

NC STATE UNIVERSITY



Technology Has No Moral Compass “The Dark Side of AI”

Steve Grobman


SVP & CTO, McAfee

@SteveGrobman


"I don't have any **regrets** ...
I feel about the **airplane** much the same
as I do in regard to **fire**.

That is, I regret all the terrible **damage** caused by fire,
but I think it is **good** for the human race
that someone discovered how to start fires,
and that we have **learned** how to put fire
to thousands of **important uses**."

Orville Wright
Interview in *The Rotarian*, April 1948


$$c = m^e \pmod n$$

$$c^d \equiv m \pmod n$$



156,000,000,000,000,000,000,000
bytes encrypted monthly



Encryption Munition

WARNING

this shirt is classified as a munition and
may not be exported from the United
States, or shown to a foreign national

RSA

encryption in perl

```
#!/bin/perl -s-- -export-a-crypto-system-sig -RSA-3-lines-PERL  
m=unpack(H.$w,$m."\\0"x$w),$_=`echo "16do$w 2+40i0$d*-^1[d2%Sa  
2/d0<X+d*Lal=z\\U$n%0]SX$k"[$m*]\\EsZlXx++p|dc`,s/^.|\\W//g,print  
pack('H*',$_)while read(STDIN,$m,($w=2*$d-1+length($n)&~1)/2)
```

Artificial Intelligence

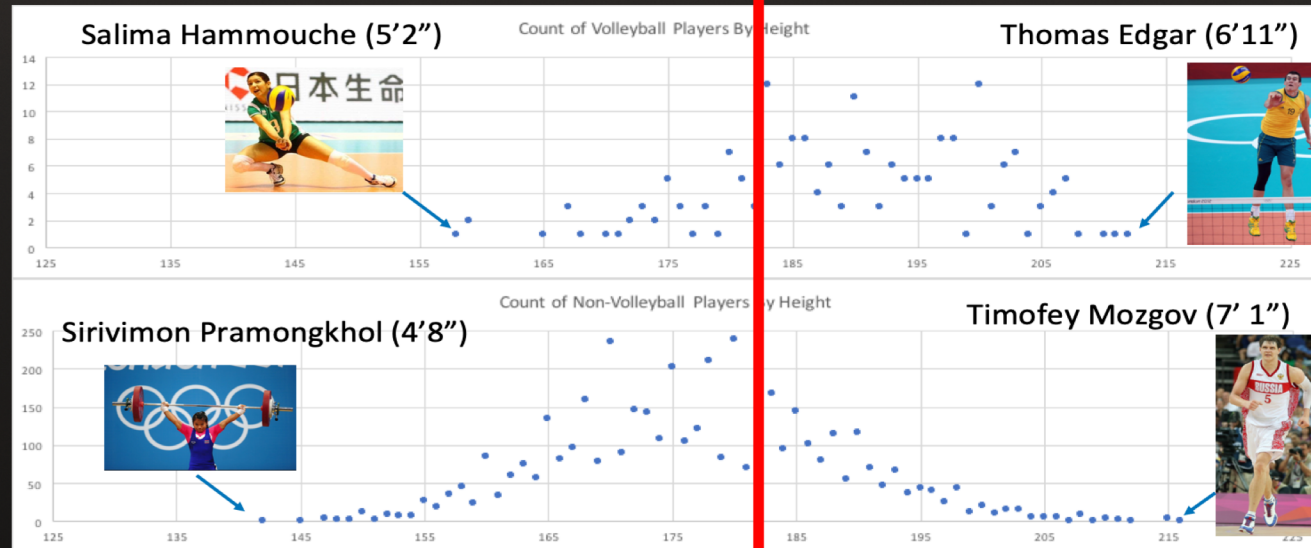
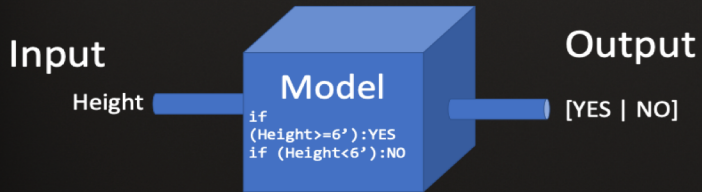
“Volleyball Players are Tall Model”

What if we say – if ≥ 6 feet (183cm) then they are a Volleyball Player

29% Missed (False Negative)

6' 0" (183cm)

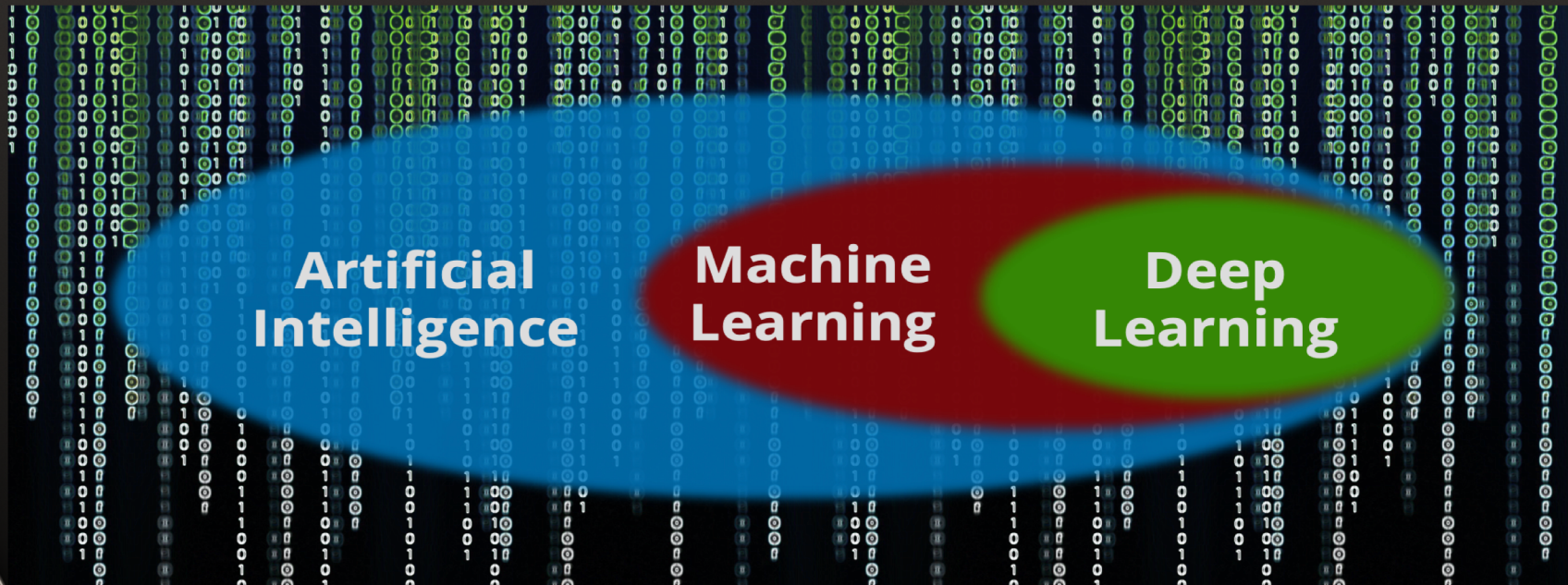
71% Correctly Identified (True Positive)



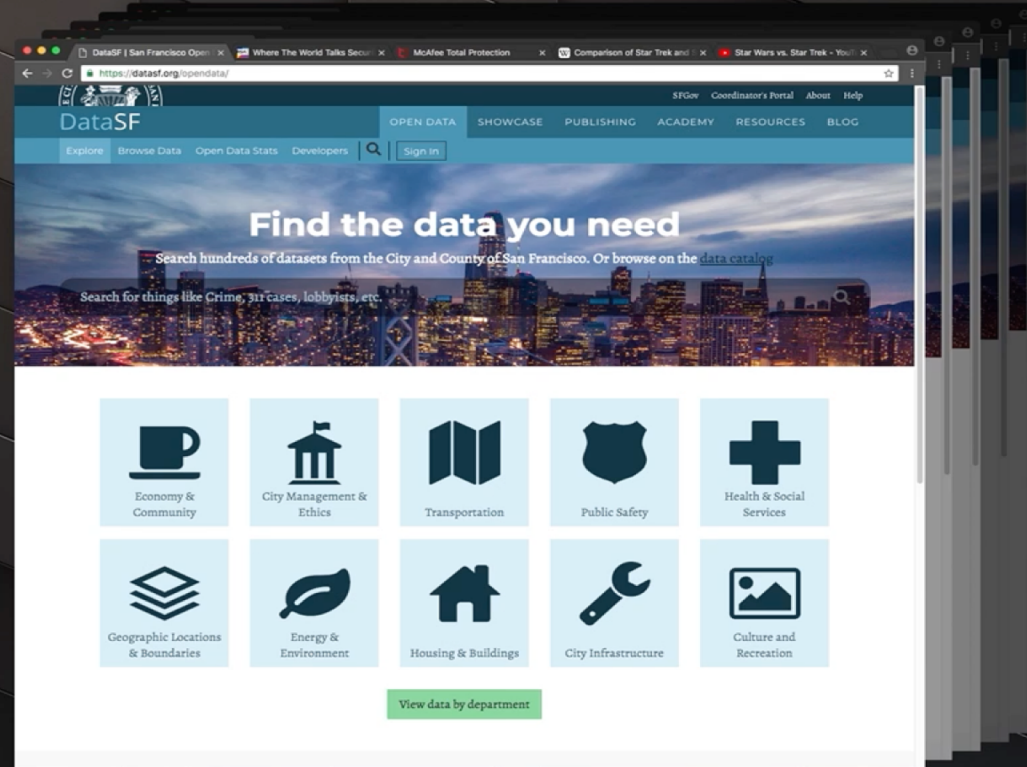
73% Correctly Identified as “Not Volleyball” (True Negative)

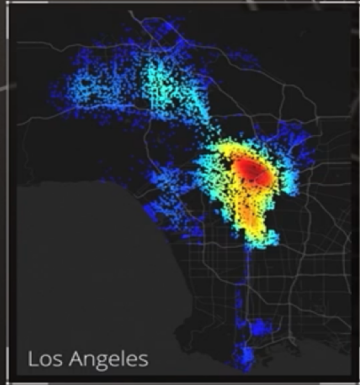
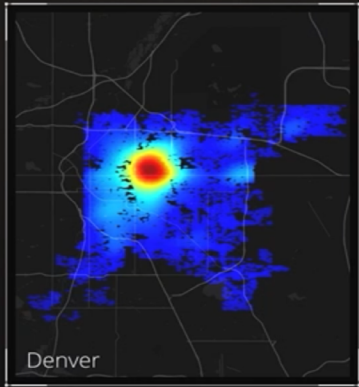
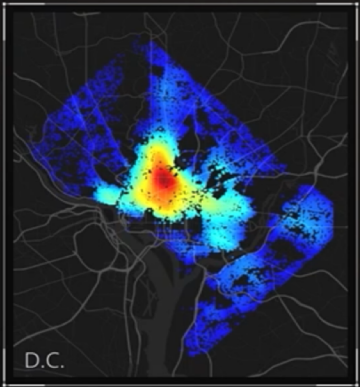
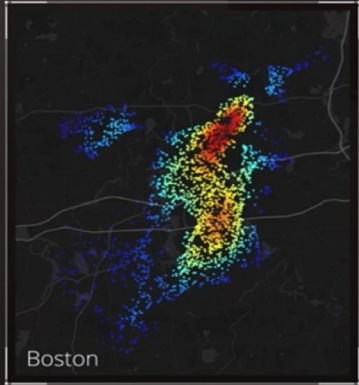
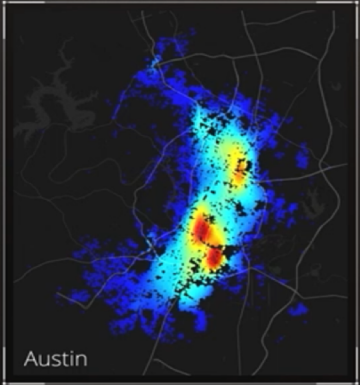
27% Wrongly Identified as “Volleyball” (False Positive)

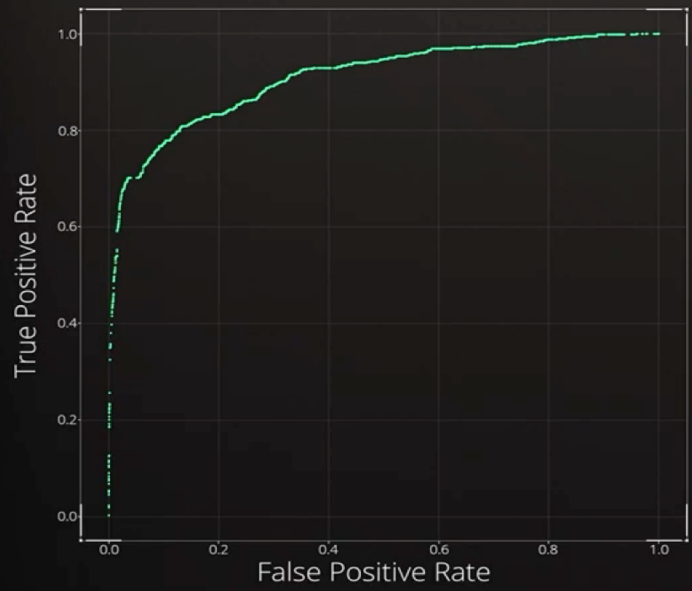
AI Categories



<https://datasf.org/opendata/>







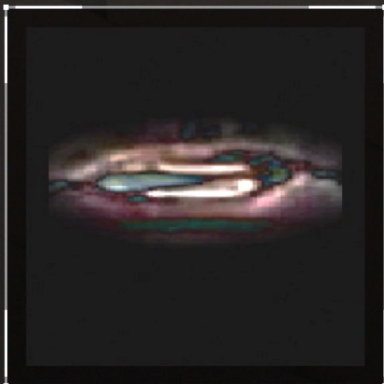
Original Source



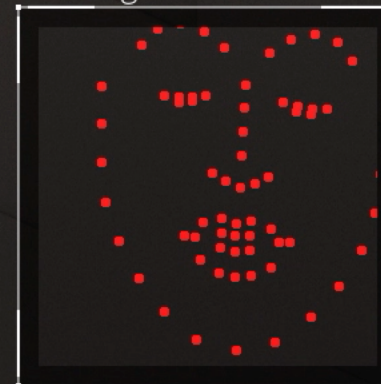
Original Mouth Capture



Mouth Calculations



Tracking Markers



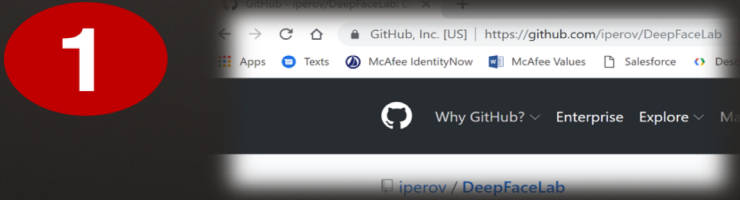
DeepFake Demo ©McAfee 2019

Rough Cut Render - Frame: **000001**

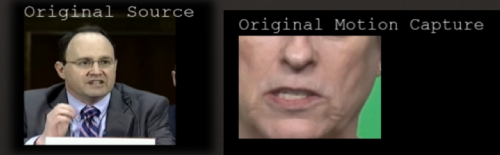
Source: Committee on Commerce, Science, and Transportation
United States Senate - One Hundred Fifteenth Congress
March 22, 2017

DeepFake "How-To"

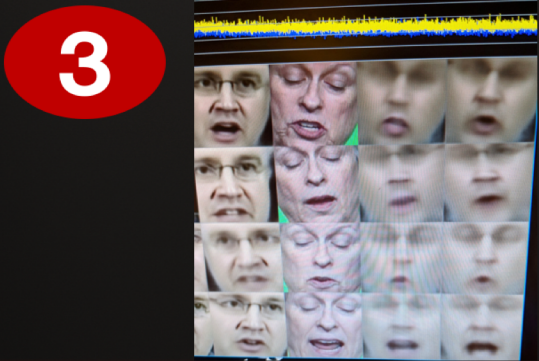
1 Download DeepFake Software



2 Find historic video and record fake dialog



3 Train machine learning model



4 Composite generated content onto historic video

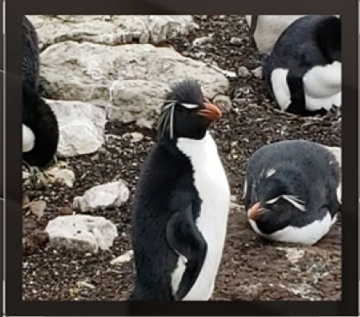


5 Render final DeepFake video



Adversarial Machine Learning

Rockhopper Penguin



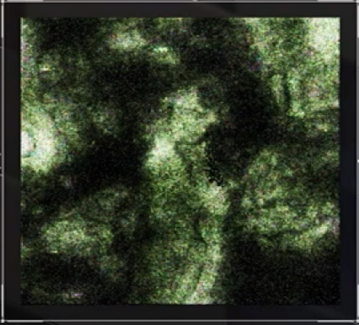
Accuracy: 0.9997957

Target Designation

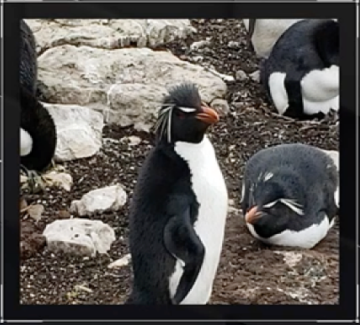
```
runfoolbox('penguin.jpg', 567)  
#567=Frying Pan
```

$$c \cdot \|X_{adv} - X\|_2^2 + loss_f(X_{adv}, L_{target})$$
$$X_{adv} \in [0, 1]^m$$

Adversarial
Machine Learning



Frying Pan



Accuracy: 0.9997547

Targeted Attack on Traffic Signs



After 20 iterations, the classifier achieved a 100% success rate of a 35 MPH speed sign for all 15 stop signs.

Physical Attack

Top Predictions	Confidence
1. stop	100.0%
2. yield	1.938594285100237e-08%



Before Attack

Top Predictions	Confidence
1. addedLane	81.44199848175049%
2. speedLimit25	9.326659888029099%



After Attack

Adversarial Example



Malware

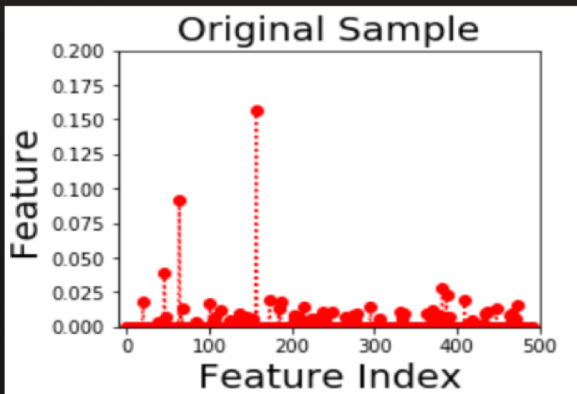
+

Perturbation

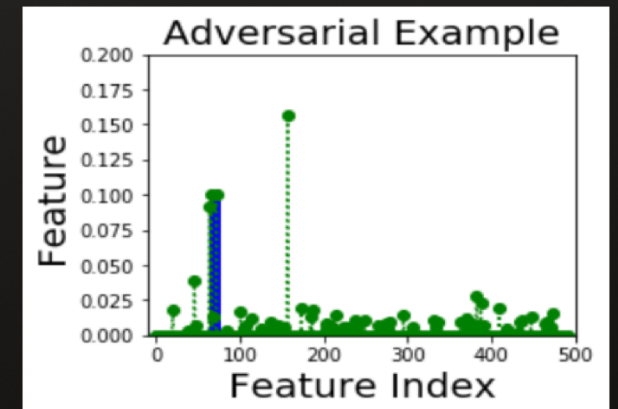
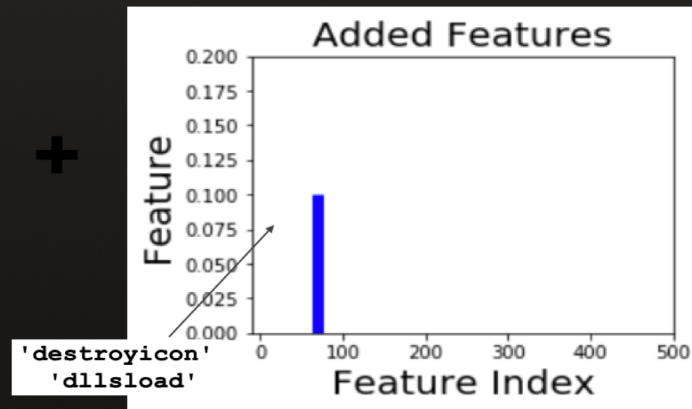
=



Benign

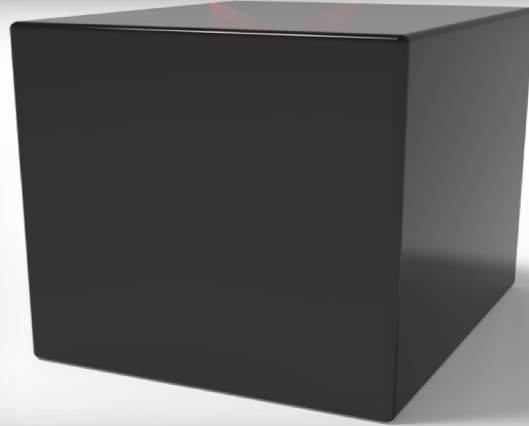


+



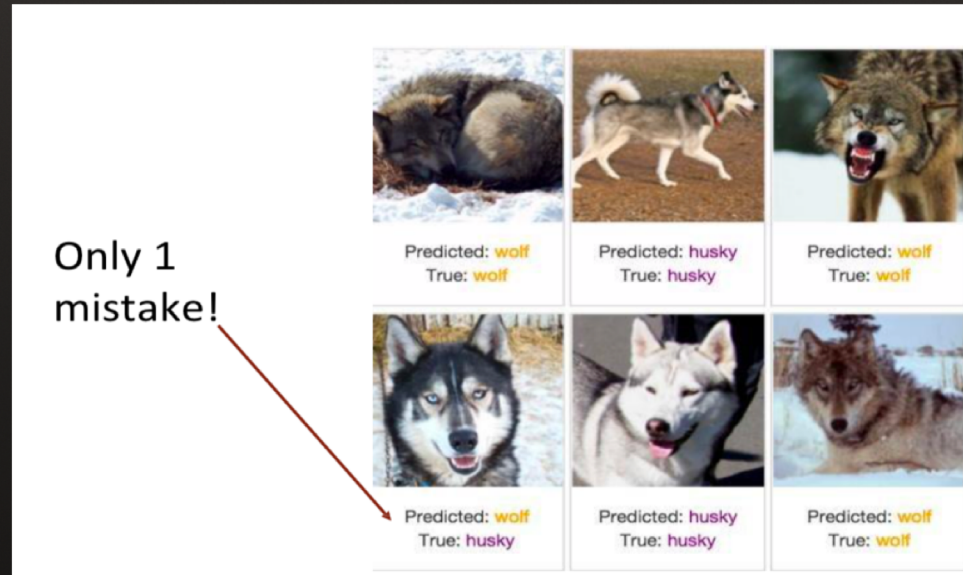
“Black Box Duality”

Solve Problems That We
Don't Understand



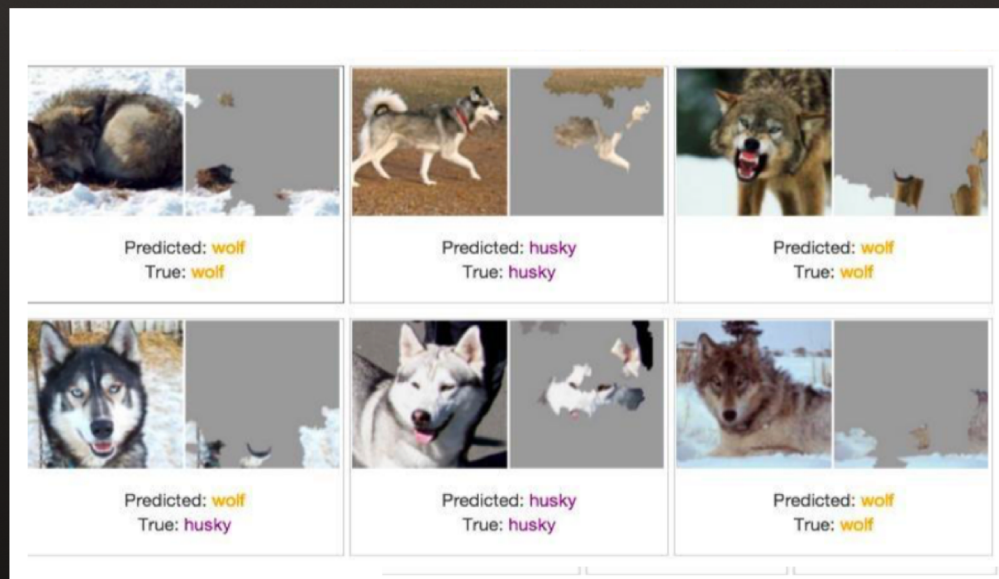
Often Unable to Explain
“Rationale for the
conclusion”

Can you build trust based on accuracy?



Source: "Why Should I Trust You?": Explaining the Predictions of Any Classifier
Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin

Can you build trust based on accuracy?



Source: "Why Should I Trust You?": Explaining the Predictions of Any Classifier
Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin

Artificial ≠ Intelligence

Thank You!

Questions??