

DeepXplore: Automated Whitebox Testing of Deep Learning Systems

by Kexin Pei, Yinzhi Cao, Junfeng Yang, Suman Jana

Eric Langlois

University of Toronto

March 11, 2019

Fast automatic generation of test inputs for a set of neural networks, where

- the networks disagree, and
- the examples have high diversity.



Contributions

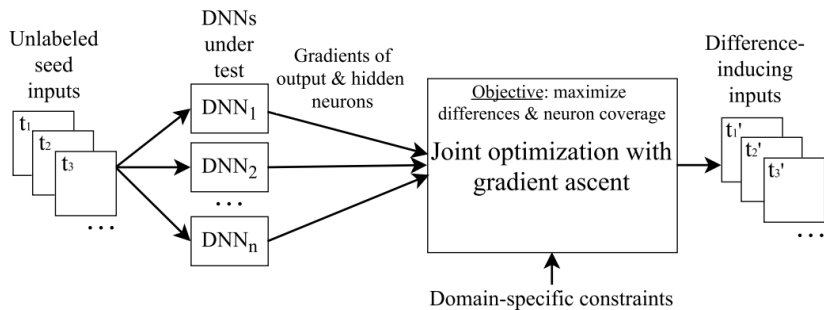
Introduced **neuron coverage** as a testing metric for DL systems.

Formulated the task of finding **behaviour differences** in a set of networks as a gradient descent optimization problem.

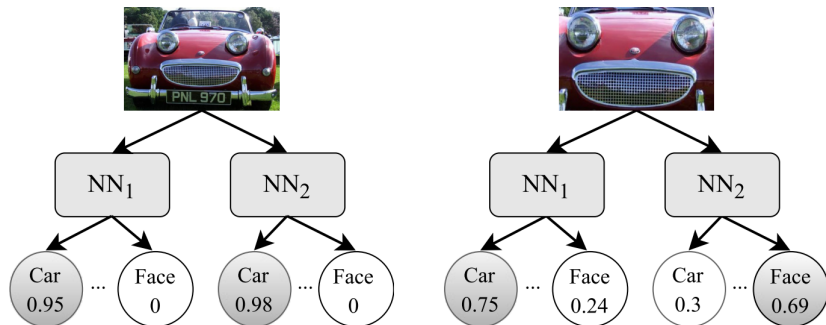
Created the **DeepXplore** open-source deep learning testing framework.

Showed that training on DeepXplore-generated tests can **increase classification accuracy**.

Algorithm



Network Prediction Differences

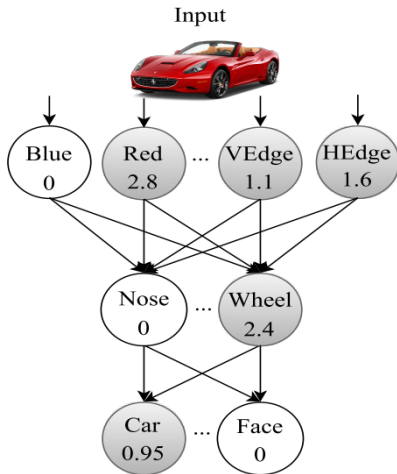
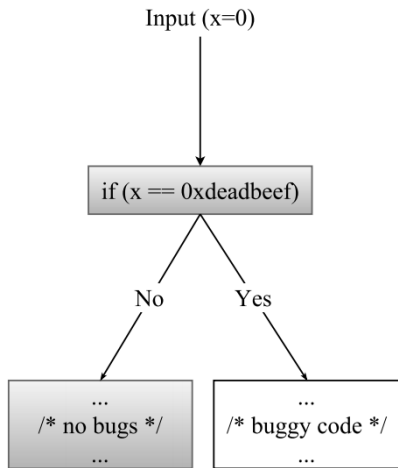


$$\text{obj}_1(\mathbf{x}) = \sum_{k \neq j} F_k(\mathbf{x})[c] - \lambda_1 F_j(\mathbf{x})[c]$$

$F_k(\mathbf{x})[c]$ = The probability according to network k that input \mathbf{x} is class c .

j is chosen randomly.

Coverage



Input Set Coverage

For a test set T , neuron set N , and threshold t

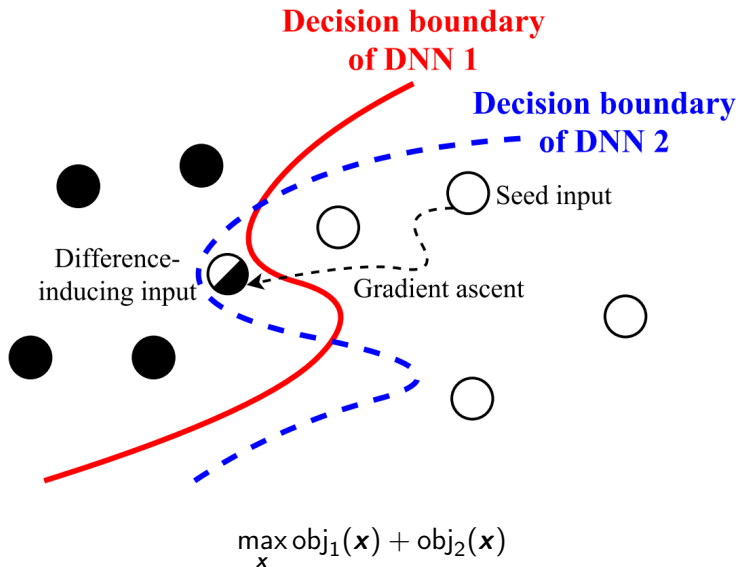
$$\text{NCov}(T, t) = \frac{|\{n \in N \mid \exists \mathbf{x} \in T, \text{out}(n, \mathbf{x}) > t\}|}{|N|}$$

Coverage Loss

$$\text{obj}_2(\mathbf{x}) = \sum_k \text{value of neuron } n_k \text{ in network } k \text{ on input } \mathbf{x}$$

n_k is chosen randomly among neurons that are not yet covered.

Test Inputs via Optimization



Constraints

Generated inputs should be realistic.
Enforce this with domain-specific constraints.

Images

- Pixel value bounds: [0, 255]
- Only modify brightness
- Only modify a small region
- Only add small black boxes

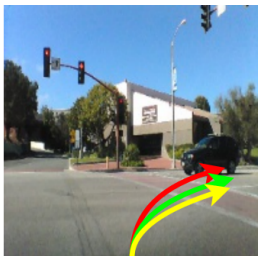
Constraints applied by modifying the gradient

$$\mathbf{x} \leftarrow \mathbf{x} + \text{step_size} \cdot \text{constrain_grad}\left(\frac{\partial \text{obj}(\mathbf{x})}{\mathbf{x}}\right)$$

Examples — Constraint: Brightness



all:right



all:right



all:right



DRV_C1:left



DRV_C2:left



DRV_C3:left

Examples — Constraint: Occlusion



all:right



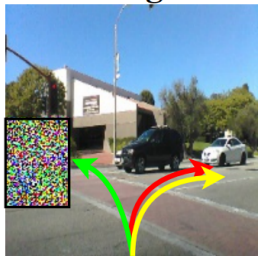
all:right



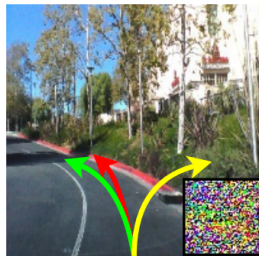
all:left



DRV_C1:left

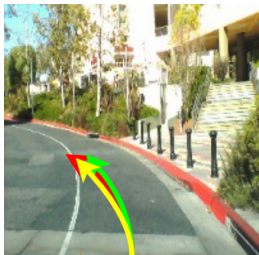


DRV_C2:left



DRV_C3:right

Examples — Constraint: Black Boxes



all:left



all:left



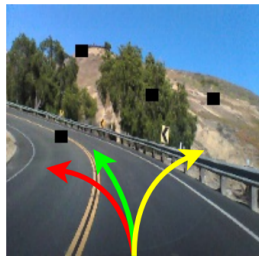
all:left



DRV_C1:right



DRV_C2:right



DRV_C3:right

Domains

- MNIST** : Classify handwritten digits
- ImageNet** : Classify images
- Driving** : Predict steering angle from images
- Contagio** : Classify PDF malware
- Drebin** : Classify android app malware

Networks

- 3 networks each
- Pre-trained networks or based on popular architecture

Coverage

- Threshold $t = 0, 0.25, \text{ or } 0.75$

Efficiency

Model differences are found for 40% – 100% of seed inputs.

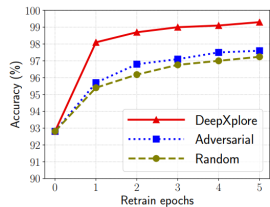
Speed

First difference-inducing input is found in seconds,
100% coverage is achieved in 6s – 200s

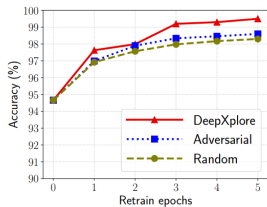
Design Validation

Using coverage moderately increases example diversity (by L1 distance)
Increasing model dissimilarity increases the number of differences found

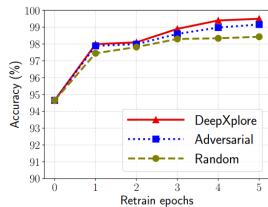
Application: Training Data Augmentation



(a) LeNet-1



(b) LeNet-4



(c) LeNet-5

Add generated inputs to the training data and retrain.

- How meaningful are the generated examples?
- Is it reasonable to use model disagreement as an objective?
- Are the constraints plausible?
- What are some alternative coverage measures?