

---

# The Odds are Odd:

## A Statistical Test for Detecting Adversarial Examples

---

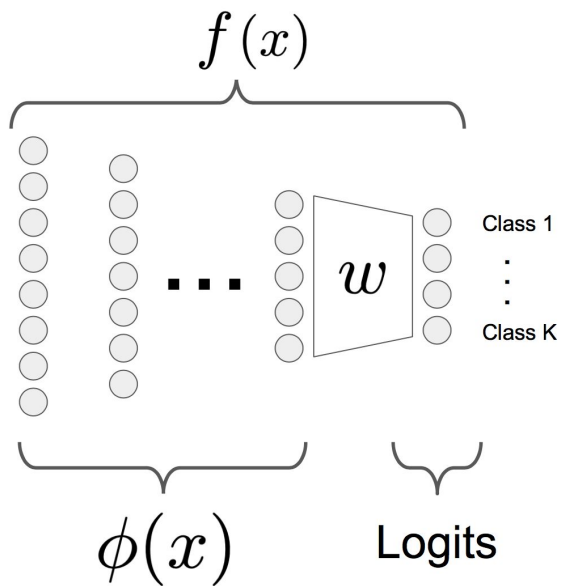
Kevin Roth\*, **Yannic Kilcher\***, Thomas Hofmann

**ETH** Zürich

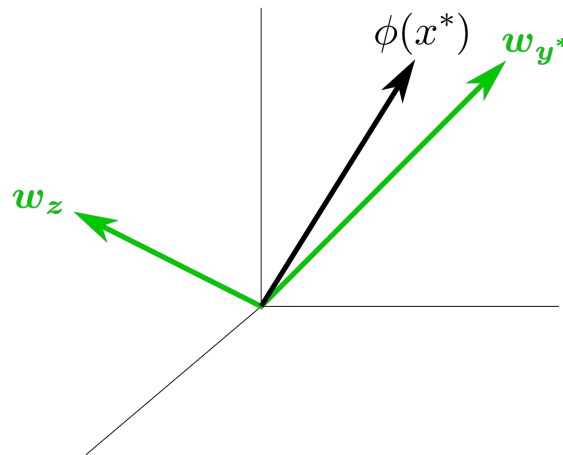


poster #62

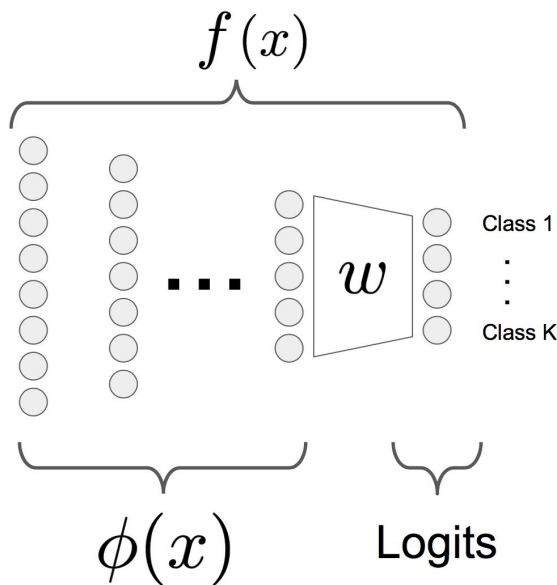
# Log-Odds & Adversarial Examples



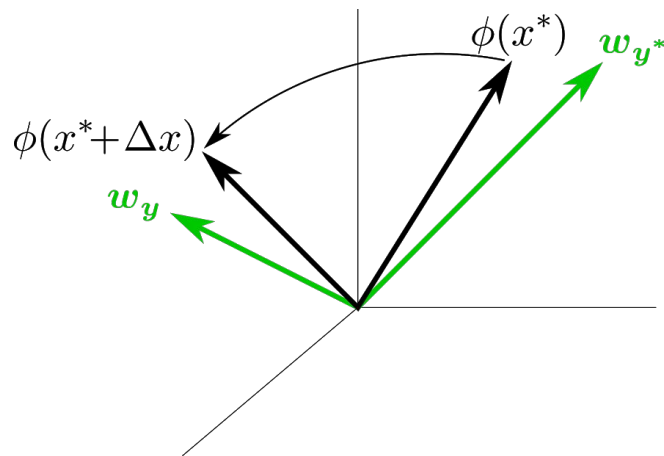
$$f_y(x) = \langle w_y, \phi(x) \rangle$$



# Log-Odds & Adversarial Examples



$$f_y(x) = \langle w_y, \phi(x) \rangle$$

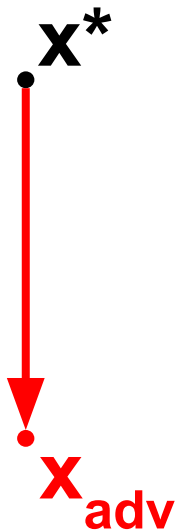


Adversarial examples cause atypically large feature space **perturbations along the weight-difference** direction

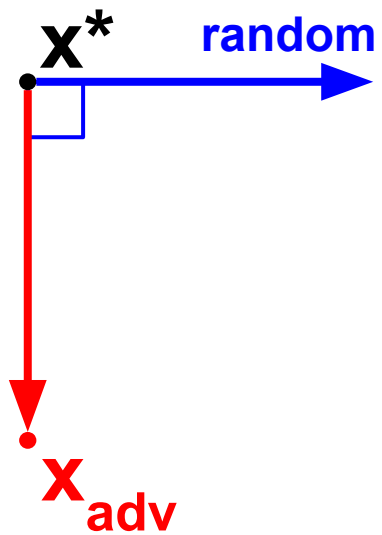
# Adversarial Cone

$\mathbf{x}^*$

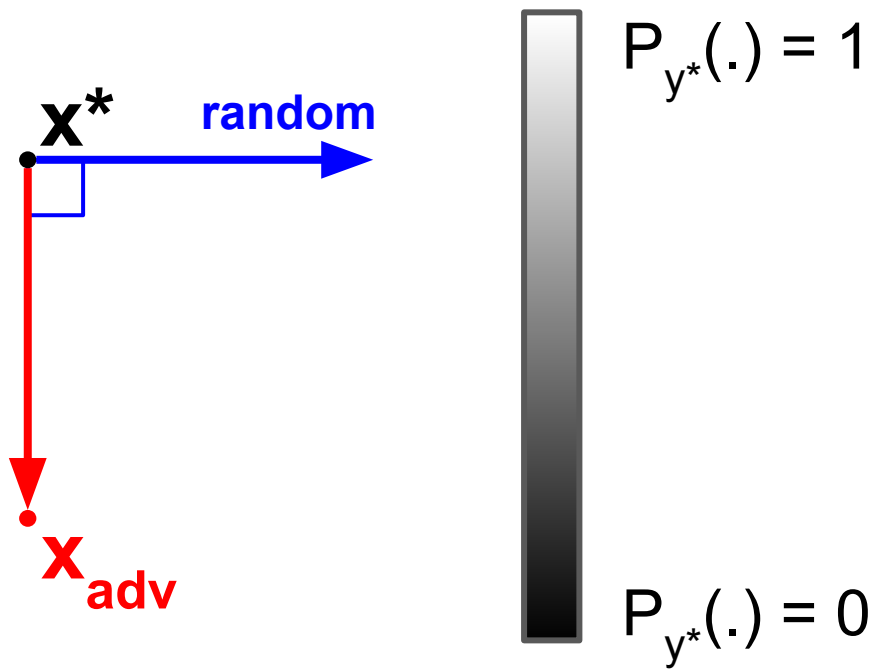
# Adversarial Cone



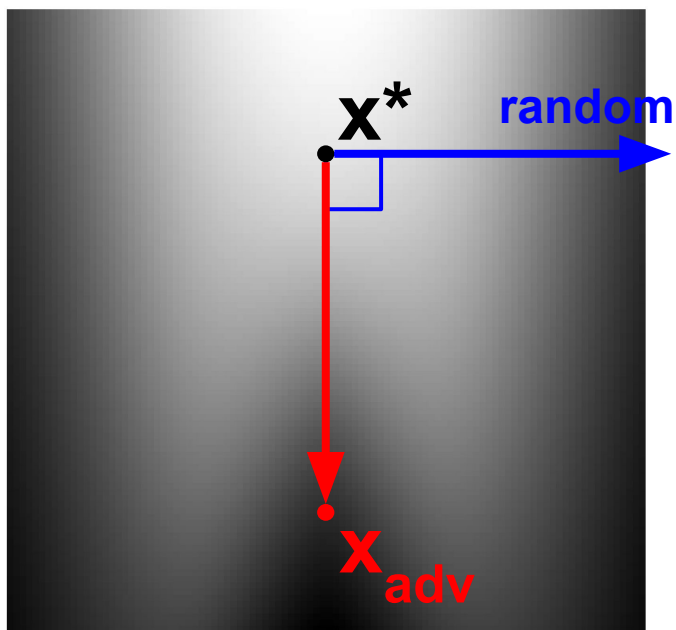
# Adversarial Cone



# Adversarial Cone

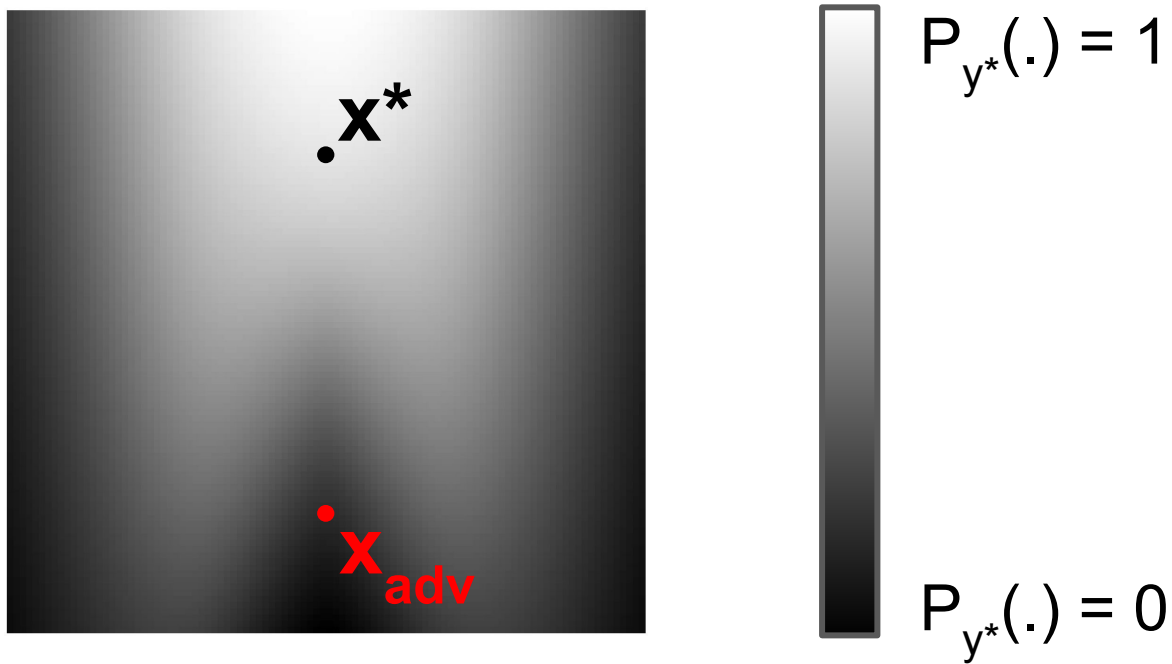


# Adversarial Cone



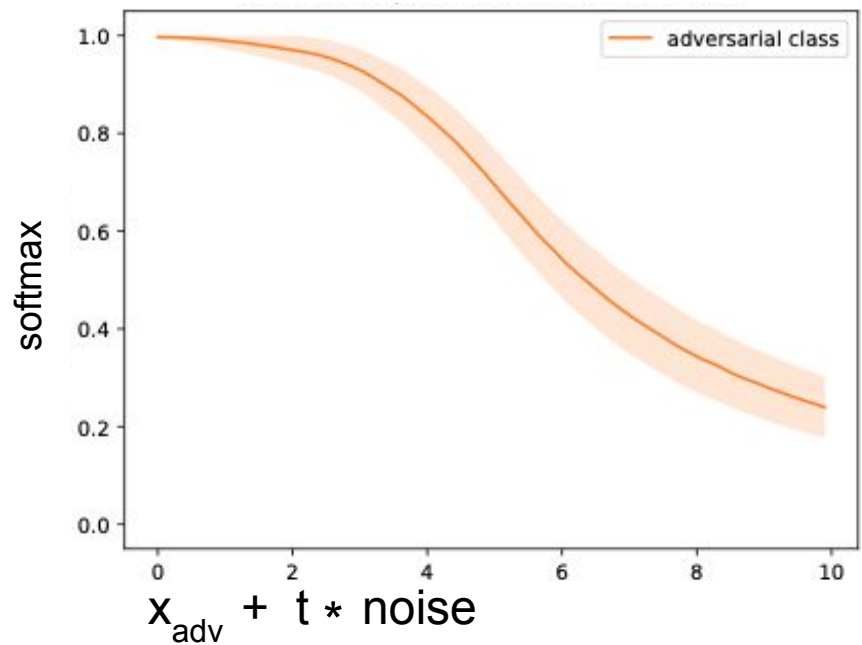


# Adversarial Cone

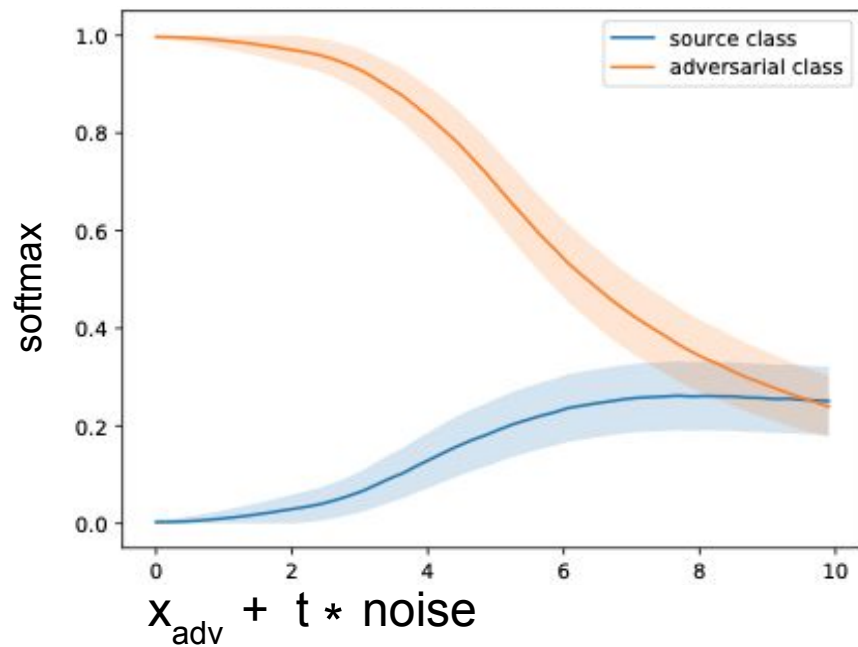


Adversarial examples are embedded in a cone-like structure

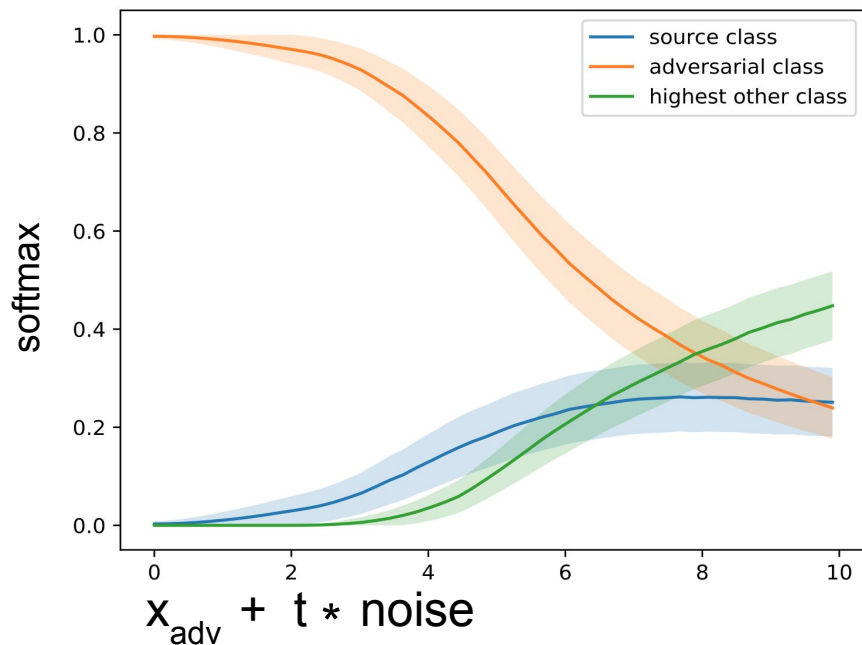
# Adversarial Cone



# Adversarial Cone

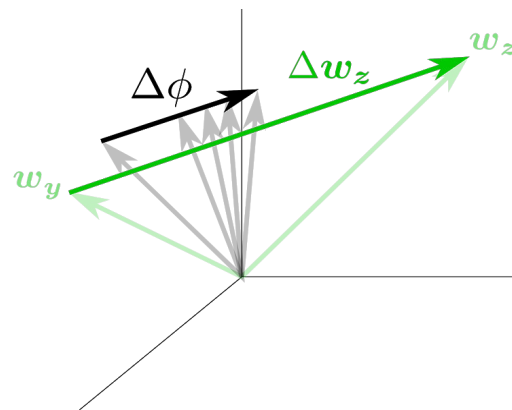
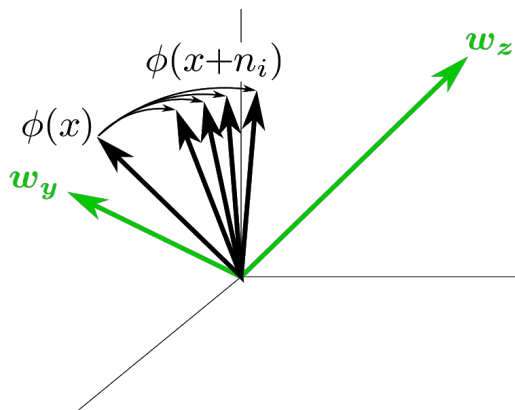


# Adversarial Cone



**Noise as a *probing instrument***

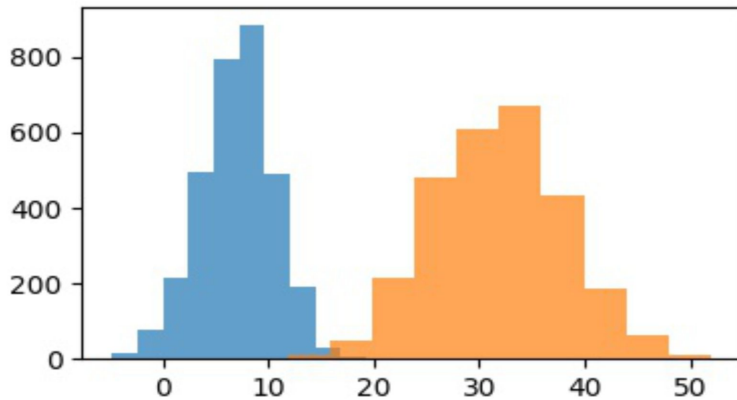
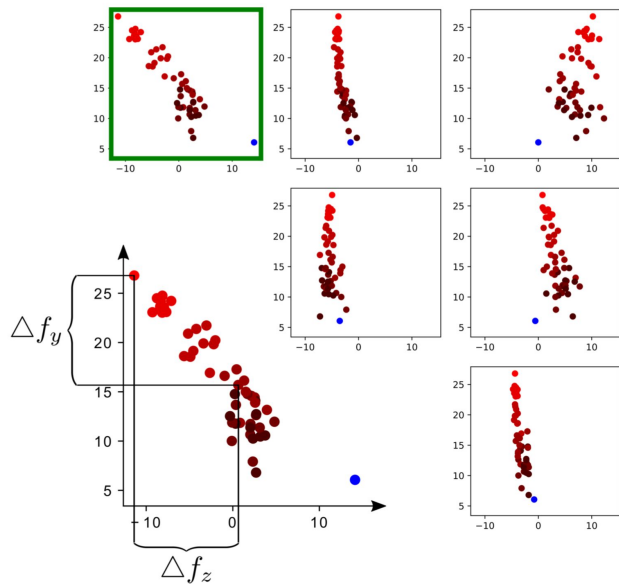
# Main Idea: Log-Odds Robustness



⇒ The robustness properties of  $\phi(x+n_i)$  are different dependent on whether  $x = x^*$  or  $x = x^* + \Delta x$

⇒  $\Delta\phi$  tends to have a **characteristic direction** if  $x = x^* + \Delta x$  whereas it tends not to have a specific direction if  $x = x^*$

# Main Idea: Log-Odds Robustness



natural adversarial



**Noise can partially undo effect of adversarial perturbation and directionally revert log-odds towards the true class  $y^*$**

# Statistical Test & Corrected Classification

⇒ We propose to use **noise-perturbed pairwise log-odds**

$$g_{y,z}(x, \eta) = \langle w_z - w_y, \phi(x + \eta) - \phi(x) \rangle$$

**to test whether  $x$  classified as  $y$  should be thought of as a manipulated example of true class  $z$  :**

$$x \text{ adversarial if } \max_{z \neq y} \{ \mathbf{E}_\eta [\bar{g}_{y,z}(x, \eta)] - \tau_{y,z} \} \geq 0$$

⇒ **Corrected classification :**  $G(x) = \arg \max_z \{ \bar{g}_{y,z}(x) - \tau_{y,z} \}$

# Detection Rates & Corrected Classification

Table 1: CIFAR10

Model	Detection rate (clean / pgd)	Corrected Accuracy (clean / pgd)
WRResNet	<b>0.2%</b> / <b>99.1%</b>	<b>96.0%</b> / <b>92.7%</b>
CNN7	<b>0.8%</b> / <b>95.0%</b>	<b>93.6%</b> / <b>89.5%</b>
CNN4	<b>1.4%</b> / <b>93.8%</b>	<b>71.0%</b> / <b>67.6%</b>

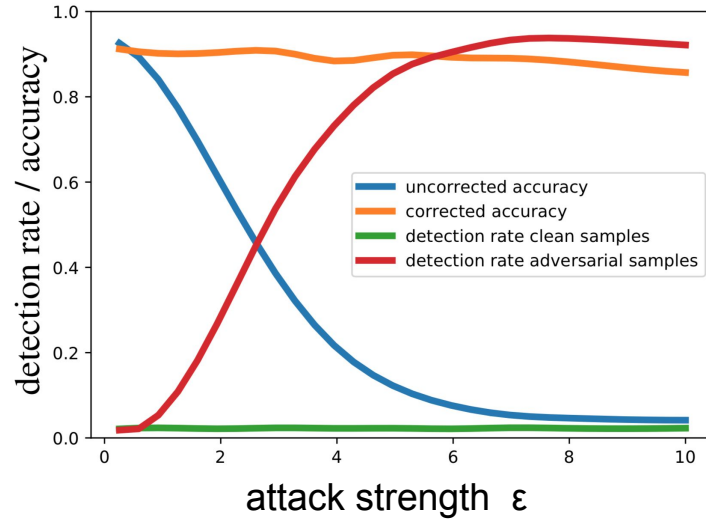
Table 2: ImageNet

Model	Detection rate (clean / pgd)
Inception V3	<b>1.9%</b> / <b>99.6%</b>
ResNet 101	<b>0.8%</b> / <b>99.8%</b>
VGG16(+BN)	<b>0.3%</b> / <b>99.9%</b>

- ⇒ Our statistical test **detects nearly all adversarial examples** with FPR ~1%
- ⇒ Our correction method **reclassifies almost all adversarial examples** successfully
- ⇒ **Drop in performance on clean samples is negligible**



# Detection Rates & Corrected Classification



**Detection rate increases with increasing attack strength**



**Corrected classification manages to compensate for decay in uncorrected accuracy due to increase in attack strength**

# Defending against Defense-Aware Attacks

Model	Detection rate (clean / attack)	Accuracy (clean / attack)
WResNet	<b>4.5% / 71.4%</b>	<b>91.7% / 56.0%</b>
CNN7	<b>2.8% / 75.5%</b>	<b>91.2% / 56.6%</b>
CNN4	<b>4.1% / 81.3%</b>	<b>69.0% / 56.5%</b>



Attacker has **full knowledge of the defense** :

perturbations that work in expectation under noise source used for detection



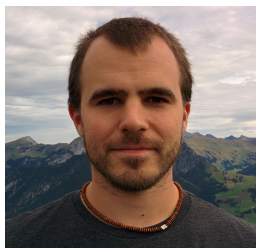
**Detection rates and corrected accuracies remain remarkably high**

# Thank You

poster #62



**Kevin Roth**



**Yannic Kilcher**



**Thomas Hofmann**



data analytics lab

**ETH** zürich

Follow-Up Work: *Adversarial Training Generalizes*  
*Data-dependent Spectral Norm Regularization*

ICML Workshop on  
Generalization (June 14)



# References

The approaches most related to our work are those that detect whether or not the input has been perturbed, either by detecting characteristic regularities in the adversarial perturbations themselves or in the network activations they induce.

- Grosse, Kathrin, et al. "On the (statistical) detection of adversarial examples." (2017).
- Metzen, Jan Hendrik, et al. "On detecting adversarial perturbations." (2017).
- Feinman, Reuben, et al. "Detecting adversarial samples from artifacts." (2017).
- Xu, Weilin, David Evans, and Yanjun Qi. "Feature squeezing: Detecting adversarial examples in deep neural networks." (2017).
- Song, Yang, et al. "Pixeldefend: Leveraging generative models to understand and defend against adversarial examples." (2017).
- Carlini, Nicholas, and David Wagner. "Adversarial examples are not easily detected: Bypassing ten detection methods." (2017).
- ... and many more