

---

---

# Generalization Comparison of Deep Neural Networks via Output Sensitivity

—— Mahsa Forouzesh, Farnood Salehi, Patrick Thiran ——

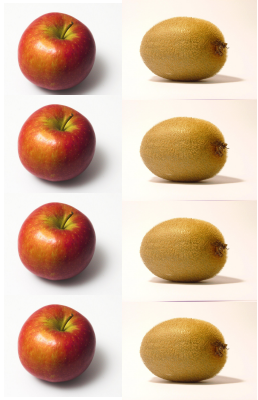
**EPFL**

---

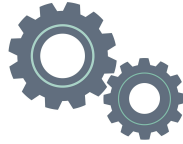
---

# Supervised Learning

Training Data



Machine Learning Model



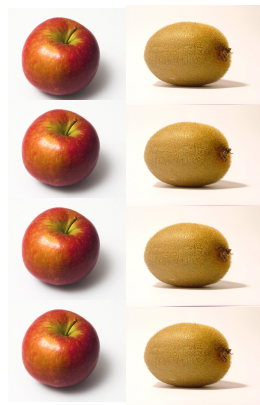
Class: Apple

Unseen Data



# Supervised Learning

Training Data



Machine Learning Model



Unseen Data

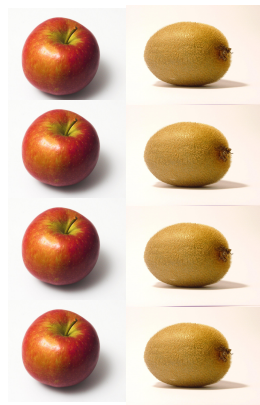


How Accurate?  $\Rightarrow$  Generalization Error

Class: Apple

# Supervised Learning

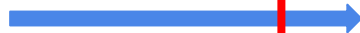
Training Data



Machine Learning Model



Unseen Data



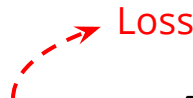

Can we compare the generalization performance of different models?

Class: Apple

# Notation

- Input:  $x$       Output:  $y$
- Neural Network:  $f_{\theta}$
- Cross Entropy Loss:  $L(f_{\theta}(x), y)$
- Output Sensitivity:  $S = \text{Var}(\overline{f_{\theta}(x + \varepsilon_x) - f_{\theta}(x)})$

# Loss vs. Sensitivity

- Cross Entropy vs. Mean Square Error:  $L \approx \sqrt{L_{\text{MSE}}/2}$   

- Bias-Variance Decomposition:  $L_{\text{MSE}} = \epsilon_{\text{bias}} + \epsilon_{\text{variance}}$
- Variance vs. Sensitivity:  $\epsilon_{\text{variance}} \approx S \cdot C + \Sigma$   


## Loss vs. Sensitivity

$$L \approx \sqrt{\frac{1}{2} (S \cdot C + \Sigma)}$$

# Loss vs. Sensitivity

$$L \approx \sqrt{\frac{1}{2} (S \cdot C + \Sigma)}$$

Requires labels

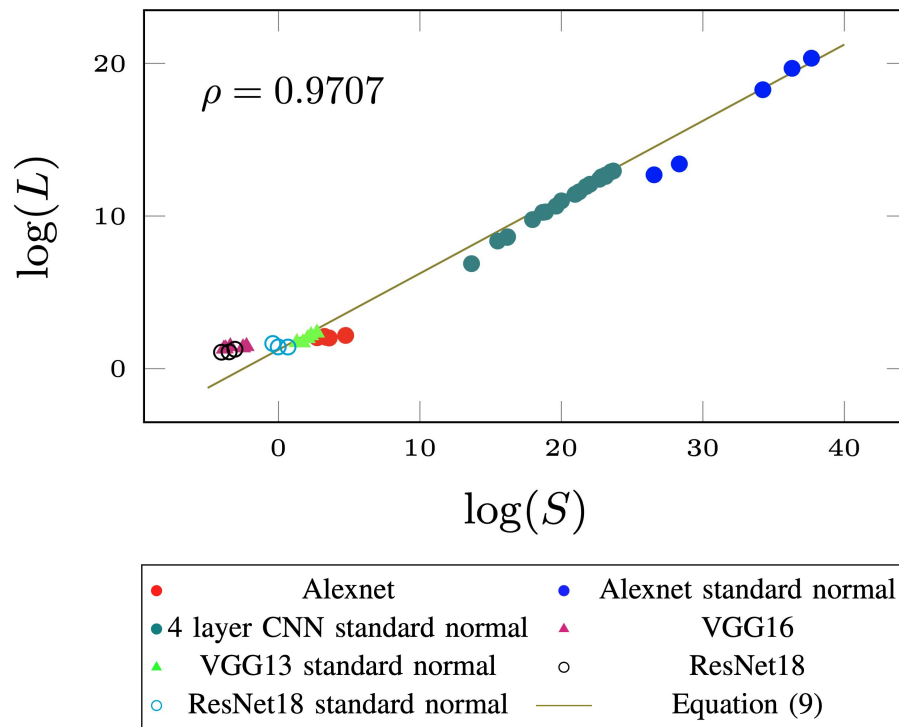


Does not require labels





# Loss vs. Sensitivity

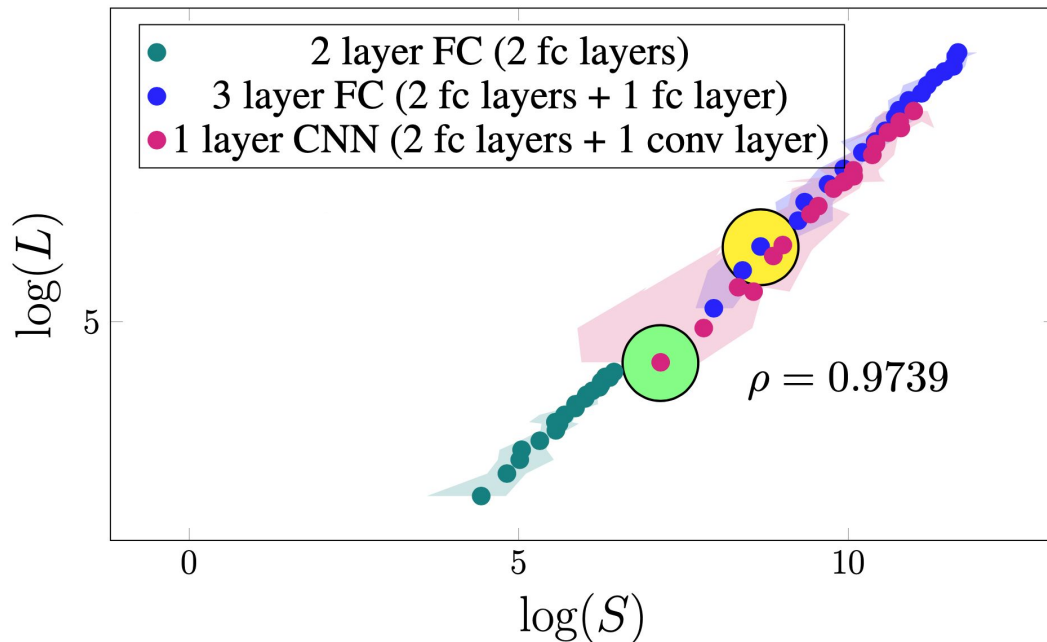


# Sensitivity as a Proxy for Loss

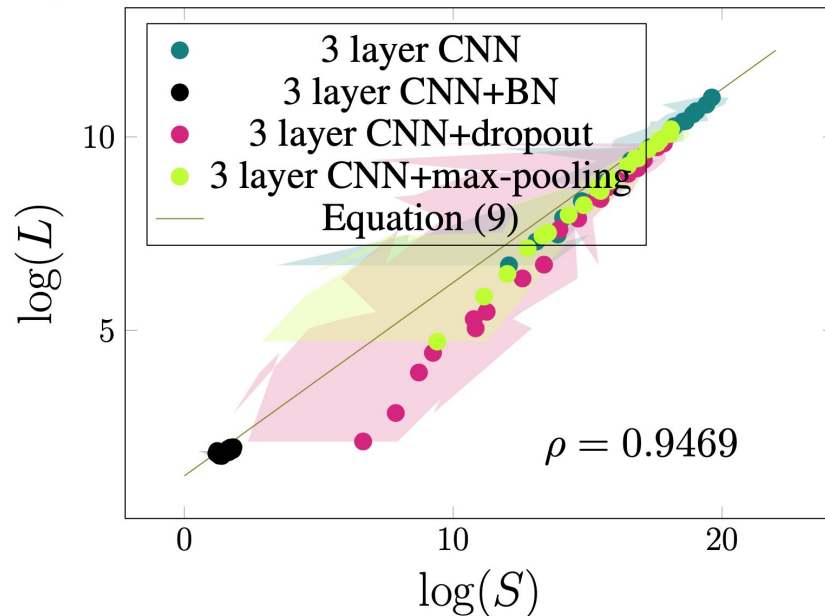
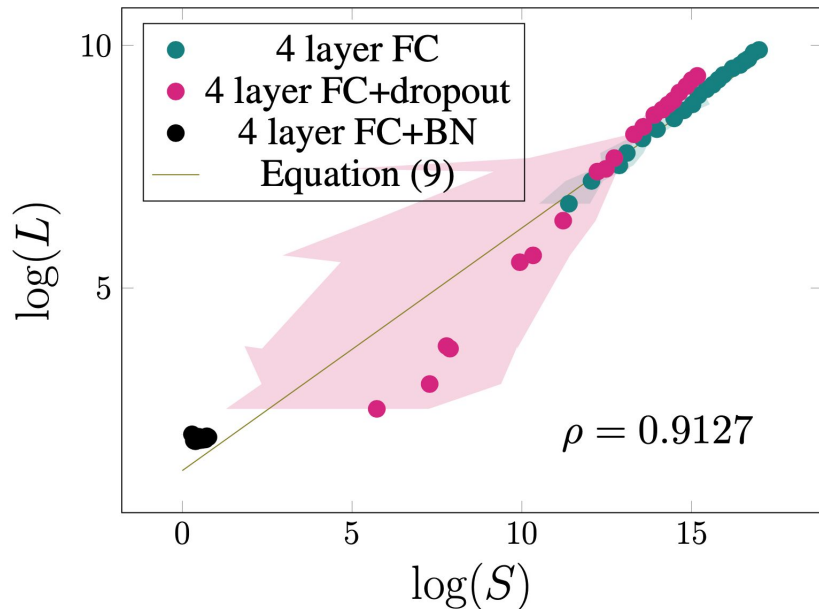
Convolutional

vs.

Fully-connected Layers

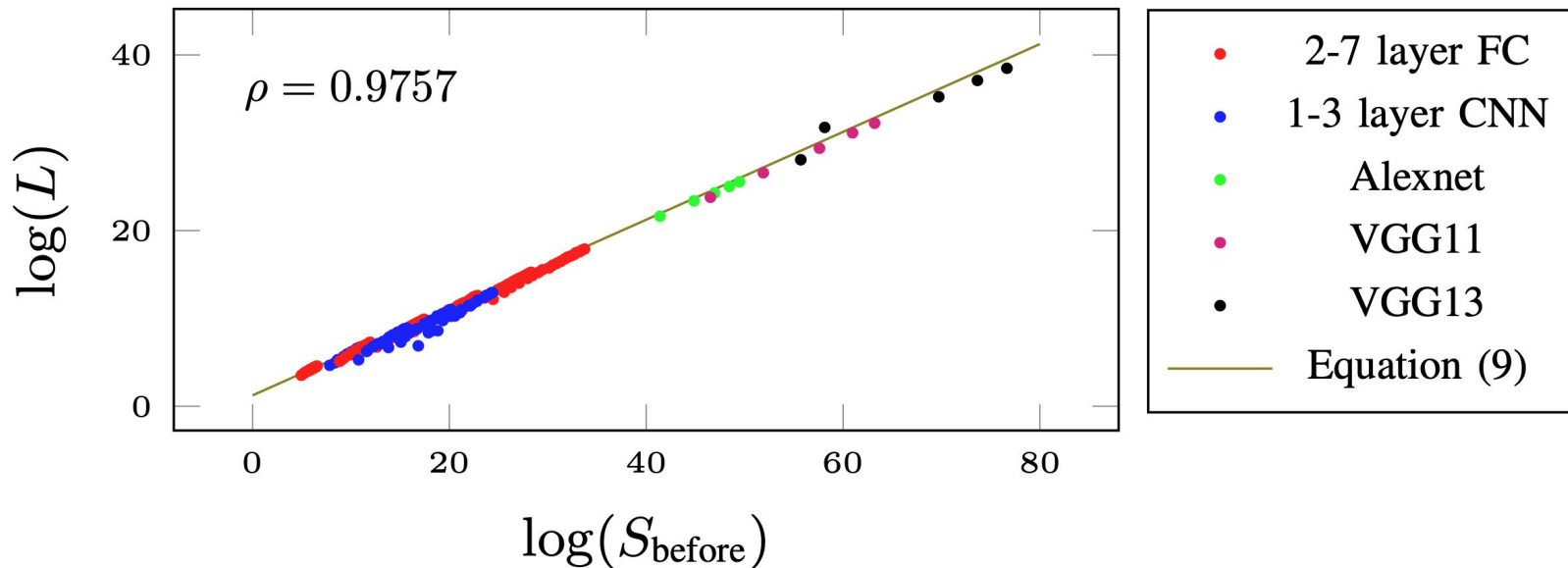


# Sensitivity as a Proxy for Loss



# Sensitivity of Untrained Networks

Output sensitivity computed *before* training:  $S_{\text{before}}$



# Future Directions

- When is the variance not the dominant term in the loss decomposition?
- Finding a link between the classification error and sensitivity
  - Bias-variance decomposition for the classification error?
  - Applying network calibration methods
- Beyond pixel-wise linear input perturbations and positive homogenous activation functions

# Thank You!

---