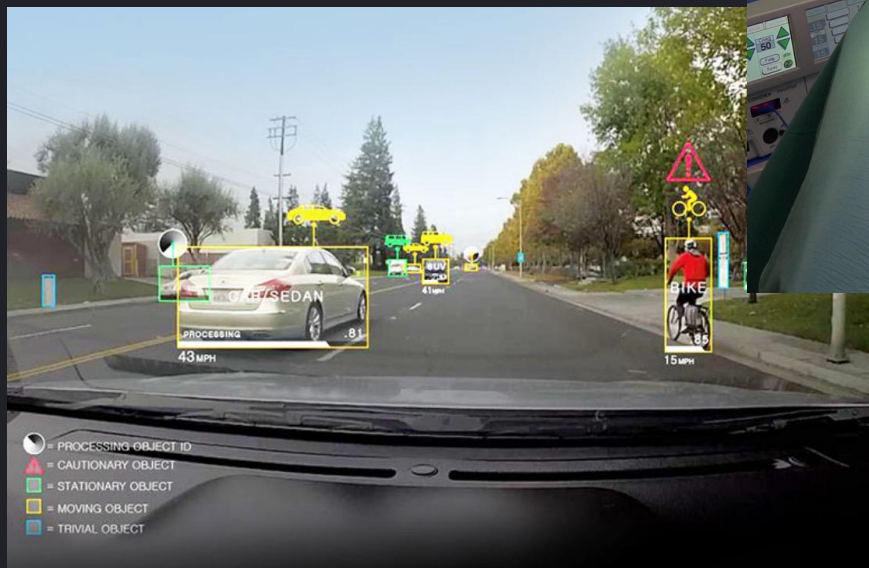


Poking holes in your deep learning vision model



Irina Vidal Migallón
2019.05.02 PyData Hamburg



Source: NVIDIA



Source: microsoftinsider.es, Hospital G. Marañón (Madrid, Spain) & Exovite

TESLA & UBER

2016: White truck on white

2018: Fire trucks

2018: Dismissed pedestrian

210,000,000 km driven by the
time of the accident
[[source](#)]



A BIT ABOUT ME

Applied research (INRIA)

Startups:

- MedTechs in Madrid, Paris & Berlin
- AR/MR in Berlin

Siemens Mobility



ROBUSTNESS?

ROBUSTNESS?



Evaluation



Debugging



Interpretability



Adversarial *





EVALUATION

EVALUATION

Model	Metric
Baseline	
New Model	

EVALUATION

Model	Metric
Baseline	
New Model	

Automated

Fast enough to iterate

Traced

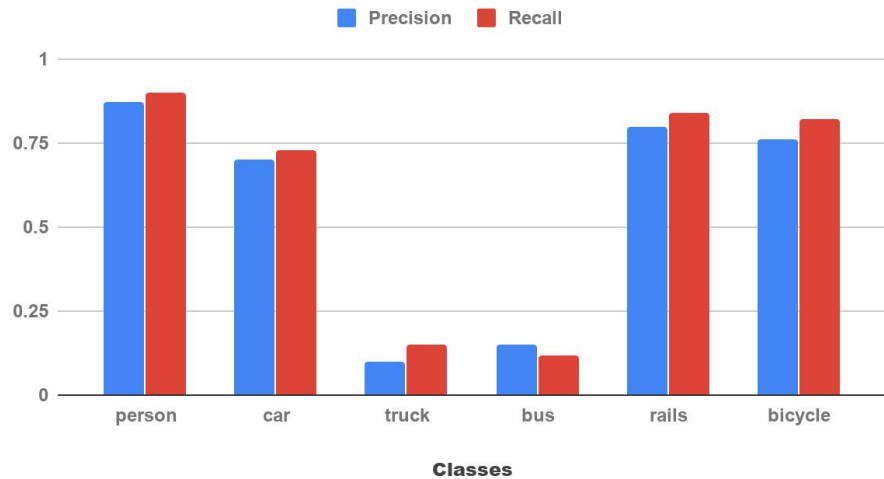
Close to production



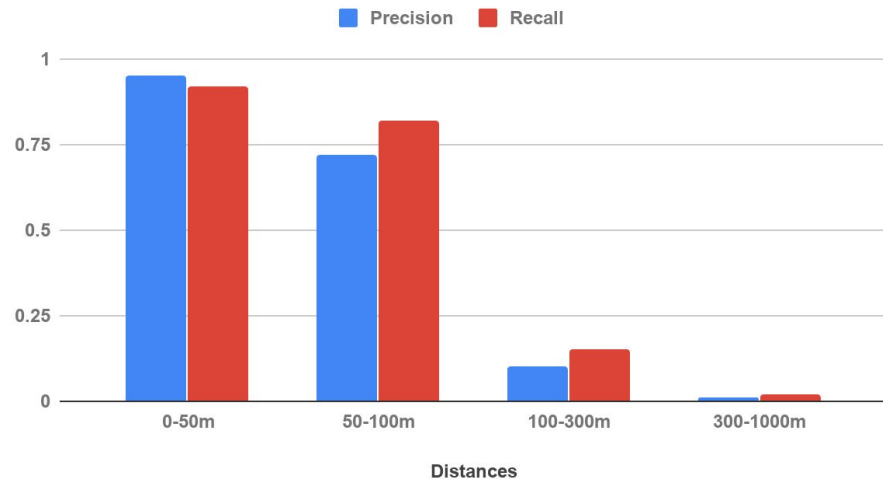
DEBUGGING

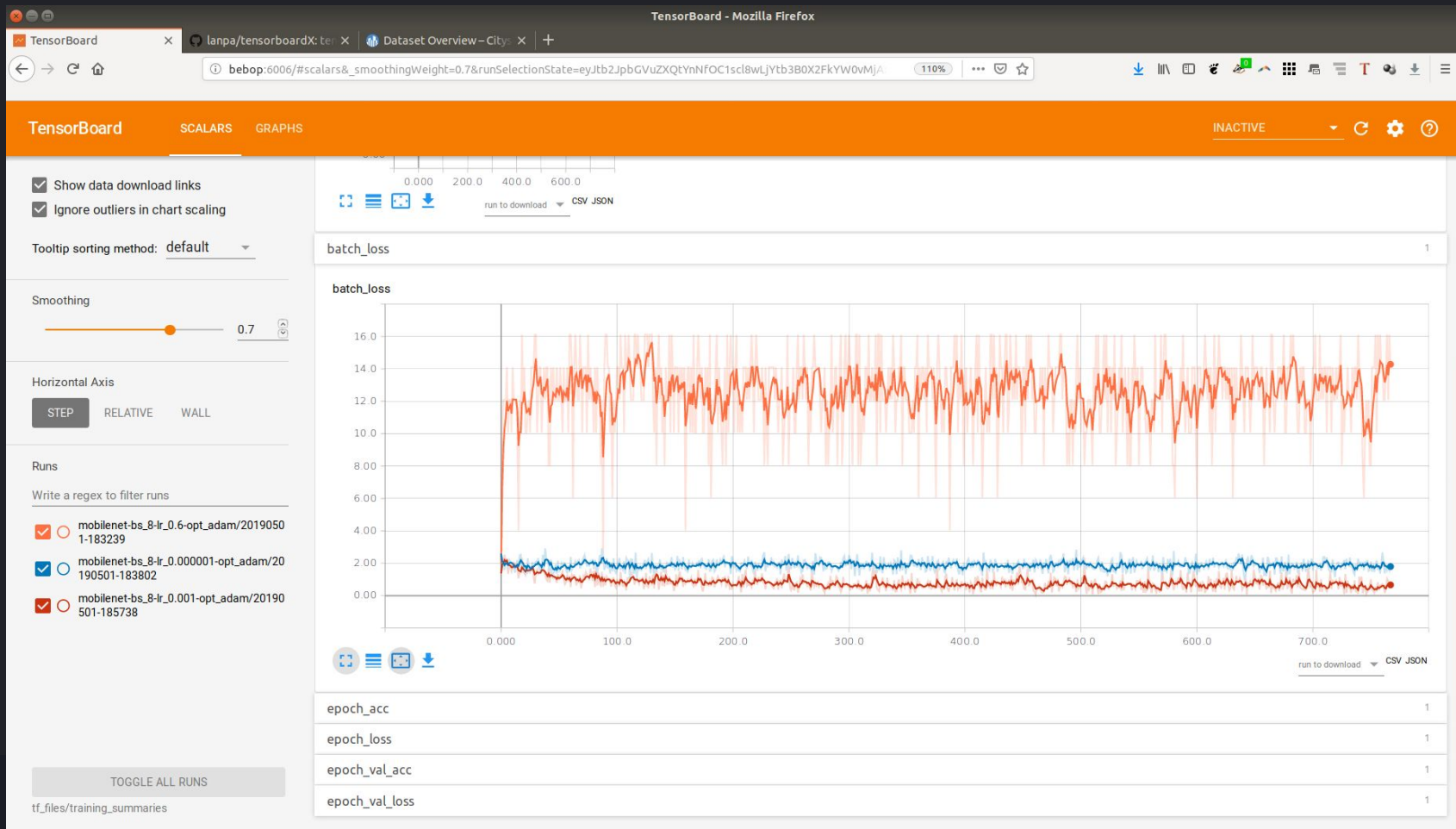
Dataset issues?

AP = 0.72



AP = 0.72





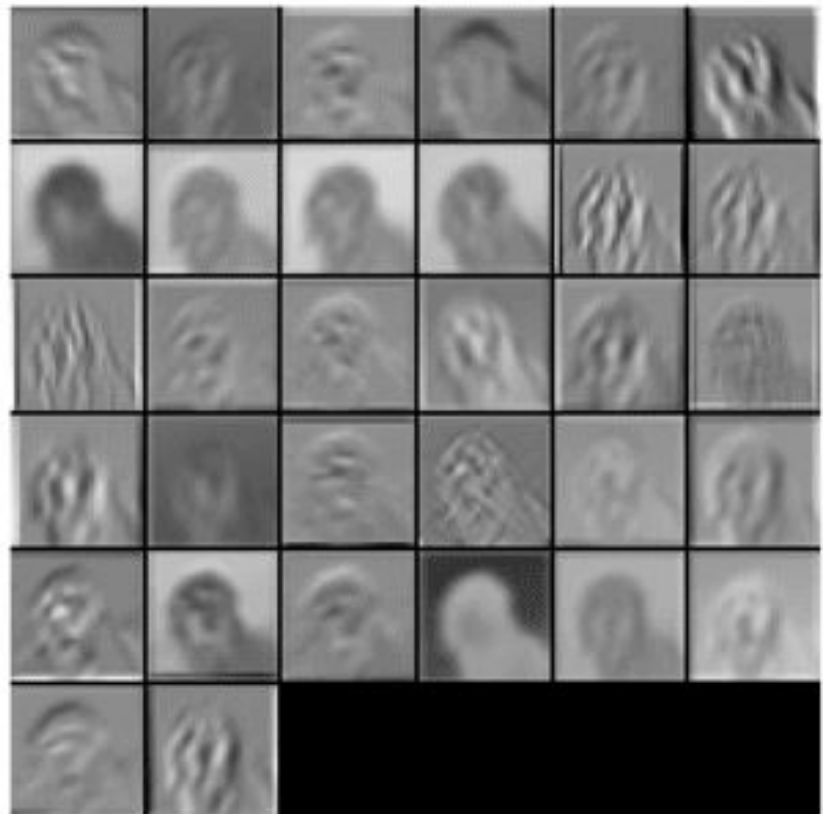
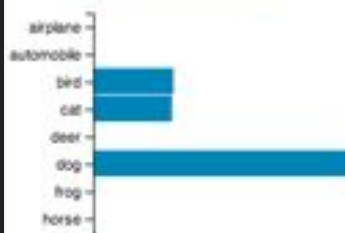
Selected filter: None

Training Images

dog



Selected image



<https://github.com/bruckner/deepViz>
<http://yosinski.com/deepvis>

<https://github.com/bruckner/deepViz>
<http://yosinski.com/deepvis>

Selected filter: None

Training Images ✕ Clear Selection

dog

Selected image

Generated images

airplane
automobile
bed
cat
deer
dog
frog
horse



INTERPRETABILITY

WHY CARE?

To...

Debug

Explain

Trust

Generalize

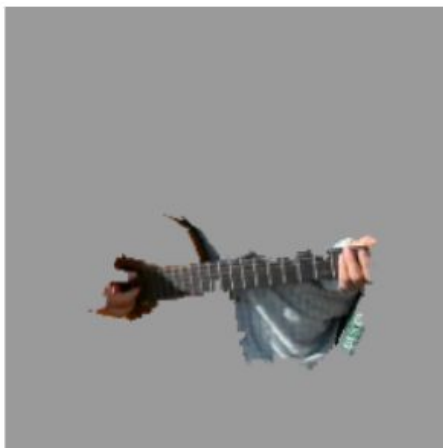
Optimize

To...

Avoid silent failure



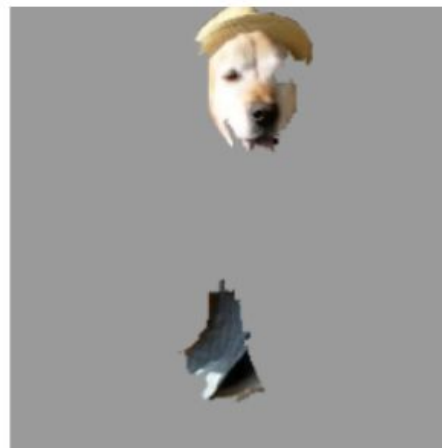
(a) Original Image



(b) Explaining *Electric guitar*



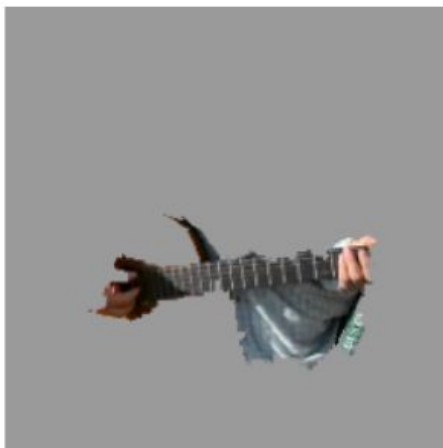
(c) Explaining *Acoustic guitar*



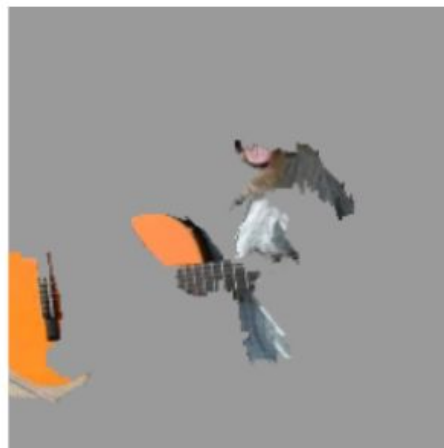
(d) Explaining *Labrador*



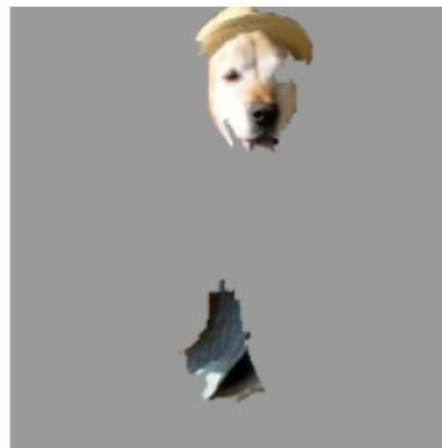
(a) Original Image



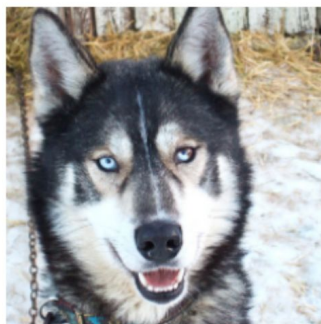
(b) Explaining *Electric guitar*



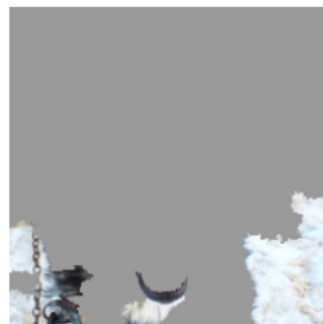
(c) Explaining *Acoustic guitar*



(d) Explaining *Labrador*



(a) Husky classified as wolf



(b) Explanation



ADVERSARIAL *

ADVERSARIAL SAMPLES 101

Collect your failures!



ADVERSARIAL SAMPLES 101





ADVERSARIAL SAMPLES



“panda”
57.7% confidence

+ .007 ×



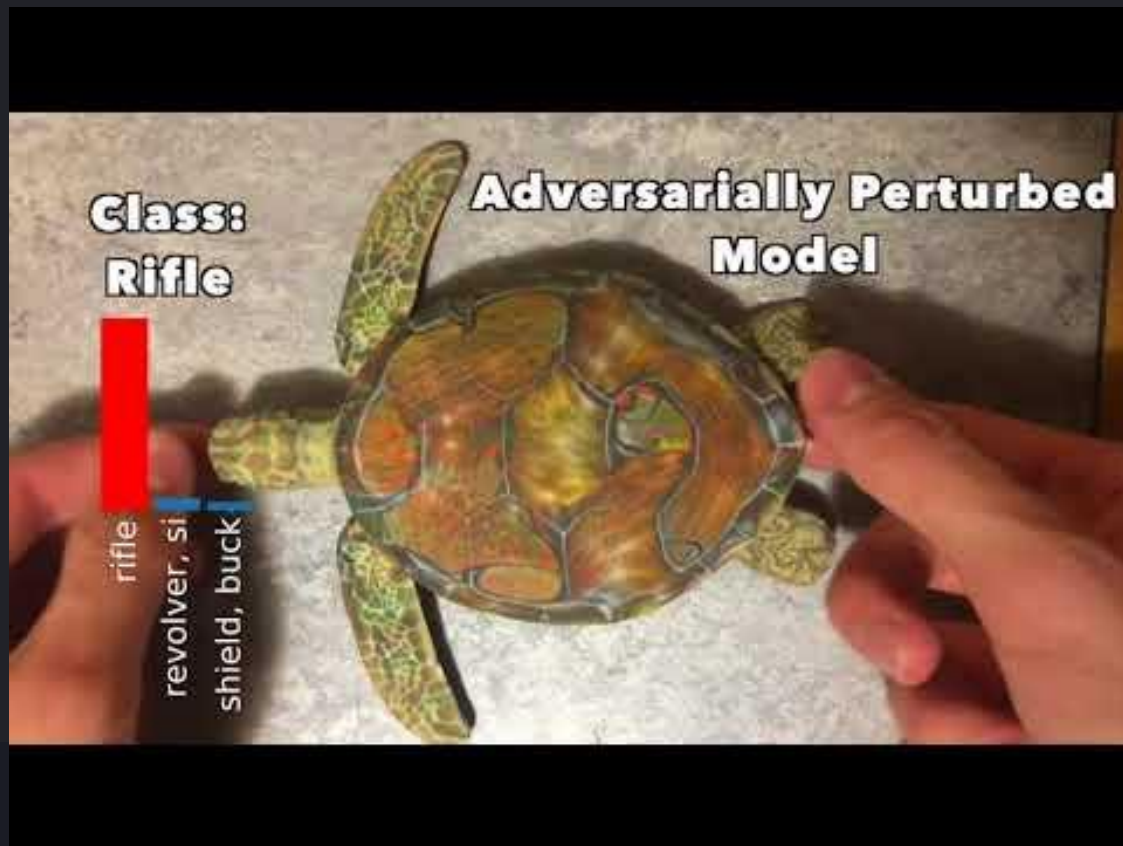
“nematode”
8.2% confidence

=



“gibbon”
99.3 % confidence

ADVERSARIAL SAMPLES



[Fooling Neural Networks in the Physical World with 3D Adversarial Objects](#)

ADVERSARIAL SAMPLES



"Robust Physical-World Attacks on Deep Learning Visual Classification", Evtimov et al. (CVPR 2018)

(a) original image



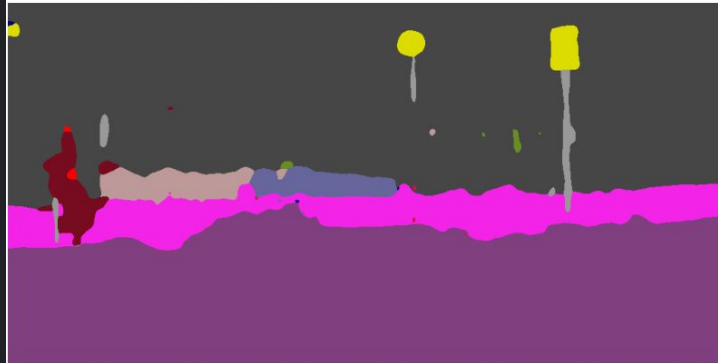
(b) adv. example



(i) prediction



(g) pred. on adv.





RECAP

ROBUSTNESS IS MANY THINGS

Evaluation

Debugging

Interpretability

Adversarial samples

ROBUSTNESS IS MANY THINGS

Evaluation

Debugging

Interpretability

Adversarial samples

Which ones are you already
using?

Thank you!

[linkedin.com/in/irinavidal/](https://www.linkedin.com/in/irinavidal/)

Thank you!

linkedin.com/in/irinavidal/