



Towards Testing of Deep Learning Systems

Jianjun Zhao
Kyushu University

iMLSE 2018 (Nara, Japan, December 4, 2018)



Pangu Research Group (Kyushu Univ.)

(知能ソフトウェア工学研究室)

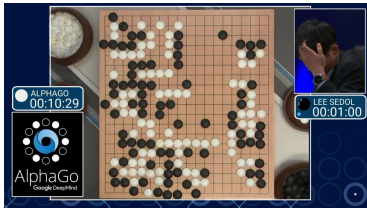
<https://pangukaitian.github.io/pangu/en/index.html>

- On-going work
 - researches on the potential symbioses between software engineering **SE** and artificial intelligence **AI**
- The overall goal
 - to obtain better software and AI systems making them **more robust, reliable, and secure**, and **easier to specify, build, maintain, or improve**
- Group members
 - 2+ faculties and 3 PhD and 7 MS students


2

Deep Learning Matches Human Intelligence

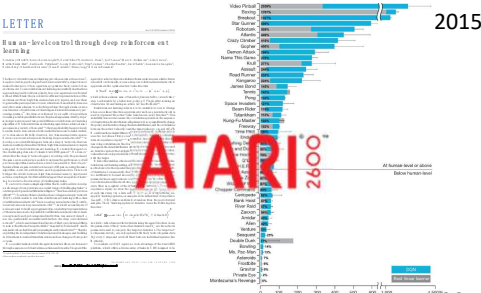
AlphaGo 4:1 Human Champion 2016




AlphaGo ZERO 100:0 AlphaGo 2017



2015



2017



3

Wide Deployment in Real World















4

Driving Force of Many Novel Technologies



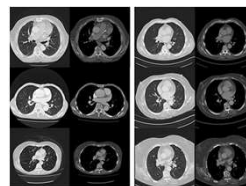
5

Deep learning is increasingly used in safety-critical systems

- Deep learning reliability and safety is crucial



Self-driving car



Medical diagnosis



Robotics



Malware detection

6

Problems and Critiques for DL Systems

- Un-safe, e.g., lack of robustness (reliability and safety)
- Hard to explain to human users (interpretability)
 - Deep neural networks are essentially black-boxes and researchers have a hard time to understand how they deduce conclusions
- Fairness, accountability, ethics, trustworthiness, etc.
 - What would human review entail if models were available for direct inspection?

7

Problems and Critiques for DL Systems

- Un-safe, e.g., lack of robustness (reliability and safety)
- Hard to explain to human users (interpretability)
 - Deep neural networks are essentially black-boxes and researchers have a hard time to understand how they deduce conclusions
- Fairness, accountability, ethics, trustworthiness, etc.
 - What would human review entail if models were available for direct inspection?

8

Current Deep Learning is Vulnerable



Classified as panda

Small adversarial noise

Classified as gibbon

Ian Goodfellow, Jon Shlens, Christian Szegedy, Explaining and Harnessing Adversarial Examples, ICLR, 2014



Accident

9

Unreliable Deep Learning



Tesla autopilot failed to recognize a white truck against bright sky leading to fatal crash

10

How to assure the quality of DL systems?

Al image recognition fooled by single pixel change
© 9 November 2017

Tesla in fatal California crash was on Autopilot
© 31 March 2018

The New York Times
Alexa and Siri Can Hear This Hidden Command. You Can't.
Researchers can now send secret audio instructions undetectable to the human ear to Apple's Siri, Amazon's Alexa and Google's Assistant.

日本経済新聞
2018年9月5日 (水)
人間は自動運転車を信頼できる？
自動運転 B P 速報

IBM
Three former managers say IBM fired them because they spoke up against cutting a black software-sales rep's commission. AP file
BUSINESS
IBM says it's reaching for the 'moon' with Watson Health. That hasn't stopped layoffs.

To design reliable systems, engineers typically engage in both testing and verification

- **By testing:** we mean evaluating the system in several conditions and observing its behavior, watching for defects.
- **By verification:** we mean producing a compelling argument that the system will not misbehave under a very broad range of circumstances.

* Ian Goodfellow and Nicolas Papernot. 2017. The Challenge of Verification and Testing of Machine Learning.

To design reliable systems, engineers typically engage in both testing and verification

- **By testing:** we mean evaluating the system in several conditions and observing its behavior, watching for defects. (this talk focuses on testing)
- **By verification:** we mean producing a compelling argument that the system will not misbehave under a very broad range of circumstances.

* Ian Goodfellow and Nicolas Papernot. 2017. The Challenge of Verification and Testing of Machine Learning.

13

Testing Issues for DL Systems

- **Test coverage criterion**
 - How to define the test coverage criteria of DL systems?
- **Test data generation**
 - How to automatically generate a mass of test data for DL systems?
- **Test data quality**
 - How to evaluate the quality of test data for DL systems?

14

Testing Issues for DL Systems

- **Test coverage criterion**
 - How to define the test coverage criteria of DL systems?
- **Test data generation**
 - How to automatically generate a mass of test data for DL systems?
- **Test data quality**
 - How to evaluate the quality of test data for DL systems?

15

Quality Assurance for Traditional Software Testing Criteria and Tools

- Line Coverage
- Branch Coverage
- Function Coverage
- Data Flow Coverage
- Combinatorial Coverage
- Mutation testing Coverage

JACOBO
Java Code Coverage

NCOVER

µJava



AgitarOne
(Continue Java To The End)



Atlassian
Clover

Cobertura

EMMA

Major

EVASUITE

K&K

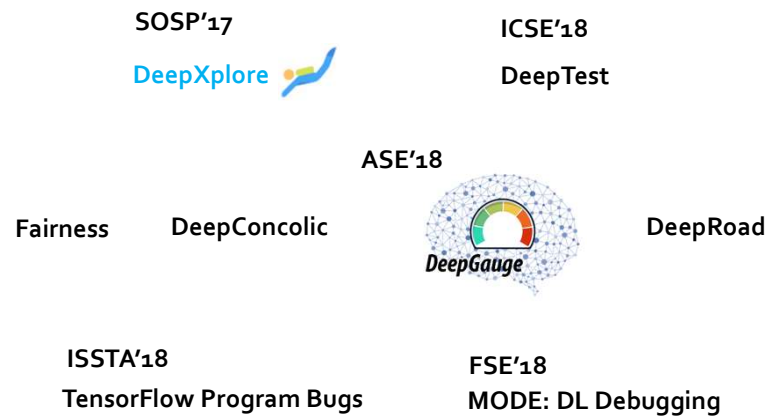


GRT
Test Apps Better

..... **OCELOT**

16

Quality Assurance for DL is at Early Stage



17

Traditional Software

SQA**Deep Learning Software**

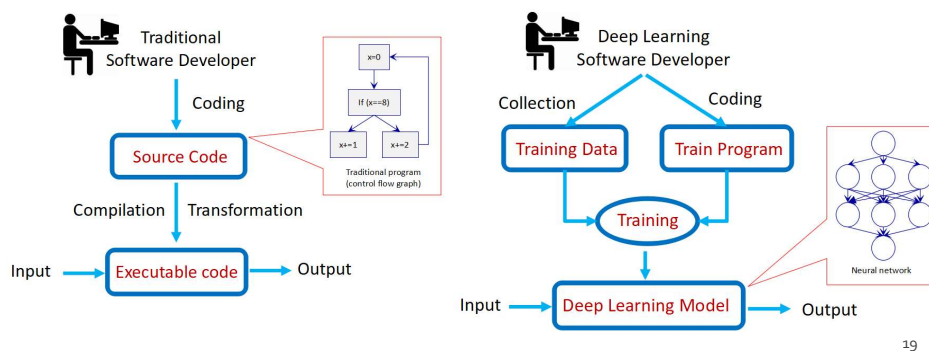
Fundamental Different Programming Paradigm

- The decision logic of a traditional software:

- In the form of code

- The decision logic of a DL system:

- The structure of DNN
- The connection weights



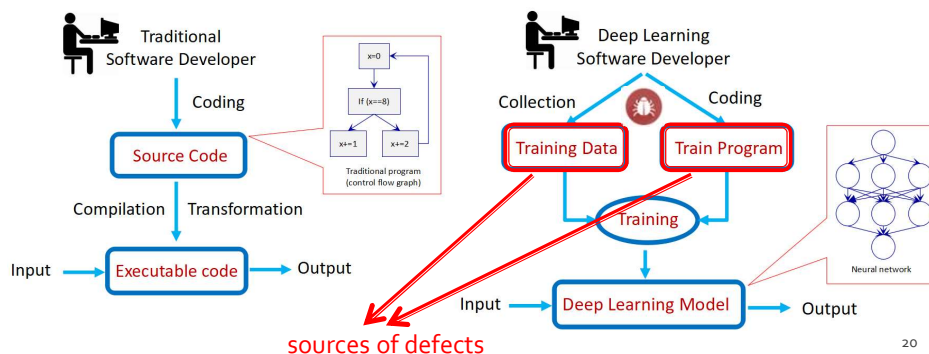
Source of defects: programming paradigm

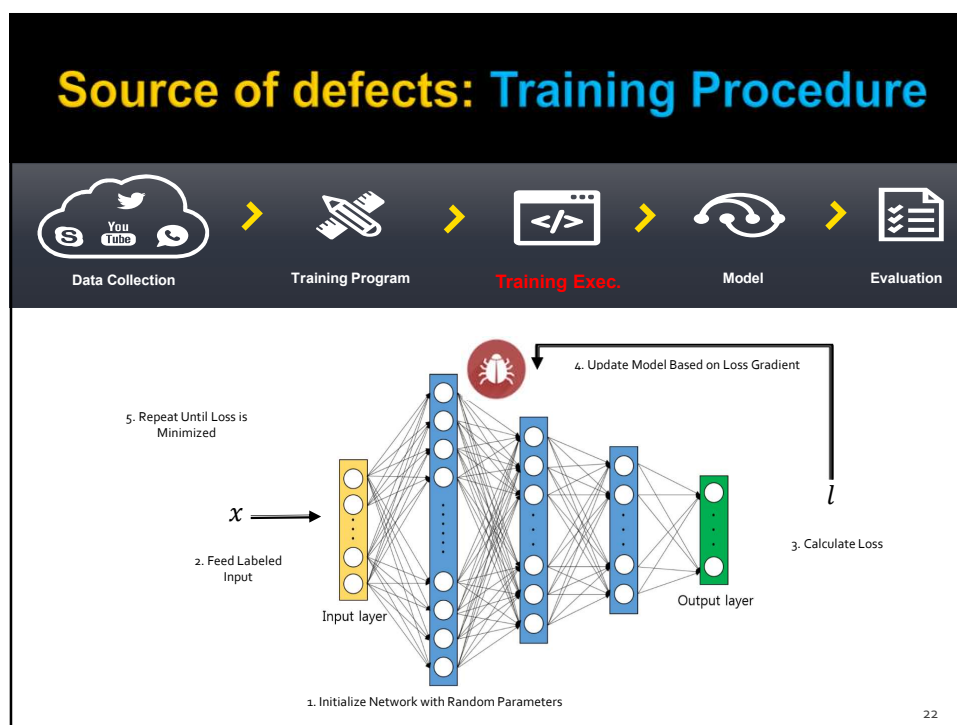
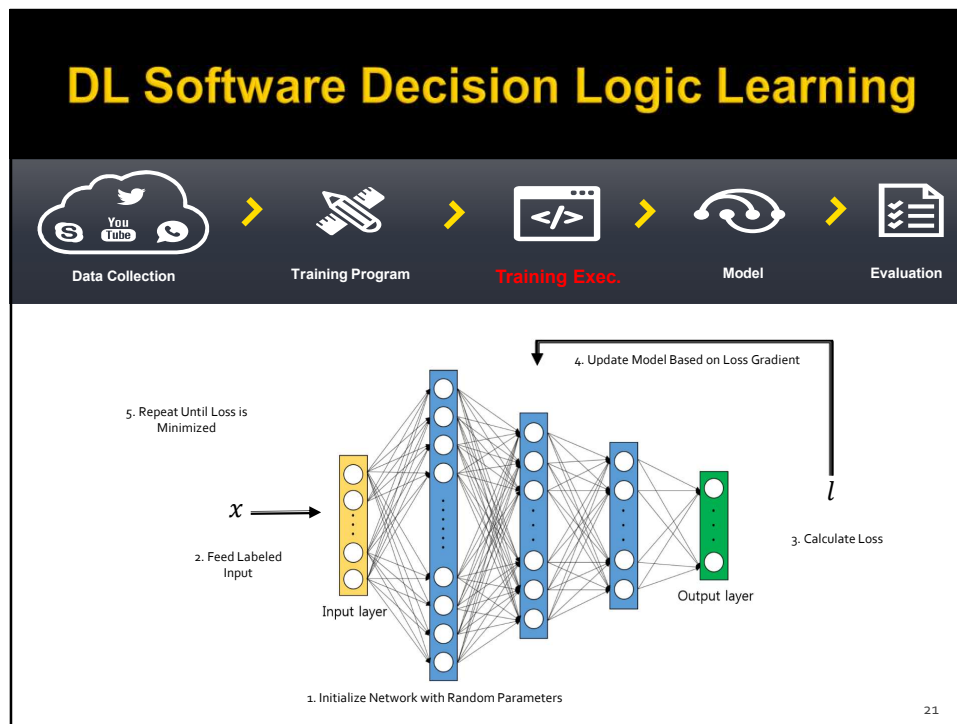
- The decision logic of a traditional software:

- In the form of code

- The decision logic of a DL system:

- The structure of DNN
- The connection weights

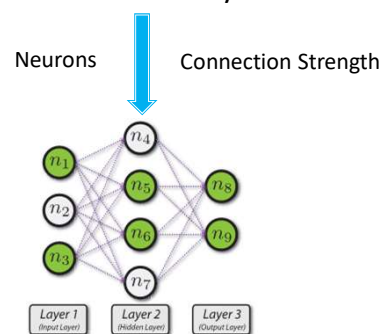




Lacking of Interpretability and Understandability

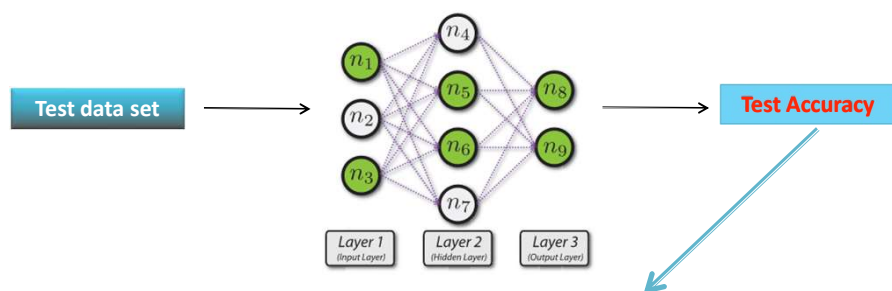
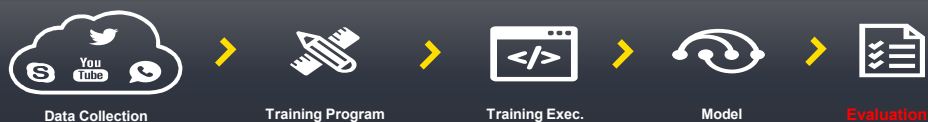


Behavior of the DL system?



23

Quality Measurement Immature



High **accuracy** \nrightarrow High **DL quality**

24

Towards Quality Assurance for DL Systems



- **Multi-granularity testing criteria for DL systems (ASE 2018)**
 - ACM SIGSOFT Distinguished Paper Award



- **Coverage-guided fuzzing testing framework**

25

Towards Quality Assurance for DL Systems



- **Multi-granularity testing criteria for DL systems (ASE 2018)**
 - ACM SIGSOFT Distinguished Paper Award

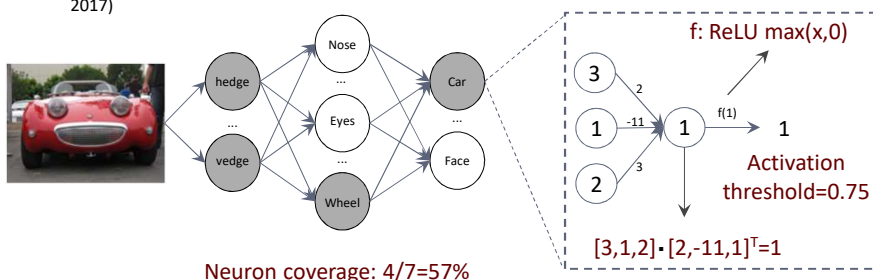


- Coverage-guided fuzzing testing framework

26

Test coverage criteria for DL systems (1)

- **Neuron coverage**: how much decision logic exercised
 - Neuron coverage = # neurons activated / # total neurons
- Kexin Pei, Yinzhi Cao, Junfeng Yang, Suman Jana. "DeepXplore: Automated Whitebox Testing of Deep Learning Systems", in Proceedings of the 26th ACM Symposium on Operating Systems Principles (SOSP 2017)



27

Problems for Neuron Coverage

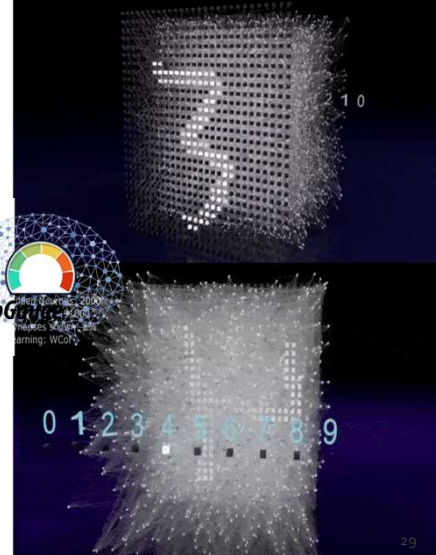
- **Neuron coverage** uses the same threshold as the activation evaluation for all the neurons
- It is straightforward to obtain a trivial test suite that has high **neuron coverage** but does not provide any adversarial example

*Testing Deep Neural Networks. Youcheng Sun, Xiaowei Huang, Daniel Kroening. arXiv:1803.04792, 2017

28

Overview of DeepGauge

- **Enable** testing quality evaluation of DL systems from multiple portrayals
- **Provide** systematic guidance of test generation for detecting defects
- **Facilitate** interpretation & understanding



29

Design of DeepGauge

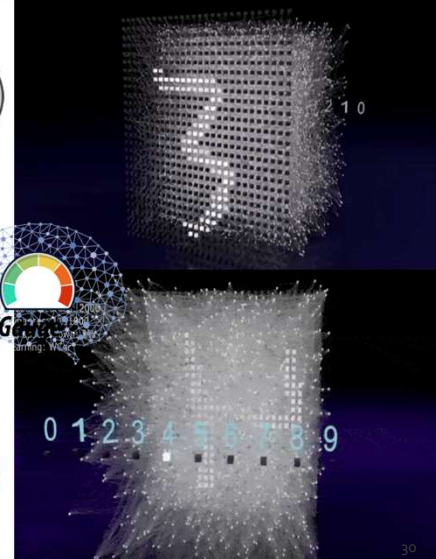
Simple to understand & use

Efficient to compute

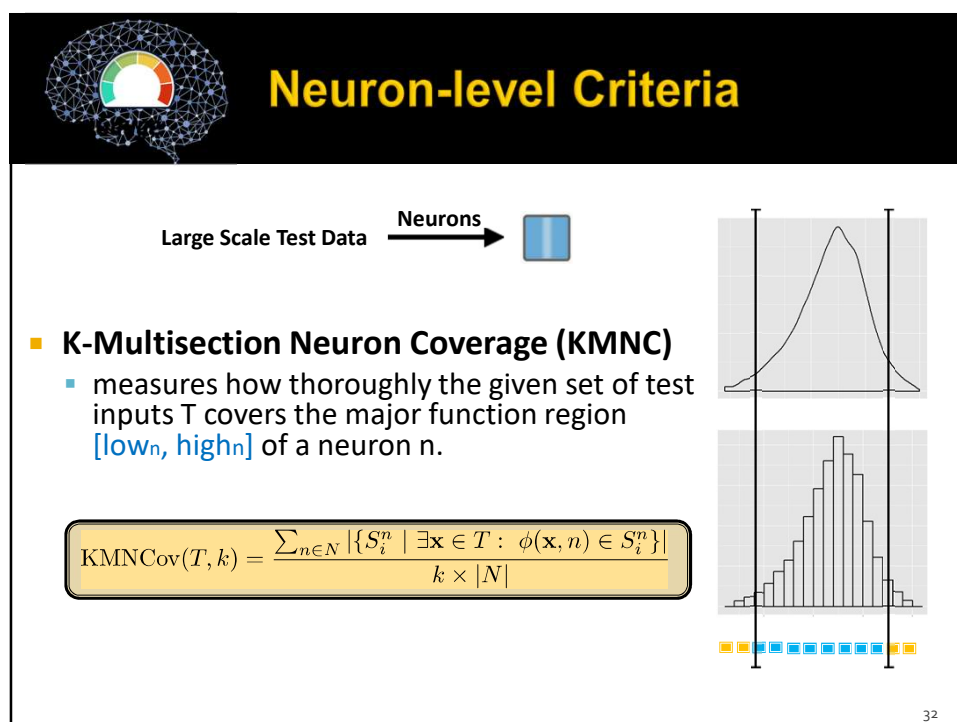
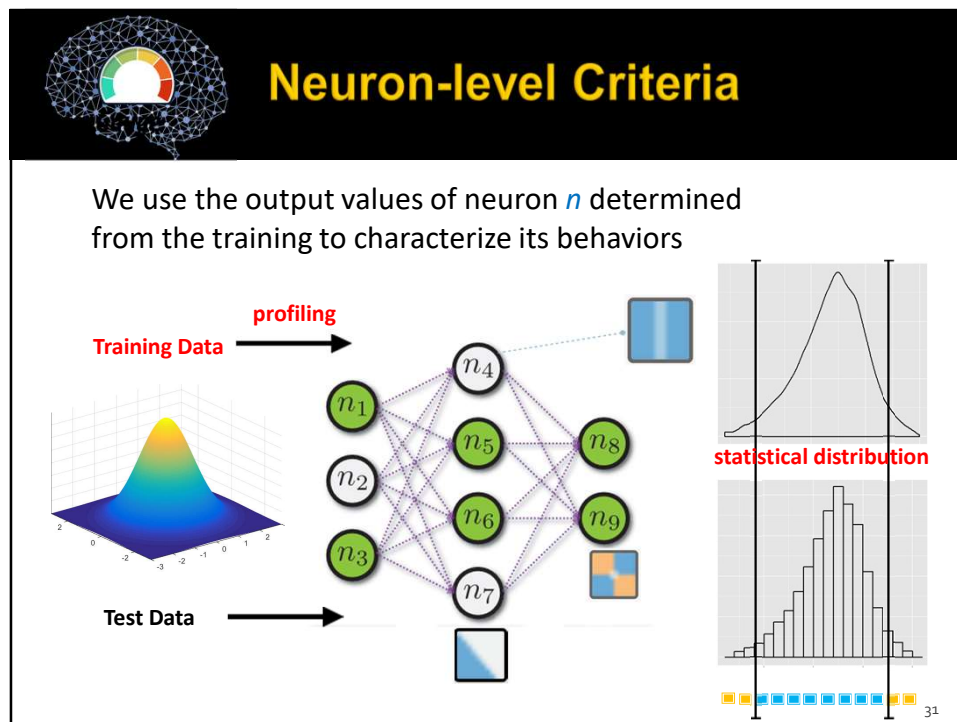
General to diverse DNNs

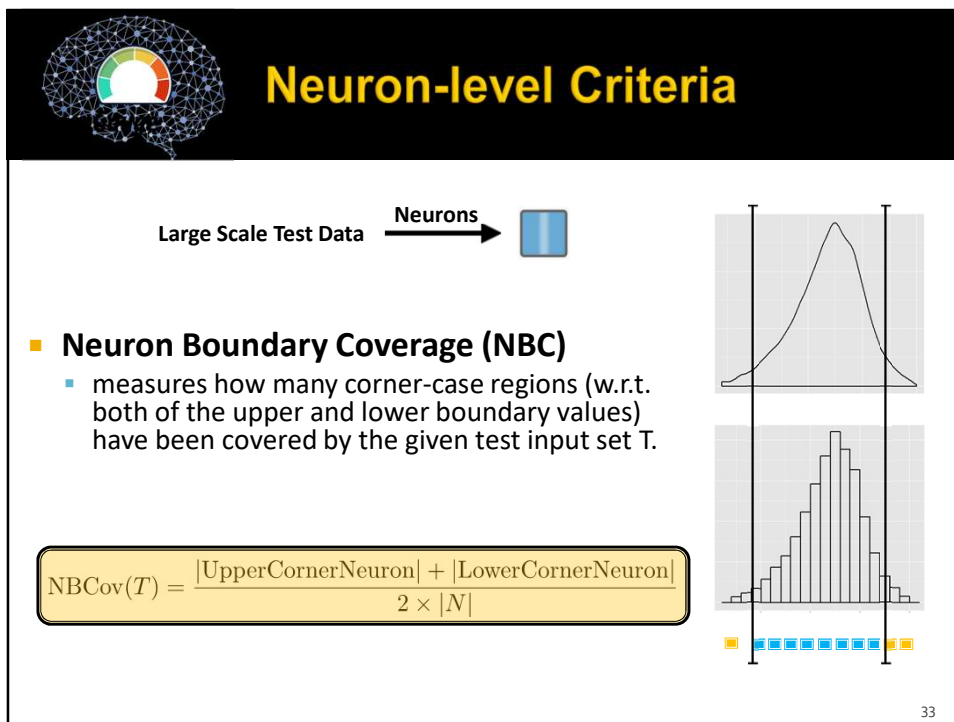
Scale to large DNNs

Adaptable by cases




30

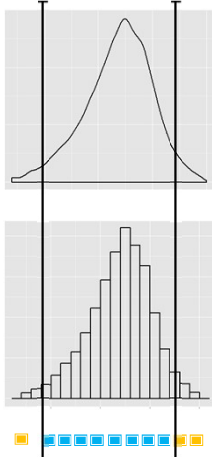




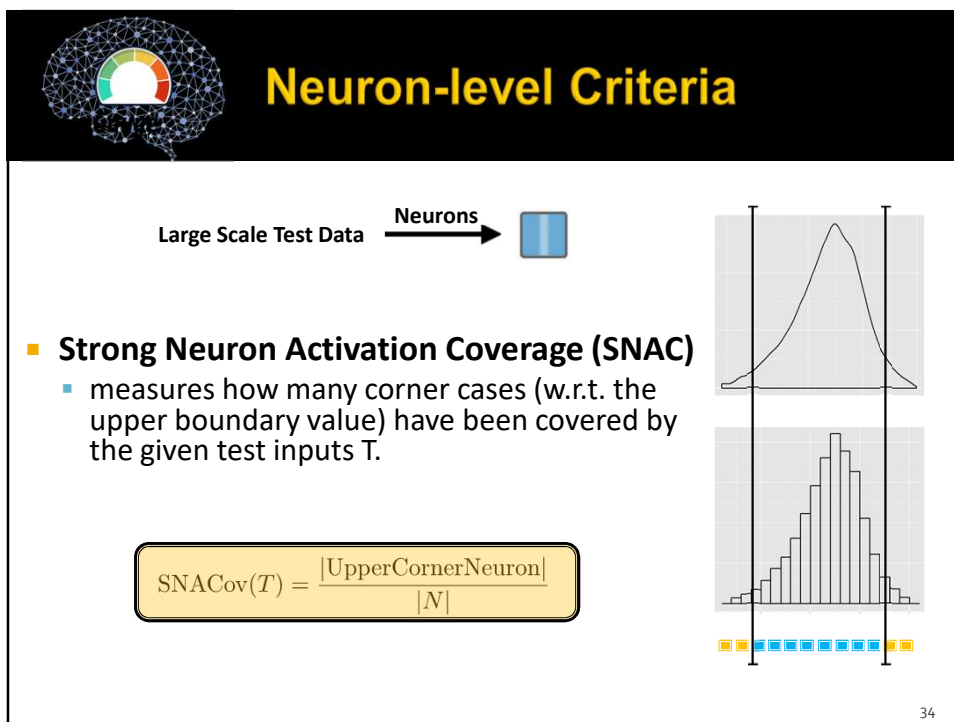
Neuron-level Criteria

Large Scale Test Data $\xrightarrow{\text{Neurons}}$ 


- **Neuron Boundary Coverage (NBC)**
 - measures how many corner-case regions (w.r.t. both of the upper and lower boundary values) have been covered by the given test input set T.

$$\text{NBCov}(T) = \frac{|\text{UpperCornerNeuron}| + |\text{LowerCornerNeuron}|}{2 \times |N|}$$


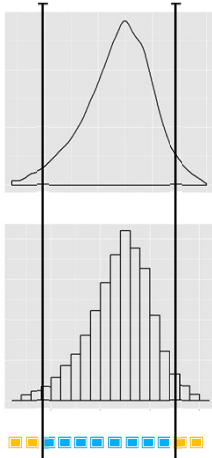
33



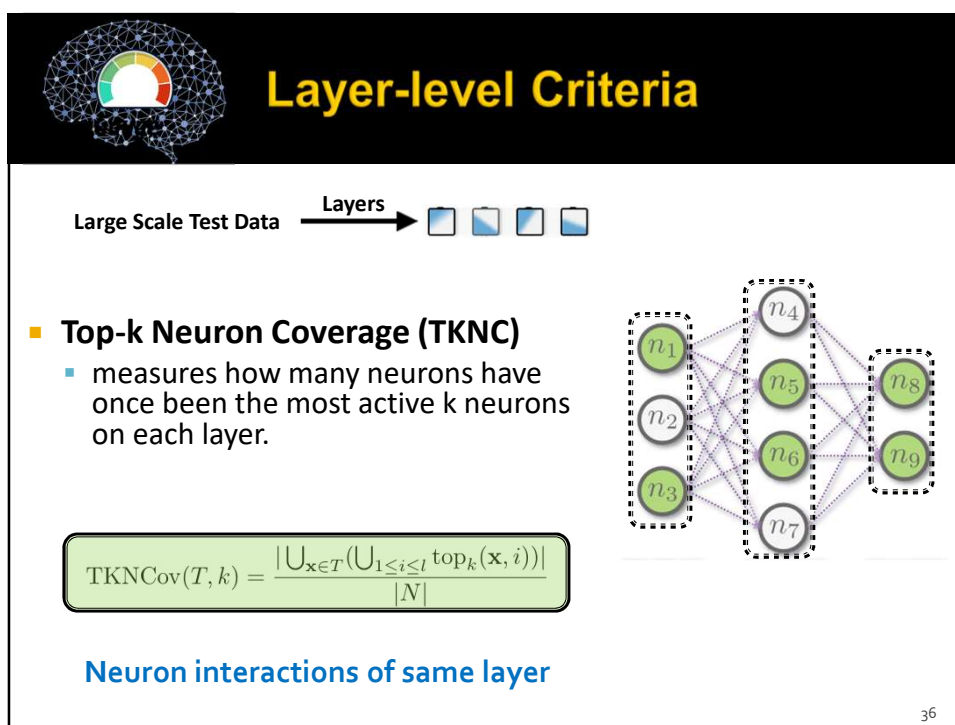
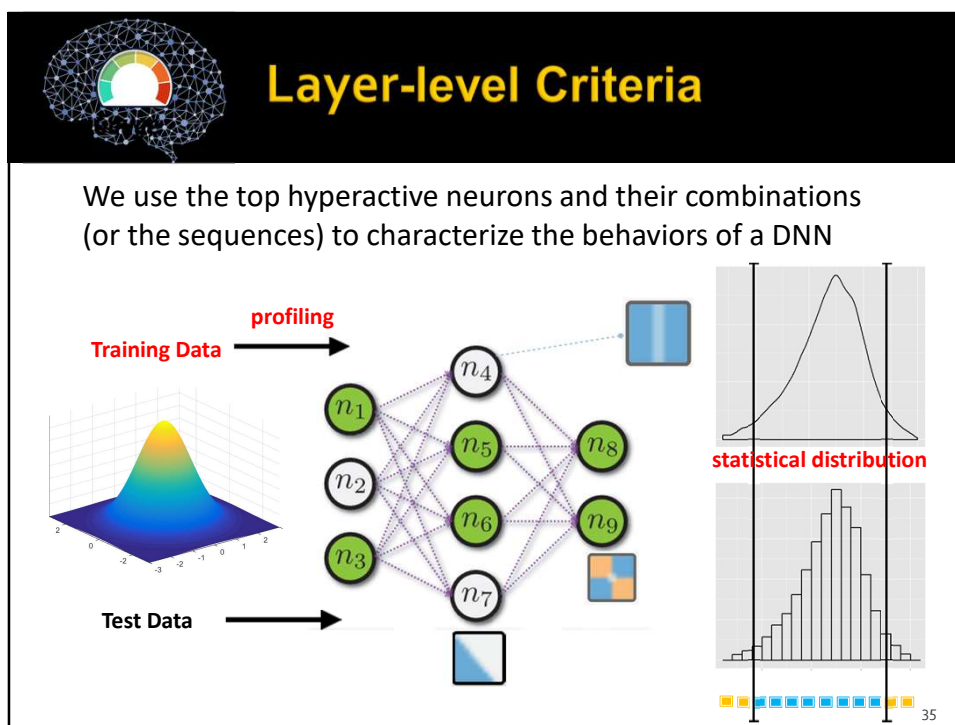
Neuron-level Criteria

Large Scale Test Data $\xrightarrow{\text{Neurons}}$ 


- **Strong Neuron Activation Coverage (SNAC)**
 - measures how many corner cases (w.r.t. the upper boundary value) have been covered by the given test inputs T.


$$\text{SNACov}(T) = \frac{|\text{UpperCornerNeuron}|}{|N|}$$


34



Layer-level Criteria

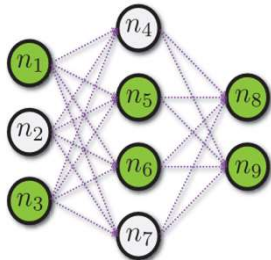


Large Scale Test Data $\xrightarrow{\text{Layers}}$ 

- **Top-k Neuron Patterns (TKNP)**
 - Given a test input x , the sequence of the top-k neurons on each layer forms a pattern

$$\text{TKNPat}(T, k) = |\{(\text{top}_k(x, 1), \dots, \text{top}_k(x, l)) \mid x \in T\}|$$

Neuron interactions across layers



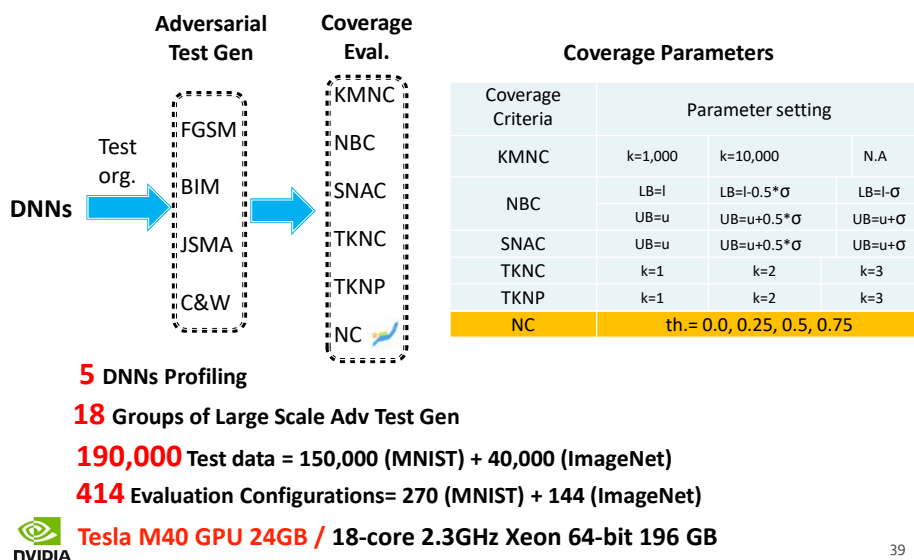
37

Large Scale Empirical Study

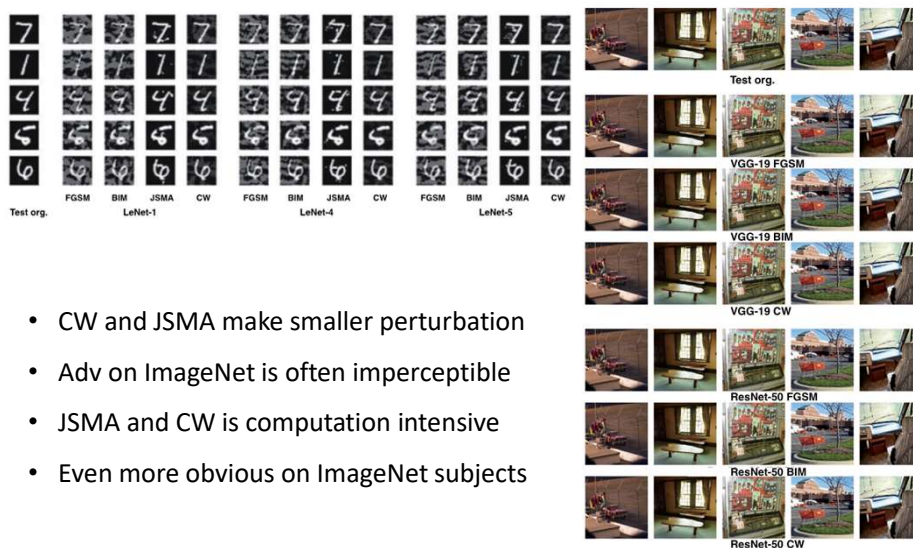
	DNNs	#Neurons	#Layers
MNIST			
60,000			
10,000	LeNet-1	52	7
(28,28,1)			
784 dim.	LeNet-4	148	8
IMAGENET			
(LSVRC-2012)	LeNet-5	268	9
1,000,000+	VGG-19	16,168	25
50,000			
(224,224,3)	ResNet-50	95,059	176
150528 dim.			

38

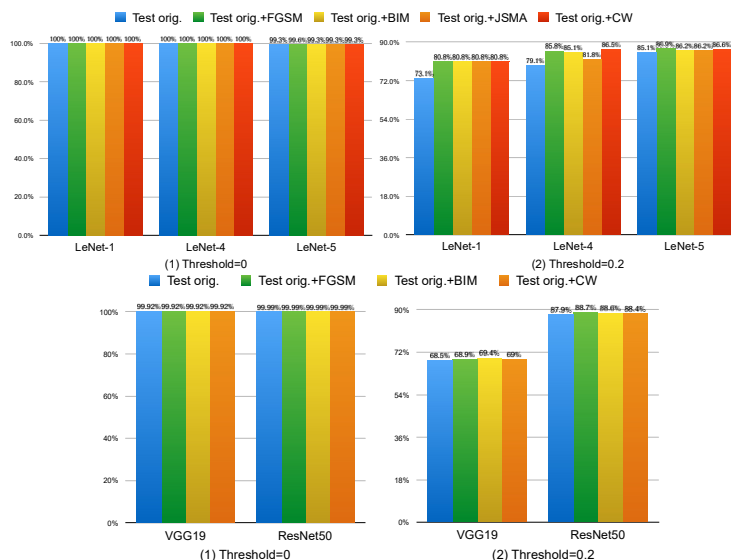
Large Scale Empirical Study



Adversarial Examples

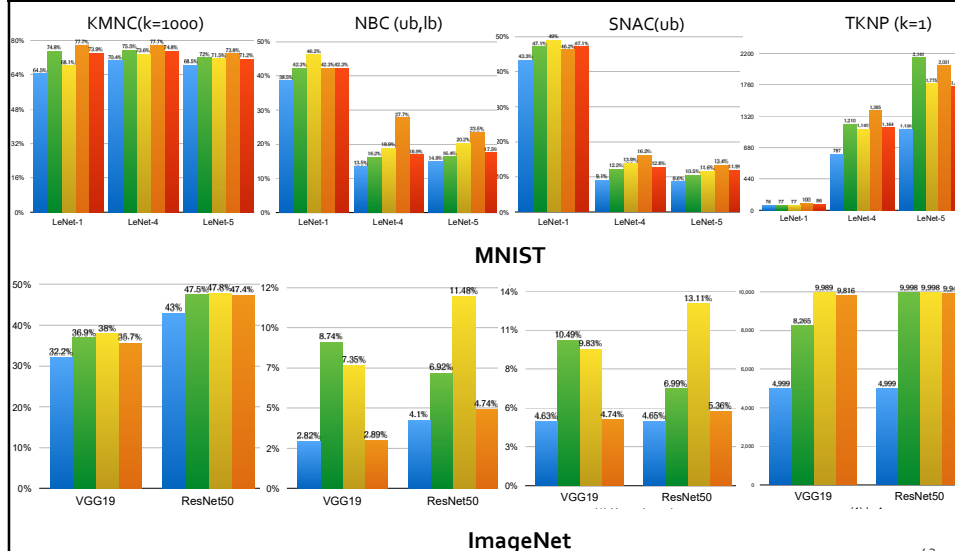


DeepXplore's Neuron Coverage



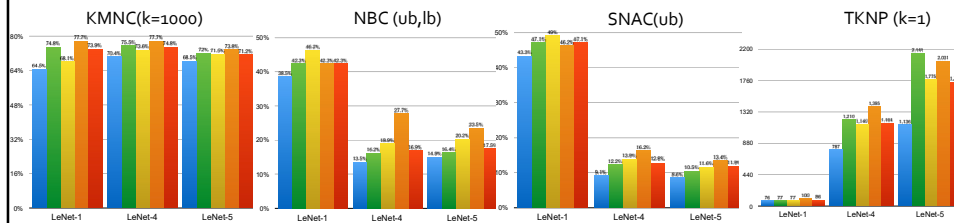
41

DeepGauge Results Glimpse



42

DeepGauge Results Summary (1)



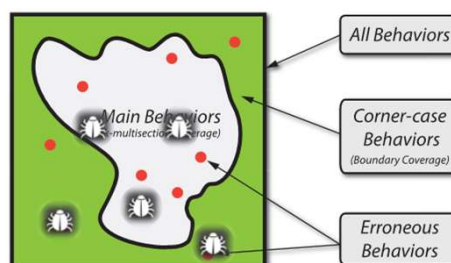
- Our criteria could distinguish benign and defects with little perturbations
- Increase the coverage could potentially increase the chance of defect detection
- TestGen guided by DeepGauge generates thousands of unique error trigger tests

43

DeepGauge Results Summary (2)

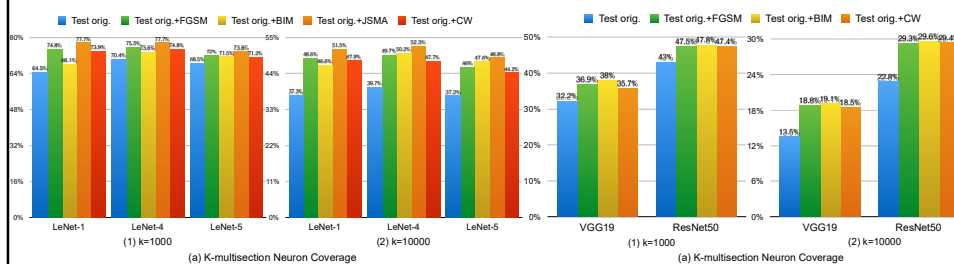


- Defects could occur in both regions, which need to be tested
- KMNC larger than NBC and SNAC, corner cases are difficult to cover



44

DeepGauge Results Summary



- Criteria show different precision with different parameters
- The accuracy could be adapted on specific DNN models
- Precision could be further enhanced efficiently with bucket

45

Towards Quality Assurance for DL Systems



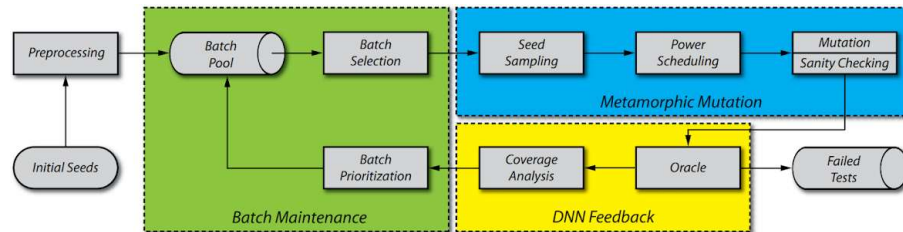
- Multi-granularity testing criteria for DL systems (ASE 2018)
 - ACM SIGSOFT Distinguished Paper Award



- Coverage-guided fuzzing testing framework

46

Overview of DeepHunter

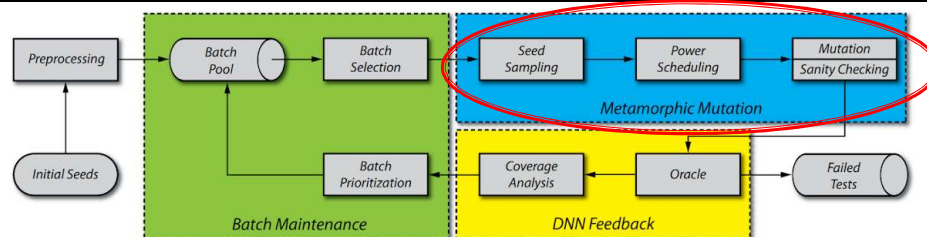


6 coverage criteria



47

Metamorphic Transformation



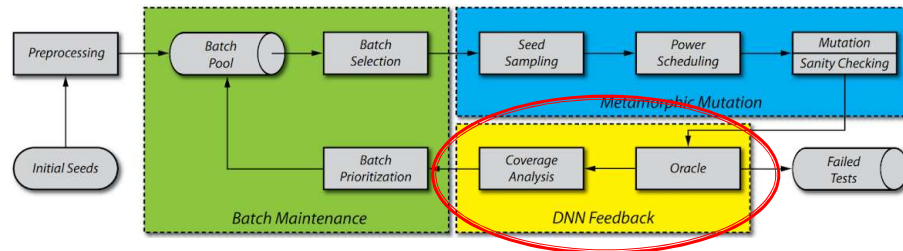
- Brightness
- Contrast
- Pixel Noise
- Blurring
- Translation
- Scaling
- Horizontal Shearing
- Rotation



Pixel value transformation

Affine transformation

48

Coverage Guidance & Oracles



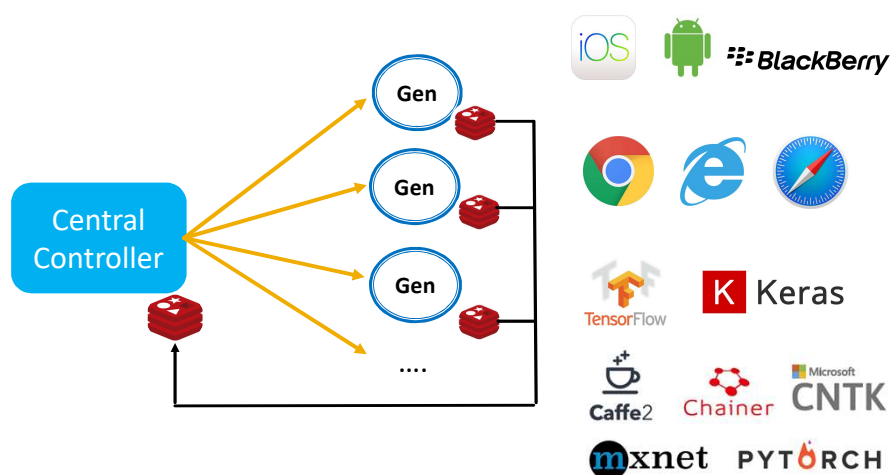
- 

- NC
 - KMNC
 - NBC
 - SNAC
 - TKNC

General Oracles

- Robustness
- Fragmentation
- Cross Platform

49

Large Scale Parallelization



50

DeepHunter In Action

- **RQ1:** What coverage can DeepHunter achieve when guided by the six testing criteria?
- **RQ2:** Can DeepHunter enable diverse erroneous behavior detection of DNNs?
- **RQ3:** Can DeepHunter detect potential defects introduced during DNN quantization?

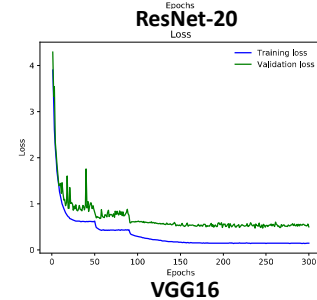
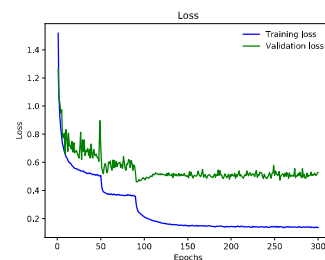
51

Subject Dataset and DNNs

DataSet	Dataset Description	DNN Model	#Neuron	#Layer	Test Acc.
MNIST	Hand written digits recog. from 0 to 9	LeNet-1	52	7	0.976
		LeNet-4	148	8	0.989
		LeNet-5	268	9	0.990
CIFAR-10	General image with 10-class	ResNet-20	2,570	70	0.917
		VGG-16	12,426	17	0.928
ImageNet	1000-class large scale image cla.	MobileNet	38,904	87	0.871*
		ResNet-50	94,059	176	0.929*

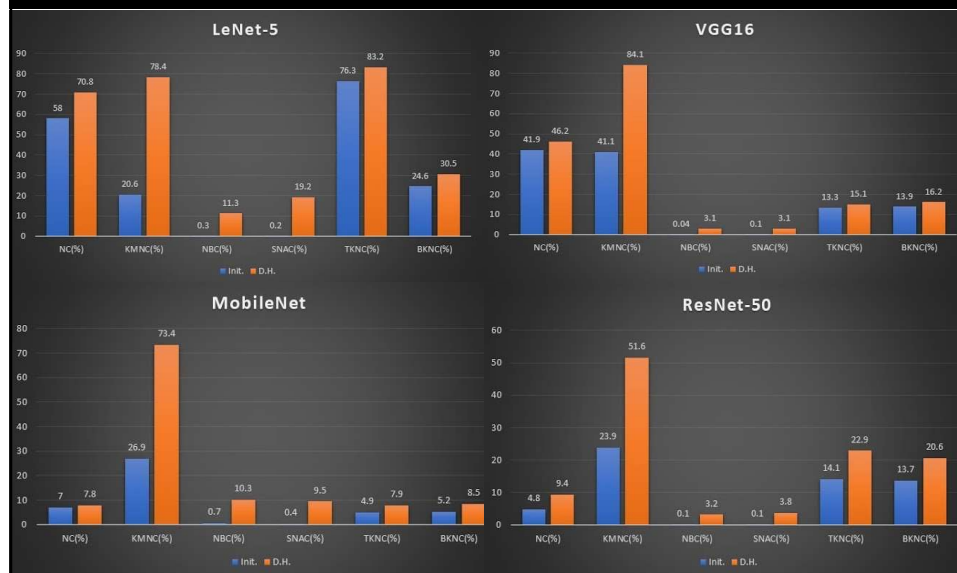
* The reported top-5 test accuracy of pretrained DNN model in [45].

DataSet	DNN	Epoch	Syno.	Train Loss	Train Acc.	Test Acc.
MNIST	LeNet-1	10	A	0.131	0.965	0.967
		30	B	0.099	0.975	0.975
		45	C	0.087	0.979	0.976
	LeNet-4	10	A	0.117	0.974	0.978
		25	B	0.077	0.986	0.986
		50	C	0.058	0.990	0.989
	LeNet-5	10	A	0.116	0.977	0.983
		30	B	0.071	0.988	0.989
		45	C	0.056	0.992	0.990
CIFAR-10	ResNet-20	40	A	0.515	0.894	0.859
		55	B	0.385	0.932	0.880
		95	C	0.239	0.977	0.917
	VGG-16	30	A	0.623	0.914	0.850
		55	B	0.443	0.965	0.900
		95	C	0.316	0.995	0.928



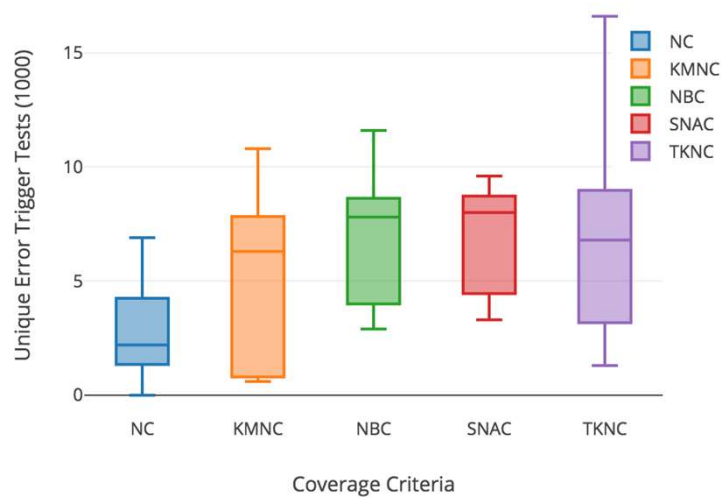
52

Significantly Improve the Coverage

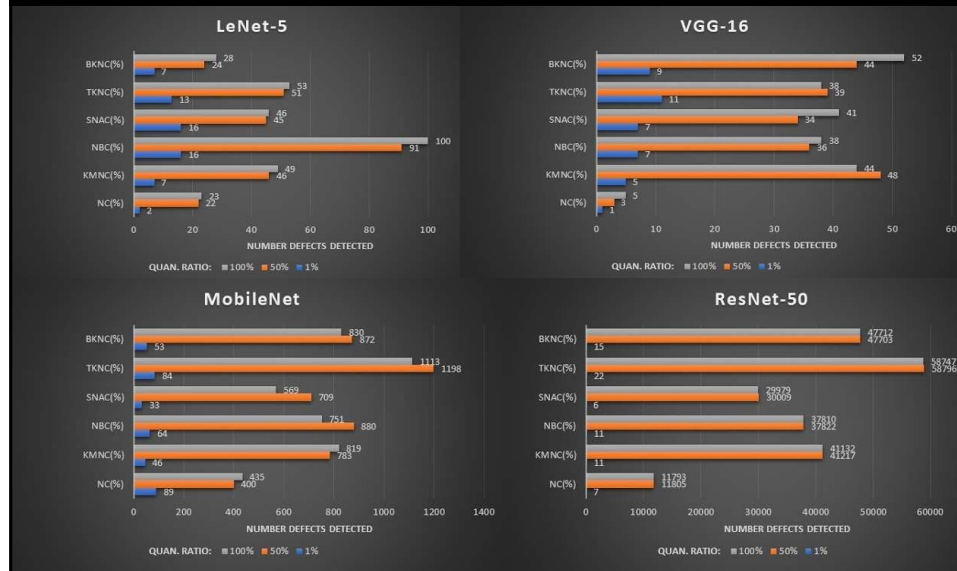


Defect Detection

LeNet-1
LeNet-4
LeNet-5
ResNet-20
VGG-16
MobileNet
ResNet-50



DNN Quantization Defects Detection



DeepGuage and DeepHunter

- **DeepGuage** defines a set of testing criteria to provide a way to gauge testing quality and guide test generation
- **DeepHunter** leverages coverage feedbacks and performs large scale fuzzing test generation for defect detection of DNN development and deployment

DeepMutation: Mutation Testing of DL Systems (ISSRE'18)

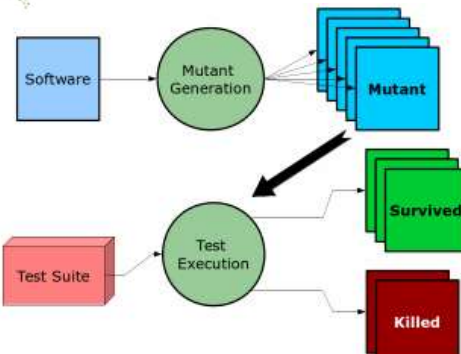
DeepMutation: Mutation Testing of Deep Learning Systems

Lei Ma^{1,2}, Fuyuan Zhang³, Fuyuan Sun⁴, Minhui Xue⁵, Bo Li⁶, Felix Jia^{1,2}
¹Harbin Institute of Technology, China ²Nanjing Technological University, Singapore ³Kyushu University, Japan
⁴University of Illinois at Urbana-Champaign, USA ⁵Carnegie Mellon University, USA ⁶Michigan State University, USA

DeepMutation

Test data Quality Assessment

Mutation Testing Workflow



Combinatorial Testing for DL Systems



Combinatorial Testing for Deep Learning Systems

Lei Ma^{1,2}, Fuyuan Zhang³, Minhui Xue⁵, Bo Li⁶,
Yong Liu⁴, Fuyuan Sun⁴, and Yong Wang¹
¹Harbin Institute of Technology, China
²Nanjing Technological University, Singapore
³New York University, USA
⁴University of Illinois at Urbana-Champaign, USA
⁵Kyushu University, Japan
Contact: minhui.xue@nyu.edu.cn

Abstract. Deep learning (DL) has achieved remarkable progress over the past decade and been widely applied to many safety-critical applications. However, the robustness of DL systems recently receives great concerns, such as adversarial examples against image classifiers, which could potentially result in severe consequences. Adversarial testing techniques could help to evaluate the robustness of a DL system and discover their vulnerabilities at an early stage. This paper challenges the testing such systems is that its runtime state space is too large if we consider each neuron as a runtime state for DL. Using a DL system, adversarial testing is a runtime state space is too large if we consider each neuron as a runtime state. In this paper, we perform an exploratory study of CT on DL systems. We adapt the concept of CT to DL systems and propose a set of coverage criteria for DL systems, as well as a CT coverage guided test generation technique. Our evaluation demonstrates that CT is a promising avenue for testing DL systems. We further pose several open questions and interesting directions for combinatorial testing of DL systems.

Keywords: Combinatorial testing; Deep learning; Adversarial attacks

1 Introduction

Deep learning (DL) systems have been widely applied to various applications due to their high accuracy, such as computer vision [25], natural language processing [19], autonomous driving [7], and natural audio processing [8]. However, recently, DL systems have been shown to be vulnerable against different attacks, such as adversarial examples in computer vision and audio systems. Given that more and more safety-critical applications start to adopt DL, deploying DL without thorough testing to safety-critical applications can lead to severe consequences, such as possible accidents in autonomous driving [12]. DL systems are

- Neuron Activation Configuration
- T-way combination sparse coverage
- T-way combination dense coverage
- (p,t)-completeness coverage

Parameters	All Combinations			2-Pair (Pairwise)		
	P1	P2	P3	P1	P2	P3
P1 : A, B, C	TC1 : A 1 X			TC1 : A 1 X		
P2 : 1, 2	TC2 : A 1 Y			TC4 : A 2 Y		
P3 : X, Y	TC3 : A 2 X			TC6 : B 1 Y		
	TC4 : A 2 Y			TC7 : B 2 X		
	TC5 : B 1 X			TC9 : C 1 X		
	TC6 : B 1 Y			TC12 : C 2 Y		
	TC7 : B 2 X					
	TC8 : B 2 Y					
	TC9 : C 1 X					
	TC10 : C 1 Y					
	TC11 : C 2 X					
	TC12 : C 2 Y					

Thank you!

Q&A

61