

Optimization for Robust Deep Learning

M. Pawan Kumar

University of Oxford

Joint work with Rudy Bunel, Alban Desmaison,
Krishnamurthy Dvijotham, Pushmeet Kohli and Philip Torr

ImageNet Challenge

~1 M images, 1000 classes

Flute



Strawberry



Traffic light



Backpack



Matchstick



Sea lion



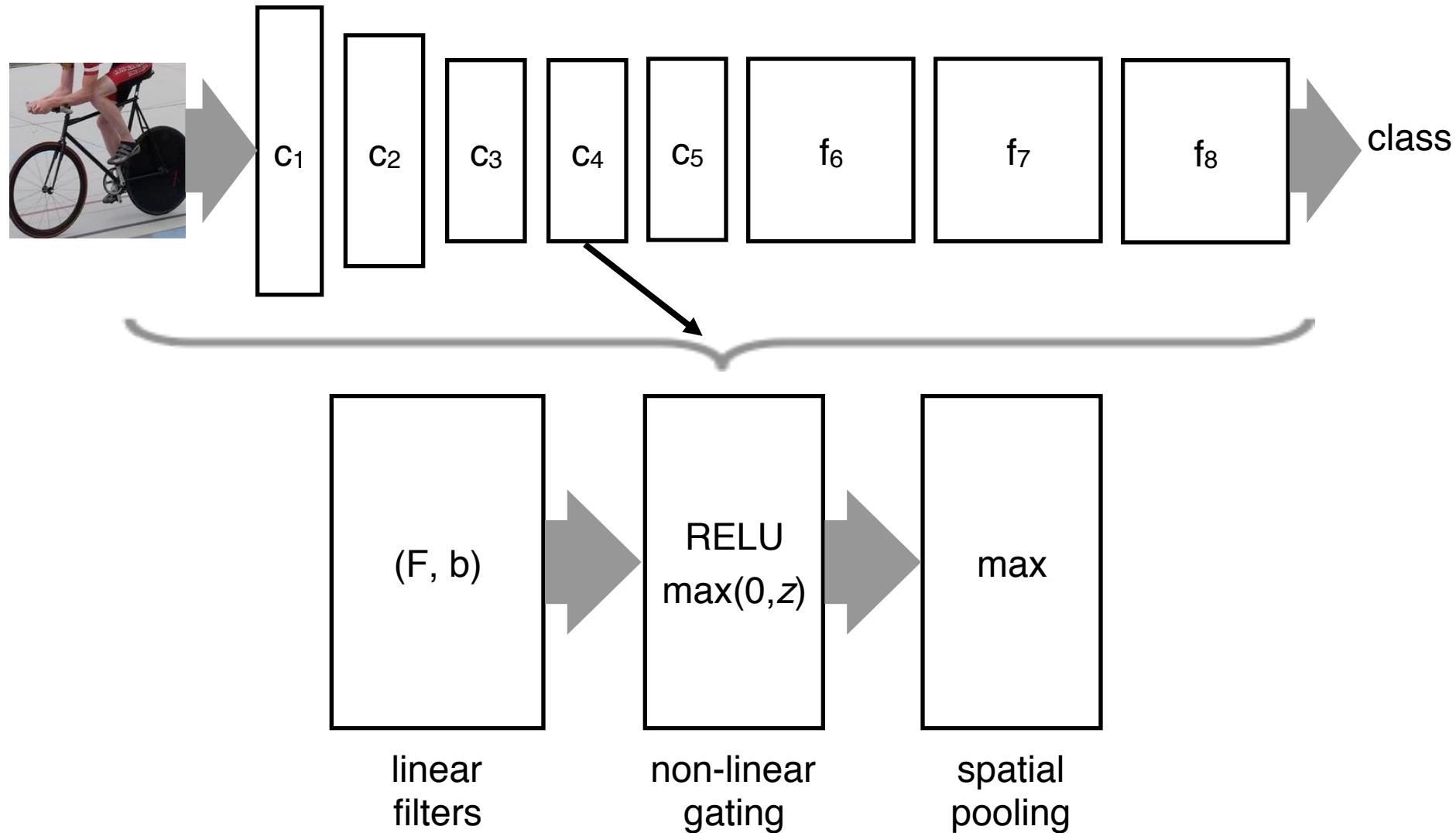
Bathing cap



Racket



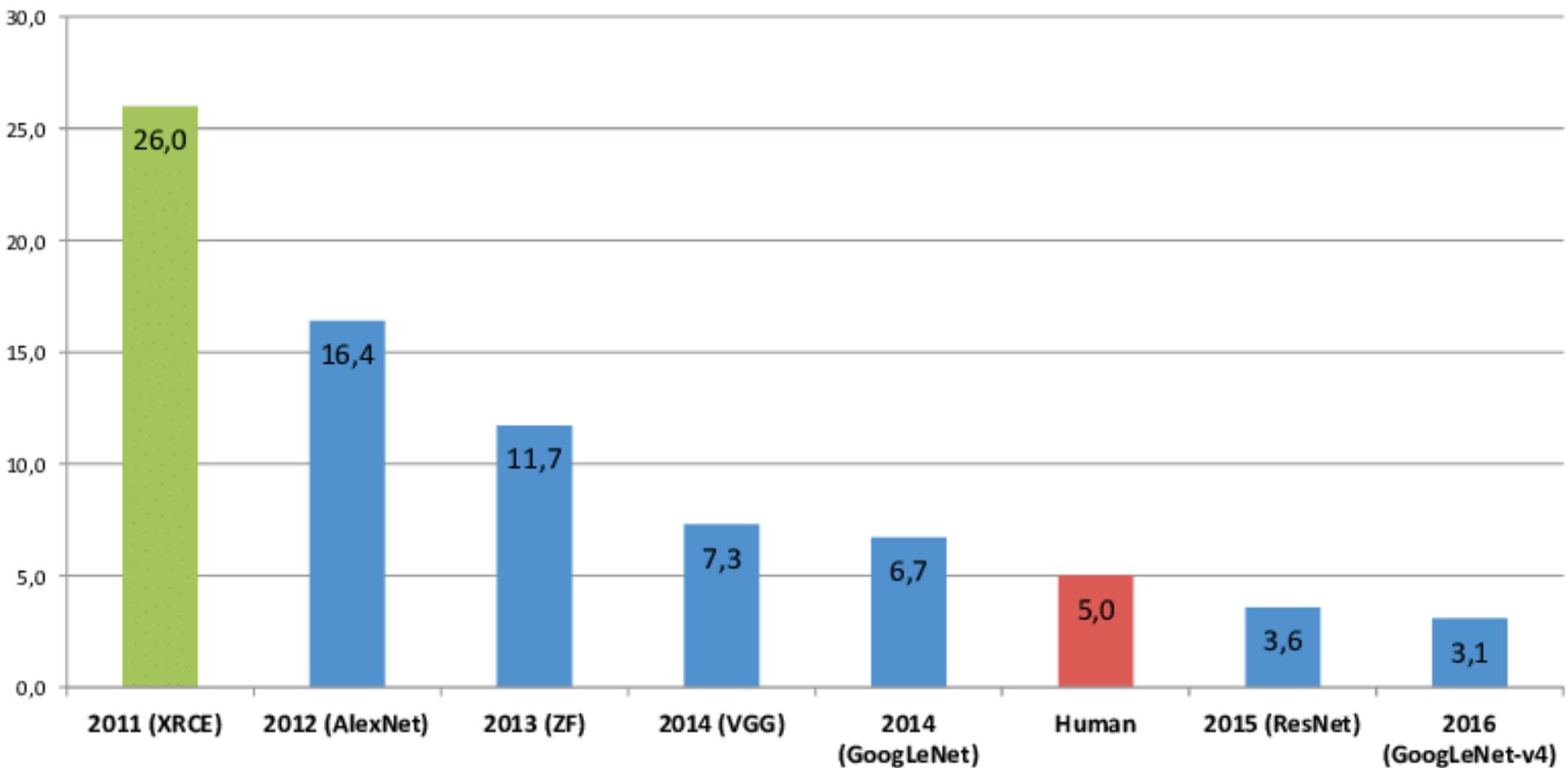
Deep Learning



Linear functions followed by piecewise linear non-linearities

ImageNet Challenge

ImageNet Classification Error (Top 5)



Self-Driving Car

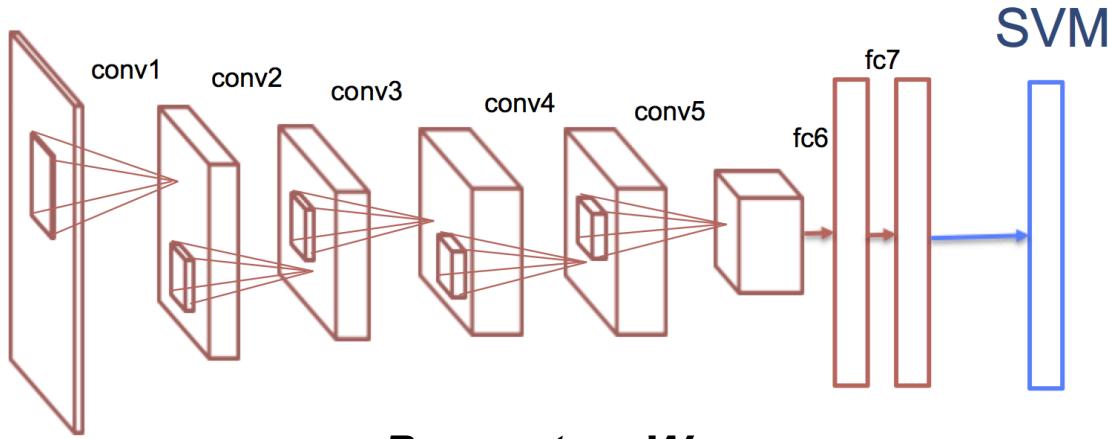


Human drivers replaced by deep neural networks

Deep Learning



Training Data



Parameters \mathbf{W}

Road sign recognition

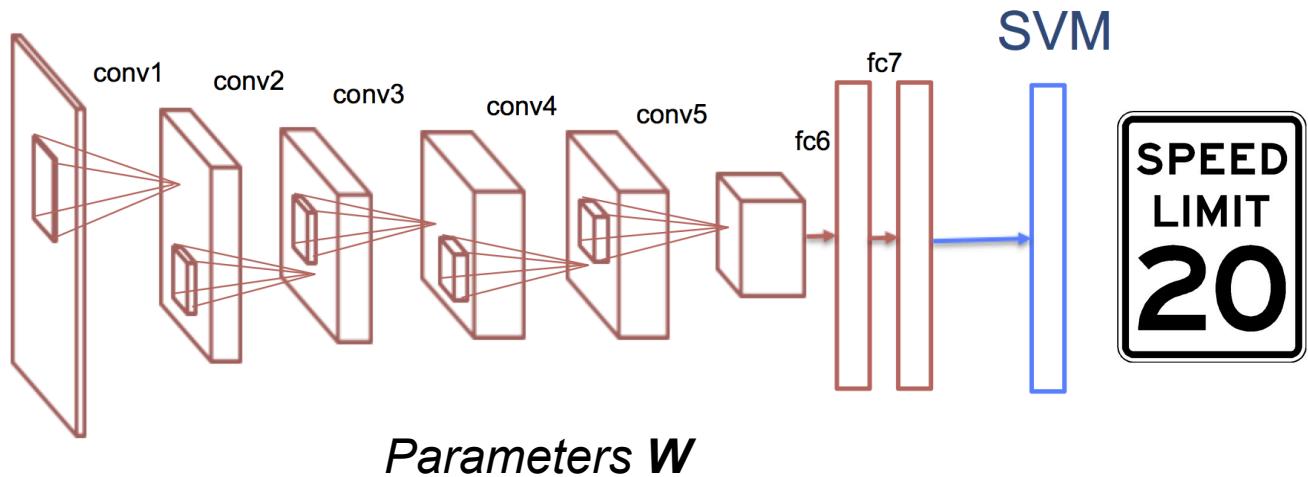
Parameters \mathbf{W} perform multiclass classification

Estimate \mathbf{W} using training data

Deep Learning



Test Image

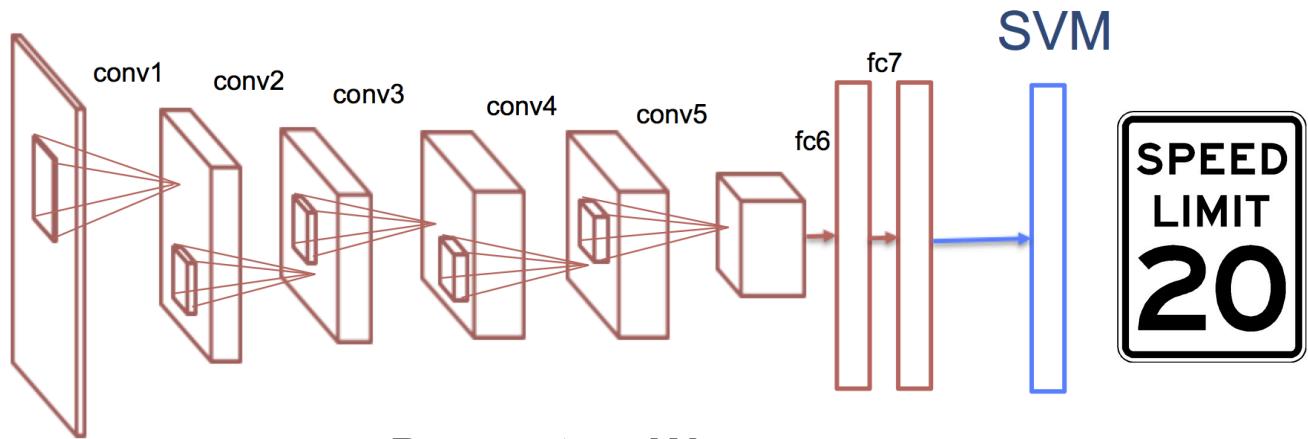


Parameters **W** select the class of a new input image

Deep Learning



Test Image



Small deformations can cause fatal errors



Blurring, Saturation



Aung et al., 2017

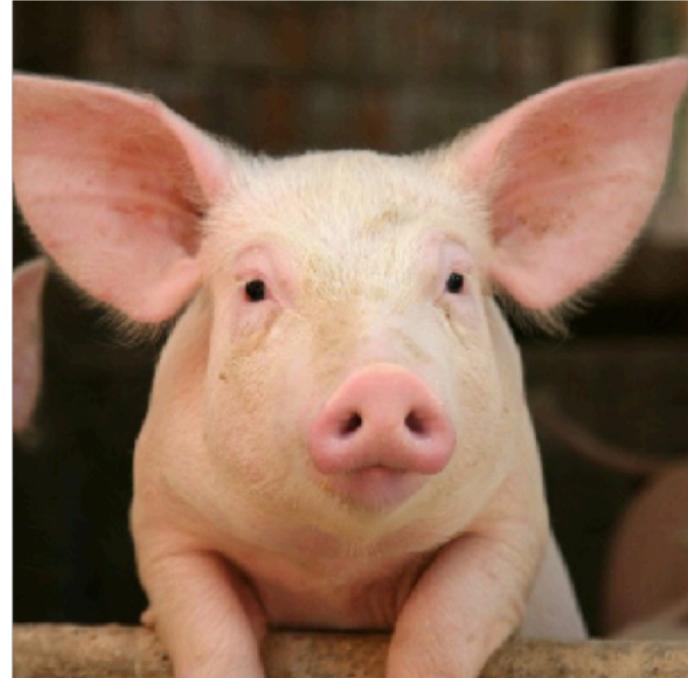
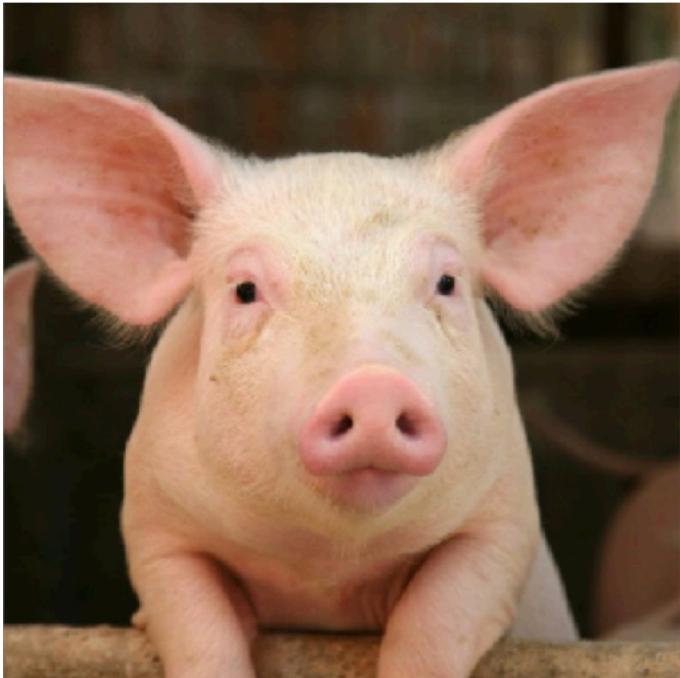


Pixel Errors



Evtimov et al., 2017

Spot the Difference

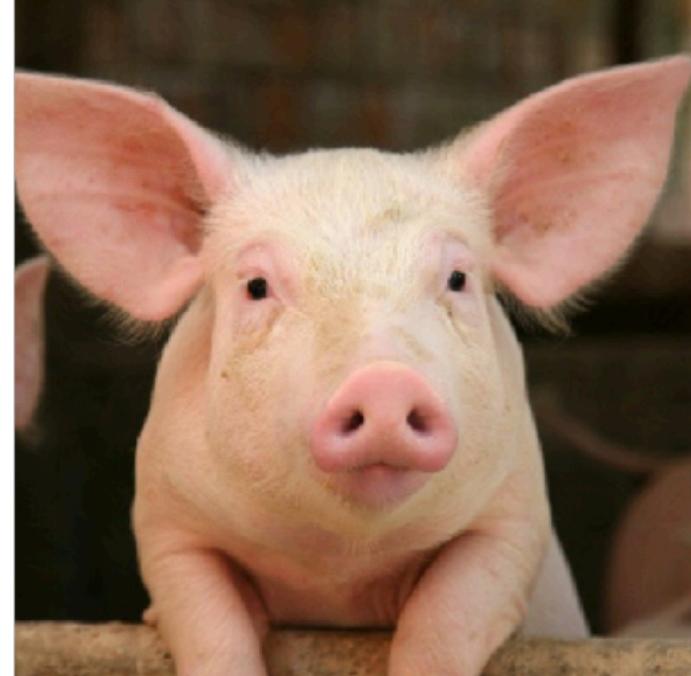


Spot the Difference

“pig” (91%)

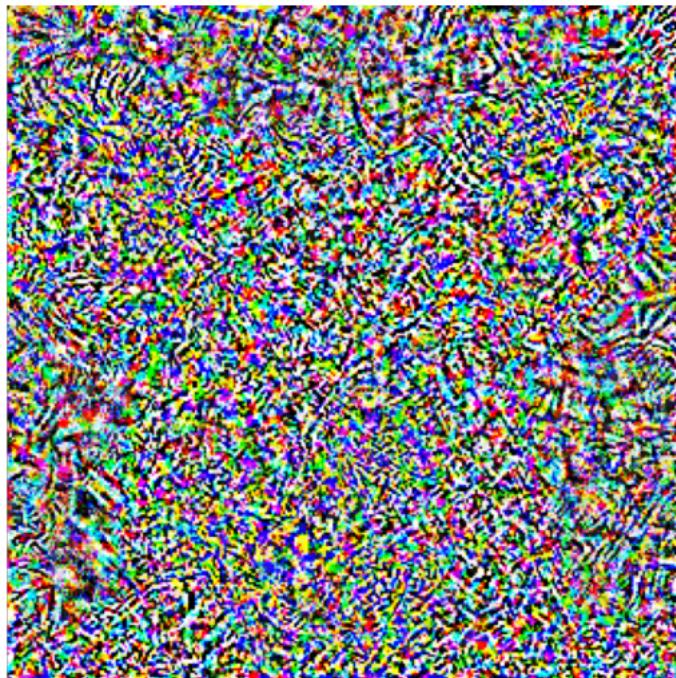


“airliner” (99%)

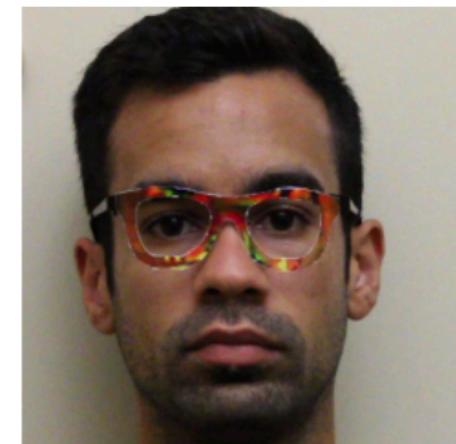


Difference

0.005 x



Glasses



Sharif, Bhagavatula, Bauer and Reiter, 2016

3D Object

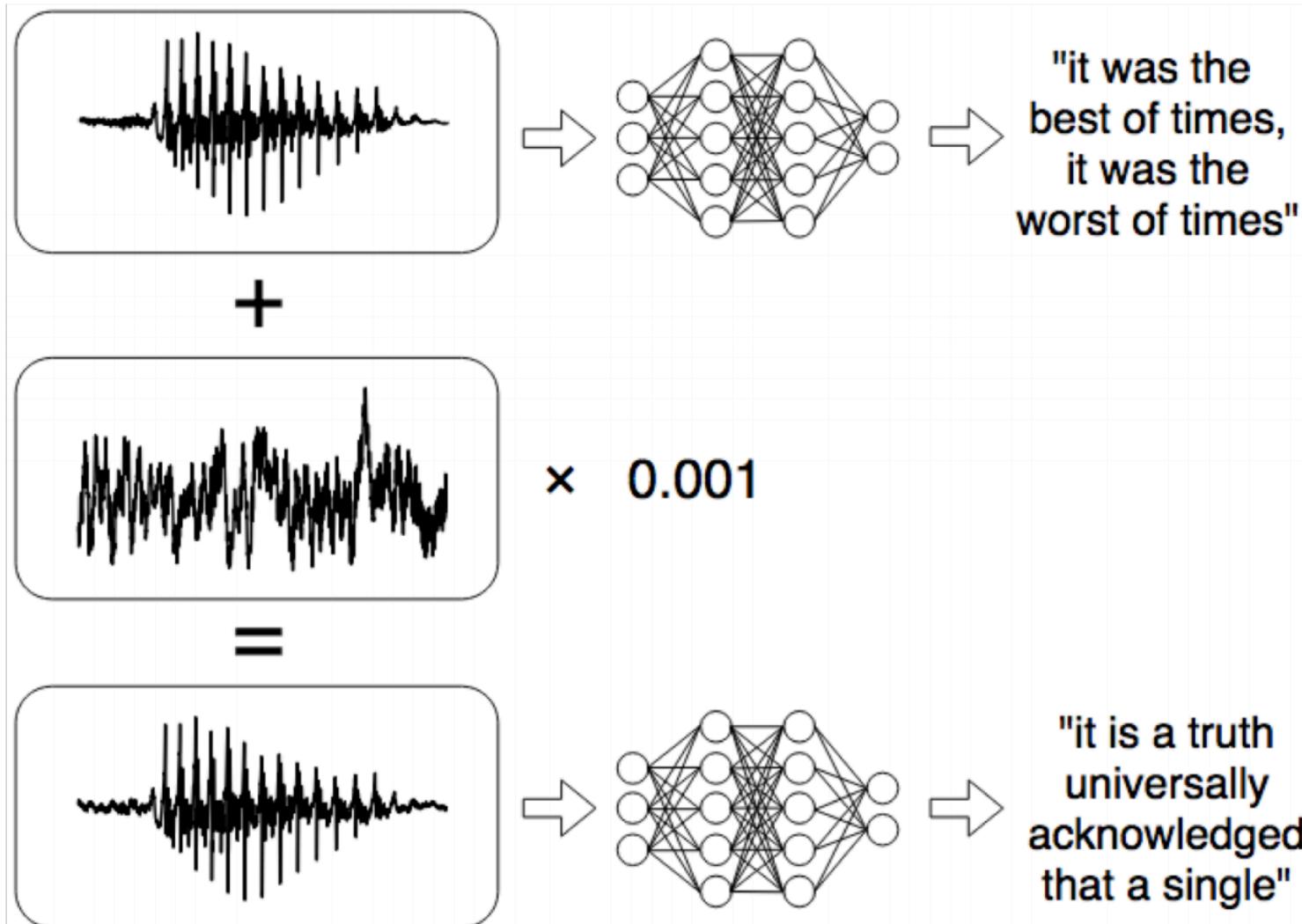


■ classified as turtle

■ classified as rifle

■ classified as other

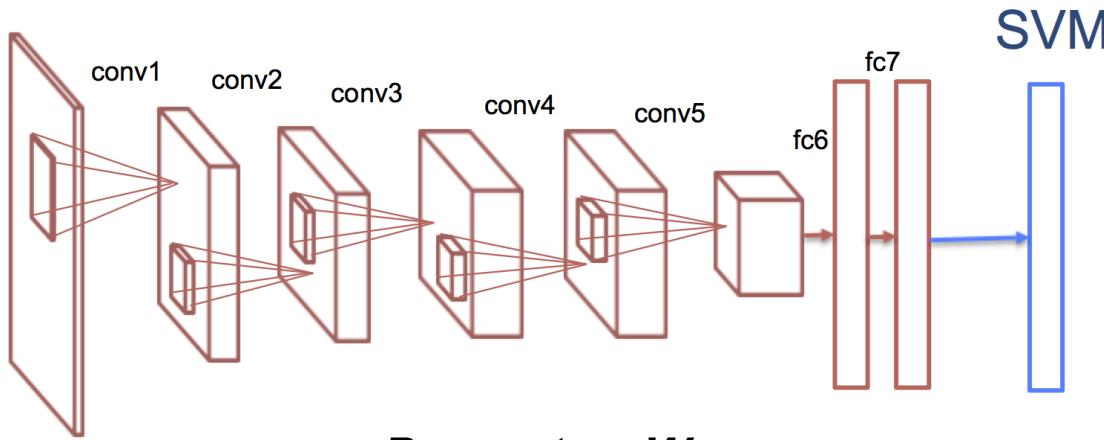
Audio



Robust Deep Learning



Training Data

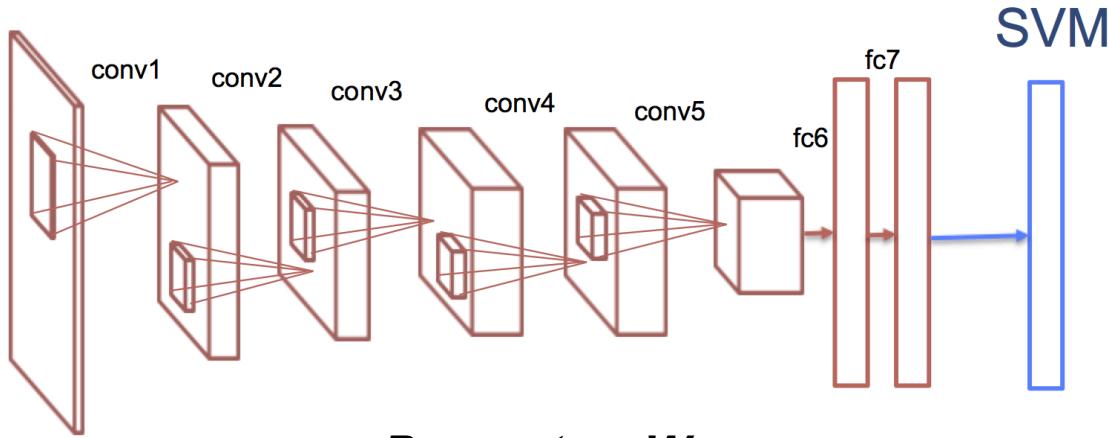


Parameters W

Robust Deep Learning



Training Data



Parameters W

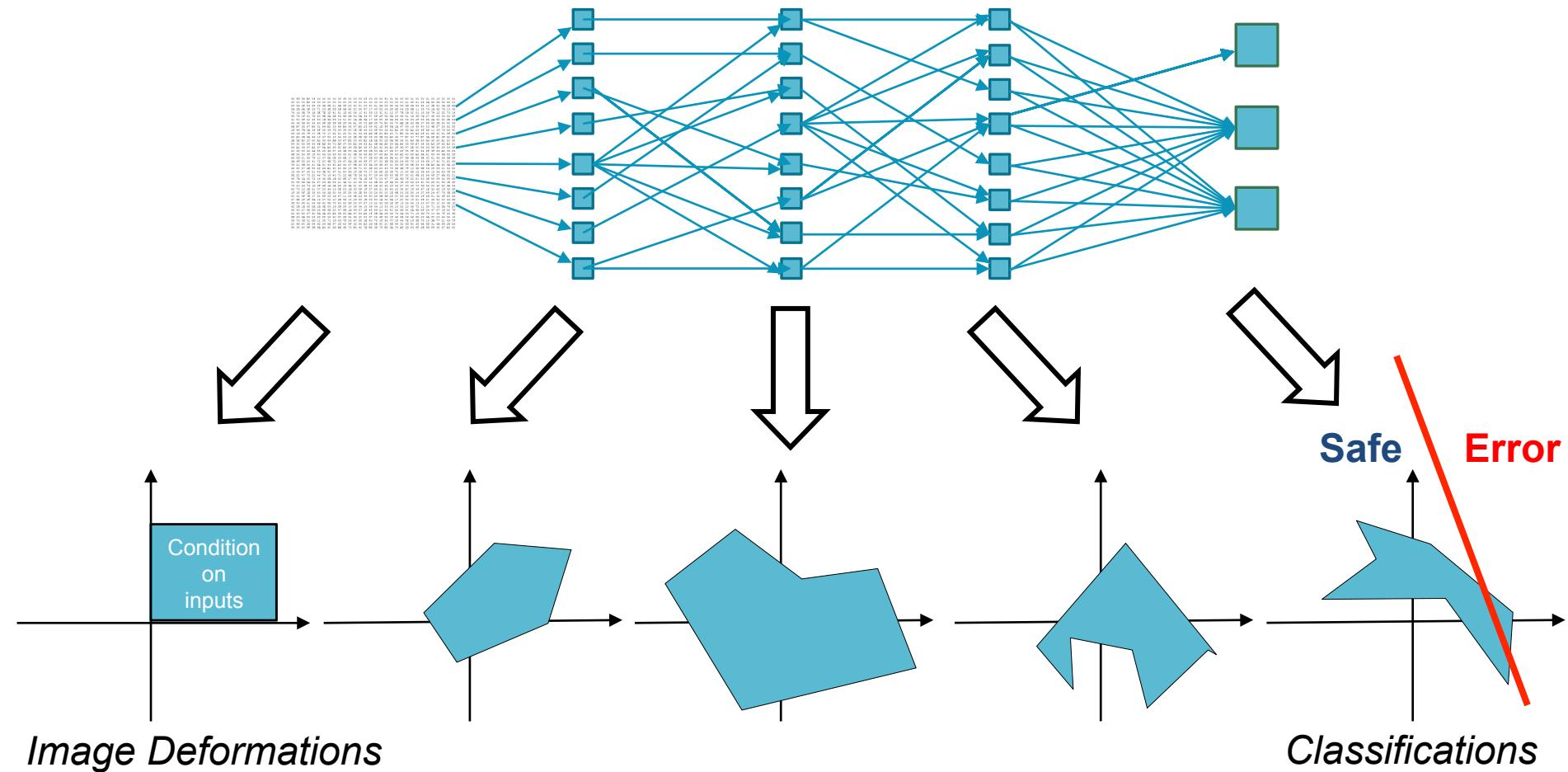
Augment
Training
Data Set

Identify
Deformation
With Error



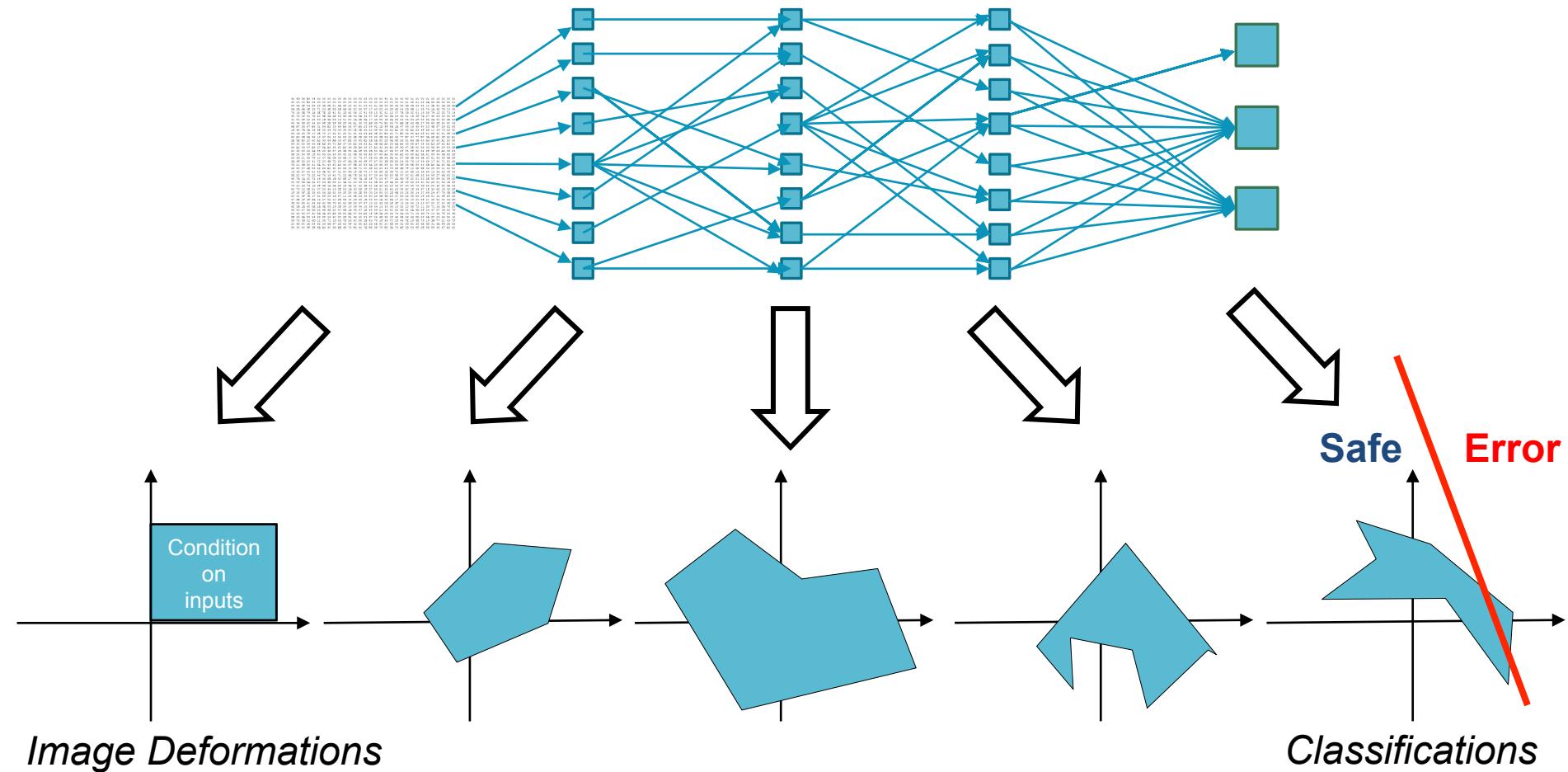
Robust Deep Learning

Is there an erroneous output?



Robust Deep Learning

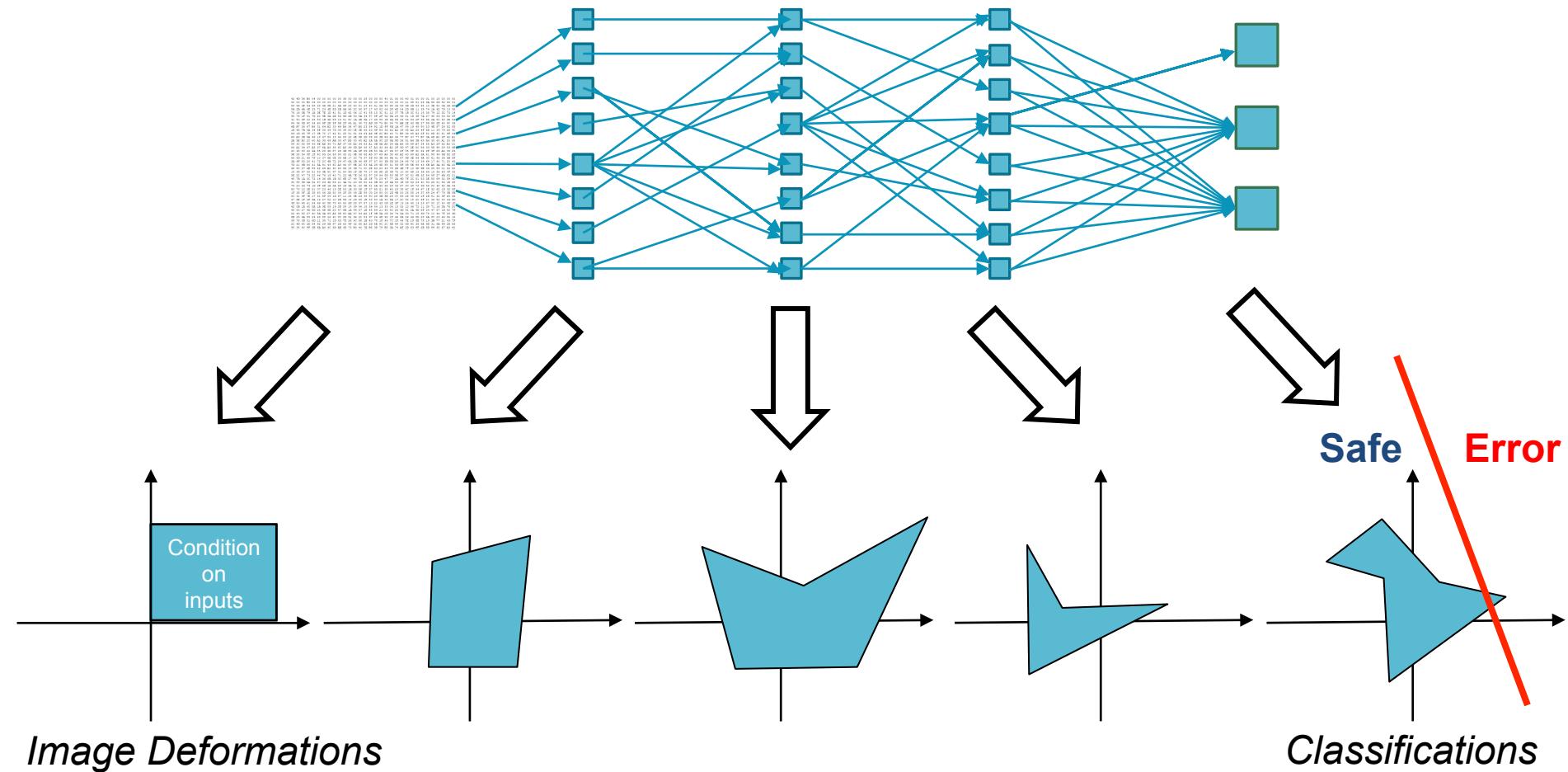
Is there an erroneous output?



Re-estimate parameters \mathbf{W}

Robust Deep Learning

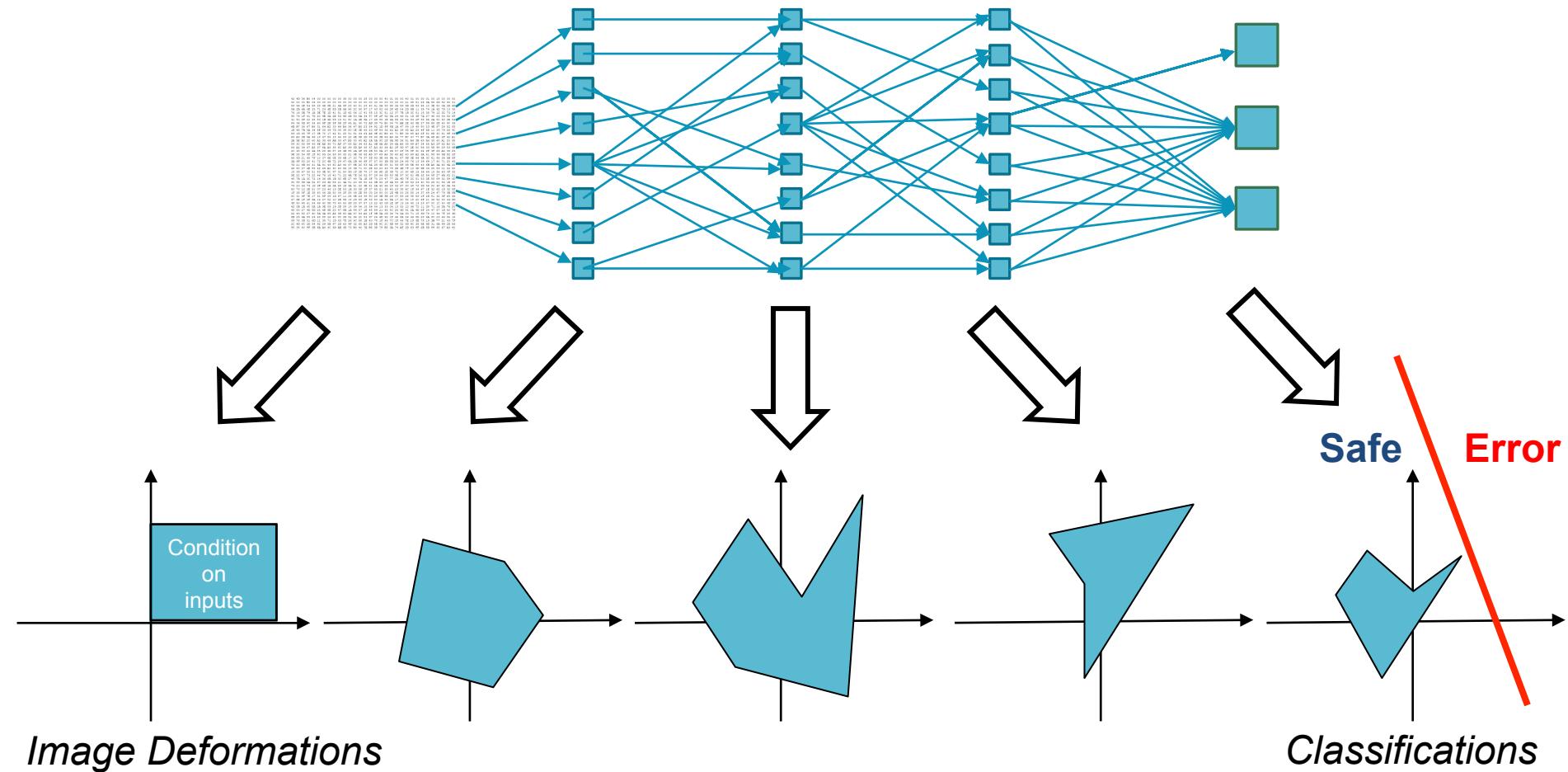
Is there an erroneous output?



Re-estimate parameters \mathbf{W}

Robust Deep Learning

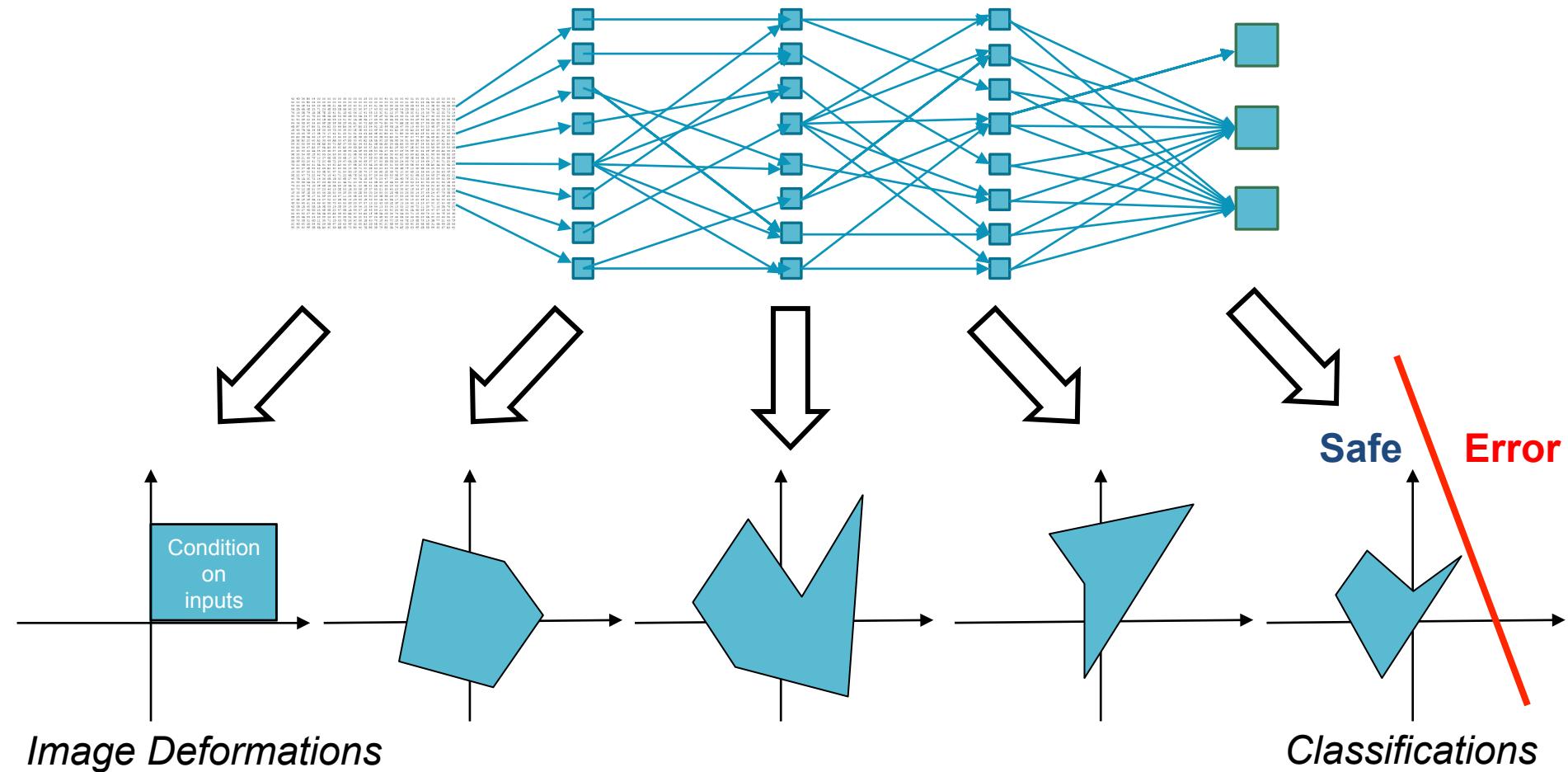
Is there an erroneous output?



Re-estimate parameters \mathbf{W}

Robust Deep Learning

Is there an erroneous output?



Terminate when all putative outputs are safe

Robust Deep Learning

Is there an erroneous output?

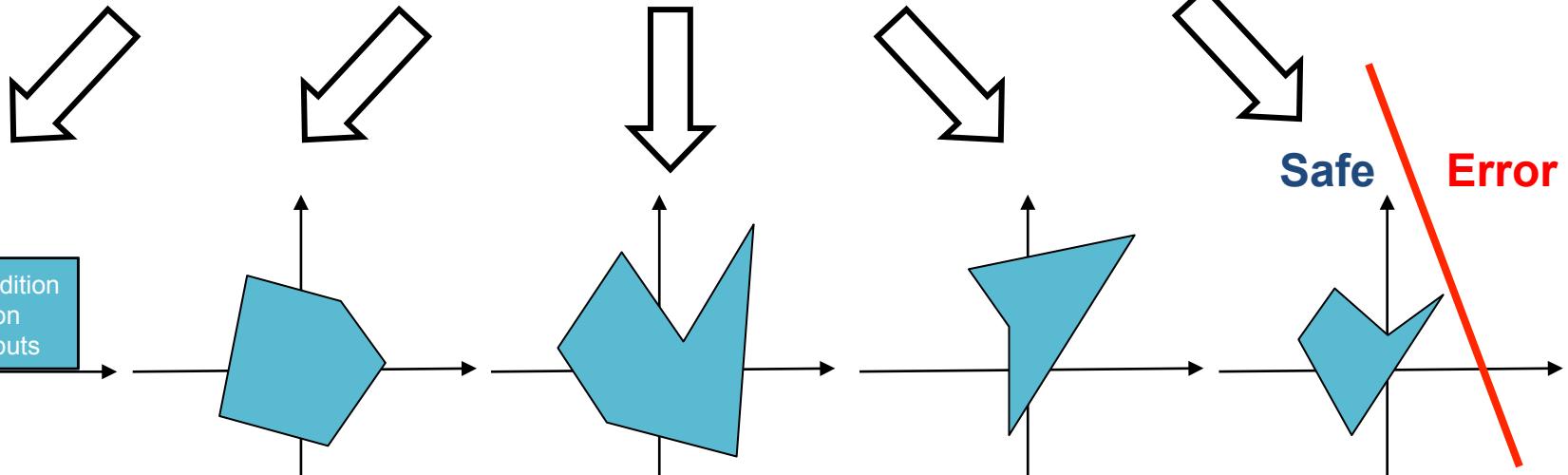
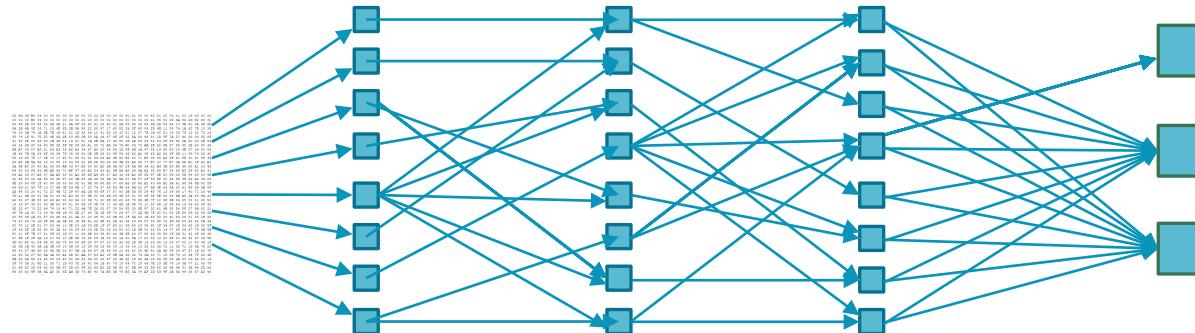


Image Deformations

Classifications

Terminate when all putative outputs are safe

Outline

- Neural Network Bounds

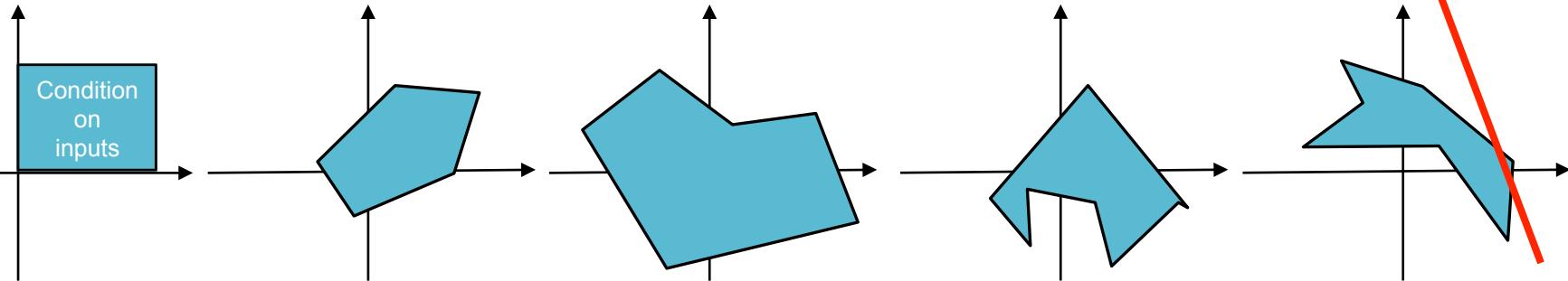
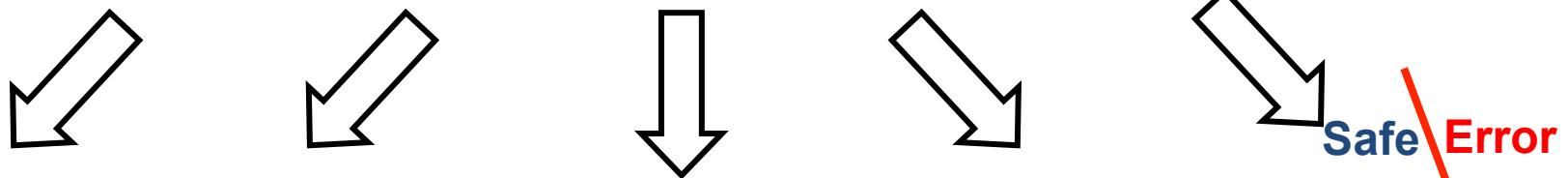
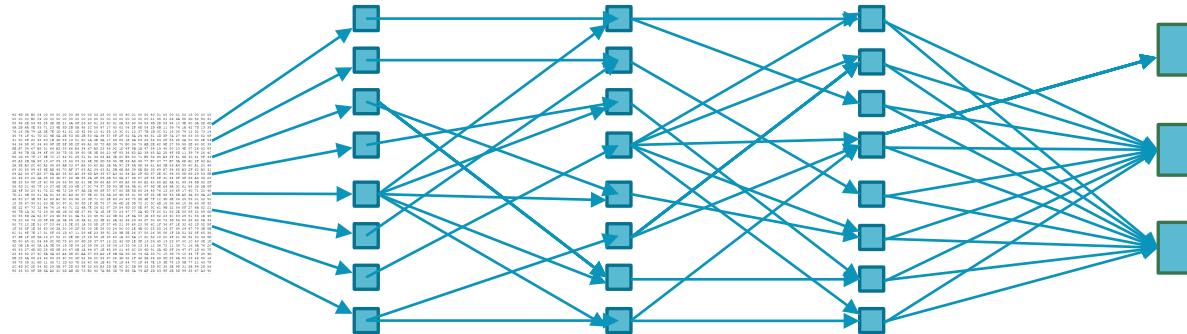
- Lagrangian Decomposition
- Proximal Minimization

Method of
Multipliers

- Results

Robust Deep Learning

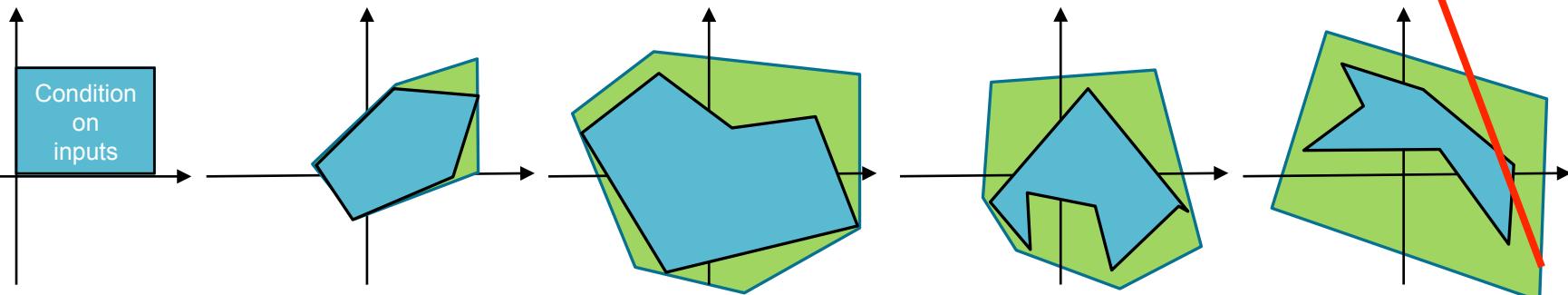
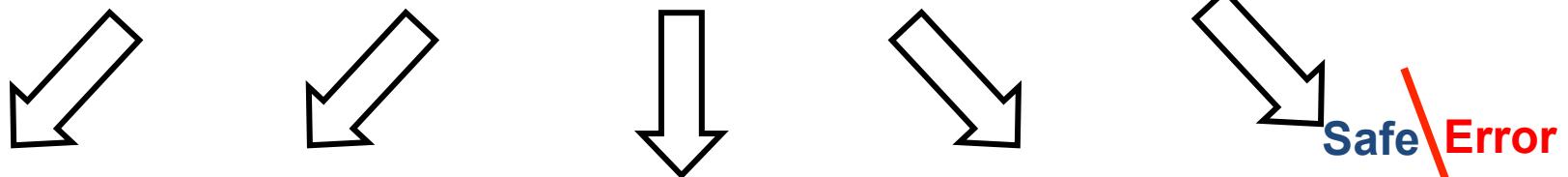
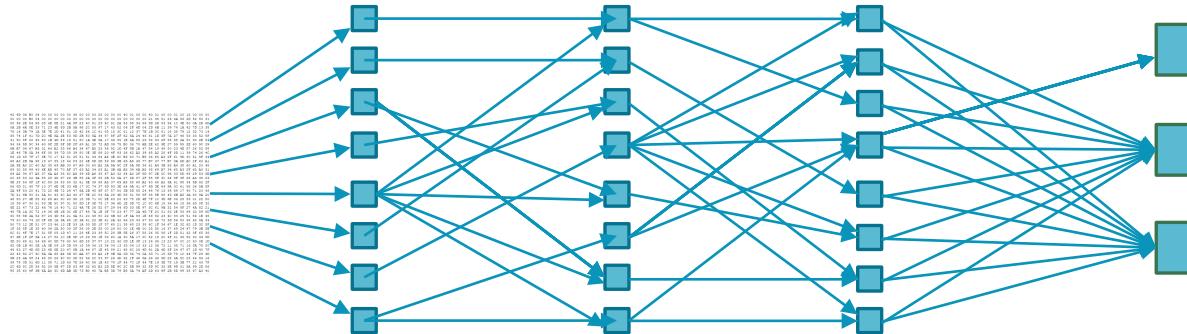
Is there an erroneous output?



Non-convexity makes the problem NP-hard

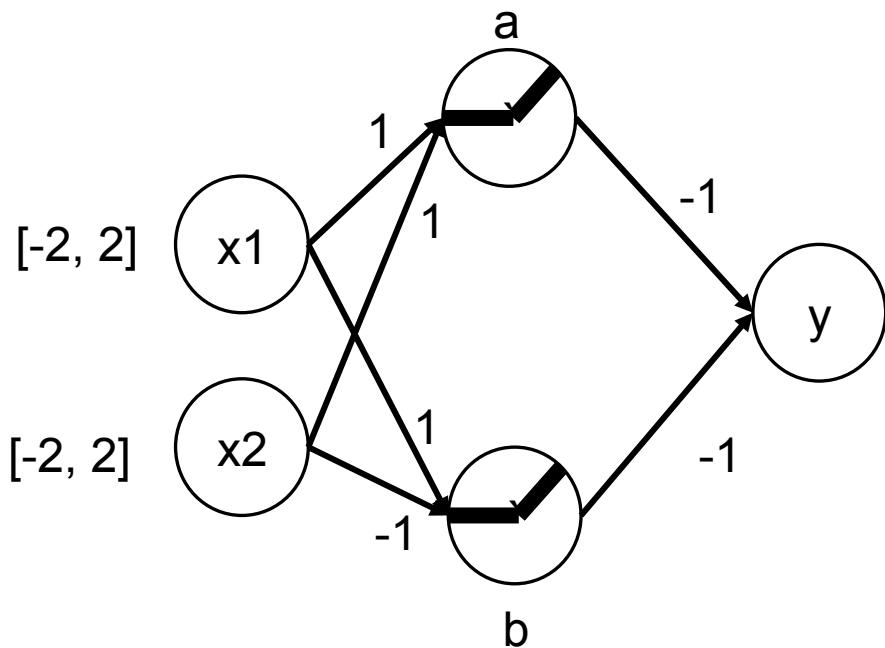
Robust Deep Learning

Is there an erroneous output?



Replace by a convex superset

Example



Find the minimum valid output

$$\begin{aligned} & \min && y \\ \text{s.t. } & -2 \leq x_1 \leq 2 \\ & -2 \leq x_2 \leq 2 \end{aligned}$$

$$a_{\text{in}} = x_1 + x_2$$

$$b_{\text{in}} = x_1 - x_2$$

$$a_{\text{out}} = \max\{a_{\text{in}}, 0\}$$

$$b_{\text{out}} = \max\{b_{\text{in}}, 0\}$$

$$y = -a_{\text{out}} - b_{\text{out}}$$

Example

Linear constraints

Easy to handle

$$\min \quad y$$

$$\text{s.t.} \quad -2 \leq x_1 \leq 2$$

$$-2 \leq x_2 \leq 2$$

$$a_{in} = x_1 + x_2$$

$$b_{in} = x_1 - x_2$$

$$a_{out} = \max\{a_{in}, 0\}$$

$$b_{out} = \max\{b_{in}, 0\}$$

$$y = -a_{out} - b_{out}$$

Example

$$\min \quad y$$

$$\text{s.t.} \quad -2 \leq x_1 \leq 2$$

$$-2 \leq x_2 \leq 2$$

$$a_{in} = x_1 + x_2$$

$$b_{in} = x_1 - x_2$$

$$a_{out} = \max\{a_{in}, 0\}$$

$$b_{out} = \max\{b_{in}, 0\}$$

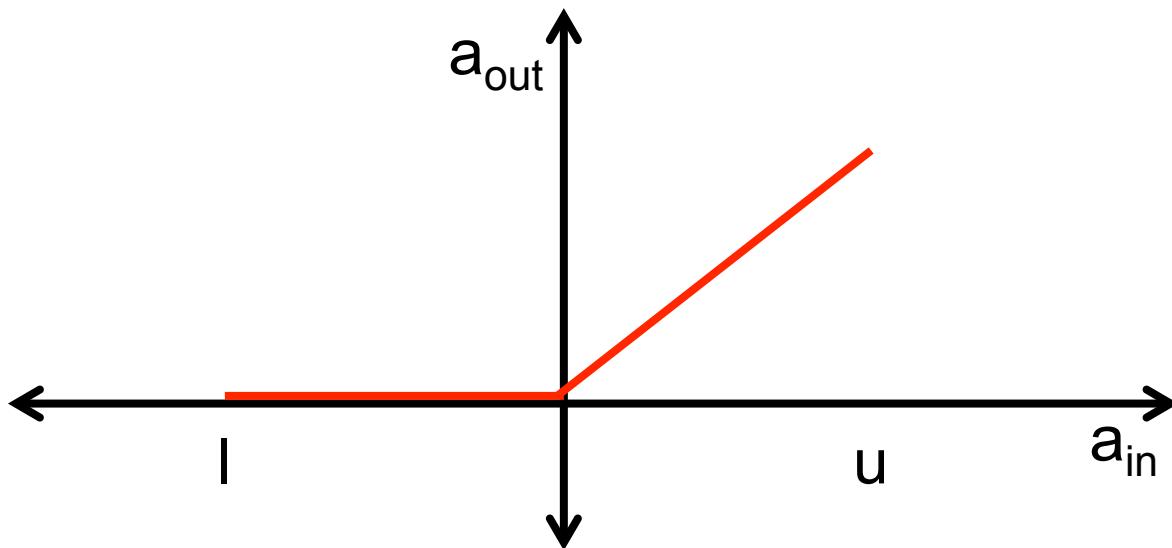
$$y = -a_{out} - b_{out}$$

Non-linear constraints

ReLU

$$a_{\text{out}} = \max\{a_{\text{in}}, 0\}$$

$$a_{\text{in}} \in [l, u]$$

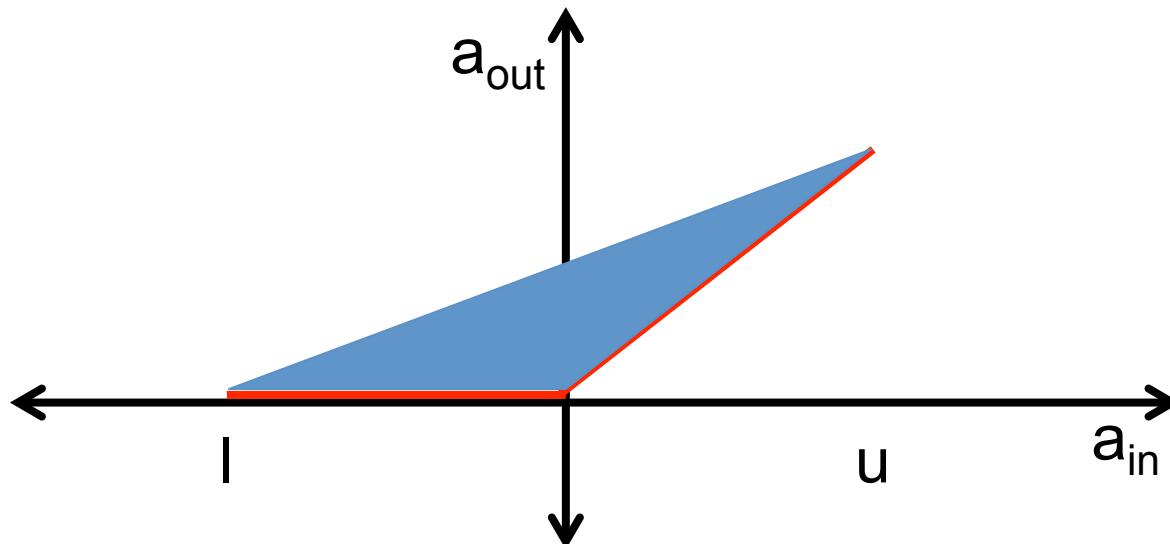


Set of feasible $(a_{\text{in}}, a_{\text{out}})$ is not convex

Convex Relaxation

$$a_{\text{out}} = \max\{a_{\text{in}}, 0\}$$

$$a_{\text{in}} \in [l, u]$$



Ehlers 2017

Replace with smallest convex superset

Original Problem

$$\min \quad y$$

$$\text{s.t.} \quad -2 \leq x_1 \leq 2$$

$$-2 \leq x_2 \leq 2$$

$$a_{in} = x_1 + x_2$$

$$b_{in} = x_1 - x_2$$

$$a_{out} = \max\{a_{in}, 0\}$$

$$b_{out} = \max\{b_{in}, 0\}$$

$$y = -a_{out} - b_{out}$$

Convex Relaxation

Linear Program

$$\min \quad y$$

$$\text{s.t.} \quad -2 \leq x_1 \leq 2$$

$$-2 \leq x_2 \leq 2$$

$$a_{in} = x_1 + x_2$$

$$b_{in} = x_1 - x_2$$

$$a_{out} \geq 0, \quad a_{out} \geq a_{in}, \quad a_{out} \leq 0.5a_{in} + 2$$

$$b_{out} \geq 0, \quad b_{out} \geq b_{in}, \quad b_{out} \leq 0.5b_{in} + 2$$

$$y = -a_{out} - b_{out}$$

General Formulation

$$\min_{\mathbf{x}} \quad \mathbf{c}^T \mathbf{x}$$

s.t.

$$A_k(\mathbf{x})$$

$$k = 1, \dots, L$$

$$B_k(\mathbf{x})$$

$$k = 1, \dots, L$$

Linear mapping of the k-th layer (“easy” constraints)

General Formulation

$$\min_{\mathbf{x}} \quad \mathbf{c}^T \mathbf{x}$$

$$\text{s.t.} \quad A_k(\mathbf{x}) \quad k = 1, \dots, L$$

$$B_k(\mathbf{x}) \quad k = 1, \dots, L$$

Linear mapping of the k-th layer (“easy” constraints)

Convex hull of activations (relaxations)

Outline

- Neural Network Bounds
- Lagrangian Decomposition
- Proximal Minimization
- Results

Additional Variables

$$\min_{\mathbf{x}, \mathbf{x}_k, \mathbf{z}_k} \mathbf{c}^T \mathbf{x}$$

$$\text{s.t.} \quad A_k(\mathbf{x}) \quad k = 1, \dots, L$$

$$B_k(\mathbf{x}) \quad k = 1, \dots, L$$

Additional Variables

$$\min_{\mathbf{x}, \mathbf{x}_k, \mathbf{z}_k} \quad \mathbf{c}^T \mathbf{x}$$

$$\text{s.t.} \quad A_k(\mathbf{x}_k) \quad k = 1, \dots, L$$

$$B_k(\mathbf{z}_k) \quad k = 1, \dots, L$$

$$\mathbf{x}_k = \mathbf{x} \quad k = 1, \dots, L$$

$$\mathbf{z}_k = \mathbf{x} \quad k = 1, \dots, L$$

Partial Lagrangian

$$\max_{\lambda, \mu} \quad \min_{x, x_k, z_k} \quad \mathbf{c}^T x + \sum_k \lambda_k^T (x_k - x) + \sum_k \mu_k^T (z_k - x)$$

$$\text{s.t.} \quad A_k(x_k) \quad k = 1, \dots, L$$

$$B_k(z_k) \quad k = 1, \dots, L$$

Eliminating x using KKT conditions

Lagrangian Decomposition

$$\max_{\lambda, \mu} \quad \min_{x_k, z_k} \quad \sum_k \lambda_k^T x_k + \sum_k \mu_k^T z_k$$

$$\text{s.t.} \quad A_k(x_k) \quad k = 1, \dots, L$$

$$B_k(z_k) \quad k = 1, \dots, L$$

$$\sum_k \lambda_k + \sum_k \mu_k = c$$

Projected subgradient descent (Dvijotham et al., 2018)

Tuning step-size schedule is hard

Outline

- Neural Network Bounds
- Lagrangian Decomposition
- **Proximal Minimization**
- Results

Proximal Minimization

$$\max_{\lambda, \mu} \min_{x_k, z_k} \sum_k \lambda_k^T x_k + \sum_k \mu_k^T z_k + \eta \|\lambda - \lambda'\|^2 + \eta \|\mu - \mu'\|^2$$

$$\text{s.t.} \quad A_k(x_k) \quad k = 1, \dots, L$$

$$B_k(z_k) \quad k = 1, \dots, L$$

$$\sum_k \lambda_k + \sum_k \mu_k = c$$

Initialize $\lambda = \lambda'$ and $\mu = \mu'$ Solve proximal problem

Proximal Minimization

$$\max_{\lambda, \mu} \min_{x_k, z_k} \sum_k \lambda_k^T x_k + \sum_k \mu_k^T z_k + \eta \|\lambda - \lambda'\|^2 + \eta \|\mu - \mu'\|^2$$

s.t. $A_k(x_k) \quad k = 1, \dots, L$

$$B_k(z_k) \quad k = 1, \dots, L$$

$$\sum_k \lambda_k + \sum_k \mu_k = c$$

Smooth quadratic dual Conditional gradient

Repeatedly solve proximal problems until convergence

Advantages

Conditional gradient of dual is the subgradient

Bach, 2015

Analytically computable optimal step-size

A single easy-to-tune parameter η

Highly parallelizable over GPUs

Outline

- Neural Network Bounds
- Lagrangian Decomposition
- Proximal Minimization
- **Results**

Network Architecture

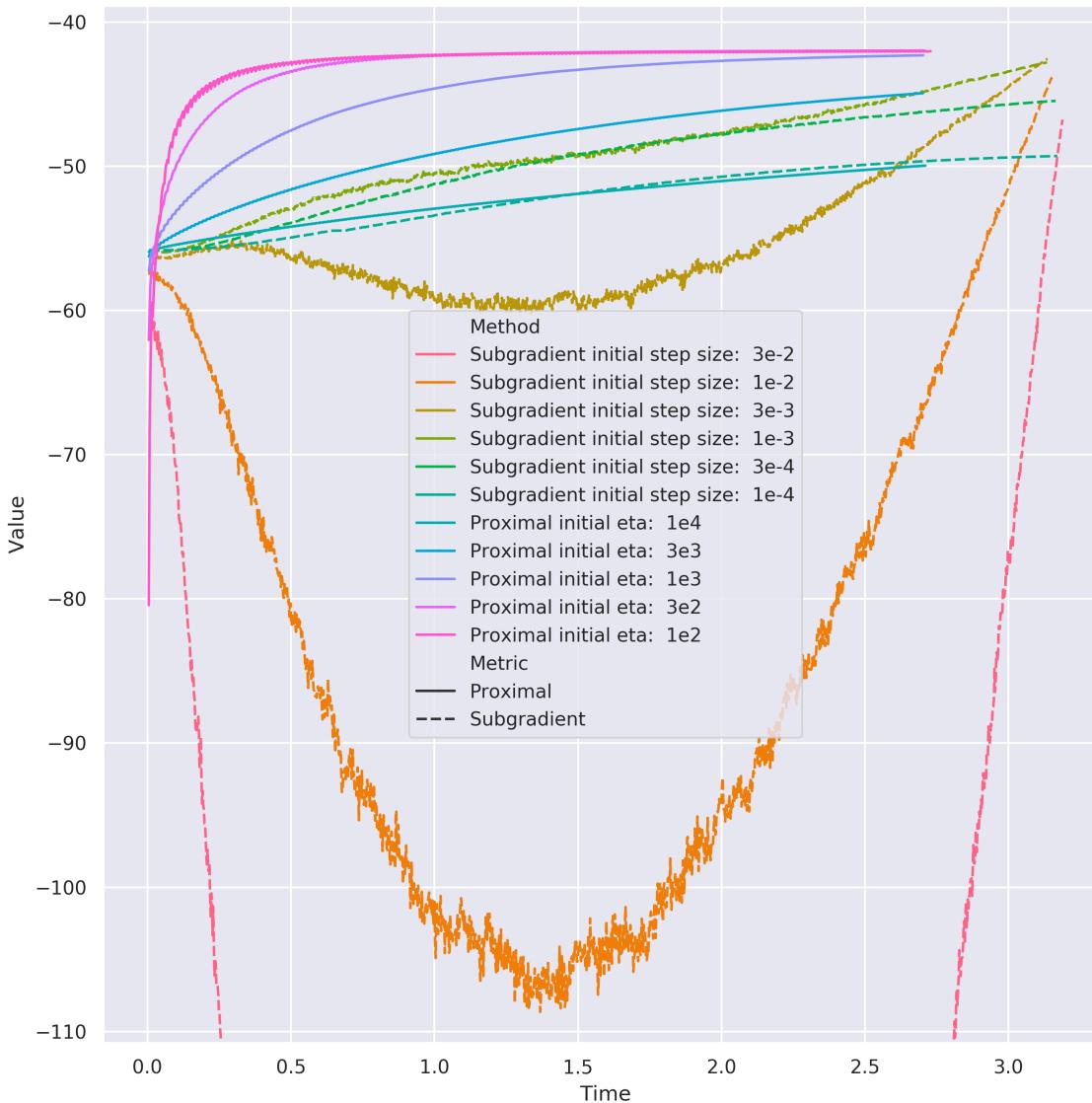
Conv2D(channels=16, kernel_size=4, stride=2)
ReLU

Conv2D(channels=32, kernel_size=4, stride=2)
ReLU

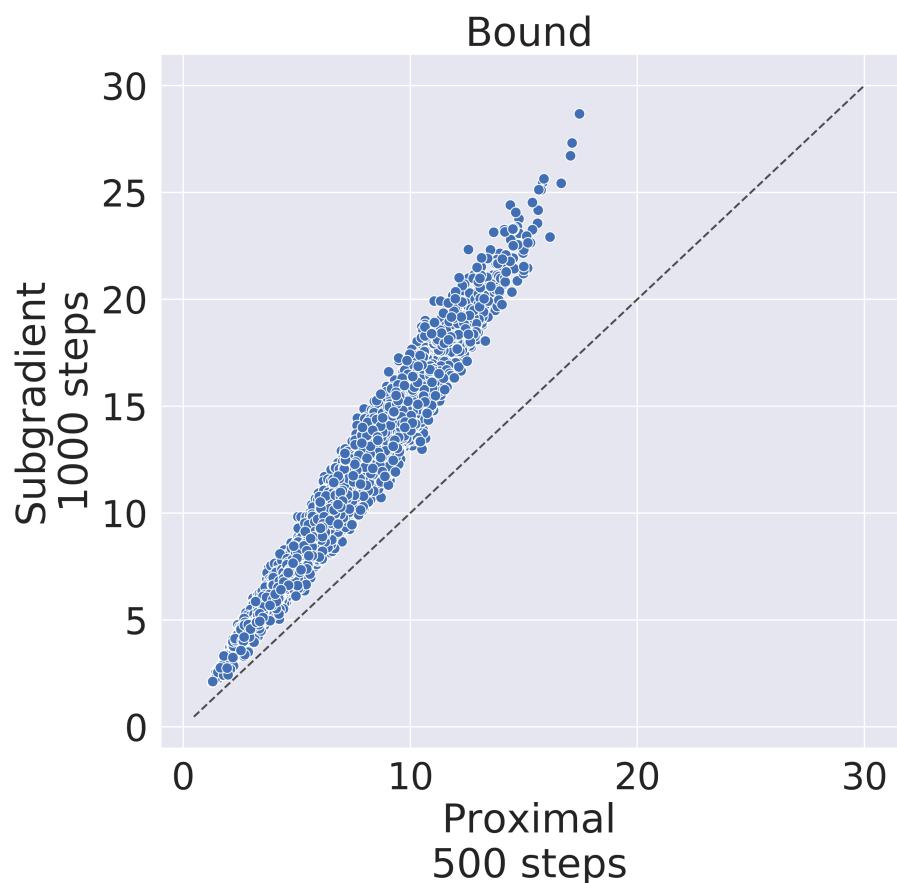
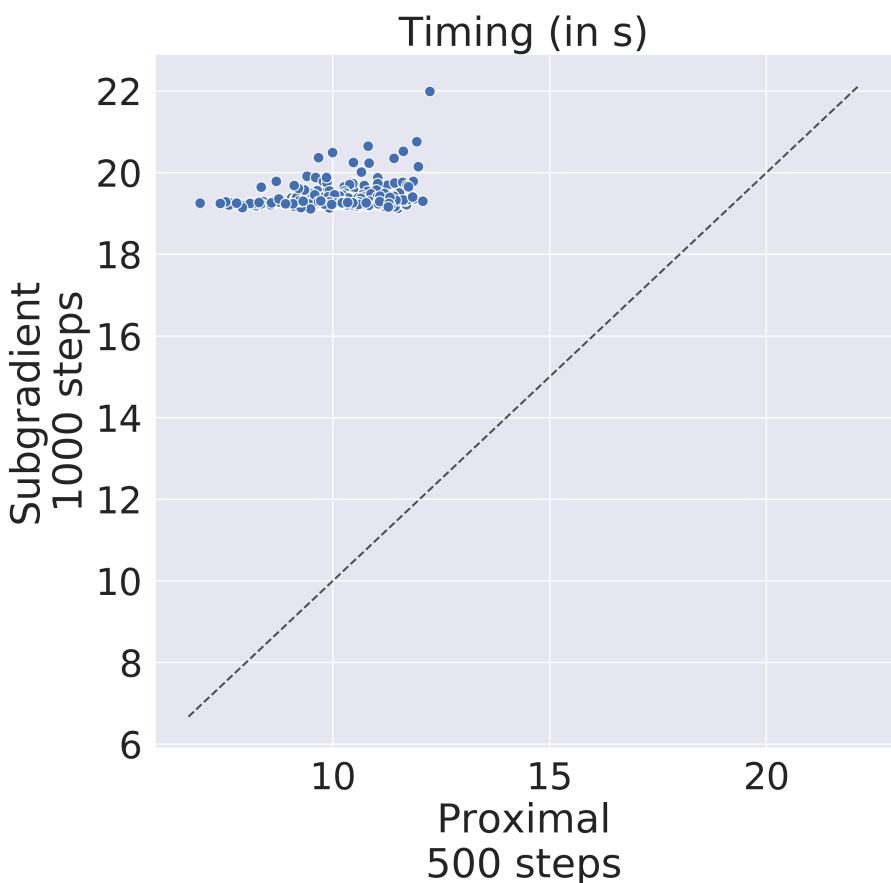
Linear(channels=100)
ReLU

Linear(channels=10)

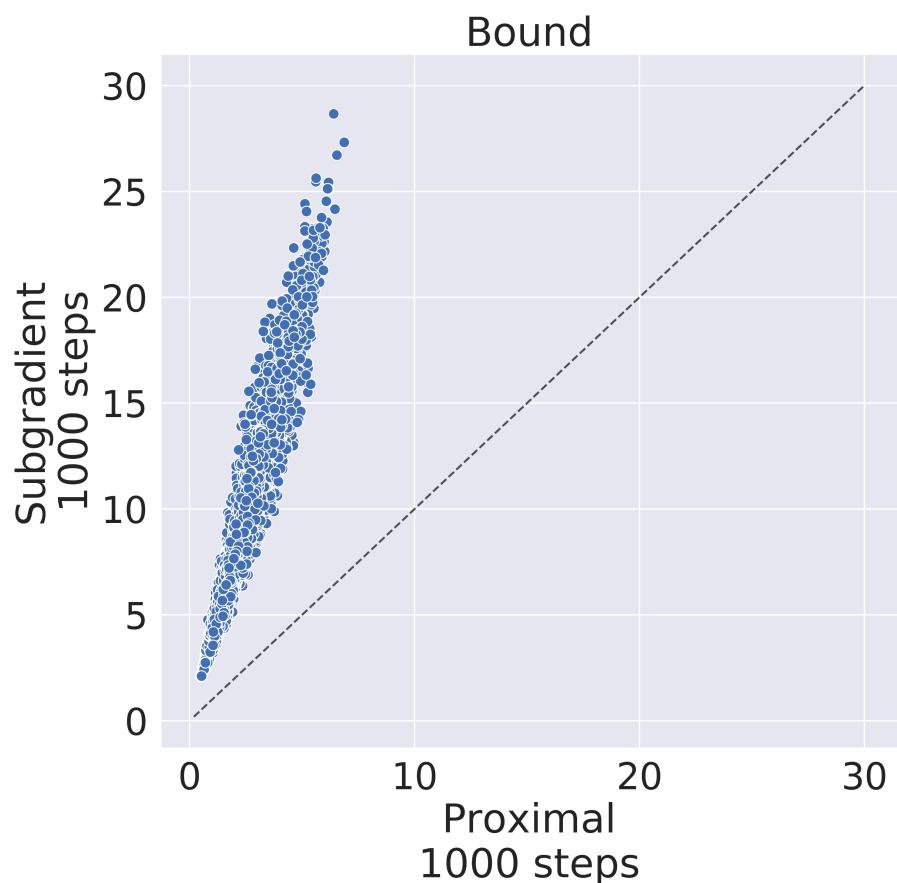
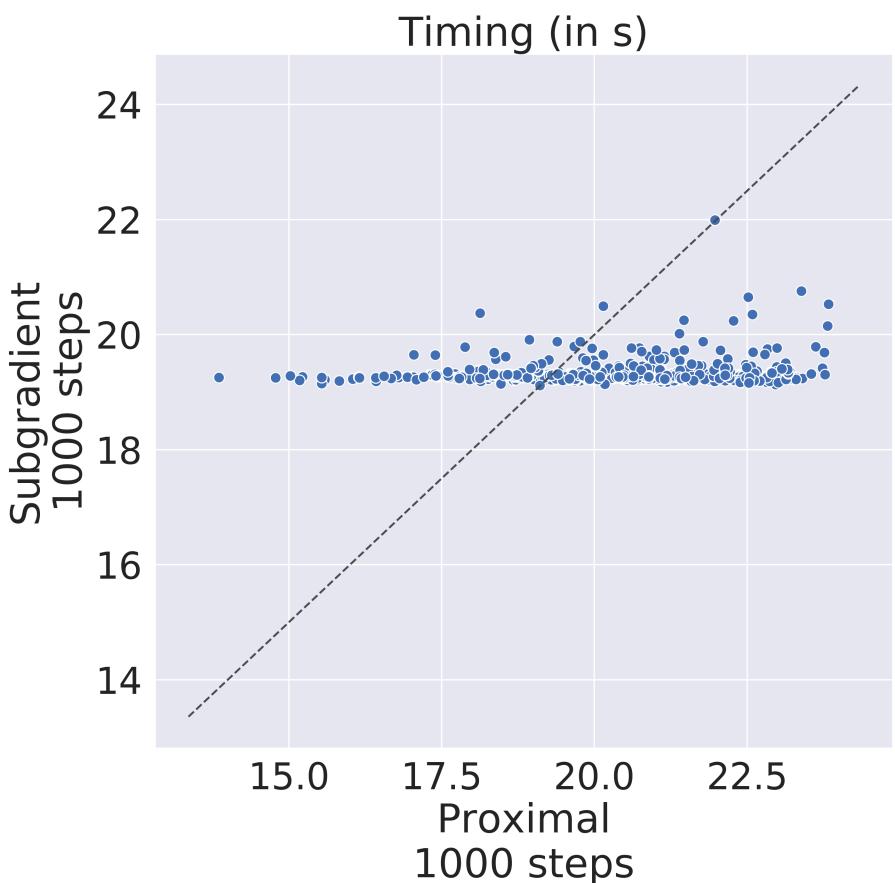
Results



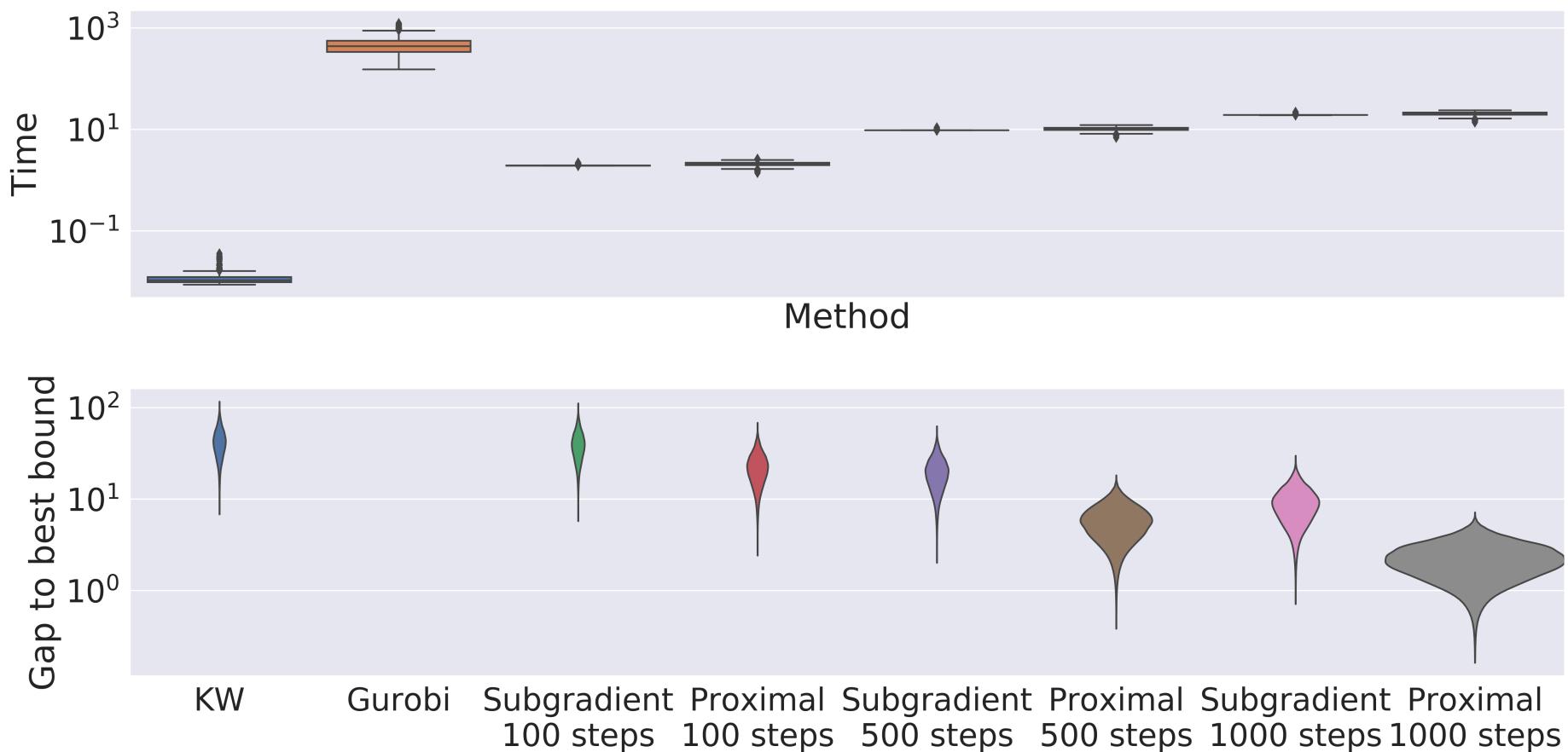
Results



Results



Results



Questions?

Paper + Code + Data Available*

<http://www.robots.ox.ac.uk/~oval>