# Towards Reliability Assurance of Deep Learning Systems

Jianjun Zhao
Kyushu University

CSBSE 2018, November 17, 2018



## Pangu Research Group
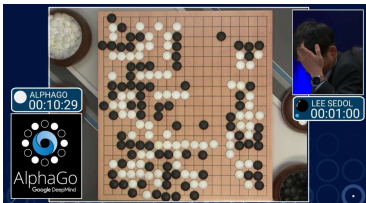(https://pangukaitian.github.io/pangu/en/index.html)

- Our group focuses on researches on the potential symbioses between software engineering (SE) and artificial intelligence (AI)

- The overall goal is to obtain better software and AI systems, making them more robust, reliable, and secure, and easier to specify, build, maintain, or improve

- Our group consists of 2+ faculty members, 3 PhD students, and 7 master students

# Deep Learning Matches Human Intelligence



# Rushing for Development and Real-world Deployment

# Driving Force of Many Novel Tech.



FINANCE  RETAIL  INSURANCE  EDUCATION

TRANSPORTATION  FOOD  LEGAL  HEALTHCARE

5

# Current Deep Learning is Vulnerable



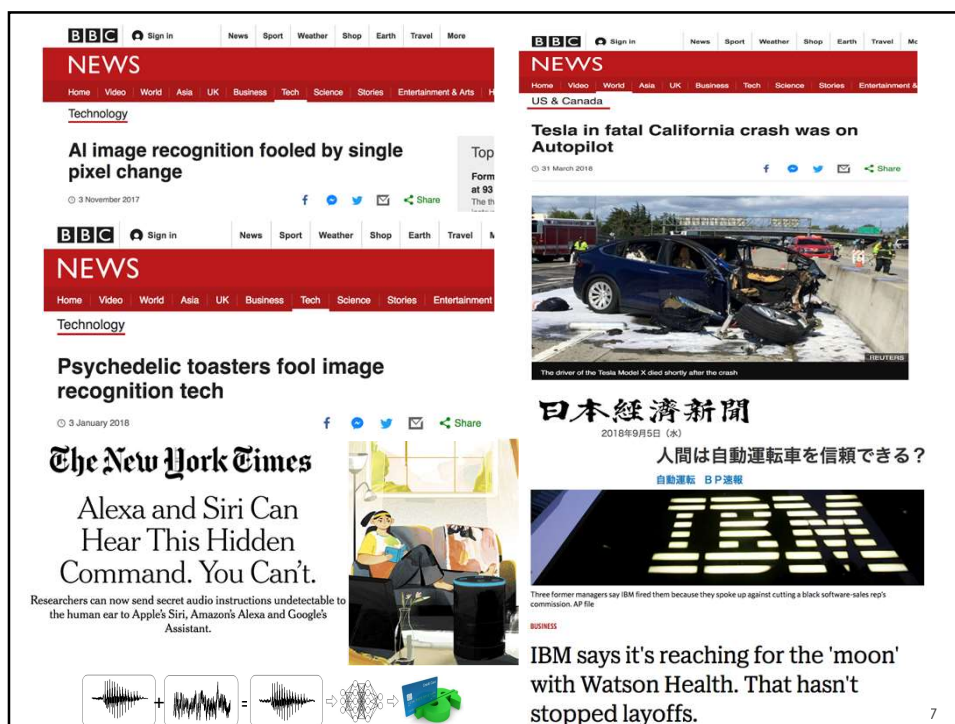Classified as panda  Small adversarial noise  Classified as gibbon

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

Accident

6

# Reliability Assurance for DL is at Early Stage

**DeepXplore**
SOSP'17

**DeepTest**
ICSE'18

**ML Fairness Testing**
ASE'18

**DeepRoad**
ASE'18

**DeepConcolic**
ASE'18

**DeepGauge**
ASE'18

**TensorFlow Program Bugs** ......
ISSTA'18

**MODE: DNN Debugging**
FSE'18

8

## Reliability Assurance for Traditional Software Testing Criteria and Tools

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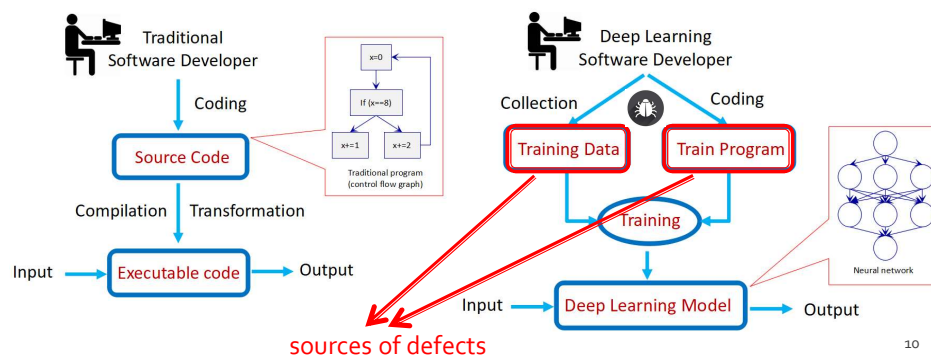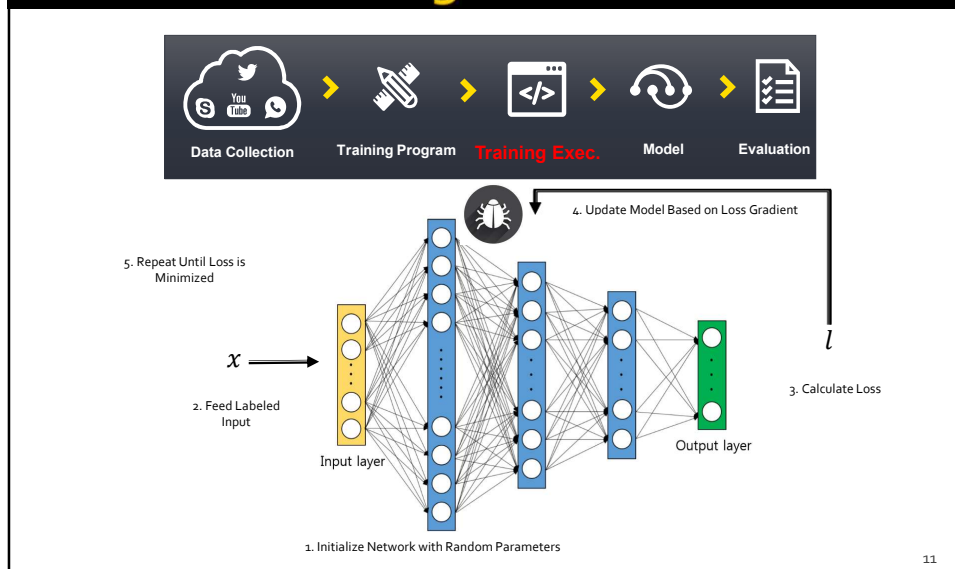- …



9

## Source of Defects Programming Paradigm

- ☐ **The decision logic of a traditional software:**
  - ◼ In the form of code
- ☐ **The decision logic of a DL system:**
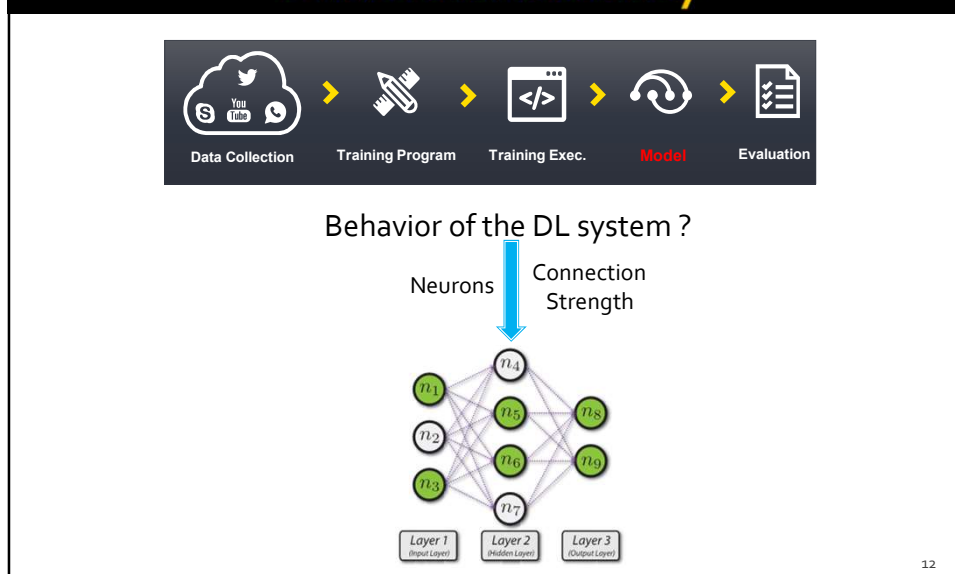  - ◼ The structure of DNN
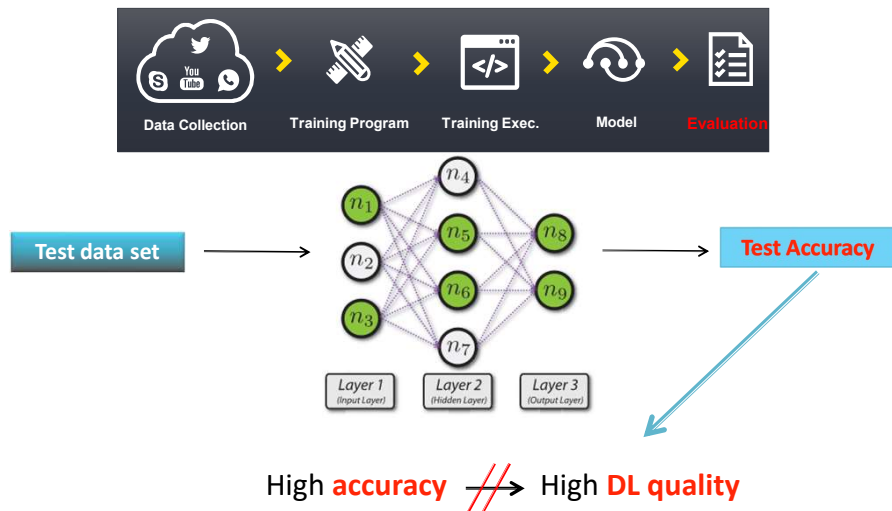  - ◼ The connection weights



sources of defects

10

Source of Defects
Training Procedure



Lacking of Interpretability and Understandability

# Reliability Measurement Immature



Data Collection → Training Program → Training Exec. → Model → **Evaluation**

Test data set → [neural network: $n_1$, $n_2$, $n_3$ (Layer 1 Input Layer); $n_4$, $n_5$, $n_6$, $n_7$ (Layer 2 Hidden Layer); $n_8$, $n_9$ (Layer 3 Output Layer)] → **Test Accuracy**

High **accuracy** ⇸ High **DL quality**

13

# Towards Reliability Assurance for Deep Learning Software

- **Testing quality & logic coverage** — *DeepGauge*

- **MT-based test data quality evaluation** — *DeepMutation*

- **Combinatorial latent space exploration** — *DeepCT*

14

## Overview of DeepGauge (ASE2018)

**Multi-Granularity Testing Criteria for Deep Learning Systems**

- **Enable** quality evaluation of DLs from multiple portrayals

- **Provide** systematic guidance of TestGen for Defects

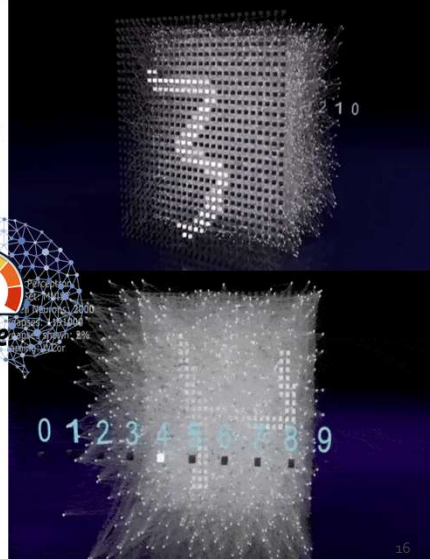- **Facilitate** Interpretation & Understanding



## Design of DeepGauge

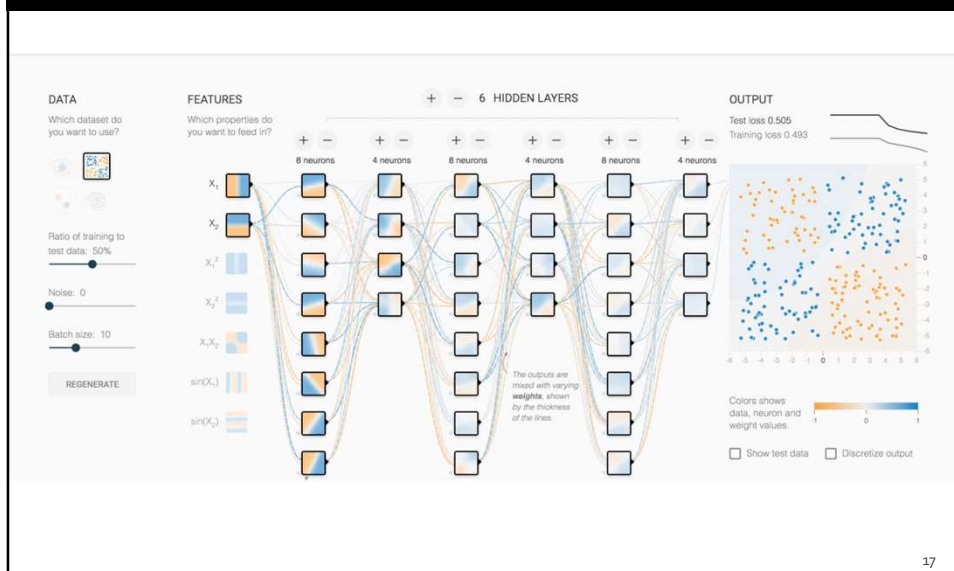**Simple** to understand & use

**Efficient** to compute

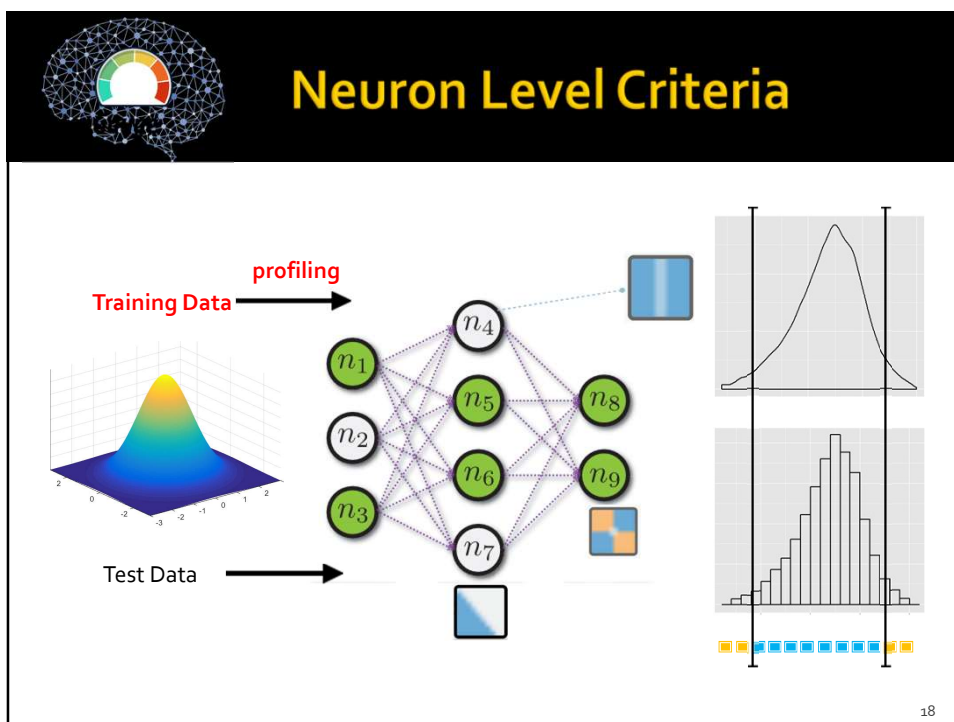**General** to diverse DNNs

**Scale** to large DNNs

**Adaptable** by cases

# Neuron Level Criteria

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

**K-Multisection Neuron Coverage**

$$\text{KMNCov}(T, k) = \frac{\sum_{n \in N} |\{S_i^n \mid \exists \mathbf{x} \in T : \ \phi(\mathbf{x}, n) \in S_i^n\}|}{k \times |N|}$$

**Neuron Boundary Coverage**

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

**Strong Neuron Activation Coverage**

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

19

# Layer Level Criteria

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

**Top-K Neuron Coverage**

$$\text{TKNCov}(T, k) = \frac{|\bigcup_{\mathbf{x} \in T} (\bigcup_{1 \le i \le l} \text{top}_k(\mathbf{x}, i))|}{|N|}$$

**Neuron interactions of same layer**

**Top-K Neuron Patterns**

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

**Neuron interactions across layers**

20

# Large Scale Empirical Study

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

21

# Large Scale Empirical Study

**Adversarial Test Gen**  **Coverage Eval.**  **Coverage Parameters**

KMNC
FGSM
Test org.  NBC
BIM
DNNs → JSMA → SNAC
TKNC
C&W
TKNP
NC

| Coverage Criteria | Parameter setting | | |
|---|---|---|---|
| KMNC | k=1,000 | k=10,000 | N.A |
| NBC | LB=l | LB=l-0.5*σ | LB=l-σ |
| | UB=u | UB=u+0.5*σ | UB=u+σ |
| SNAC | UB=u | UB=u+0.5*σ | UB=u+σ |
| TKNC | k=1 | k=2 | k=3 |
| TKNP | k=1 | k=2 | k=3 |
| NC | th.= 0.0, 0.25, 0.5, 0.75 | | |

**5** DNN Profiling

**18** Group Large Scale Adv Test Gen

**190,000** Test data= 150,000 (MNIST) + 40,000 (ImageNet)

**414** Evaluation Configurations= 270 (MNIST) + 144 (ImageNet)

**Tesla M40 GPU 24GB /** 18-core 2.3GHz Xeon 64-bit 196 GB

22

# Adversarial Examples



- CW and JSMA make smaller perturbation
- Adv on ImageNet is often imperceptible
- JSMA and CW is computation intensive
- Even more obvious on ImageNet subjects

# DeepMutation: Mutation Testing of Deep Learning System (ISSRE'18)