

# Adversarial Training and Security in ML

Lili Mou, Priyank Jaini  
[doublepower.mou@gmail.com](mailto:doublepower.mou@gmail.com)

David R. Cheriton School of Computer Science  
University of Waterloo

UROC 22 Sep 2017

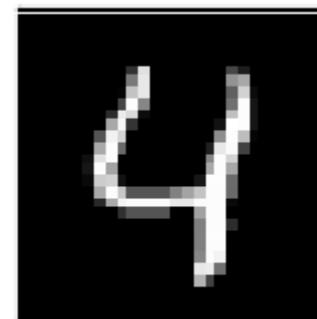
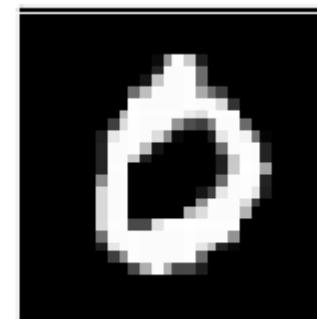
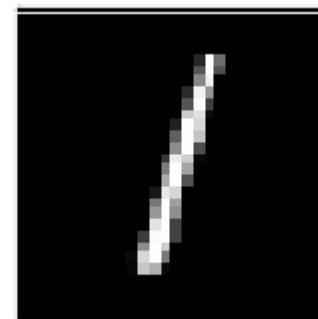


# Agenda

- **Background of neural networks**
  - Miniproject: CNN and its visualization
- Adversarial samples
  - Miniproject: Crafting adversarial data
- Open research

# Philosophy of Deep Learning

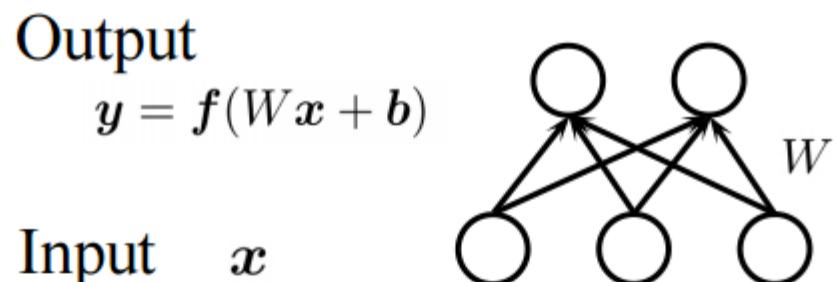
- Consider hand-written digit recognition



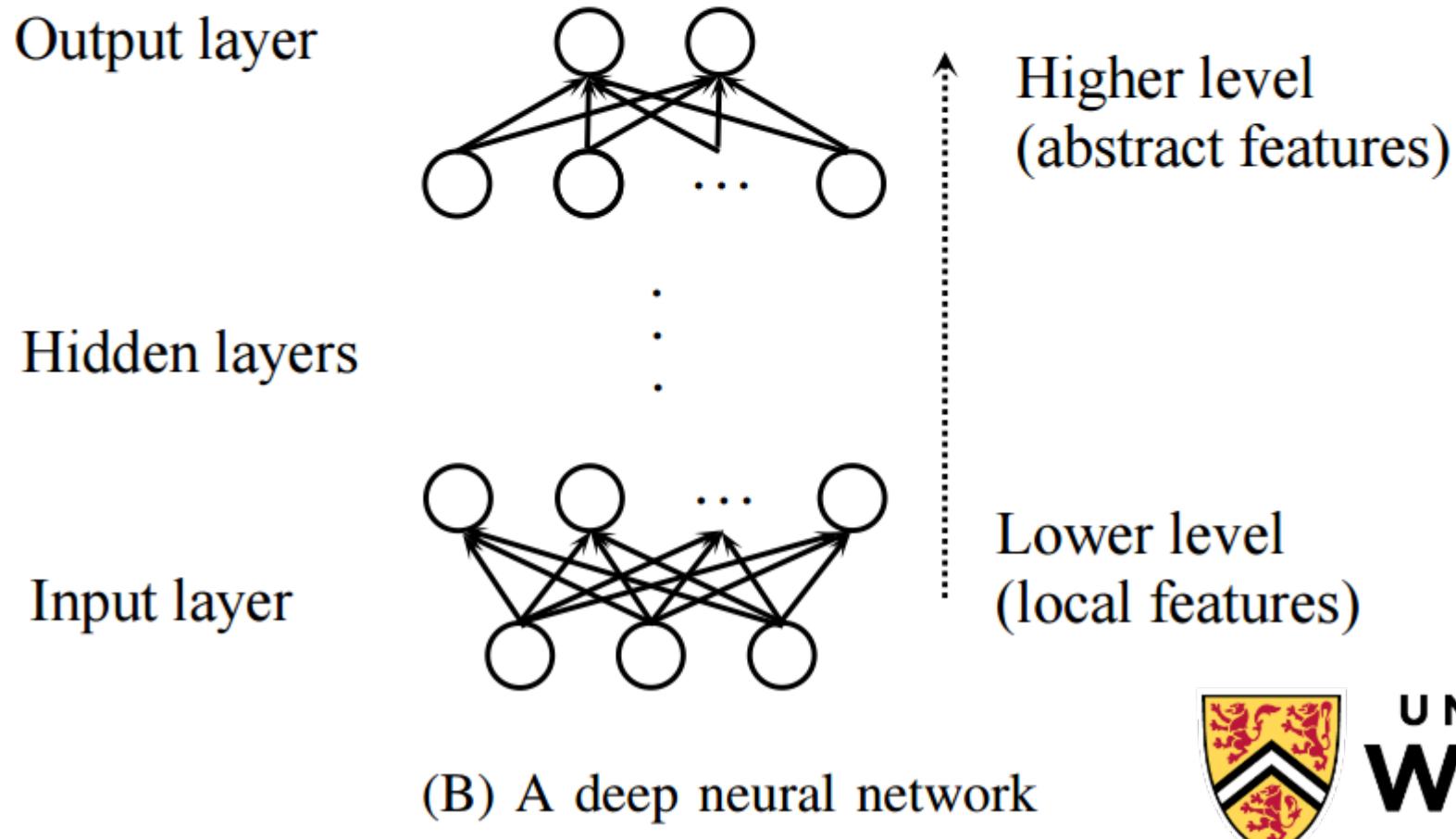
- Deep learning: End-to-end training
  - Input: raw signal (28 x 28 pixels)
  - Output: the target labels "7," "2," "1," "0," and "4."



UNIVERSITY OF  
**WATERLOO**

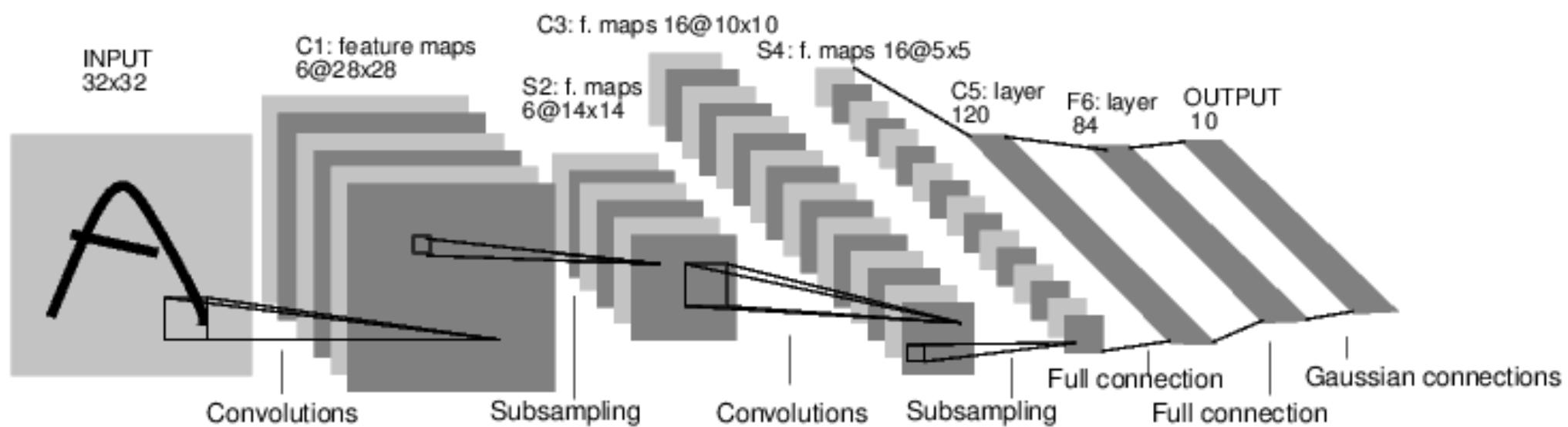


(A) A single layer of neurons



UNIVERSITY OF  
**WATERLOO**

# A Convolutional Neural Network



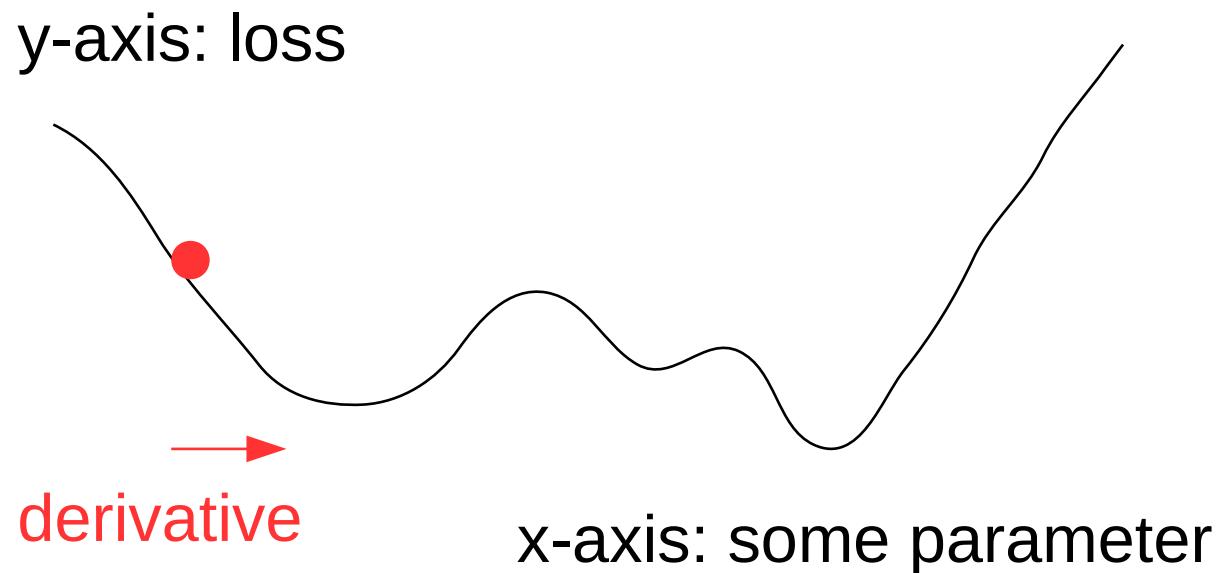
LeNet5



UNIVERSITY OF  
**WATERLOO**

# Training

- How do we learn weights?
  - Backpropagation
  - Compute the partial derivative of a "loss" w.r.t. each parameter
  - Take a small step towards the derivative



UNIVERSITY OF  
**WATERLOO**

# What are these features?

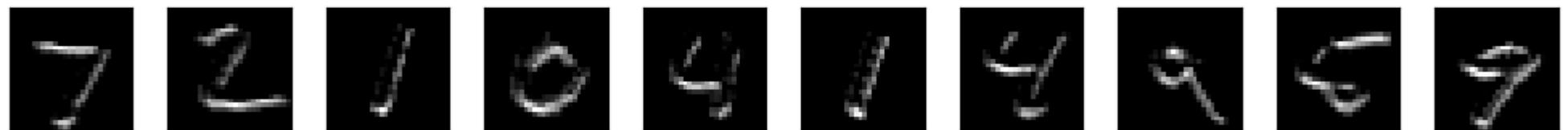
- Visualizing weights



Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *NIPS*. 2012.

# What are these features?

- Visualizing the activation functions



UNIVERSITY OF  
**WATERLOO**

# Mini-Project

- Code

<https://www.dropbox.com/s/iafhbi87mtk67gk/Adversarial.zip?dl=0>

- Cached data

[https://www.dropbox.com/s/m2qn92q5b8ky4mo/adv\\_data.zip?dl=0](https://www.dropbox.com/s/m2qn92q5b8ky4mo/adv_data.zip?dl=0)

- Slides available at

<http://sei.pku.edu.cn/~moull12>



UNIVERSITY OF  
**WATERLOO**

# Agenda

- Background of neural networks
  - Miniproject: CNN and its visualization
- **Adversarial samples**
  - Miniproject: Crafting adversarial data
- Open research



UNIVERSITY OF  
**WATERLOO**

# Think of the Training Process

- Loss:  $L = f(x; w)$
- Training objective: minimize  $L$
- Compute:  $\nabla_w f(x; w)$

# Think of the Training Process

- Loss:  $L = f(x; w)$
- Training objective: minimize  $L$
- Compute:  $\nabla_w f(x; w)$
- What if we compute:  $\nabla_{\textcolor{red}{x}} f(x; w)$

# Adversarial Samples from Random

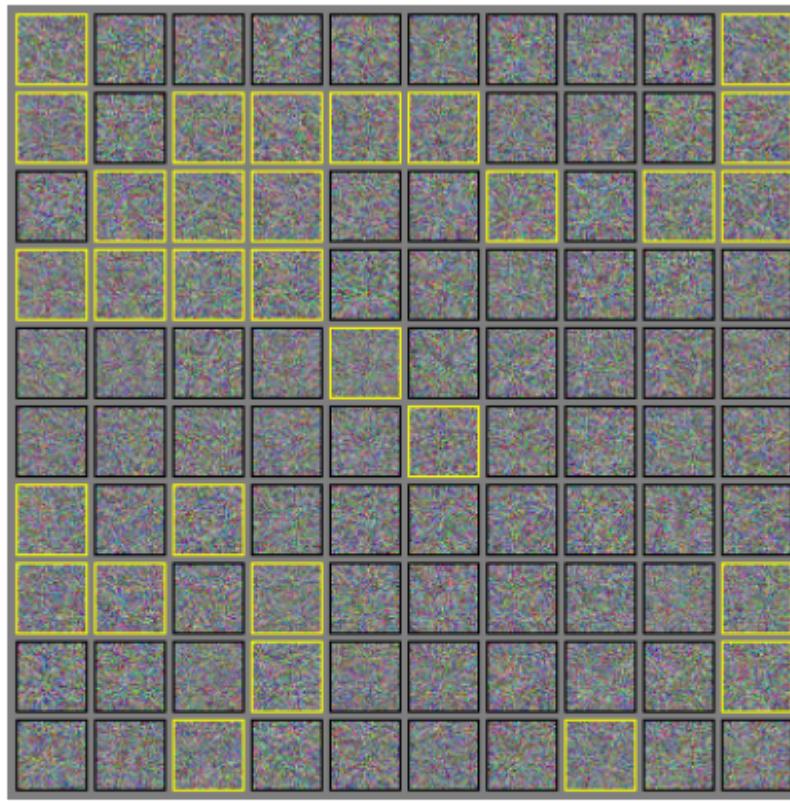


Figure 5: Randomly generated fooling images for a convolutional network trained on CIFAR-10. These examples were generated by drawing a sample from an isotropic Gaussian, then taking a gradient sign step in the direction that increases the probability of the “airplane” class. Yellow boxes indicate samples that successfully fool the model into believing an airplane is present with at least 50% confidence. “Airplane” is the hardest class to construct fooling images for on CIFAR-10, so this figure represents the worst case in terms of success rate.

Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." *ICLR*, 2015.



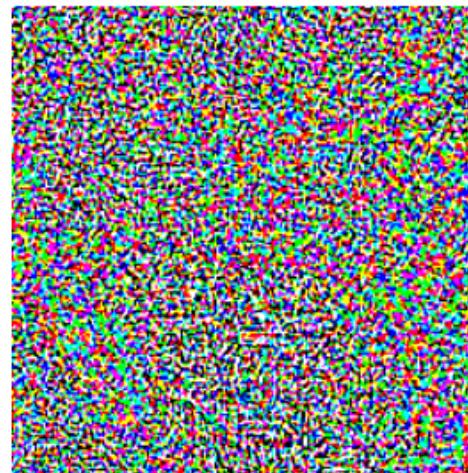
UNIVERSITY OF  
**WATERLOO**

# Adversarial Samples from Real Data



$x$   
“panda”  
57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$   
“nematode”  
8.2% confidence

$=$



$x +$   
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”  
99.3 % confidence



UNIVERSITY OF  
**WATERLOO**

# Approach

$$X_{\text{adv}} = X - \epsilon \cdot \text{sign}(\nabla_X J(X, y_{\text{target}}))$$

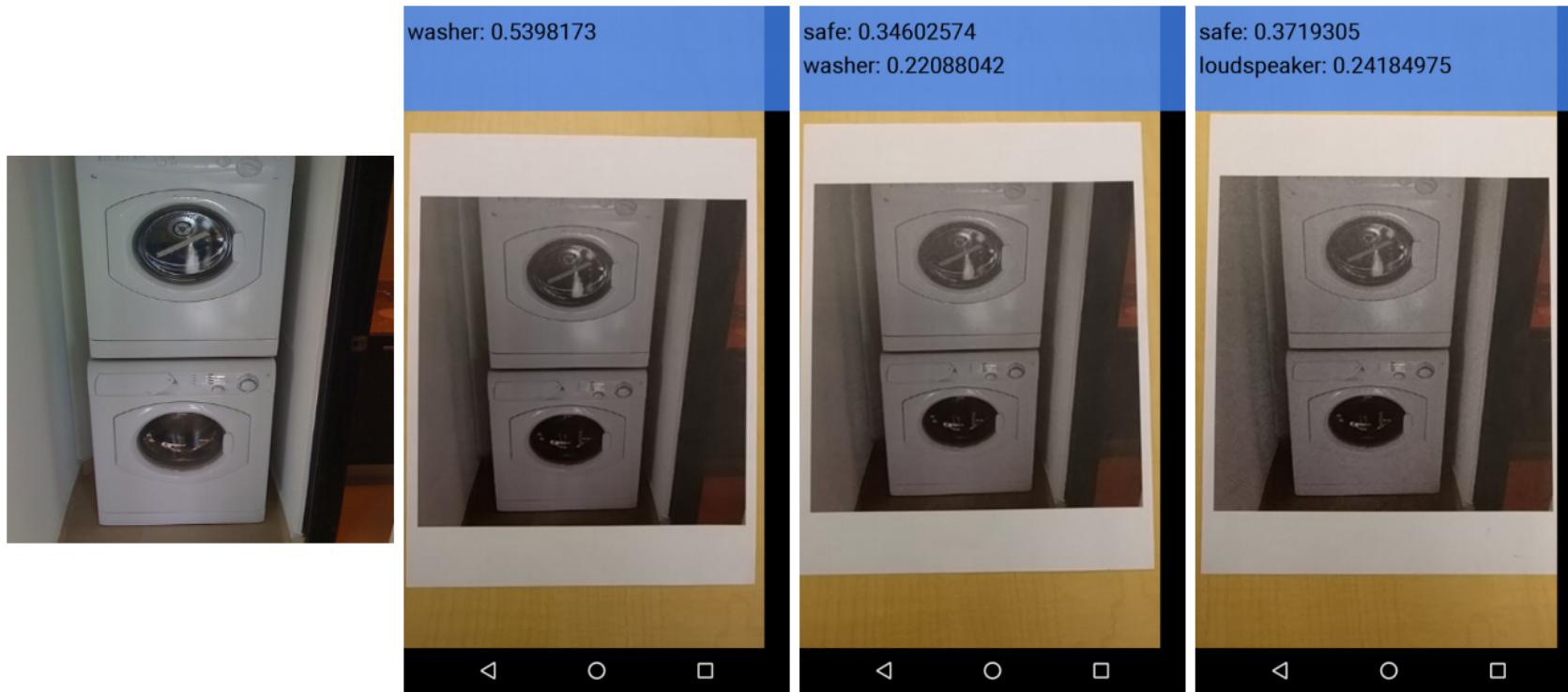
- $y_{\text{target}}$  : whatever target you want
- Take the sign of the partial derivative
  - Alternatively, we can truncate the gradient
  - So that the adversarial image is not too far away
- Can iterate several times if necessary



UNIVERSITY OF  
**WATERLOO**

# Ubiquity of Adversarial Samples

- A same adversarial sample works:
  - For different networks (models)
  - Even after further perturbation with noise



Kurakin, Alexey, Ian Goodfellow, and Samy Bengio.  
"Adversarial examples in the physical world." *ICLR*, 2017

# Mini-Project

- Generate Adversarial Samples by Yourselves
- Results



6 (0.93)



4 (0.98)



3 (0.99)

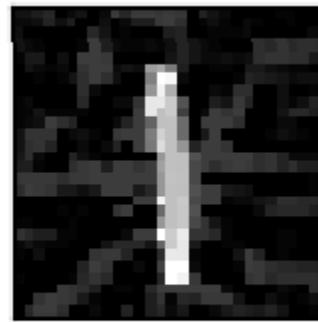


5 (0.98)



9 (1.00)

- Interpretation, predicted as "4" w.p. 98%



4 (0.98)



UNIVERSITY OF  
**WATERLOO**

# Agenda

- Background of neural networks
  - Miniproject: CNN and its visualization
- Adversarial samples
  - Miniproject: Crafting adversarial data
- Open research

# Open Topics

- Test the robustness of NNs
- Further confirm the ubiquity of adv samples
- Crafting more deceptive adversarial samples
- Training more robust machine learning models

# Thanks!

## QA