# Robust Convolutional Neural Networks under Adversarial Noise

Jonghoon Jin, Aysegul Dundar and Eugenio Culurciello
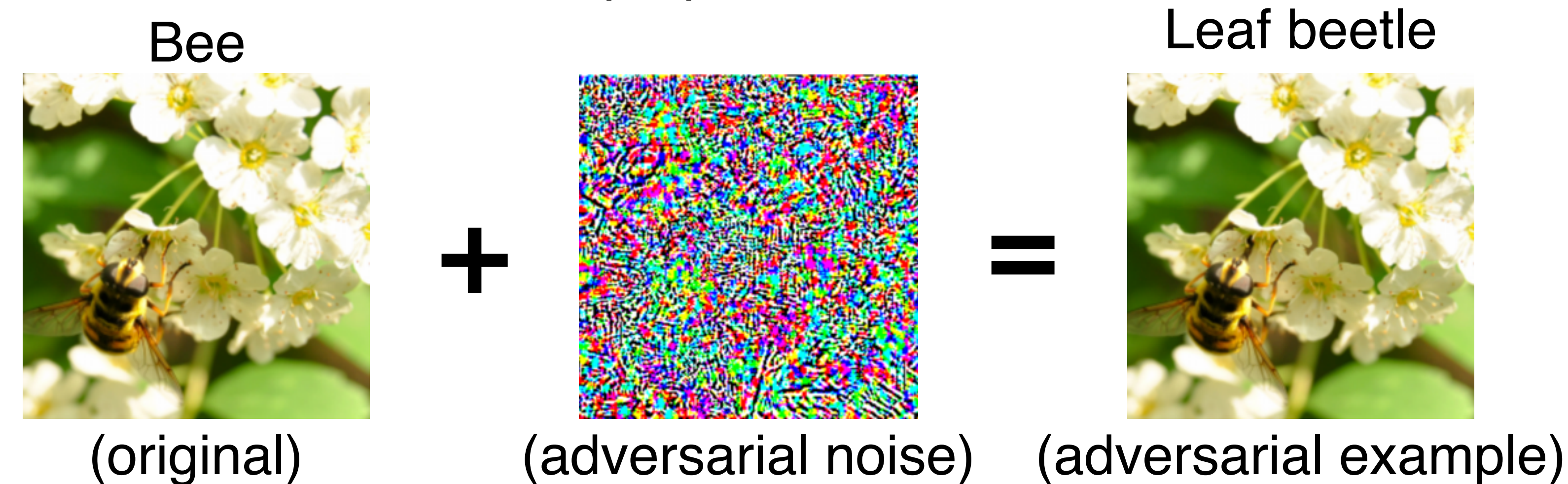
e-Lab
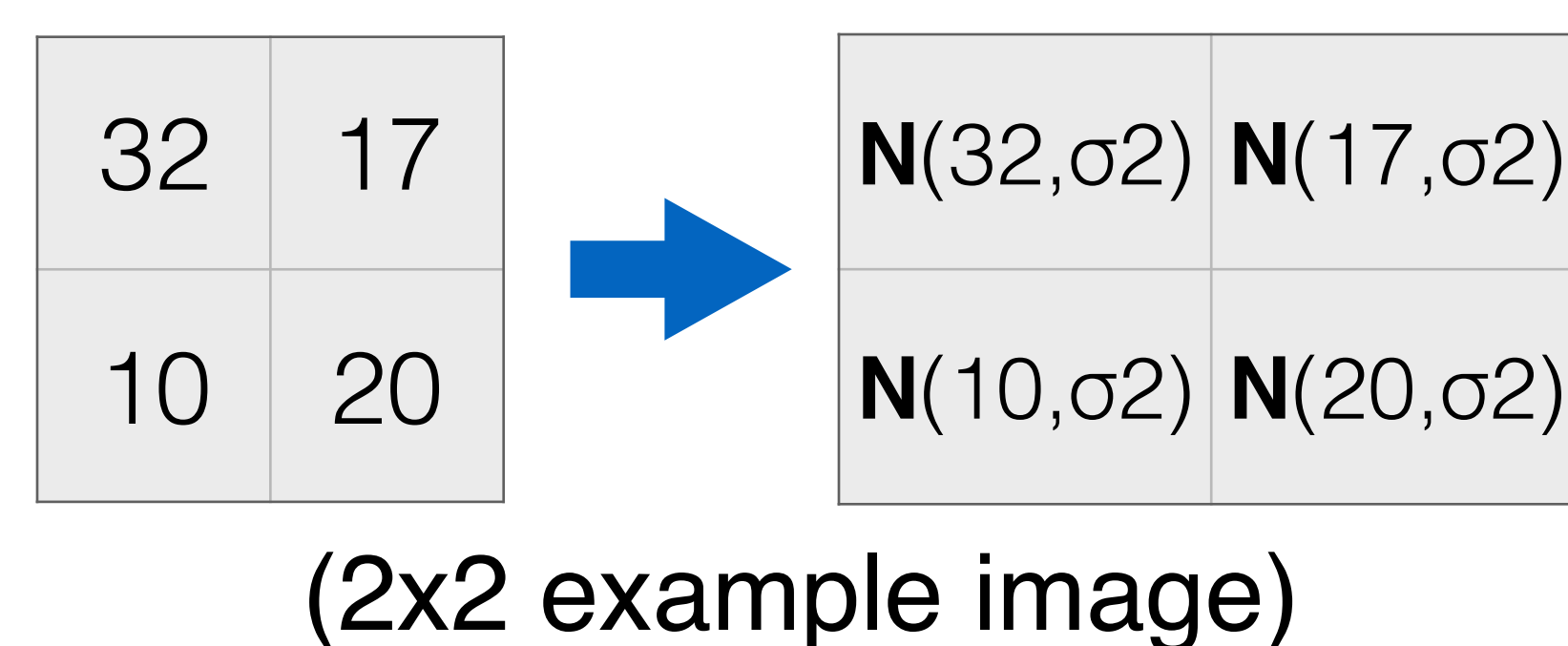
PURDUE UNIVERSITY

## Abstract

We propose a new feedforward CNN that is **robust to adversarial noise**. With uncertainty noise added to input, all operators in CNNs are modified to benefit from the noise. The model is parameterized by mean and variance per pixel and **successfully applied to deep architecture** like ResNet-101.

## Adversarial example

- generated to fool CNNs on purpose

Bee

Leaf beetle

(original) + (adversarial noise) = (adversarial example)

## Input with uncertainty

| 32 | 17 |
| 10 | 20 |

→

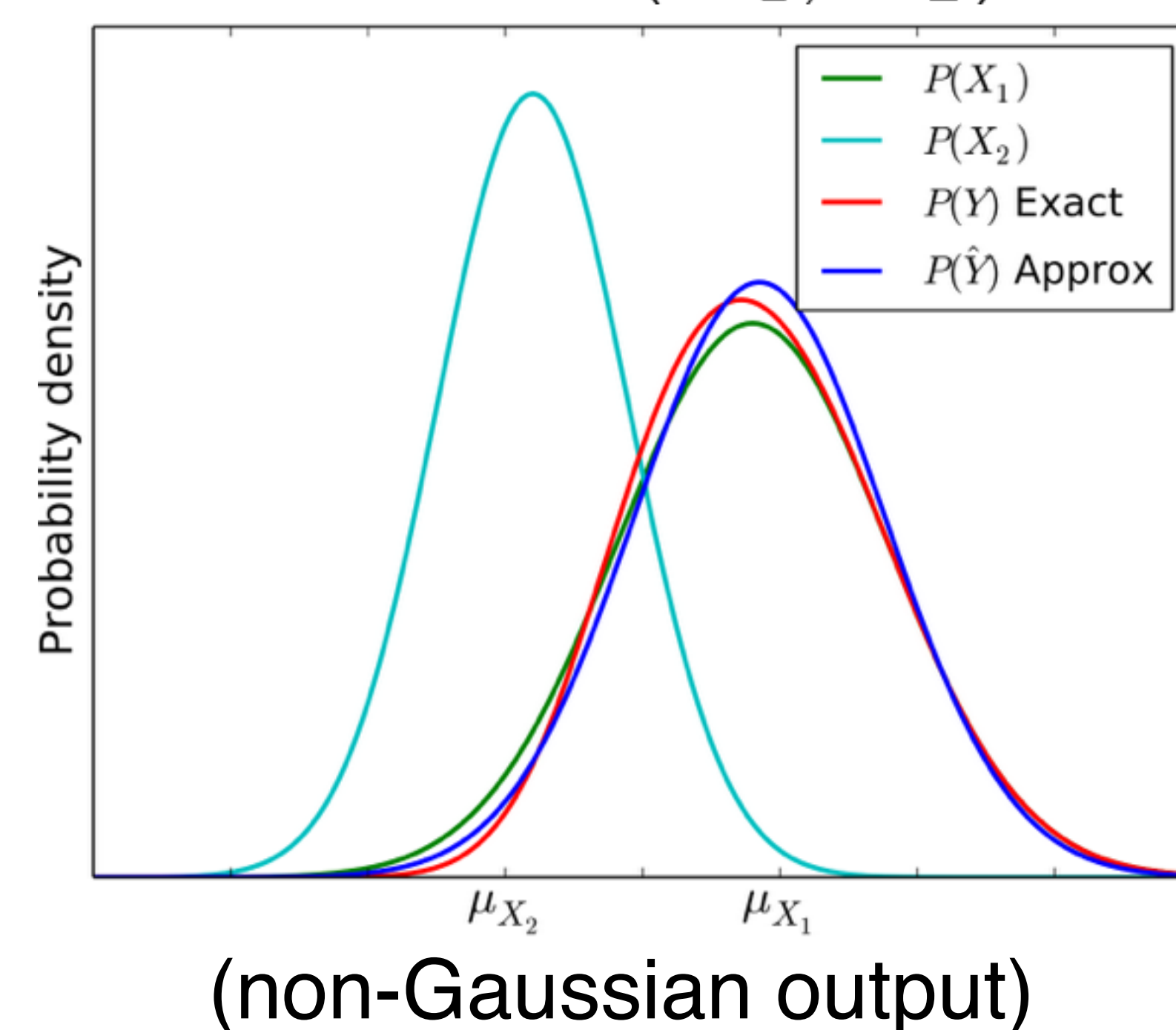| $N(32,\sigma2)$ | $N(17,\sigma2)$ |
| $N(10,\sigma2)$ | $N(20,\sigma2)$ |

(2x2 example image)

## Convolution

$$E[Y] = \sum \omega E[X] + b$$

$$Var[Y] = \sum \omega^2 Var[X]$$

(Gaussian output)

## Max-pooling

$$Y = max(X_1, X_2)$$



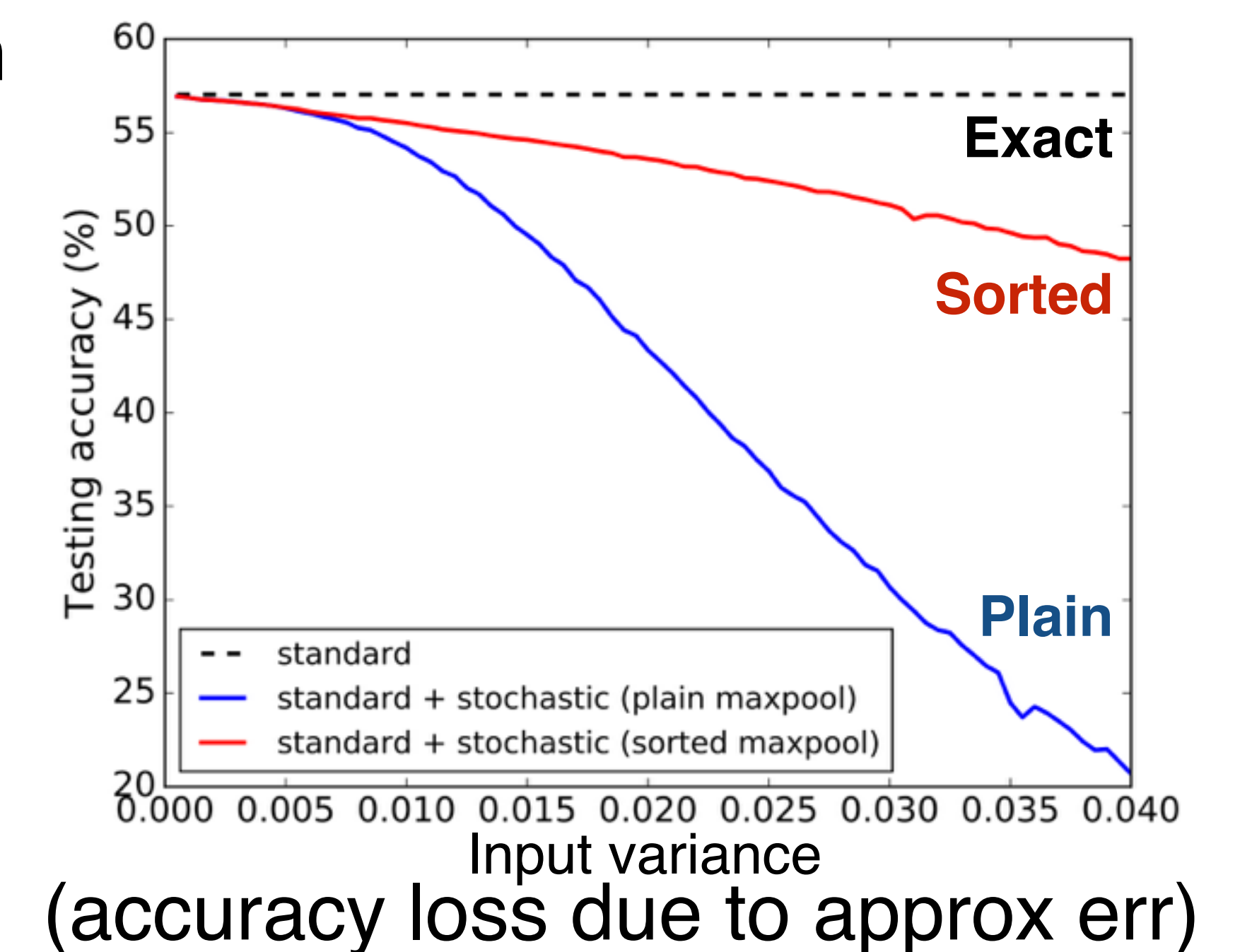(non-Gaussian output)

## ReLU

$$Y = max(X, \theta)$$



(non-Gaussian output)
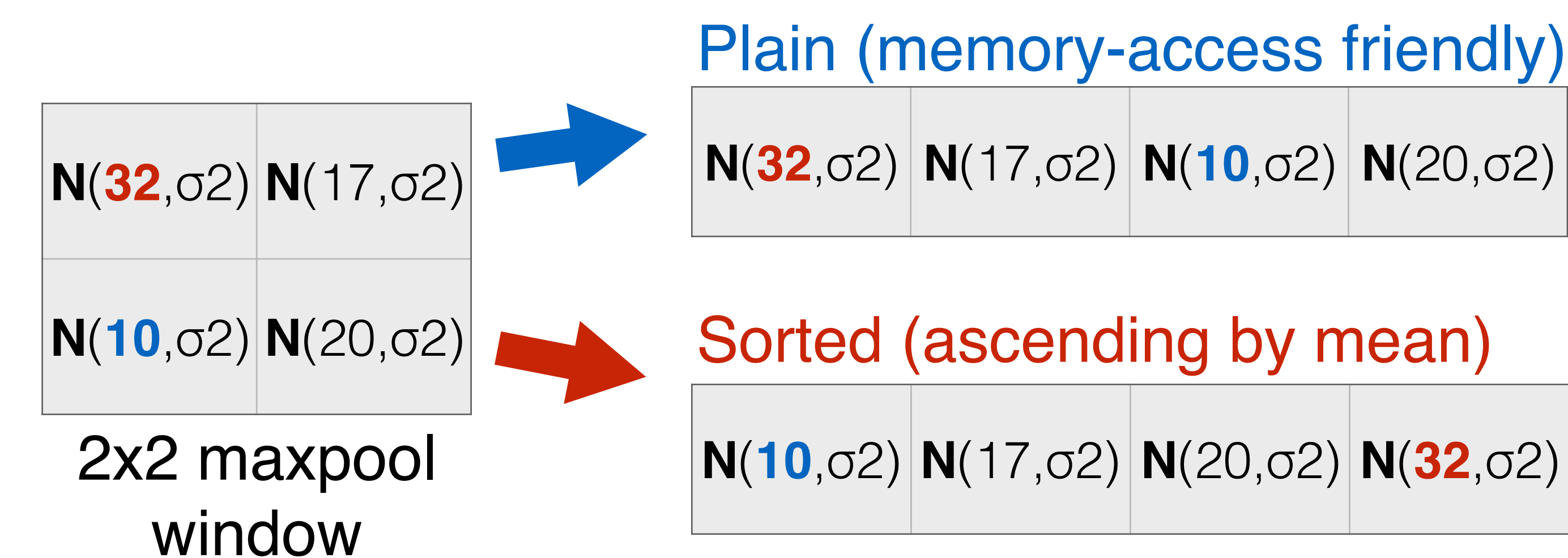
## Gaussian approximation for max-pooling (and ReLU)

- low error when max-pooled in ascending order by mean

| $N(\mathbf{32},\sigma2)$ | $N(17,\sigma2)$ |
| $N(\mathbf{10},\sigma2)$ | $N(20,\sigma2)$ |

2x2 maxpool window

**Plain (memory-access friendly)**

| $N(\mathbf{32},\sigma2)$ | $N(17,\sigma2)$ | $N(\mathbf{10},\sigma2)$ | $N(20,\sigma2)$ |

Sorted (ascending by mean)

| $N(\mathbf{10},\sigma2)$ | $N(17,\sigma2)$ | $N(20,\sigma2)$ | $N(\mathbf{32},\sigma2)$ |



(accuracy loss due to approx err)

## Parameter tuning (σ)

- trade-off (small classification loss <—> large gain with noise)



(fixed noise intensity)



(fixed input variance)

## Demo



## Code

jhjin/stochastic-cnn

## Classification accuracy under noise

(higher is better)

| Dataset | | CIFAR-10 | | ImageNet | | |
|---|---|---|---|---|---|---|
| Model | | NIN | | AlexNet | | ResNet-101 |
| Adversarial noise intensity [px] | | 0 | 0.5 | 0 | 0.01 | 0.5 | 1 |
| Standard training | | 90.1 | 72.3 | 57.0 | 56.1 | 24.8 | 17.50 |
| Standard training + stochastic (this work) | | 88.9 | **78.1** | 57.0 | **56.2** | **33.4** | **39.24** |
| LWA + BN (Huang et al. 2016) | | 89.0 | 82.3 | — | — | — | — |
| Adversarial training (Goodfellow et al. 2015) | | 88.7 | 82.1 | 43.0 | 42.9 | * | * |
| Adversarial training + stochastic (this work) | | 88.7 | **82.9** | 43.0 | 42.9 | * | * |

( * : failed to converge)