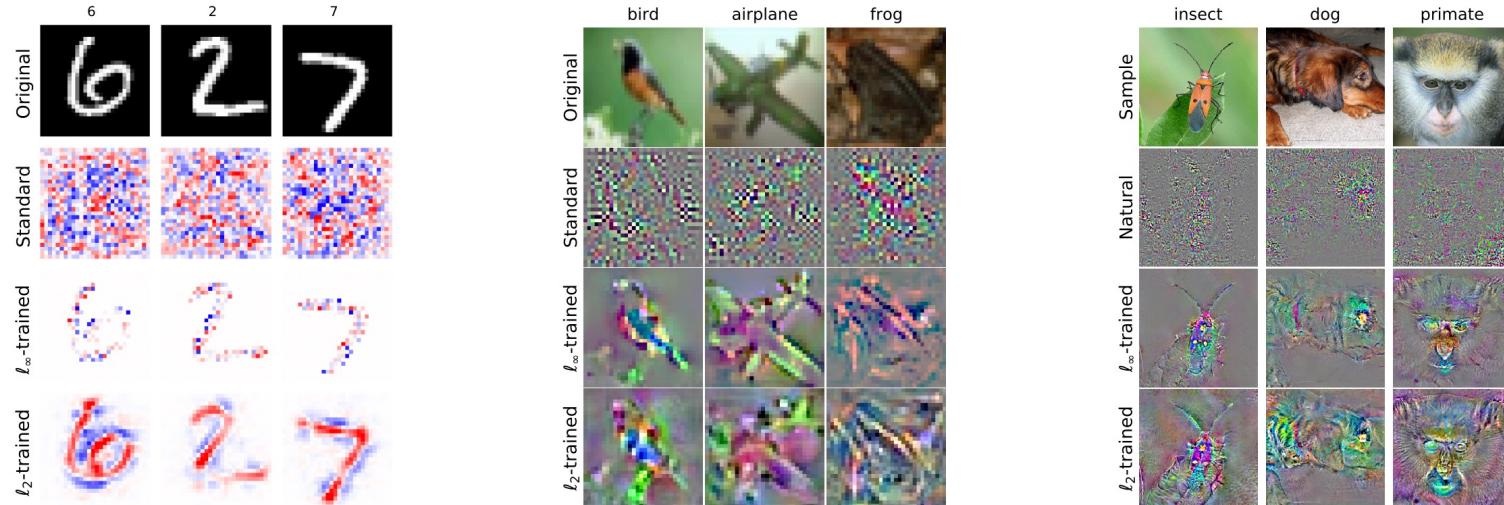
- **Key Idea:** “bugs” in accurate explanations are evidence of model quality issues
- On well-behaved models, we shouldn’t see these anomalies
- How do we improve model quality?

# Recall: obtaining robust models



# Robust models are more explainable

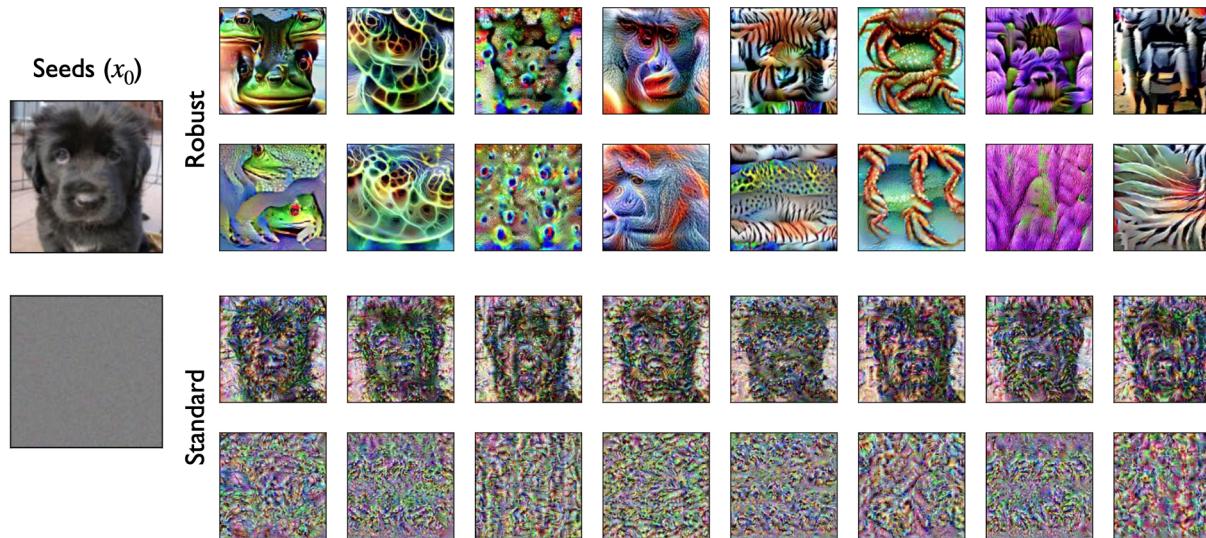
- Input gradients on robust models better align with the salient objects



**Explanations on Robust Models**  
*Tsipras et al. ICLR 2019\**  
*Etmann et al. ICML 2019*

# Robust models are more explainable

- Feature visualization on robust models yields more recognizable results



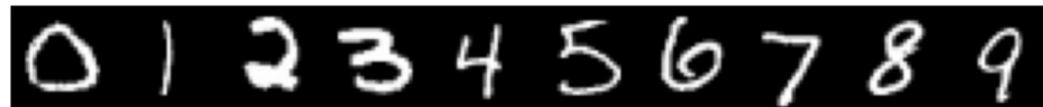
## Feature Visualization

For classifier,  $f$ , and class ,  $c$ , find  $\delta$  that maximizes  $f_c(x_0 + \delta)$

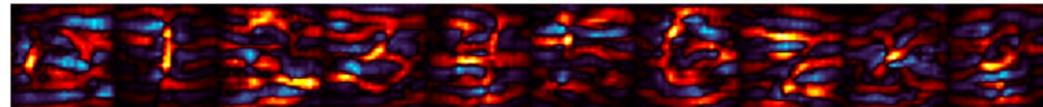
## Visualizations on Robust Models

*Tsipras et al. ICLR 2019*

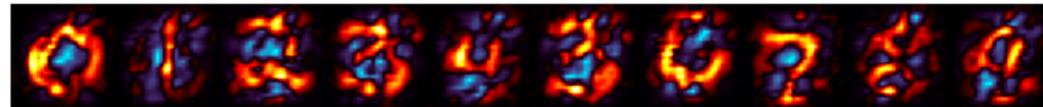
# Explanations for Provably Robust Models



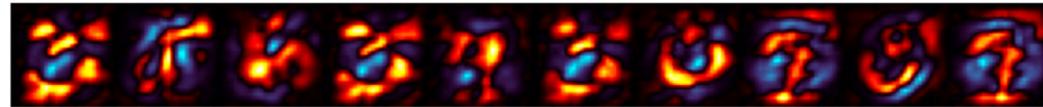
standard training



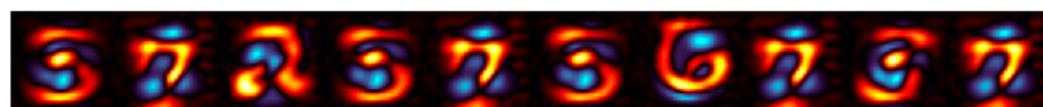
GloroNet  
[Leino et al. ICML 21']



[Lee et al. NuerIPS 20']



[Wong. et al. ICML 18']

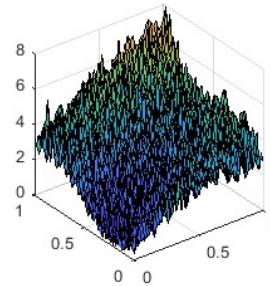


## Saliency Map

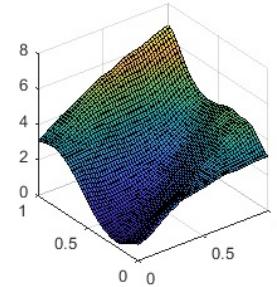
For classifier,  $f$ , input,  $x$ , and class,  $c$ , compute  $\partial f_c(x)/\partial x$

# Explanations on Robust Models are More Robust

Non-robust  
Model



Robust  
Model



Non-robust      Robust  
[Etman et al., ICML'19]

**Theorem (Wang et al.  
NeurIPS 20')**

If  $f$  is  $(K, \delta)$  locally-Lipschitz,  
then  $\nabla_x f$  is  $(O(K), \delta)$  locally-  
Lipschitz

# Proposition



## Key Idea

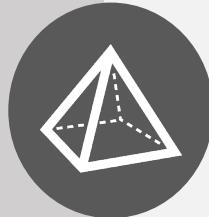
robustness is a prerequisite for human-interpretability

# Why are robust models more interpretable?



## Feature-based perspective

adversarial examples are the result of a model using non-robust features, which are inherently not human-interpretable



## Geometric perspective

the geometry of robust models leads to gradients that better convey information about the decision boundaries between classes

# Why are robust models more explainable?

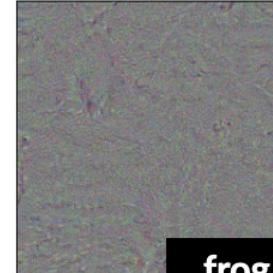
## A feature-based perspective



**Hypothesis** (*Ilyas et al. ICLR 2019*)

standard-trained models use *non-robust features* that are nonetheless predictive on the data distribution

example of non-robust  
features contained in an  
instance labeled “frog”



non-robust features only

**Non-robust Features**  
*Ilyas et al. 2019*

# Non-robust features

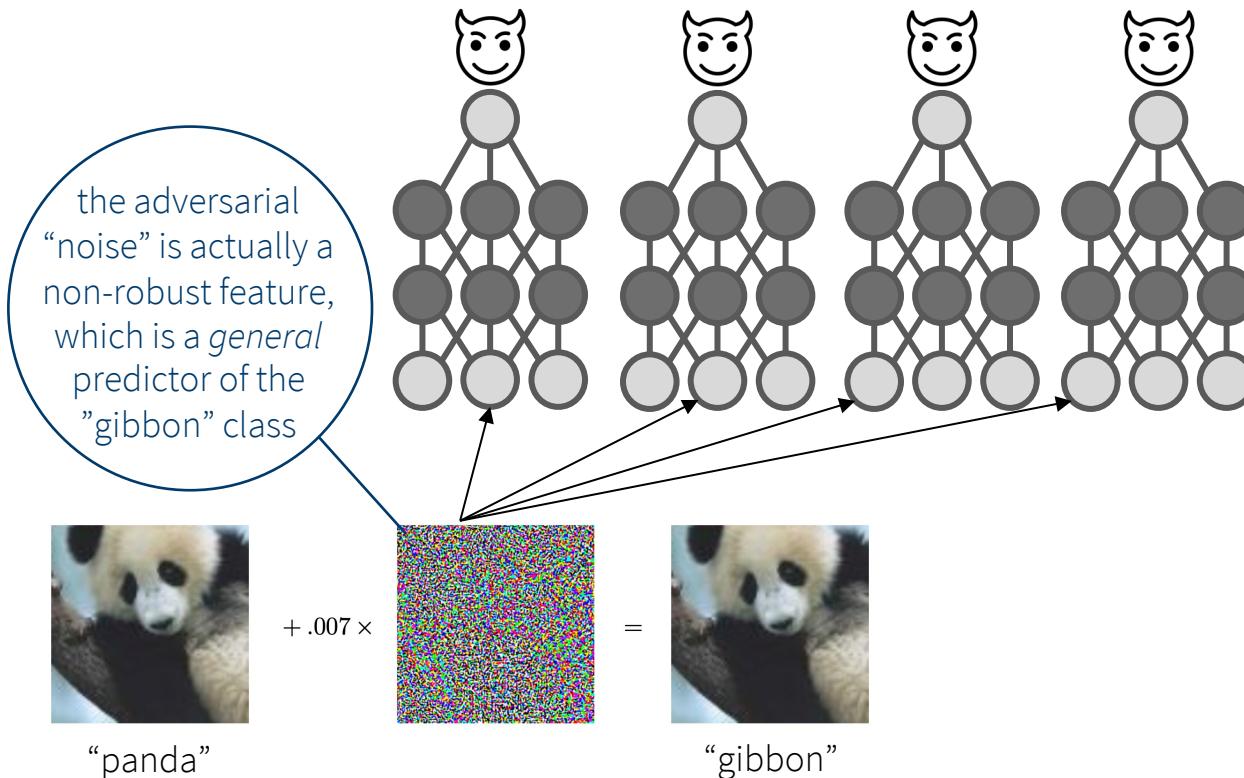
## Definition

A *feature* is a neuron in a neural network, which is a function,  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

## Definition

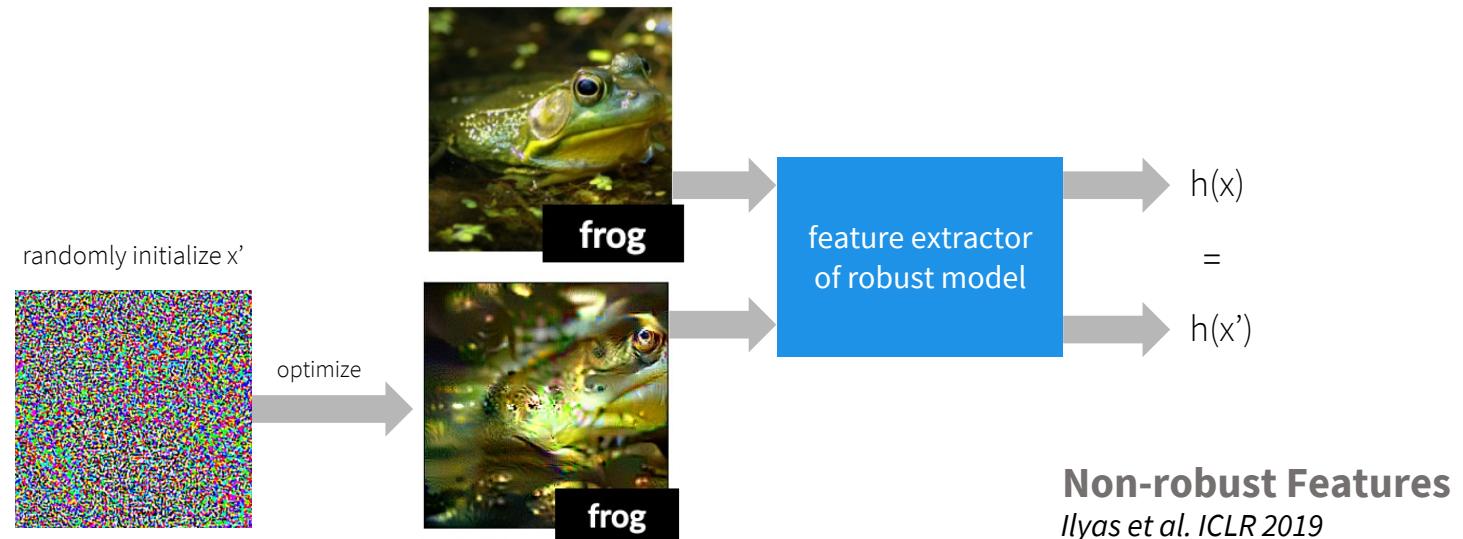
A feature is *non-robust* on data point,  $(x, y)$ , if  $f(x)$  correlates with  $y$ , but  $f(x + \delta)$  does not correlate with  $y$  for  $\|\delta\| \leq \epsilon$

# Recall: adversarial examples generalize



# Isolating robust features

Non-robust features are not useful for a robust objective, thus we do not expect robust models to learn them (i.e., robust models should only learn robust features)



# Why are robust models more explainable?

- Standard-trained models use *non-robust features* that are nonetheless predictive
- Non-robust features are not useful for a robust objective, thus we do not expect robust models to learn them
- Non-robust features are inherently less interpretable



**Non-robust Features**  
Ilyas et al. ICLR 2019

# Why are robust models more explainable?



## Feature-based perspective

adversarial examples are the result of a model using non-robust features, which are inherently not human-interpretable



## Geometric perspective

the geometry of robust models leads to gradients that better convey information about the decision boundaries between classes

# Why are robust models more explainable?

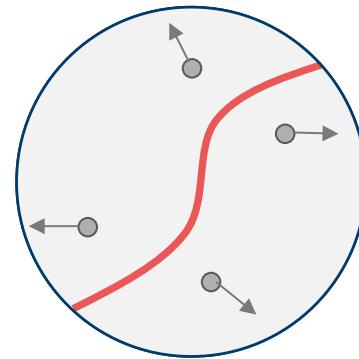
## A geometric perspective



Explanations usually focus on the gradient on the target inputs



Is the input gradient sufficient to explain the classification?



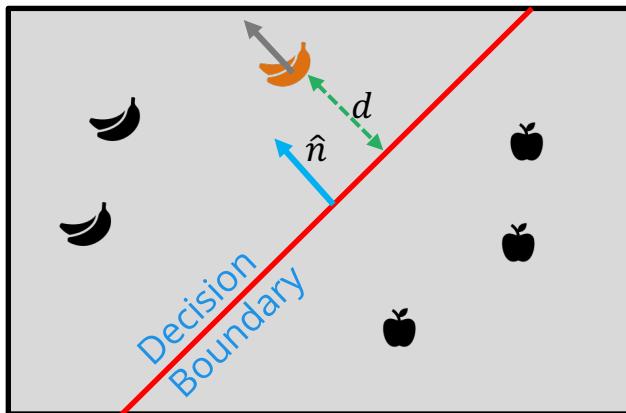
Input gradients may not accurately describe the decision boundary.

### Theorem (Informal, Wang et al. 2021)

Input gradients of a model with smaller Lipschitz Constants of gradients better capture the information of decision boundary

# How to Capture Decision Boundary ?

## Linear Classifier



	Why <i>Banana</i> instead of <i>Apple</i> ?
$\hat{n}$	Normal Vector $\hat{n}$ explains features importance
$d$	Greater projection distance corresponds to high confidence

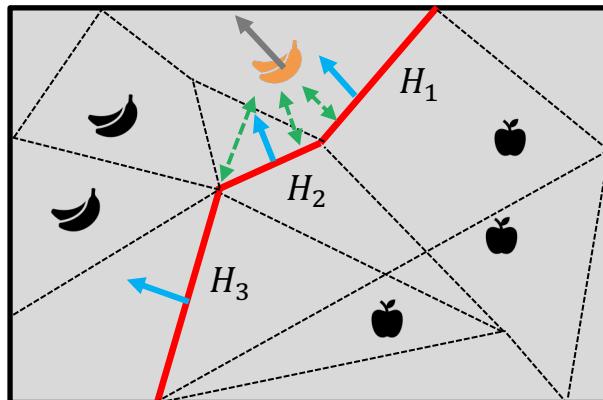
Found by axiomatic explanations

Calculated by projection

→ The input gradient has the same direction as the normal vector →

# Gradient Explanations May Not Capture Decision Boundary

## Non-linear Classifier



-  Why **Banana** is classified as *Banana* instead of *Apple*?
-  Which normal vector should be the explanation?
-  Which decision should be the confidence?

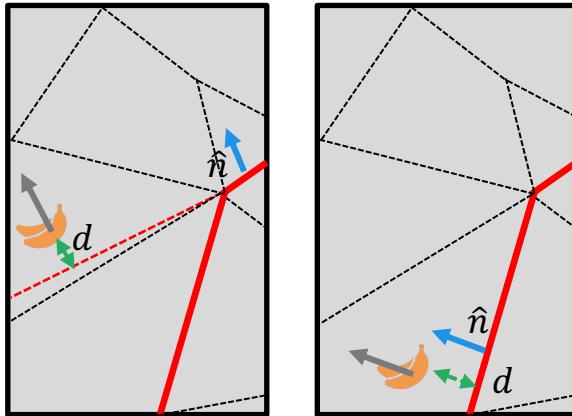


Closest / Nearby boundaries

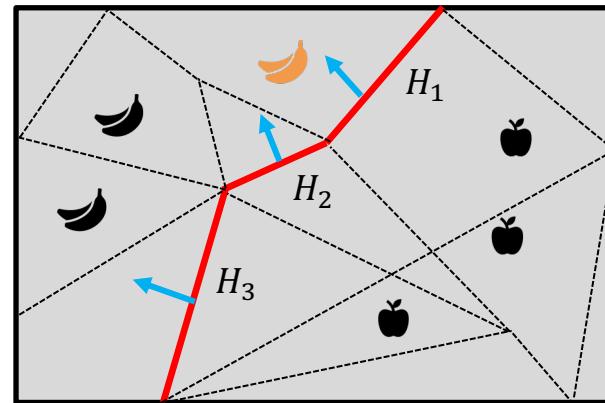
Which boundary is the input gradient explaining? ↗

# Computing Normal Vectors of Nearby Boundary

Issue 1



Issue 2

 $\hat{n}$ 

Gradient on the input may  
NOT be the  $\hat{n}$  from closest boundary

 $d$ 

Projection distance may NOT be the  
actual minimal distance

 $\hat{n}$ 

One normal vector is NOT enough  
to capture the model's behavior

# Boundary Attributions Provide Normal (Vector) Explanations

1

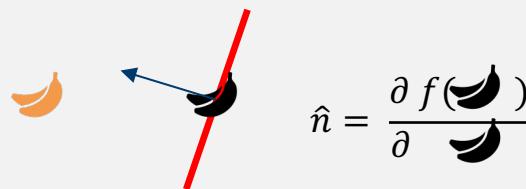
**Find the closest decision boundary**

- Projection is not reliable.
- Exact solution is NP-hard

**Approximation:** the closest adversarial example stands on the closest decision boundary

2

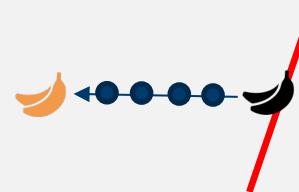
**Calculate the normal vector**



Gradient on the closest adversarial example is the normal vector

3

**Boundary-based Integrated Gradient (BIG)**



$$\int \frac{\partial f(x)}{\partial x} dx$$

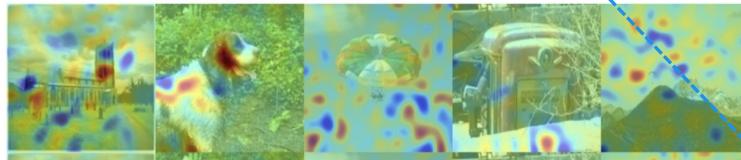
**Boundary Explanations**  
Zifan, Fredrikson, Datta, '21

# Boundary Attributions Provide Normal (Vector) Explanations

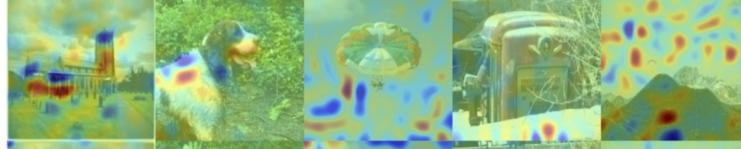
Correctly  
classified  
images



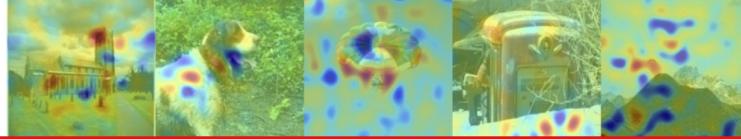
Saliency  
Map



Integrated  
Gradient



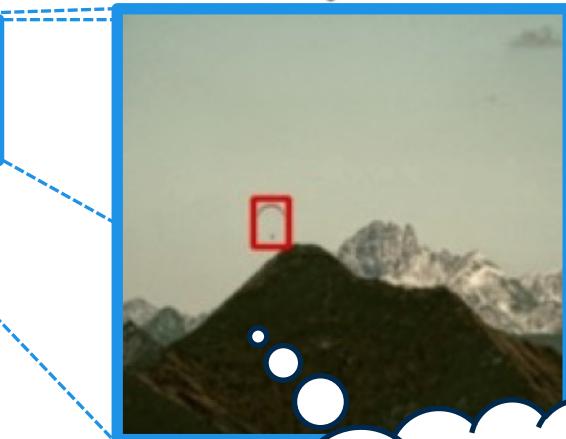
Single-boundary  
attribution



**BIG:**  
Multi-boundary  
attribution



Better  
explanation

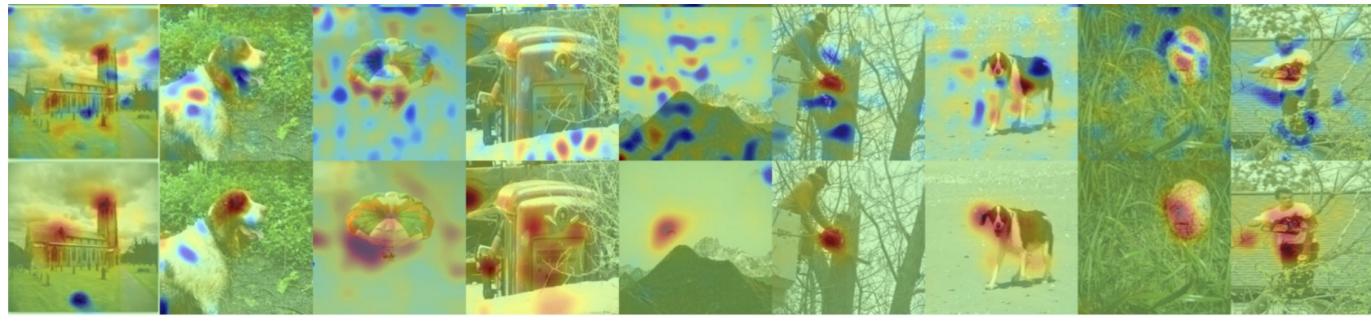


Why does the  
model think  
this is an image  
of parachute ?

# Robust Model Are More Explainable

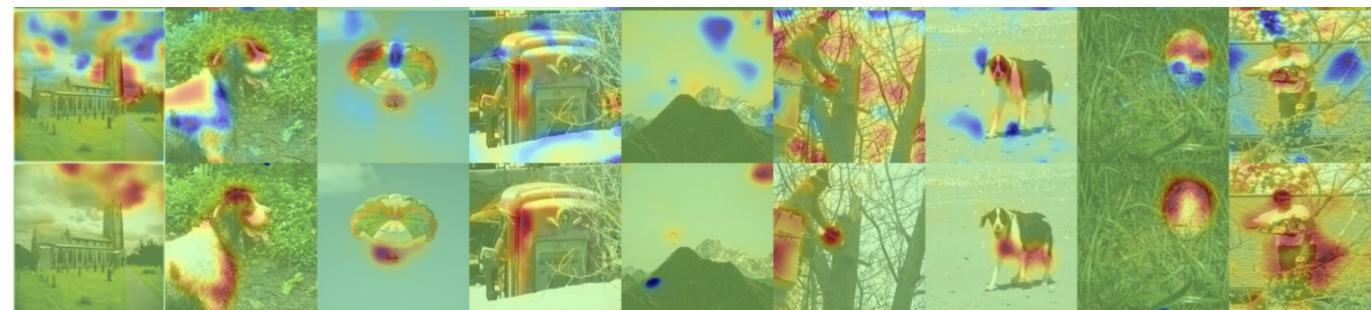
because gradient explanations better mimic boundary explanations

Gradient Explanation (IG)



**Non-Robust Model**

Boundary Explanations (BIG)



**Robust Model**

Gradient Explanation (IG)

**Boundary Explanations**  
Zifan, Fredrikson, Datta, '21

# Demo Boundary Attribution

## with TruLens

# Summary

“Bugs” in accurate explanations  
are evidence of model quality  
issues

Quality explanations require  
quality models

Robustness may be one way to  
achieve better model quality

# Q & A



truera

SIGKDD 2021

Tutorial | SIGKDD 2021

# Machine Learning *Explainability* and *Robustness* → Connected at the Hip

Anupam Datta  
Matt Fredrikson  
Klas Leino  
Kaiji Lu  
Shayak Sen  
Zifan Wang

We appreciate your participation in this tutorial.

For More Resources:

- [Tutorial Website](#)
- [Accountable Systems Lab](#)
- [TruLens and Demos](#)
- [Truera's Blog Posts on Explainability](#)

Contact Us: [shayak@truera.com](mailto:shayak@truera.com), [zifan@cmu.edu](mailto:zifan@cmu.edu)