

IOWA STATE UNIVERSITY

Department of Computer Science

DeepDiagnosis: Automatically Diagnosing Faults and Recommending Actionable Fixes in Deep Learning Programs



Mohammad Wardat, Breno Dantas Cruz, Wei Le, and Hridesh Rajan

wardat@iastate.edu, bdantasc@iastate.edu, weile@iastate.edu, and hridesh@iastate.edu

Dept. of Computer Science, Iowa State University

Deep Learning Is Being Used In Critical Systems

DNNs are used in autonomous vehicles to medical devices.



Automated Driving



Aerospace and Defense



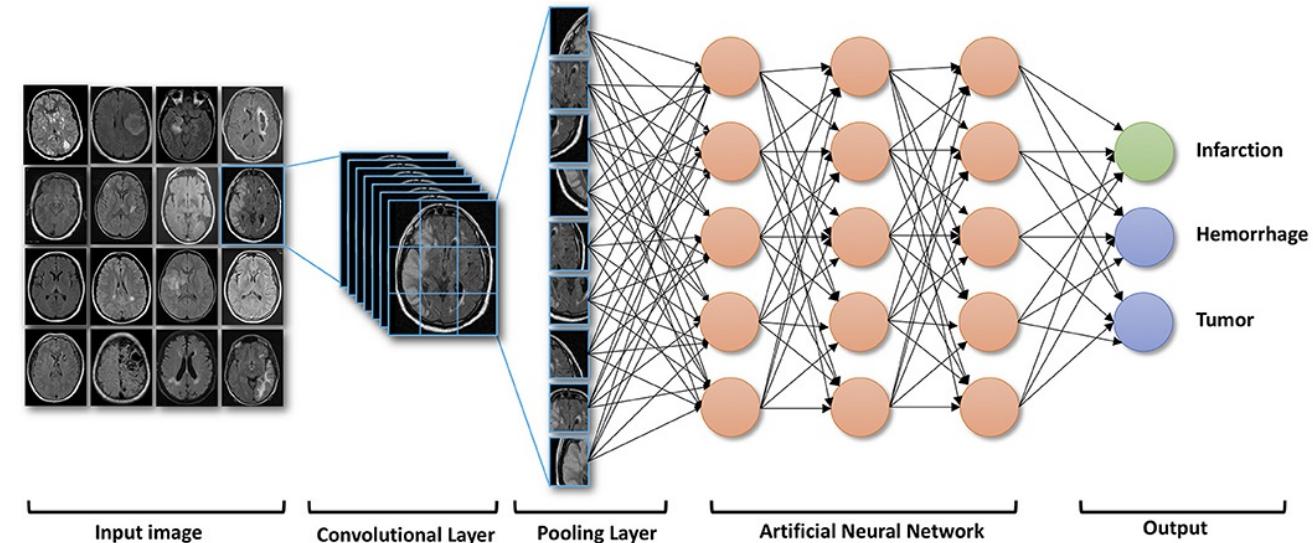
Medical Research



Industrial Automation



Medical Electronics



Problem Statement

N/A

Existing software techniques are still not applicable to identify symptoms and their root causes.



Existing methodologies do not effectively identify bugs and provide actionable fixes to developers.



There is no existing technique that finds the relation between symptoms and root causes in the structure DNNs.

Prior Work

1 DeepLocalize: fault localization for deep neural networks

Performs dynamic analysis during training to localize bugs by monitoring values.

2 UMLAUT: Debugging deep learning programs

Checks the structure of deep learning programs and model behavior during training.

3 AUTOTRAINER: An automatic DNN training problem detection and repair

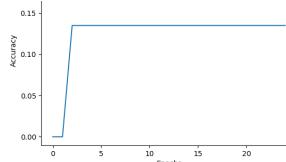
Detects, repairs five symptoms, and provides a set of possible fixes.

4 DEBAR: Detecting numerical bugs in neural network architectures

Detects numerical errors in the deep neural networks using static analysis.

- [1] M. Wardat, W. Le, and H. Rajan, “Deeplocalize: fault localization for deep neural networks,” in ICSE’21, 2021.
- [2] E. Schoop, F. Huang, and B. Hartmann, “UMLAUT: Debugging deep learning programs using program structure and model behavior,” in CHI’21, 2021.
- [3] X. Zhang, J. Zhai, S. Ma, and C. Shen, “Autotrainer: An automatic dnn training problem detection and repair system,” in ICSE’21, 2021.
- [4] Y. Zhang, L. Ren, L. Chen, Y. Xiong, S.-C. Cheung, and T. Xie, “Detecting numerical bugs in neural network architectures,” in FSE’21, 2021.

Motivation

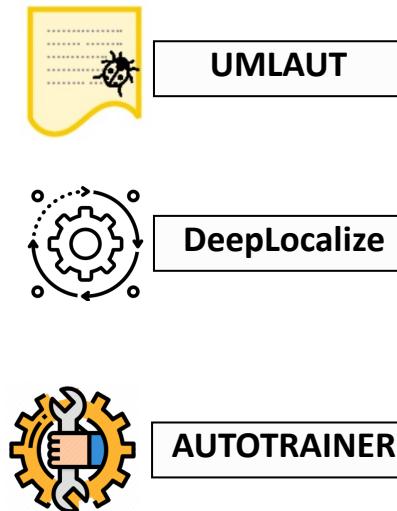
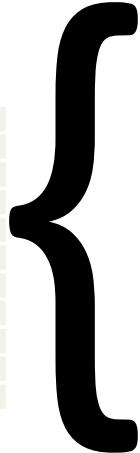


- Q1) Why do I get such a bad result for such a simple dataset?
Q2) Is my dataset malformed?

NO ANSWER

```
1 model = Sequential()
2 model.add(Dense(128, 50))
3 model.add(Activation('relu'))
4 model.add(Dropout(0.2))
5 model.add(Dense(50, 50))
6 model.add(Activation('relu'))
7 model.add(Dropout(0.2))
8 model.add(Dense(50, 1))
9 model.add(Activation('softmax'))
10 model.compile(loss='binary_crossentropy', optimizer=RMSprop())
11 model.fit(X,Y,batch_size ,epoch ,validation_data=(X_test ,Y_test))
```

<https://stackoverflow.com/questions/31880720/>



UMLAUT
No Output
DeepLocalize
layer 7: Numerical Error
Using 'SelU' each layers
Using 'BatchNormalization'
...
Unsolved.



104.65 Seconds



2.14 Seconds



495.83 Seconds

Contributions



We study the symptoms and their root causes to propose a mapping algorithm of symptoms to a fix.



Benchmark

We provide a benchmark with 444 models that practitioners can use to evaluate their techniques.



Our prototype is more comprehensive and efficient than:



UMLAUT



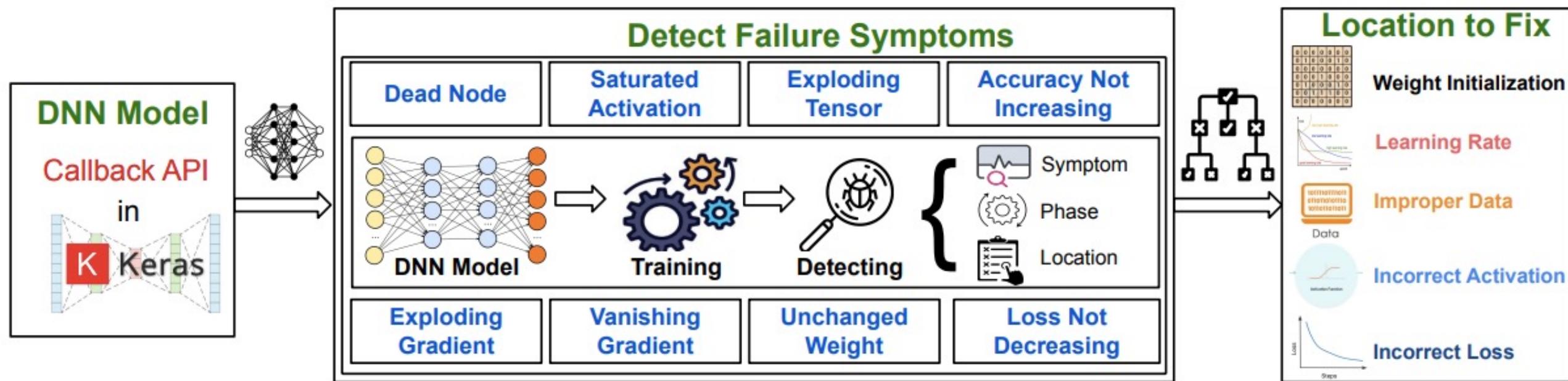
DeepLocalize



AUTOTRAINER

8

DeepDiagnosis Overview



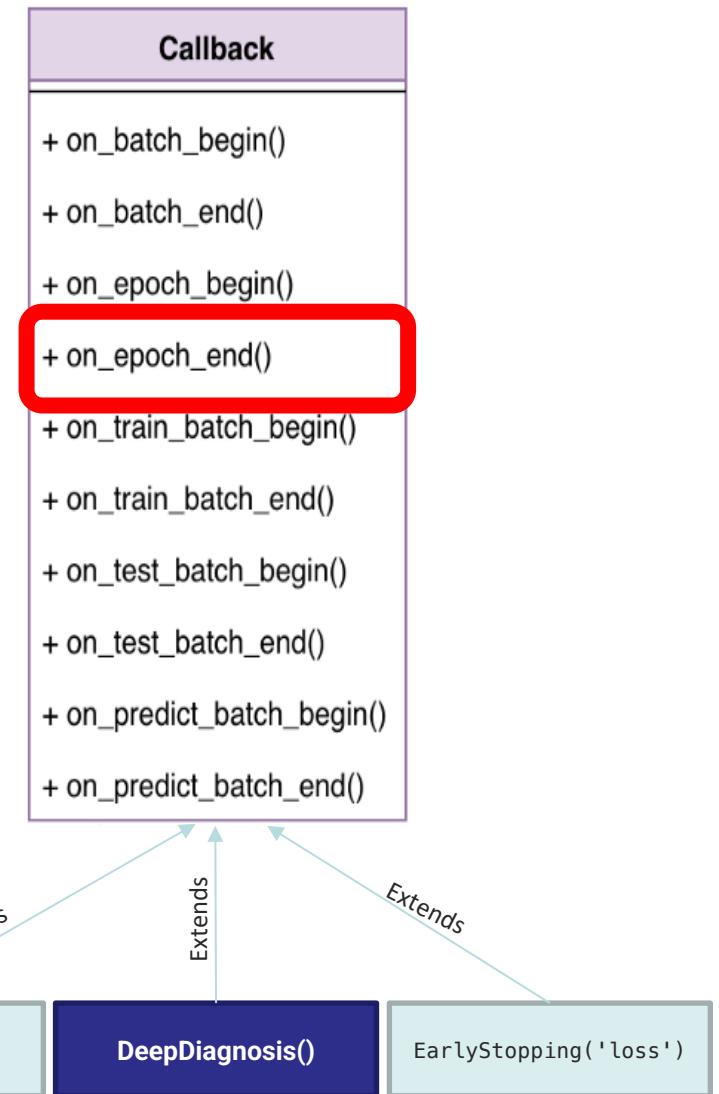
Callback Method



Our work focuses on monitoring model training.

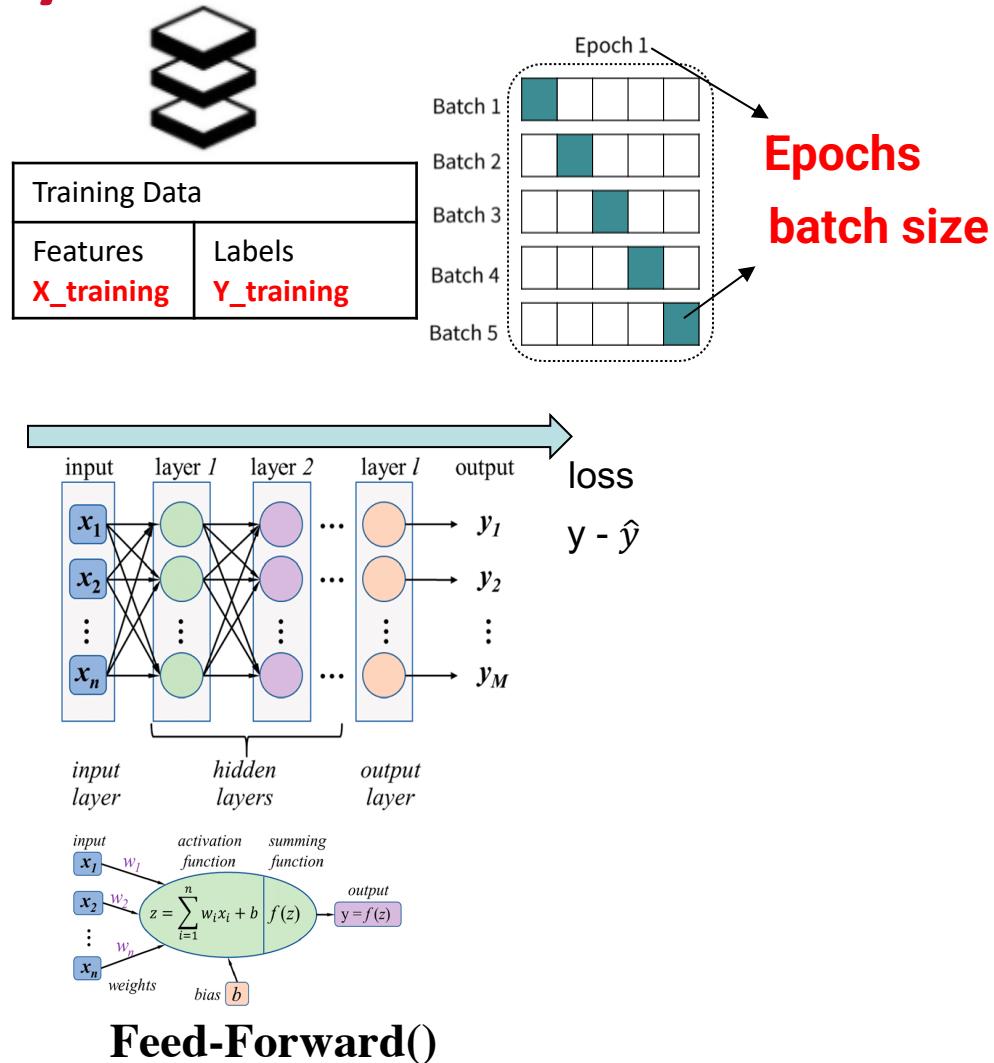
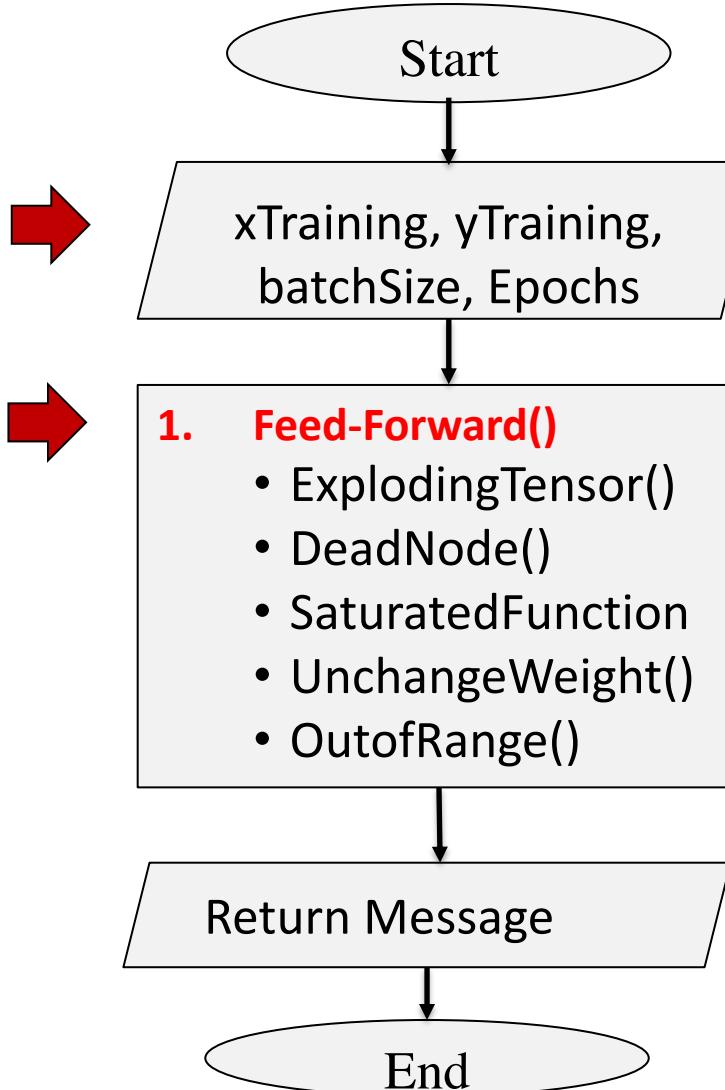


To implement DeepDiagnosis, we override the ***on_epoch_end()*** method from the ***keras.callbacks.Callback*** class

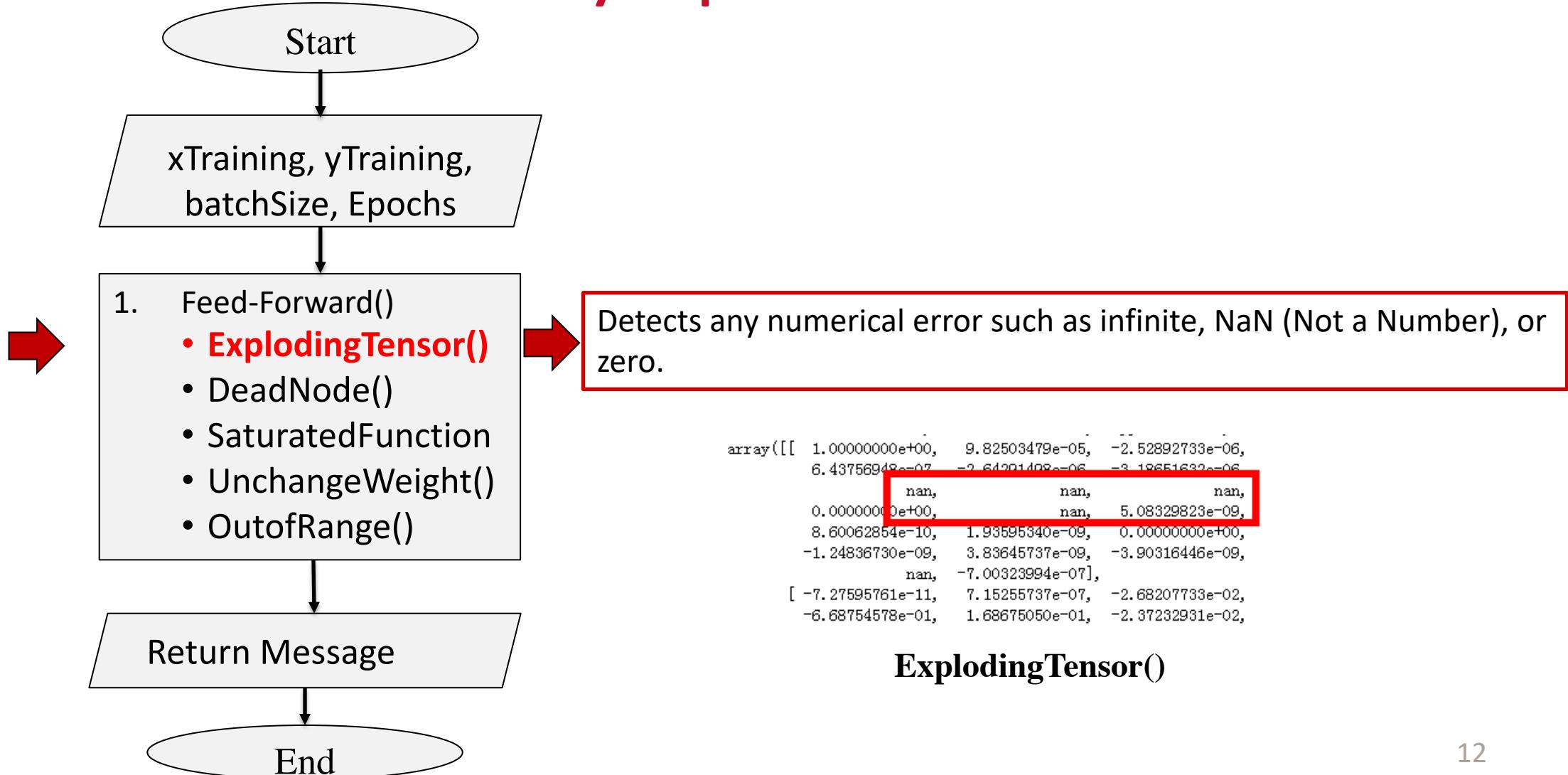


Class Diagram for Keras Callback

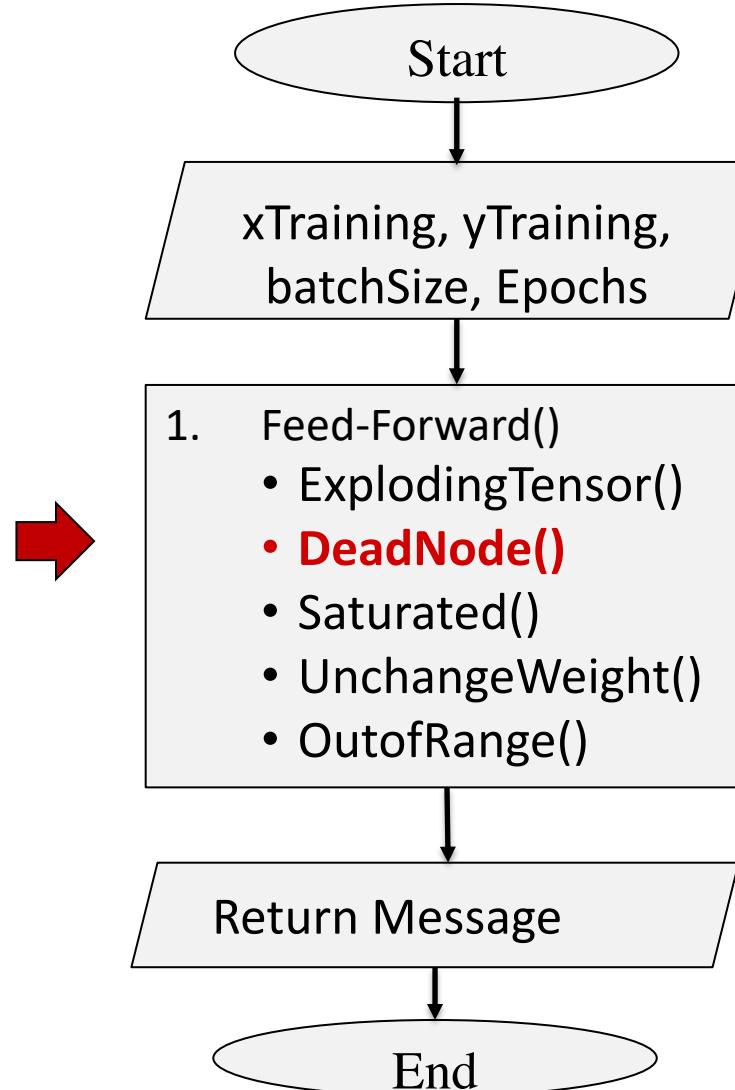
Algorithm: Failure Symptoms Detection



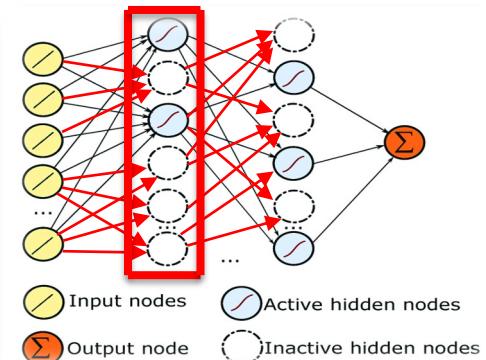
Algorithm: Failure Symptoms Detection



Algorithm: Failure Symptoms Detection

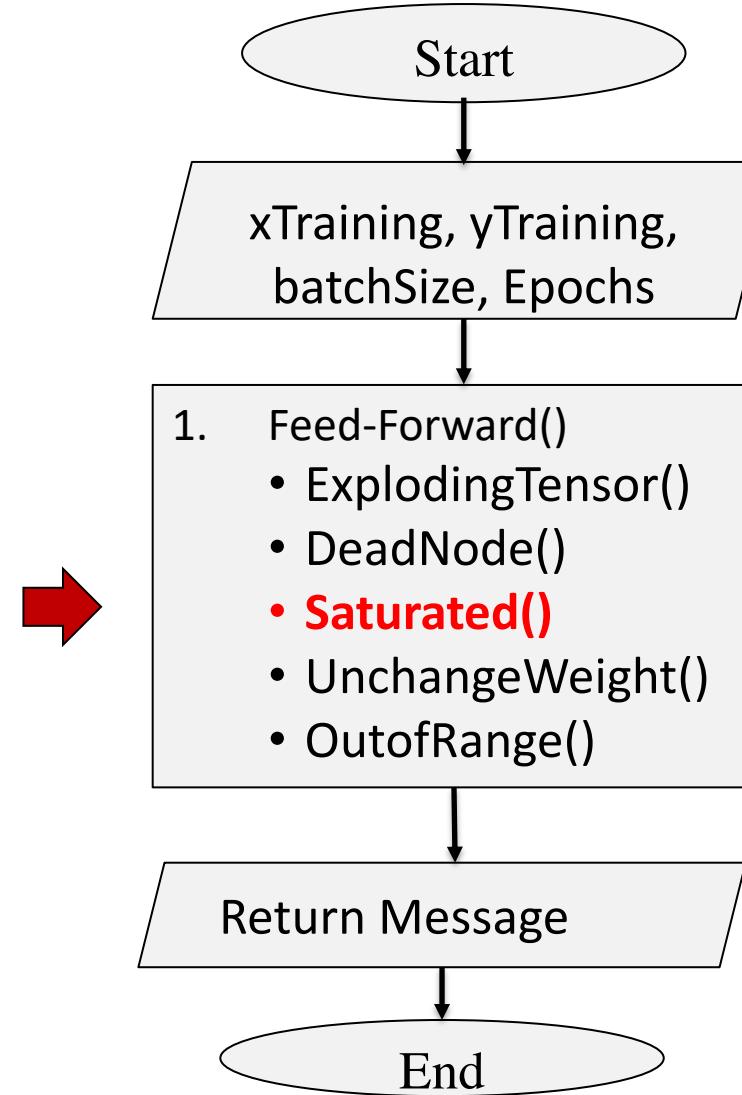


Takes the output of ReLU activation function as input, then computes how many inactive nodes dropped below Threshold.

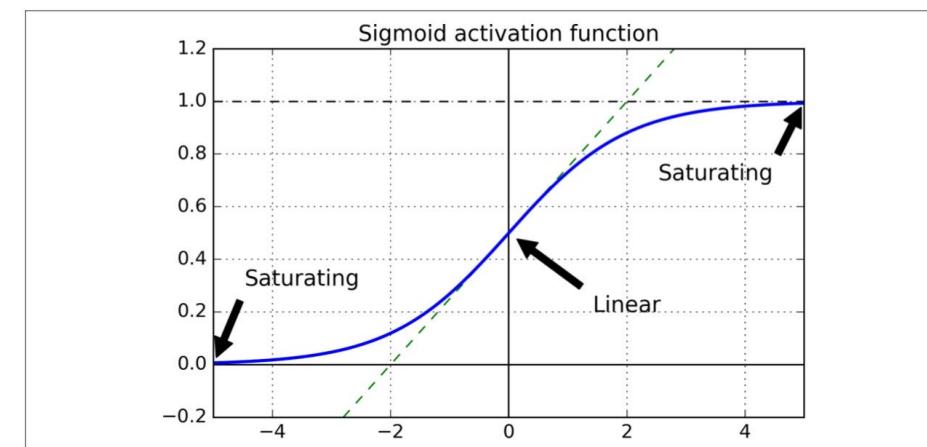


DeadNode()

Algorithm: Failure Symptoms Detection



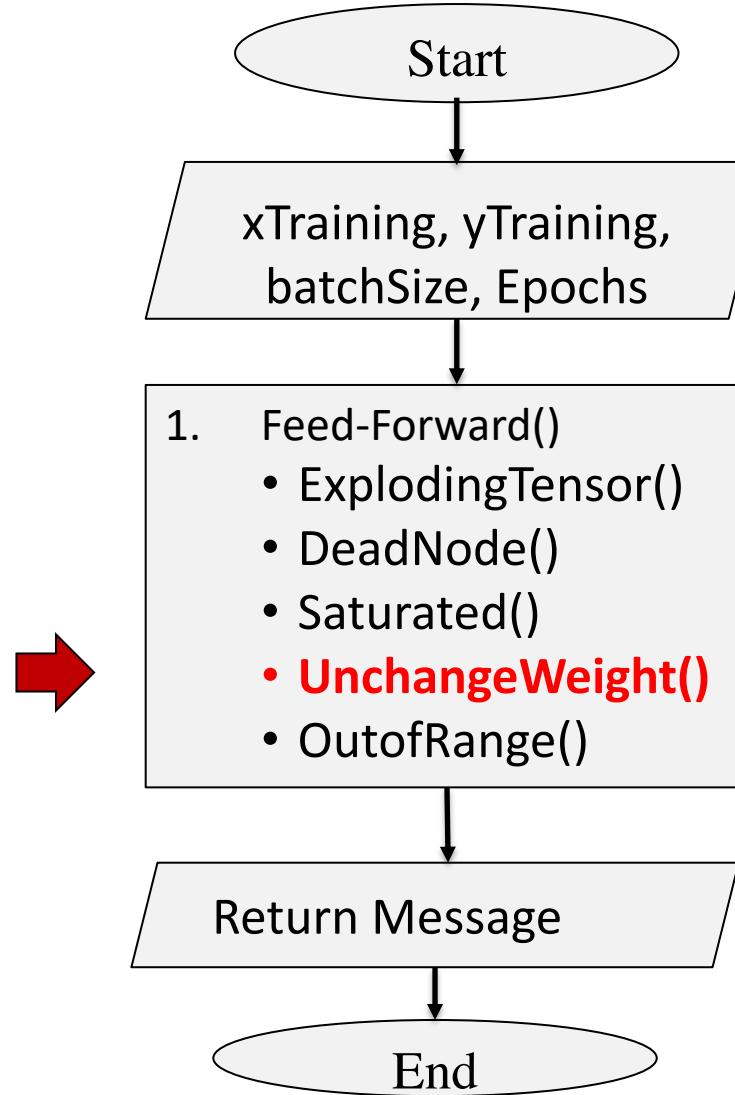
Detects if any logistic function are becoming saturated by checking when the input of function has reached either a max or min value.



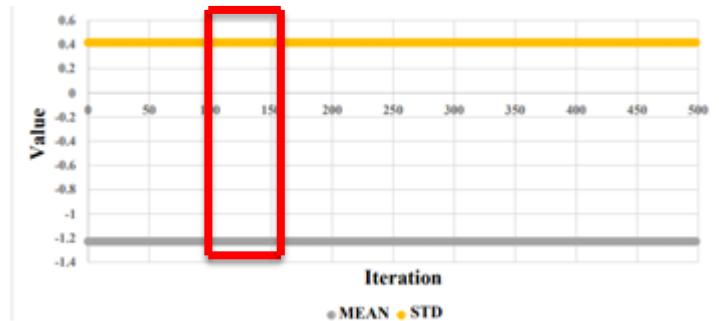
Saturated Activation Function()

14

Algorithm: Failure Symptoms Detection

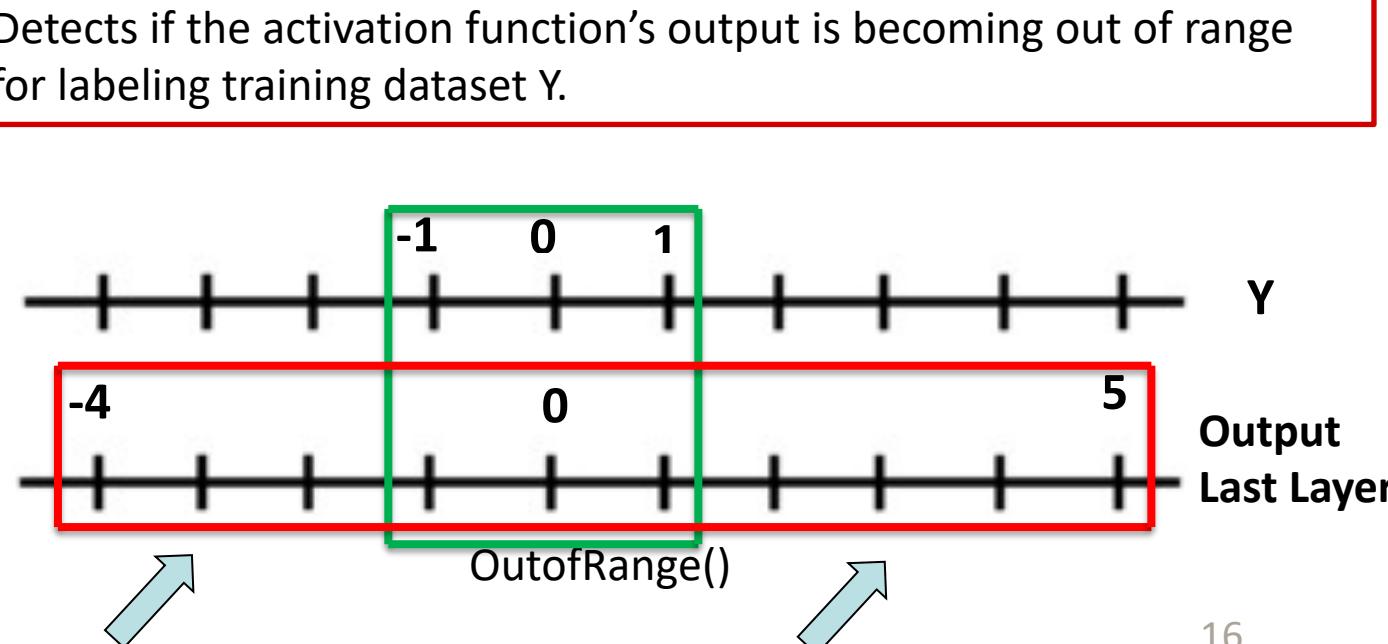
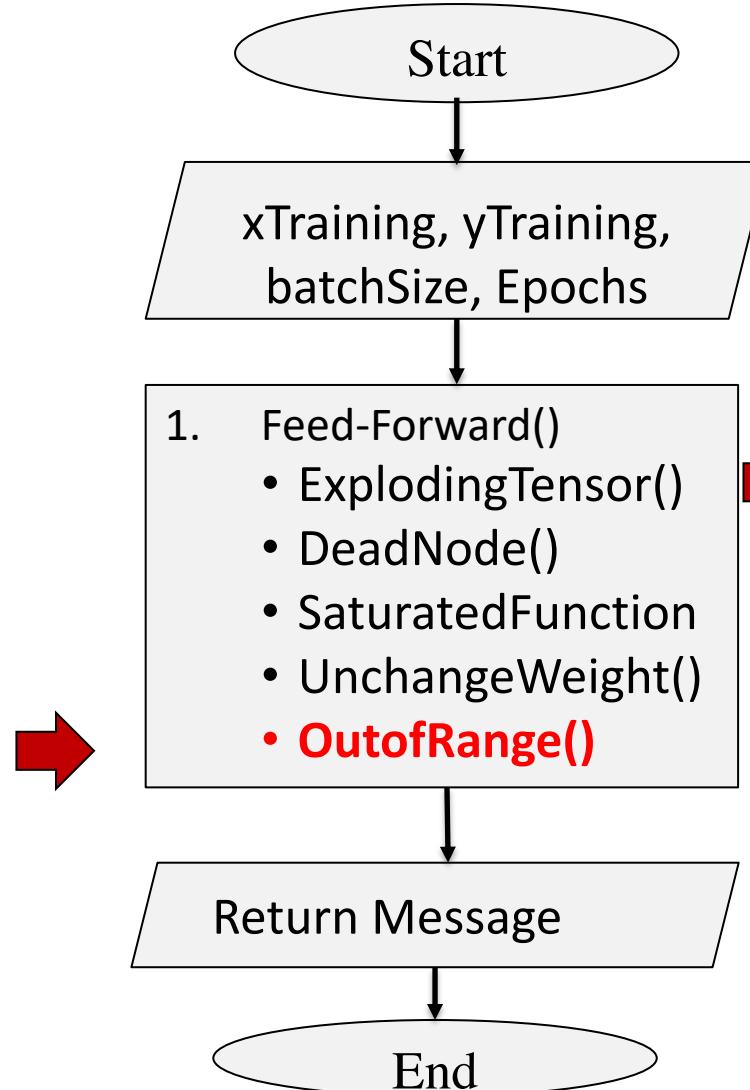


Detects if value is not changing by comparing the value for the current step with the mean value that stores in previous (N = 5) steps.

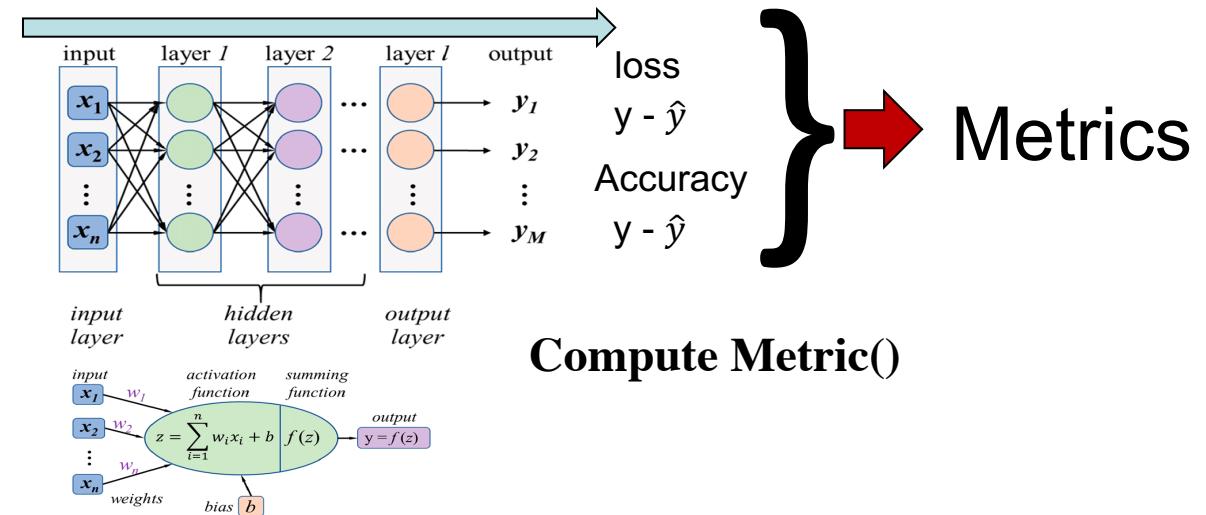
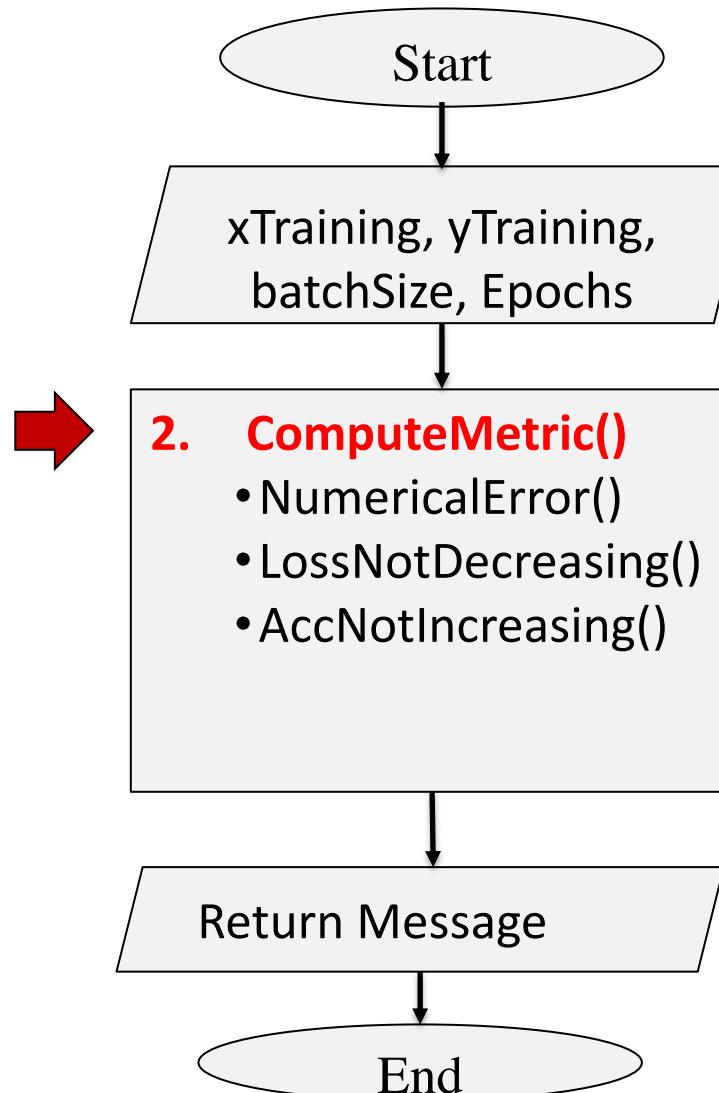


UnchangeWeight()

Algorithm: Failure Symptoms Detection

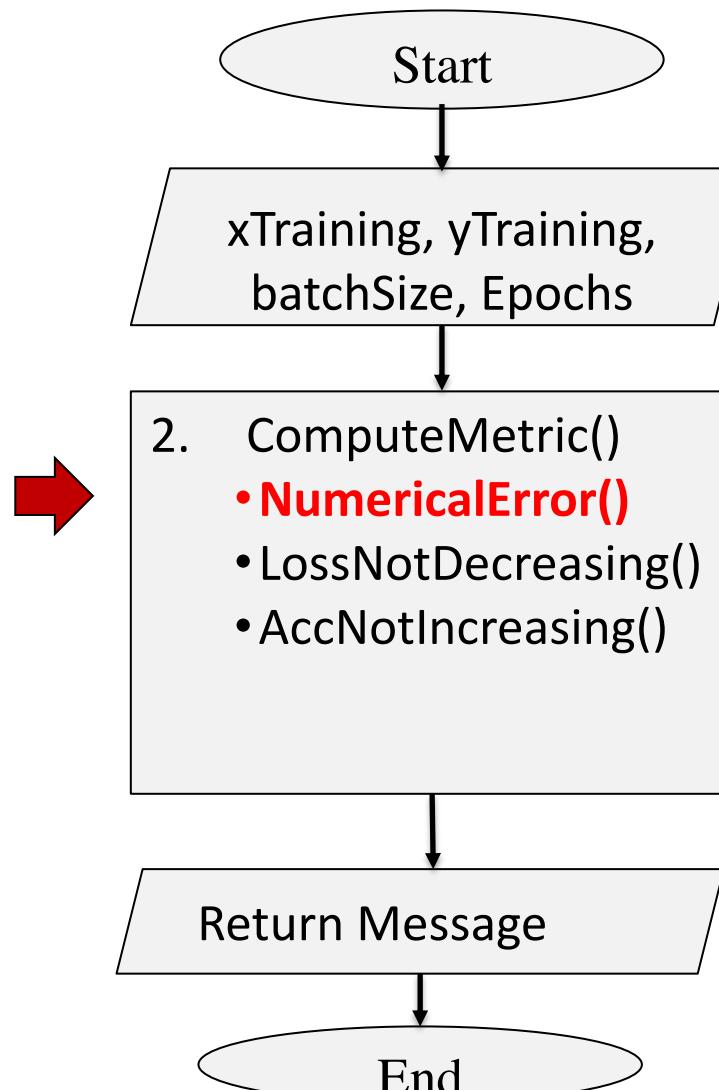


Algorithm: Failure Symptoms Detection



Compute Metric()

Algorithm: Failure Symptoms Detection

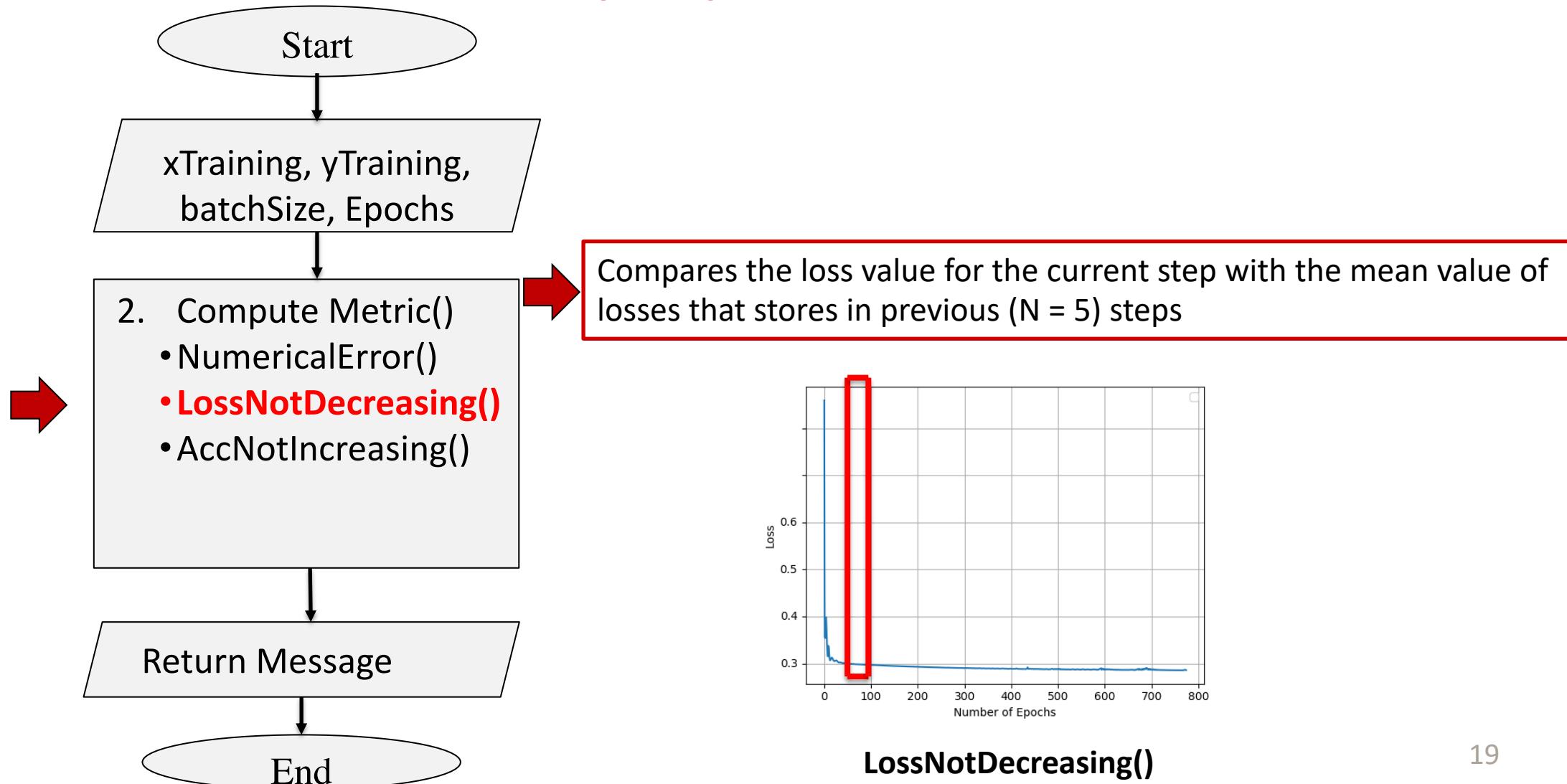


```
Epoch 1/10  
1875/1875 [=====] - 6s 3ms/step - loss: nan accuracy: 0.0987  
Epoch 2/10  
1875/1875 [=====] - 6s 3ms/step - loss: nan accuracy: 0.0987  
Epoch 3/10  
1875/1875 [=====] - 6s 3ms/step - loss: nan accuracy: 0.0987  
Epoch 4/10  
1875/1875 [=====] - 6s 3ms/step - loss: nan accuracy: 0.0987  
Epoch 5/10  
1875/1875 [=====] - 6s 3ms/step - loss: nan accuracy: 0.0987  
Epoch 6/10  
1875/1875 [=====] - 6s 3ms/step - loss: nan accuracy: 0.0987  
Epoch 7/10  
1875/1875 [=====] - 6s 3ms/step - loss: nan accuracy: 0.0987
```

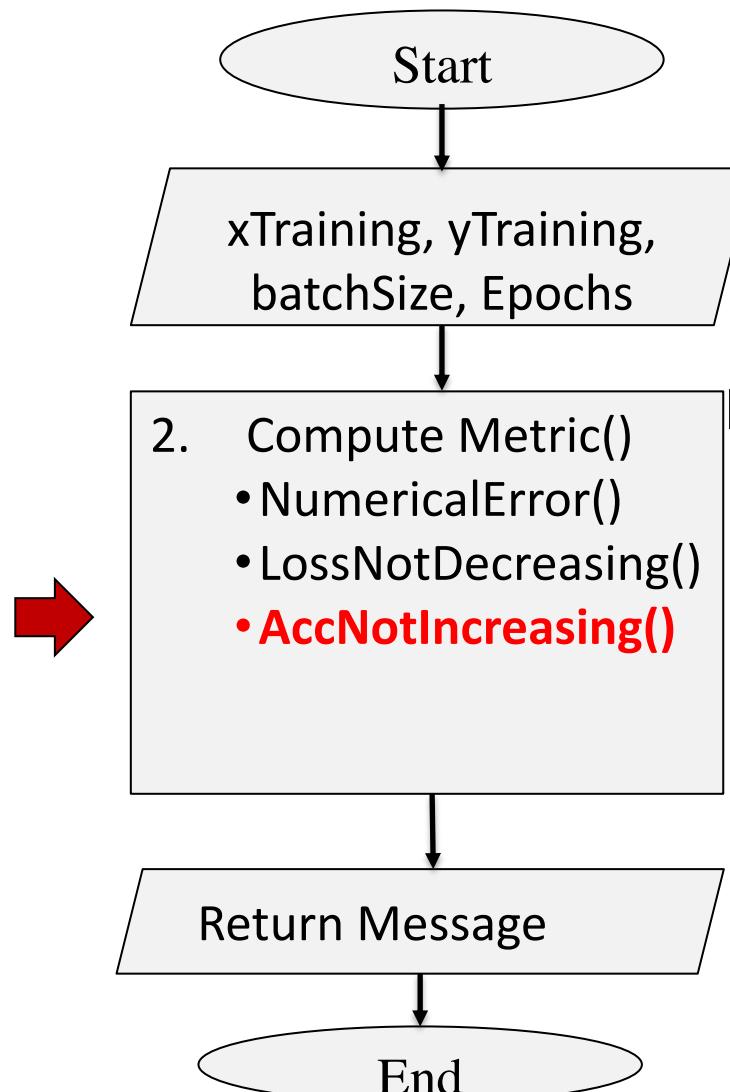
NumericalError()

18

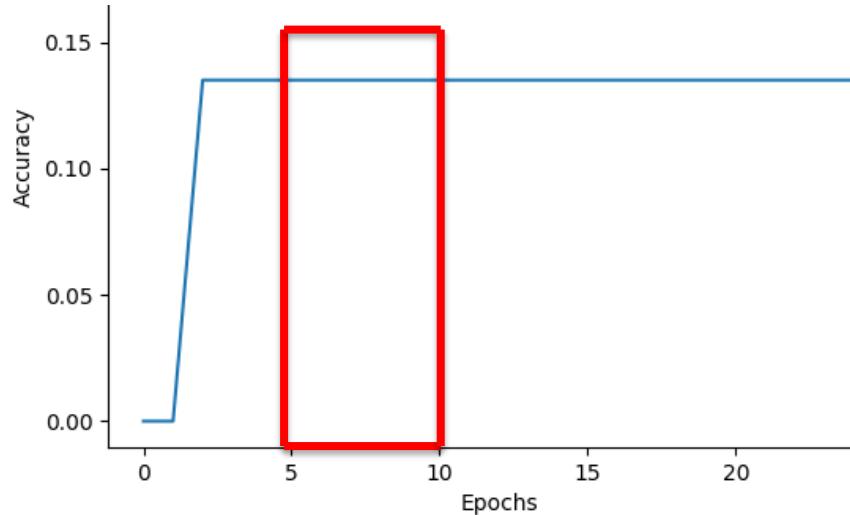
Algorithm: Failure Symptoms Detection



Algorithm: Failure Symptoms Detection



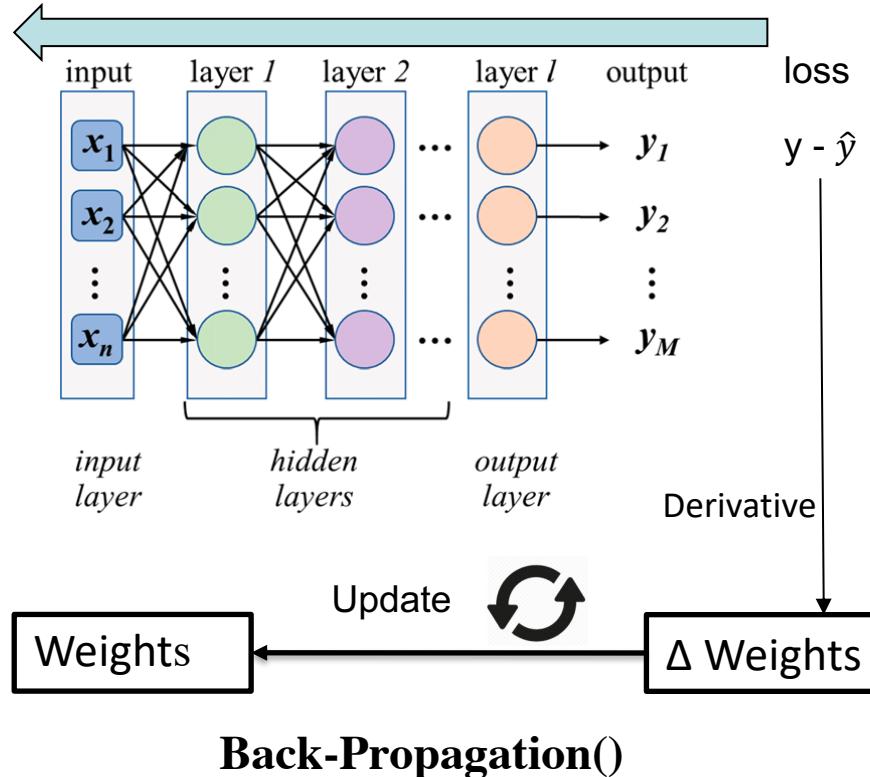
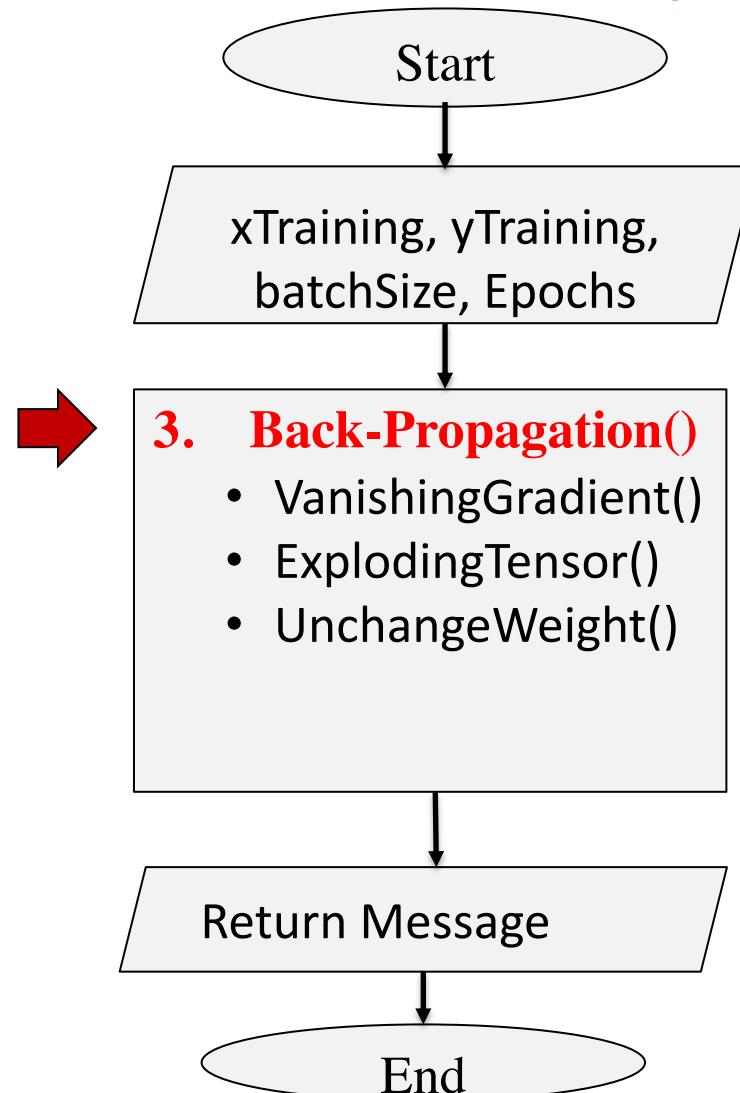
Compares the accuracy value for the current step with the mean value of accuracy that stores in previous (N = 5) steps.



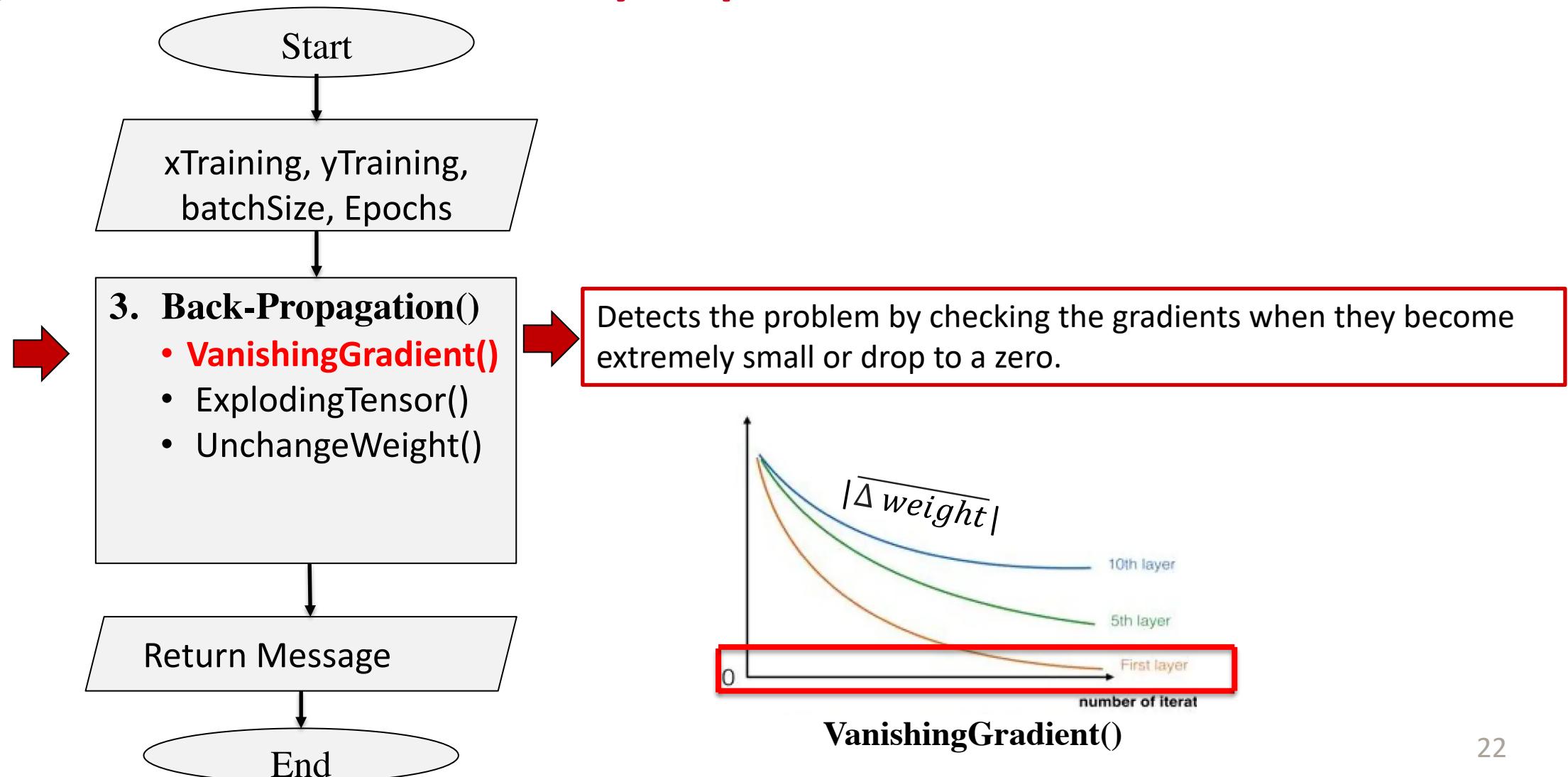
AccuracyNotIncreasing()

20

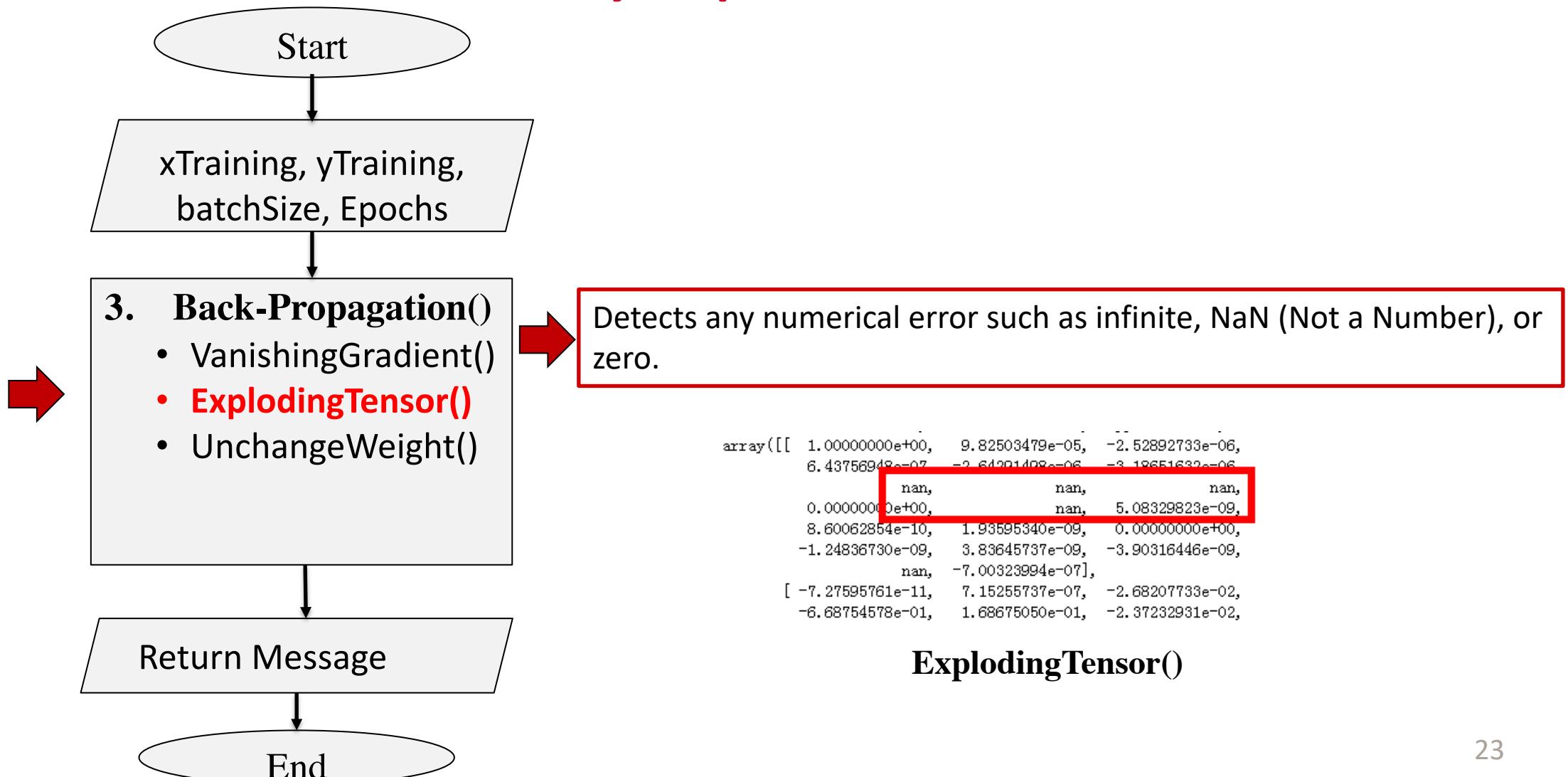
Algorithm: Failure Symptoms Detection



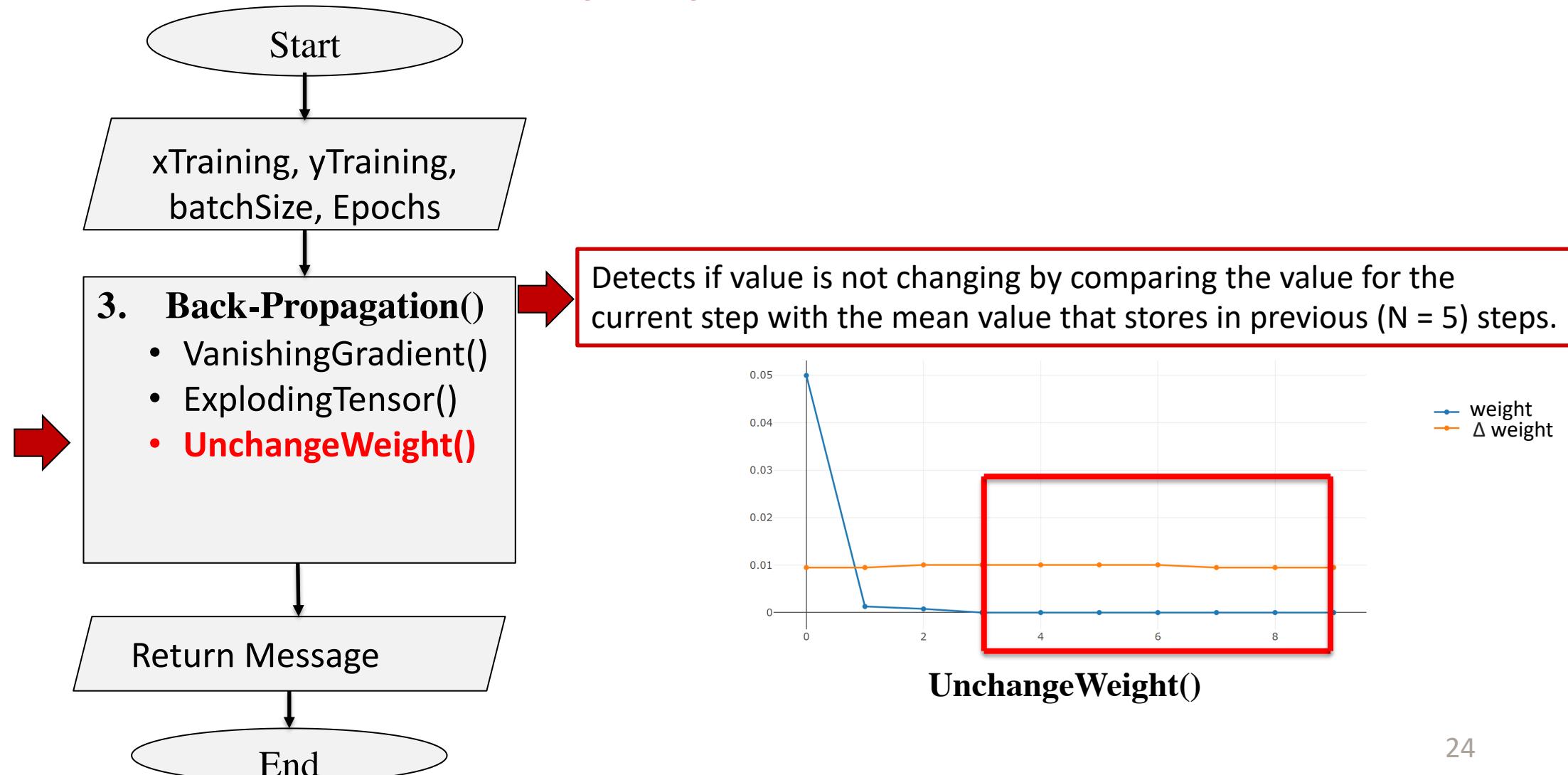
Algorithm: Failure Symptoms Detection



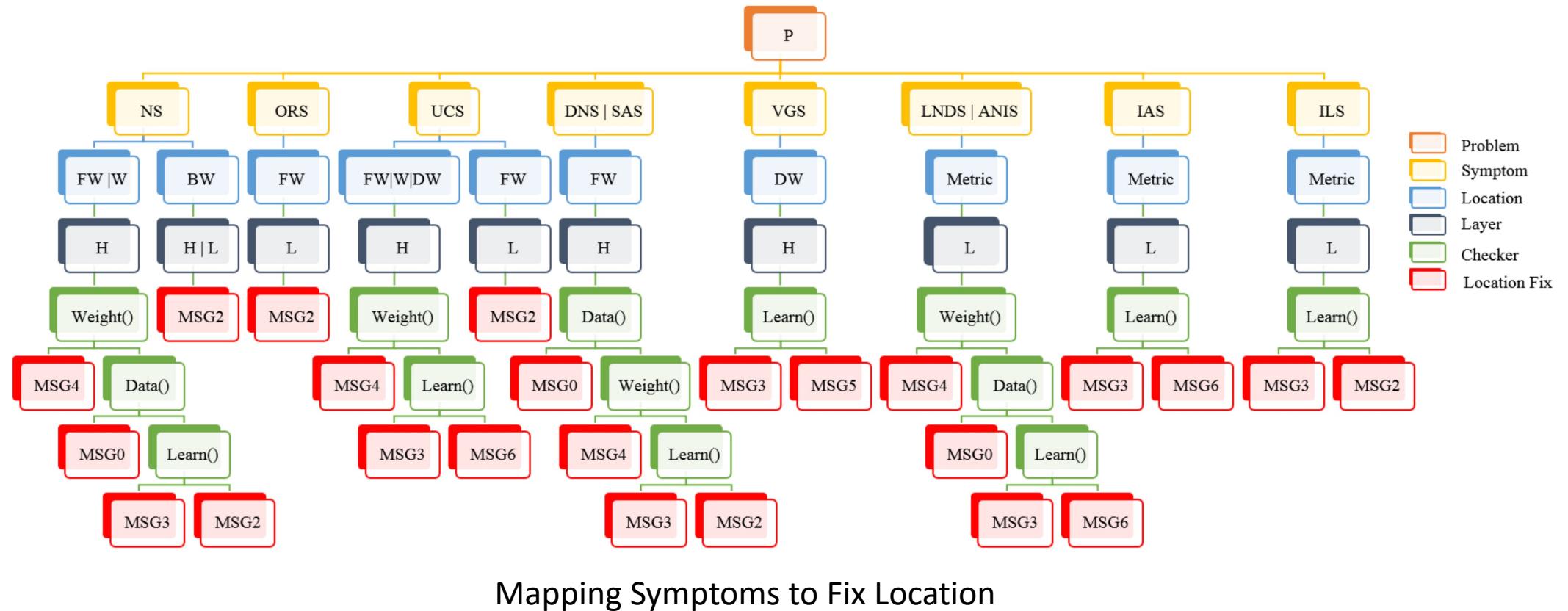
Algorithm: Failure Symptoms Detection



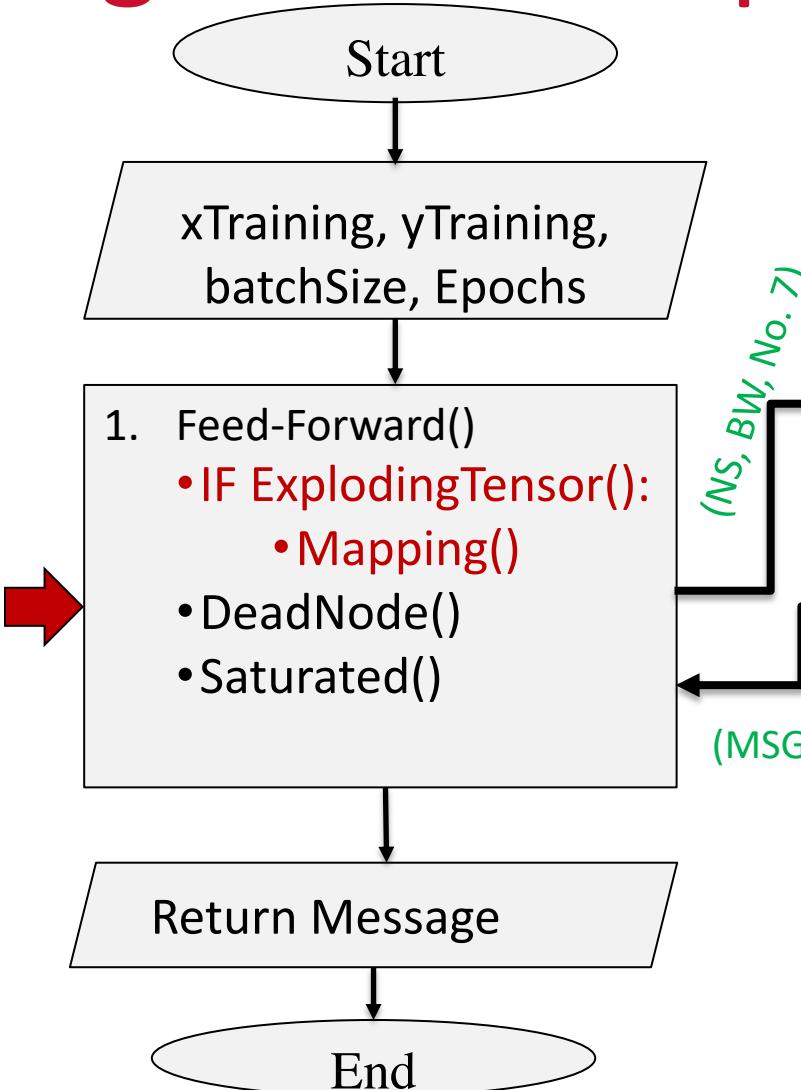
Algorithm: Failure Symptoms Detection



Decision Tree



Algorithm: Mapping Symptoms to Locate Fix



Mapping(symptom, location, LayerNo)

- SatisfiedCondition()
- ImproperData()
- WeightInitialization()
- TuneLearn()

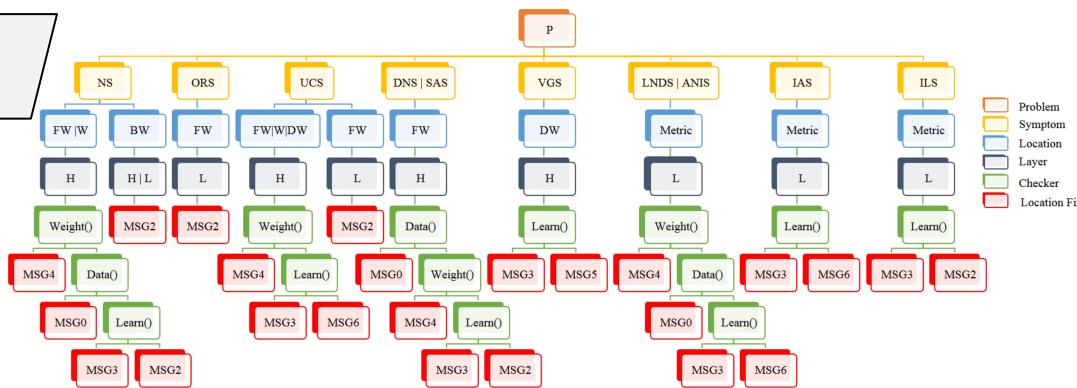
if data is not normalized

if weights aren't initialized correctly

if learning rate is very high or low

Return Message

End



Mapping Symptoms to Fix Location

Algorithm: Failure Symptoms Detection

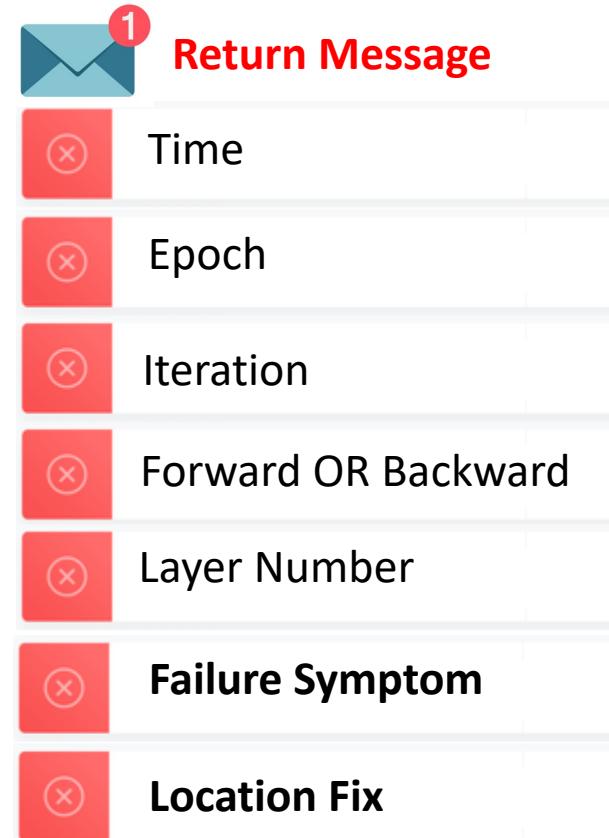
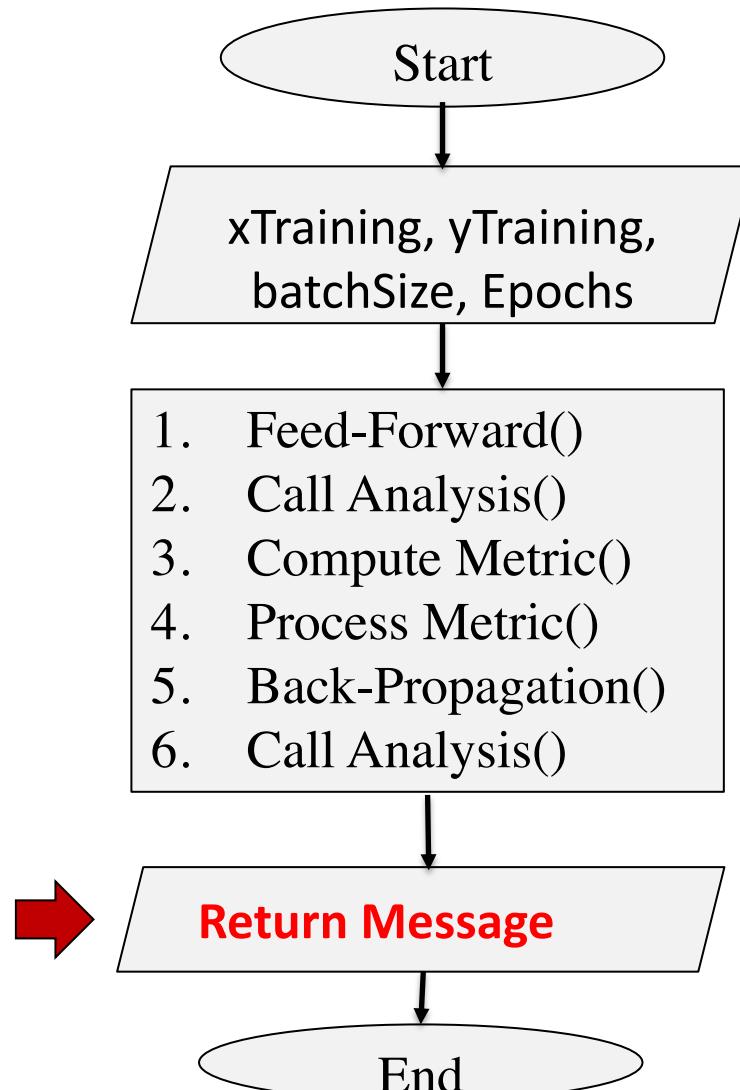


TABLE III: Abbreviation of Failure Symptoms

No	Symptoms	Abbreviation
1	Numerical Errors	NS
2	Unchanged weight	UCS
3	Saturated Activation	SAS
4	Dead Node	DNS
5	Out of Range	ORS
6	Loss Not Decreasing	LNDS
7	Accuracy Not Increasing	ANIS
8	Vanishing Gradient	VGS
9	Invalid Loss	ILS
10	Invalid Accuracy	IAS
11	Correct Model	CM

TABLE IV: Abbreviation of Actionable Changes

No	Message Guideline	Abbreviation
1	Improper Data	MSG0
2	Change the loss function	MSG1
3	Change the activation function	MSG2
4	Change the learning rate	MSG3
5	Change the initialization of weight	MSG4
6	Change the layer number	MSG5

Research Questions



Validation



Comparison



Limitation



Ablation

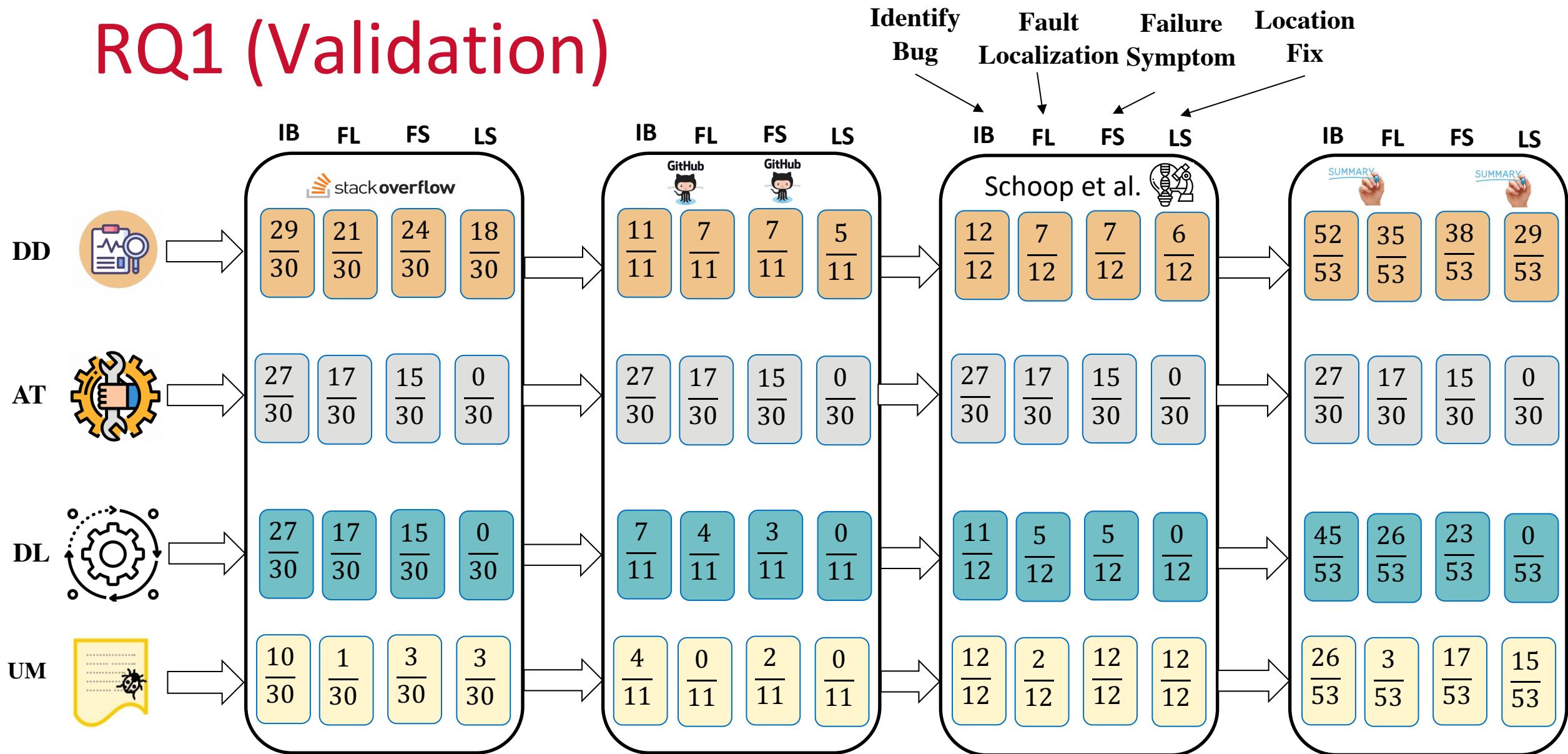
RQ1: Can our technique **localize faults and report the symptoms** in deep learning programs effectively?

RQ2: How **effectively** and **how fast** can our technique localize the faults and report the symptoms compared to existing methodologies?

RQ3: In which cases does our technique **fail to localize** the faults and report the symptoms?

RQ4: To what extent does each component of our approach **contribute to the overall performance**?

RQ1 (Validation)

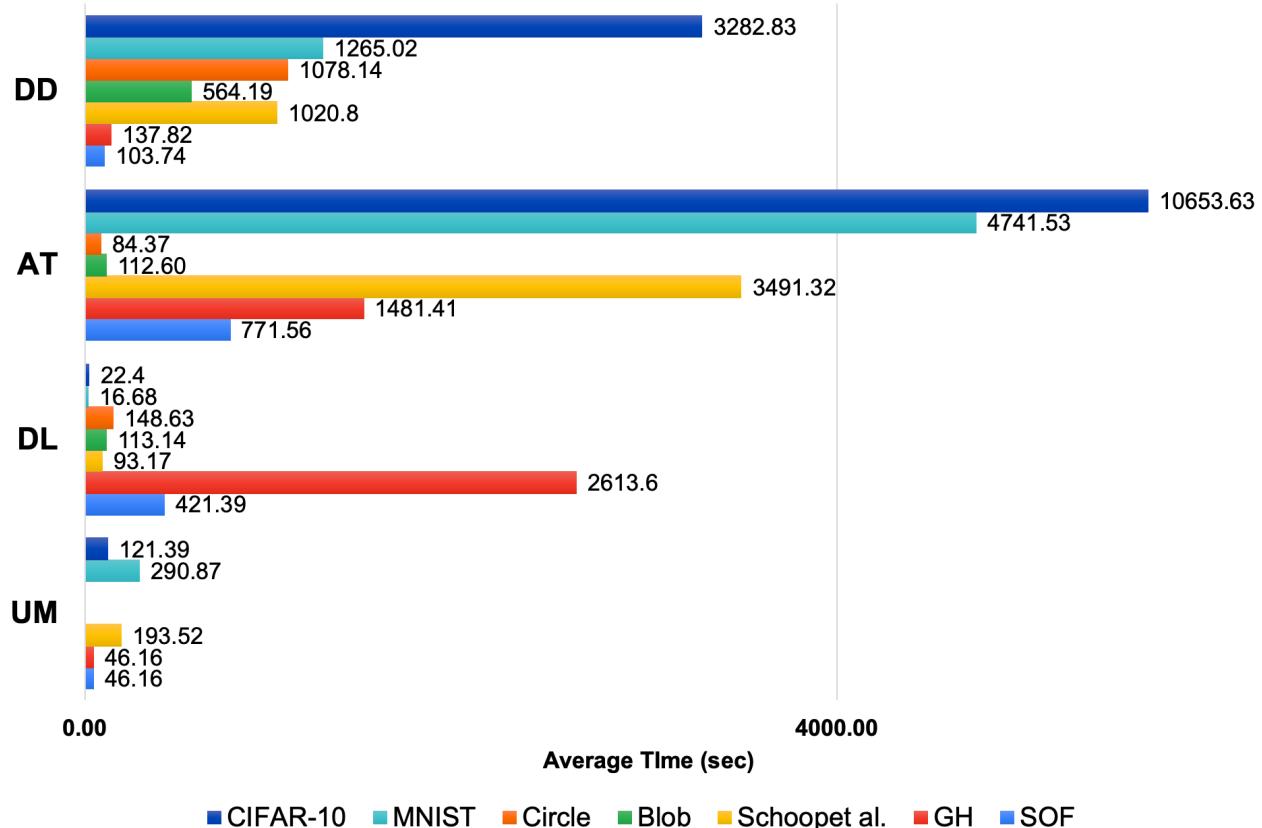


RQ2 (Comparison)



We measured the analysis time for:

- **UMLAUT (UM), DeepLocalize (DL), AUTOTRAINER (AT), and DeepDiagnosis (DD)**
- Using Stack Overflow (SOF), GitHub (GH), UM, and AT benchmarks.



DeepDiagnosis is the most efficient for Stack Overflow and *Schoop et al.*'s model.

- **Slower** than UMLAUT on the GitHub models.
- **Faster** than AT on all benchmarks except for the Blob and Circle datasets.



Our runtime overhead mainly is the online dynamic analysis performed on the internal parameters.

RQ3 (Limitation)

- • Our prototype does not support all Keras APIs such as *fit_generator()* instead of *fit()* function.
- • Our technique does not support Recurrent Neural Networks (RNNs) models.
- • Our approach does not support some hyperparameters such as the batch size, epoch, and dropout rate.

RQ4 (Ablation)

1. The result shows which **procedure** report the **symptom** in buggy model:

Dataset	Exploding	Saturated	DeadNode	OutofRange	UnchangeWeight	InvalidAcc	AccNotIncreasing	VanishingGradient
Stack Overflow, GitHub, Schoop et al.	23	7	5	5	2	2	1	1
AUTOTRAINER	50	5	8	0	23	0	21	31
Total	73	12	13	5	25	2	22	32

2. Our study found that **data preparation** is a frequently occurred issue.
3. **SGS** benchmark does not have a very deep model that contains **many layers**.
4. **VanishingGradient()** is invoked very **frequently** in **AUTOTRAINER** models, because this dataset has **many layers** using **sigmoid** and **tanh** as activation functions.

32

Conclusions and Future Work



We put forward an approach dynamically detecting bugs in DNNs that uses a map between symptoms to root causes.



We compared DeepDiagnosis against SoTA using several benchmarks. The results show that DeepDiagnosis is both more comprehensive and efficient than prior work.



Future work includes adding support to automatically:



- Detecting additional bug types in other model architectures
- Repairing models without developer intervention

33

DeepDiagnosis: Automatically Diagnosing Faults and Recommending Actionable Fixes in Deep Learning Programs

Mohammad Wardat

wardat@iastate.edu

Dept. of Computer Science, Iowa State University
226 Atanasoff Hall, Ames, IA, USA

Wei Le

weile@iastate.edu

Dept. of Computer Science, Iowa State University
226 Atanasoff Hall, Ames, IA, USA

Breno Dantas Cruz

bdantasc@iastate.edu

Dept. of Computer Science, Iowa State University
226 Atanasoff Hall, Ames, IA, USA

Hridesh Rajan

hridesh@iastate.edu

Dept. of Computer Science, Iowa State University
226 Atanasoff Hall, Ames, IA, USA

Motivation



Q1) What is the network architecture that's preventing the model from learning?
Q2) What is the root cause, and how to fix the problem?

UMLUAT: No Output, STOP, 104.65 Seconds
DeepLocalize: layer 7: Numerical Error, STOP, 2.14 Seconds
AUTOTRAINER: Using 'SelU' each layers, Using 'BatchNormalization', ... Unsolved, STOP, 495.83 Seconds

Code snippet:

```
i: model = Sequential()
2: model.add(Dense(128, 50))
3: model.add(Activation('relu'))
4: model.add(Dropout(0.2))
5: model.add(Dense(50, 50))
6: model.add(Activation('relu'))
7: model.add(Dropout(0.2))
8: model.add(Dense(1))
9: model.add(Activation('softmax'))
10: model.compile(loss='binary_crossentropy', optimizer=RMSprop())
11: model.fit(X,Y,batch_size=epoch,validation_data=(X_test,Y_test))
```

<https://stackoverflow.com/questions/31880720/>

Contributions

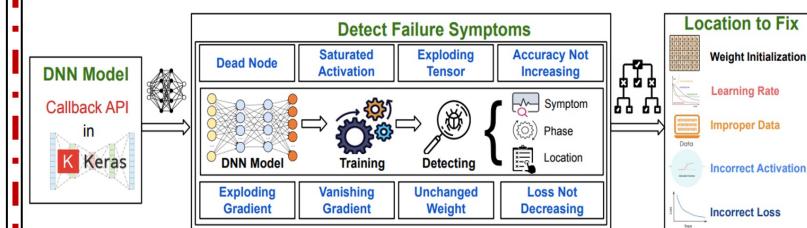
We study the symptoms and their root causes to propose a mapping algorithm of symptoms to fix.

We provide a benchmark with 444 models that practitioners can use to evaluate their techniques.

Our prototype is more comprehensive and efficient than:

- UMLUAT
- DeepLocalize
- AUTOTRAINER

DeepDiagnosis Overview



Preprint: arxiv.org/abs/2112.04036



Artifacts: github.com/DeepDiagnosis/ICSE2022

34