Deep Neural Network Test Coverage: How Far Are We?

Junjie Chen, Ming Yan, Zan Wang, Yuning Kang, Zhuo Wu

Abstract—DNN testing is one of the most effective methods to guarantee the quality of DNN. In DNN testing, many test coverage metrics have been proposed to measure test effectiveness, including structural coverage and non-structural coverage (which are classified according to whether considering which structural elements are covered during testing). Those test coverage metrics are proposed based on the assumption: they are correlated with test effectiveness (i.e., the generation of adversarial test inputs or the error-revealing capability of test inputs in DNN testing studies). However, it is still unknown whether the assumption is tenable. In this work, we conducted the first extensive study to systematically validate the assumption by controlling for the size of test sets. In the study, we studied seven typical test coverage metrics based on 9 pairs of datasets and models with great diversity (including four pairs that have never been used to evaluate these test coverage metrics before). The results demonstrate that the assumption fails for structural coverage in general but holds for non-structural coverage on more than half of subjects, indicating that measuring the difference of DNN behaviors between test inputs and training data is more promising than measuring which structural elements are covered by test inputs for measuring test effectiveness. Even so, the current non-structural coverage metrics still can be improved from several aspects such as unfriendly parameters and unstable performance. That indicates that although a lot of test coverage metrics have been proposed before, there is still a lot of room for improvement of measuring test effectiveness in DNN testing, and our study has pointed out some promising directions.

Index Terms—DNN Test Coverage, Empirical Study, Test Effectiveness

1 Introduction

DEEP neural network (DNN) has been widely studied in recent years and has achieved great success in many domains, e.g., autonomous driving cars [1], speech recognition [2], face recognition [3], machine translation [4], medical diagnosis [5], and code analysis [6], [7]. However, like traditional software systems, DNN also contains bugs [8], [9], [10]. DNN bugs could lead to many unexpected behaviors, even disasters in safety-critical domains. For example, an Uber autonomous driving car killed a pedestrian in Tempe, Arizona in 2018¹. Also, a Tesla Model S in autopilot mode crashed into a parked fire trunk with light flashing on a California freeway in 2018². Therefore, it is very critical to ensure the DNN quality.

DNN testing is one of the most effective methods to ensure the DNN quality. As with traditional software testing [11], [12], [13], [14], [15], [16], [17], [18], one important aspect in DNN testing is to measure test effectiveness. Indeed, many test coverage metrics have been proposed to measure test effectiveness in recent years, e.g., neuron coverage measuring which DNN elements are covered during the testing process [8], [9] and surprise coverage measuring

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

- Junjie Chen, Ming Yan, Zan Wang, Yuning Kang are with College of Intelligence and Computing, Tianjin University, Tianjin, 300350, China. E-mail: [junjiechen,yanming,wangzan,kangyuning]@tju.edu.cn.
- Zhuo Wu is with Tianjin International Engineering Institute, Tianjin University, Tianjin, 300350, China.
 E-main:wuzhuo@tju.edu.cn.
- 1. https://www.vice.com/en_us/article/9kga85/uber-is-giving-up-on-self-driving-cars-in-california-after-deadly-crash
- 2. https://www.newsweek.com/autonomous-tesla-crashes-parked-fire-truck-california-freeway-789177

the difference of DNN behaviors between test inputs and training data [19]. When evaluating these test coverage metrics to be effective, these work adopted the number of generated adversarial test inputs [8], [9], [19] or the number of error-revealing test inputs [10], [20] to represent test effectiveness. That is, these test coverage metrics are based on the assumption: they are correlated with either the generation of adversarial test inputs or the error-revealing capability of test inputs.

However, it is still unknown whether the assumption is tenable. Although some experiments have been conducted when these test coverage metrics were proposed, these experiments cannot validate the assumption due to ignoring the influence of test-set sizes. For example, in the existing experiments [9], [19], they first generated the same number of adversarial test inputs as the number of original test inputs (also called natural test inputs), and then compared the test coverage of natural test inputs and that of both natural and adversarial test inputs. They found that the latter is larger than the former, and thus concluded that the test coverage metrics are correlated with the generation of adversarial test inputs. However, the experiments do not control for the number of test inputs when comparing test coverage, and thus the increased coverage may be due to the generation of adversarial test inputs or just increasing the number of test inputs.

Further, we conducted a preliminary study to investigate whether adding natural test inputs can also increase test coverage by taking the DNN model *ResNet-20* based on the dataset *CIFAR-10* and the test coverage *KMNC* (to be introduced in Section 2) as an example. In this study, we first randomly selected 5,000 test inputs from *CIFAR-10* as the original test set and regarded the remaining 5,000 test inputs

in CIFAR-10 as the newly added natural test inputs, and then generated 5,000 adversarial test inputs via C&W [21] (to be introduced in Section 3.4), a widely-used adversarial test input generation method. We found that the KMNC values of the original test set, the original test set and added natural test inputs, and the original test set and adversarial test inputs are 58.53%, 64.50%, and 63.83%, respectively. The results demonstrate that adding natural test inputs can also increase test coverage, even the increment is larger than that achieved by the same number of adversarial test inputs. Therefore, it is still unclear whether these test coverage metrics are correlated with the generation of adversarial test inputs or the error-revealing capability of test inputs independently of the size of test sets.

In this work, we conducted the first empirical study to systematically investigate whether the assumption is tenable by controlling for the size of test sets. There are three main reasons for the necessity of this study. First, as discussed above, there is no study that has validated the assumption. Second, test coverage is an important aspect in DNN testing and many test coverage metrics have already been proposed, and thus it is time to investigate whether the current test coverage metrics are good enough or there is still a long way to go. Third, these test coverage metrics were evaluated based on a small set of datasets and DNN models, and there is no study systematically comparing these test coverage metrics based on the same set of datasets and DNN models, and thus it is necessary to revisit/compare them based on the same and larger set of datasets and DNN models to tell us which test coverage is better.

In the study, we used 9 pairs of datasets and DNN models with great diversity as subjects by considering 1) both small and large datasets, 2) both classification and regression models, 3) both image and speech domains, 4) both CNN and RNN models, and 5) both the subjects that have been used in the work proposing the corresponding test coverage and the newly added subjects. Moreover, we considered five structural coverage metrics (i.e., DNC, TKNC, KMNC, NBC, and SNAC) [8], [9] and two non-structural coverage metrics (i.e., LSC and DSC) [19] as the studied coverage metrics since they are indeed representative. More specifically, DNC is the first neuron coverage and won the best paper award in SOSP'17 [8]. TKNC, KMNC, NBC, and SNAC further complement DNC and won the distinguished paper award in ASE'18 [9]. LSC and DSC are state-of-theart test coverage metrics in DNN testing. Besides, we also studied the influence of many factors, including different adversarial test input generation methods, different sizes of test sets, and different test-effectiveness measurements (i.e., the ratio of adversarial test inputs and the ratio of errorrevealing test inputs).

From our study, we got the following major findings:

 In general, structural coverage does not have moderate to extremely strong positive correlation with both the generation of adversarial test inputs and the errorrevealing capability of test inputs when controlling for the size of test sets. Relatively, SNAC performs the best among the five studied structural coverage metrics, indicating the importance of considering the upper corner-case regions of neurons.

- Different from structural coverage, non-structural coverage has different degrees of correlation for different subjects and LSC has moderate to extremely strong positive correlation for more than half of subjects, indicating that measuring the difference of DNN behaviors between test inputs and training data is more promising than measuring which structural elements are covered by test inputs for measuring test effectiveness.
- Although non-structural coverage (especially LSC) is more promising than structural coverage, the current non-structural coverage metrics still need to be improved from several aspects, e.g., unfriendly parameters and unstable performance.
- When evaluating a DNN testing metric or a DNN testing technique, it is not enough to only use the adversarial test input generation method FGSM [22] and the simple datasets (i.e., MNIST, CIFAR-10, Driving). More adversarial test input generation methods (e.g., C&W, BIM, and JSMA) and more complex datasets (e.g., CIFAR-100 and ImageNet) should be considered due to their different characteristics.

This work makes the following main contributions:

- The first extensive study investigating the correlation between test coverage and test effectiveness for DNN. Our study used 9 pairs of datasets and DNN models with great diversity as subjects by considering different DNN structures (i.e., CNN and RNN), different domains (i.e., image and speech), different tasks (i.e., classification and regression), and both the subjects that have been used before and newly added subjects.
- An extensive analysis considering various factors, including different test coverage metrics (i.e., five structural coverage metrics), different adversarial test input generation methods, different sizes of test sets, and different test-effectiveness measurements (i.e., the ratio of adversarial test inputs and the ratio of error-revealing test inputs).
- Many empirical conclusions, for example, measuring the difference of DNN behaviors between test inputs and training data (non-structural coverage) is more promising than measuring which structural elements are covered by test inputs (structural coverage) for measuring test effectiveness but there is still a lot of room for further improvement. Moreover, our study has pointed out some promising directions.

2 BACKGROUND

In the literature, many test coverage metrics have been proposed to measure test effectiveness in the field of DNN testing [8], [9], [19]. According to whether considering which structural elements are covered during testing, we classify test coverage into two categories: *structural coverage* and *non-structural coverage*. In the following, we will introduce several typical test coverage metrics in detail, which are also the studied test coverage in our work. More details about other test coverage metrics will be presented in Section 7.

2.1 Structural Coverage

Structural coverage considers which structural elements are covered during testing, mainly referring to various neuron

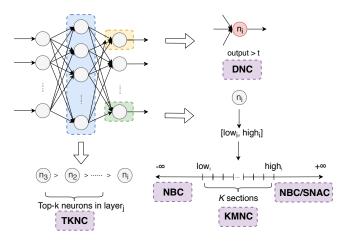


Fig. 1: Overview of Structural Coverage

coverage. Here, we study five typical structural coverage metrics, and Figure 1 shows the overview of the five structural coverage. The first neuron coverage was proposed by Pei et al. [8], which is called **DNC** (DeepXplore's Neuron Coverage). DNC divides the state of a neuron into *activated* and *non-activated*, and measures the ratio of activated neurons in a DNN. If the output value of a neuron is larger than a pre-defined threshold t after executing a test input, DNC regards the neuron as an activated neuron.

Ma et al. [9] further proposed a set of neuron coverage metrics at several levels of granularities, including TKNC (Top-K Neuron Coverage), KMNC (K-Multisection Neuron Coverage), NBC (Neuron Boundary Coverage), and SNAC (Strong Neuron Activation Coverage). TKNC is a layer-level neuron coverage metric, which considers top-k neurons in a layer, which are ranked in the descending order of their output values after executing a test input, to be covered. It measures the ratio of neurons that have once been topk neurons in the corresponding layer. KMNC is a more fine-grained neuron coverage metric. It first obtains the range of its output values for each neuron from training data, denoted as $[low_i, high_i]$ for the neuron n_i , and then partitions the range into k sections. If the output value of a neuron falls into a section after executing a test input, KMNC regards the section as a covered section. KMNC measures the ratio of covered sections of all neurons in a DNN. The insight behind KMNC is that, different sections may represent different functional modules of a DNN and more fine-grained neuron coverage can increase discrimination of different test inputs. NBC also obtains the range of its output values for each neuron from training data like KMNC. The difference is that, NBC considers whether the corner-case region, i.e., $(-\infty, low_i)$ and $(high_i, +\infty)$, of a neuron is covered after executing a test input. SNAC is similar to NBC, but the difference is that SNAC only considers whether the upper corner-case region, i.e., $(high_i)$, $+\infty$), is covered after executing a test input.

Besides, the state-of-the-art structural coverage is MCDC-inspired neuron coverage [23], which is tailored to the distinct features of a DNN. It instantiates the decision change and condition change in DNNs by defining the *sign change* and *value change* of neurons, which show the degree of the change of neuron activation values under two different test inputs. Based on the instantiation, MCDC-

inspired neuron coverage contains a set of specific metrics that compute the percentage of neuron pairs within adjacent layers that are covered by test inputs with respect to the corresponding coverage (e.g. **Sign-Sign Coverage**, which measures the independent impact of the sign change of a neuron on the sign of a neuron in the next layer). As collecting MCDC-inspired neuron coverage is very time-consuming (due to high time complexity and constraint solving) as demonstrated by the existing work [23], it is hard to apply it to large datasets and models and thus we did not use it in our extensive study. Instead, we discuss whether our conclusions can be generalized to this structural coverage by conducting a small experiment (in Section 6.1).

2.2 Non-Structural Coverage

Non-structural coverage does not consider whether structural elements of a DNN are covered during testing. The most typical non-structural coverage is surprise coverage [19]. It regards the difference of DNN behaviors between a test input and training data as surprise adequacy (SA) of the test input. Surprise coverage measures the range of SA values that a set of test inputs cover. More specifically, Kim et al. [19] proposed two kinds of SA: Likelihood-based Surprise Adequacy (LSA) and Distance-based Surprise Adequacy (DSA). LSA refers to the surprise of a test input with respect to the estimated density of each activation value in a set of activation traces observed over neurons of a selected layer for training data via KDE (Kernel Density Estimation) [24]. LSA only considers the training data whose label is the same as the predicted class of the test input. DSA is defined using the Euclidean distance between the activation trace observed over neurons of a selected layer for a test input and a set of activation traces for training data. More specifically, DSA needs to find the closest neighbor of the test input that shares the same class, denoted as x_a , and also the closest neighbor of x_a in a different class.

According to LSA and DSA, the corresponding surprise coverage metrics are called LSC and DSC, respectively. More specifically, given the upper and lower bounds of LSA/DSA denoted as *U* and *L*, LSC/DCS first divides [*L*,*U*] into *n* sections and then measures the ratio of covered sections. If the LSA/DSA value of a test input falls into a section, LSC/DCS regards the section as a covered section. In particular, a test input with a quite large SA value may be irrelevant to the problem domain (e.g., an image of an apple is irrelevant to the testing of traffic sign classifiers), and thus the parameter *U* can be used to filter out those irrelevant test inputs. Please note that DCS cannot be applicable to regression models [19].

3 STUDY DESIGN

In this section, we present the design of our empirical study. Our study aims to validate the assumption: *DNN test coverage is correlated with the generation of adversarial test inputs or the error-revealing capability of test inputs*. Here, we controlled for the size of test sets to validate the assumption. Therefore, the generation of adversarial test inputs can be measured by the ratio of adversarial test inputs in a test set, and the error-revealing capability of test inputs can be measured by the ratio of error-revealing test inputs in a test

set. That is, we validate the assumption by addressing the following two research questions:

- RQ1: Is test coverage correlated with the ratio of adversarial test inputs?
- **RQ2**: Is test coverage correlated with the ratio of errorrevealing test inputs?

Please note that adversarial test inputs are intentionally designed to cause a DNN model to make wrong predictions, but not all adversarial inputs can successfully fool the DNN model. Error-revealing test inputs are the inputs for which the DNN model makes wrong predictions. Both natural test inputs and adversarial test inputs can be error-revealing. That is, adversarial test inputs and error-revealing test inputs are not completely overlapped, and thus our study investigates the correlation between test coverage and both adversarial test inputs and error-revealing test inputs.

3.1 Datasets and DNN models

In our study, we used 9 pairs of datasets and DNN models as subjects in total. The DNN models include *LeNet-5*, *VGG-16*, *VGG-19*, *ResNet-20*, *ResNet-32*, *ResNet-50*, *Dave*, *Chauffeur*, and *DeepSpeech*. These models are trained based on 6 popular datasets respectively, including *MNIST*, *CIFAR-10*, *CIFAR-10*, *ImageNet*, *Driving*, and *Speech-Commands*. MNIST is a handwritten digit dataset³, CIFAR-10 and CIFAR-100 are 10-class and 100-class ubiquitous object datasets respectively⁴, Driving is an Udacity-provided autonompus driving dataset⁵, ImageNet is a more realistic and complex image dataset⁶, and Speech-Commands is a sequential dataset containing a set of one-second .wav audio files, each of which contains a single spoken English word⁷.

To more extensively study these test coverage, we carefully considered the diversity of subjects, including 1) both small datasets (e.g., MNIST) and large datasets (e.g., ImageNet), 2) normal classification models, multi-label classification models (i.e., DeepSpeech), and regression models, 3) both image and speech domains, 4) both CNN and RNN models, and 5) the subjects collected from the existing work proposing these test coverage metrics and the subjects that have never been used to evaluate these test coverage metrics before.

Table 1 shows the basic information of our subjects, where the last seven columns represent the model size, the number of training instances, the number of original test inputs, the model accuracy, the model task, the domain of the model, and the network type, respectively. Table 2 presents whether the subjects have been used to evaluate the studied test coverage metrics in the work proposing them before. Here, five subjects have been used to evaluate at least one test coverage metric before and four subjects have never been used to evaluate any test coverage metric.

3.2 Implementations

Since the code for collecting DNC, LSC, and DSC is released by the existing work [8], [19], we directly adopted the code

- 3. http://yann.lecun.com/exdb/mnist/.
- 4. http://www.cs.toronto.edu/~kriz/cifar.html.
- 5. https://udacity.com/self-driving-car.
- 6. http://www.image-net.org.
- $7. \quad https://github.com/bjtommychen/Keras_DeepSpeech2_Speech \\ Recognition.$

to collect them in our study. As the code for collecting TKNC, KMNC, NBC, and SNAC is not available, we reimplemented the collection code according to the descriptions in the work [9]. To check whether our implementations are correct, we applied our implementations to collect the corresponding coverage of the subjects used in the existing work [9] and then compared whether our results are consistent with those reported in the paper [9]. We found that the results are indeed consistent within the margin of error. For many parameters of these studied coverage metrics, we set them following the existing work [8], [9], [19]. More specifically, the t value of DNC is set to be 0.5, the k value of TKNC is set to be 3, the k value of KMNC is set to be 1,000, and the n value of LSC and DSC is set to be 1,000. LSC and DSC also have some parameters that tend to be specifically set for different subjects [19], i.e., the selected layer to calculate LSA/DSA, the upper bound U and lower bound L of LSC and DSC. However, there is no guide provided to help set them for different subjects. Therefore, by communicating with the authors of LSC and DSC, for each subject we set the upper bound U and the lower bound L of LSC/DSC to be the largest and smallest LSA/DSA values of all the test inputs respectively as they suggested. Also, we set the selected layer required by LSA/DSA to be the last-hidden layer of a DNN model according to the existing work [25], [26], since these work have demonstrated that deeper layers tend to perform better in DNN testing.

3.3 Test Set Creation

In our study, we controlled for the size of test sets to answer these research questions. That is, we created the test sets with the same size. First, we used an adversarial input generation method to generate the same number of adversarial test inputs as the number of original test inputs for each subject. Second, we *randomly* selected $N^*\alpha$ original test inputs and $N^*(1-\alpha)$ adversarial test inputs to form a test set, where N is the test-set size and α is a *random* number between 0 to 1 to represent the mixing ratio of natural test inputs and adversarial test inputs. In our study, we considered N to be 100, 500, and 1000, respectively. Third, we created 500 test sets for each subject under each test-set size with each adversarial input generation method.

Besides, when answering RQ2, we also considered the test sets without adversarial test inputs since some original test inputs are able to reveal errors. That is, we also *randomly* selected N natural test inputs to form a test set, and created 500 test sets for each subject under each test-set size. Here, we also considered N to be 100, 500, and 1000, respectively.

3.4 Adversarial Test Input Generation Methods

In our study, we adopted four advanced adversarial test input generation methods to generate adversarial test inputs for general image datasets, i.e., MNIST, CIFAR-10, CIFAR-100, and ImageNet. They are FGSM (Fast Gradient Sign Method) [22], C&W (Carlini&Wagner) [21], BIM (Basic Iterative Methods) [27], and JSMA (Jacobian-based Saliency Map Attack) [28]. These methods have been widely used in the existing work [9], [19]. For the dataset Driving, we did not use the four methods but used the Patch and Light methods to generate adversarial test inputs since most of the

TABLE 1: Datasets and DNN models

ID	Dataset	Model	Size(KB)	#Train	#Test	Acc(%)	Task	Domain	Network
1	MNIST	LeNet-5	1,093	60,000	10,000	98.68	classification	image	CNN
2	CIFAR-10	VGG-16	21,814	50,000	10,000	96.07	classification	image	CNN
3	CIFAR-10	ResNet-20	3,507	50,000	10,000	91.45	classification	image	CNN
4	CIFAR-100	ResNet-32	10,615	50,000	10,000	71.42	classification	image	CNN
5	ImageNet	VGG-19	562,176	1,400,000	50,000	64.73	classification	image	CNN
6	ImageNet	ResNet-50	100,352	1,400,000	50,000	68.27	classification	image	CNN
7	Driving	Dave	8,306	101,396	5,614	90.34	regression	image	CNN
8	Driving	Chauffeur	76,834	101,396	5,614	99.94	regression	image	CNN&RNN
9	Speech-Commands	DeepSpeech	6,734	51,776	6,471	94.53	classification	speech	RNN

TABLE 2: Whether the subjects have been used to evaluate the studied test coverage metrics before?

ID	DNC	TKNC	KMNC	NBC	SNAC	LSC	DSC
1	/	✓	✓	/	✓	0	0
2	X	Х	×	X	X	0	0
3	X	Х	×	X	Х	0	0
4	X	Х	×	X	Х	X	X
5	✓	✓	✓	✓	✓	X	X
6	✓	✓	✓	✓	✓	X	X
7	✓	Х	×	X	Х	1	/
8	X	×	×	Х	Х	/	✓
9	X	X	Х	Х	Х	Х	X

* V/X represents that the subject was/wasn't used to evaluate the corresponding test coverage metric. O represents special cases: LSC and DSC were evaluated based on MNIST and CIFAR-10, but the used DNN models in the existing study [19] are the designed DNN by themselves rather than the widely-used models for the two datasets (i.e., LeNet-5 for MNIST and VGG-16/ResNet-20 for CIFAR-10). In our study, we adopted the widely-used models for the two datasets.

existing work using this dataset adopted the two methods for it [19], [25], [26]. For the dataset Speech-Commands, all the above-mentioned methods cannot be applicable to it since this dataset is sequential data (i.e., audio) rather than images. Therefore, we adopted the widely-used CTCM (CTC loss based Method) [29] to generate adversarial test inputs for this dataset. In the following, we introduce the seven adversarial test input generation methods in detail.

- FGSM is a gradient-based adversarial test input generation method and its key idea is to increase loss. It first obtains the derivative for an input, then leverages symbolic function to determine its gradient direction, and finally generates an adversarial test input by perturbing the original input based on a step size and the direction.
- BIM is an improved version of FGSM, which divides the perturbation into multiple small steps and then performs multiple iterations to generate an adversarial test input.
- JSMA is a targeted adversarial test input generation method, which is based on L₀ distance metric. Its key insight is that different input features have different degrees of influence on different outputs produced by the classifier. It aims to find the features corresponding to a specific output of the classifier, and then produce an adversarial test input with the specific output by enhancing the corresponding features in the original test input.
- C&W is an optimization-based adversarial test input generation method. It aims to minimize the distance between
 the adversarial test input and the original test input and
 is able to generate an adversarial test input changing the
 output class of the original test input. It would be better

if the probability of the wrong class is larger. We adopted the version of C&W based on the L_2 distance metric.

- Patch is to randomly block some parts of a test input to generate an adversarial test input for the dataset Driving, in order to simulate block some parts of a camera.
- Light is to randomly change the intensities of lights for a test input to generate an adversarial test input for the dataset Driving.
- CTCM is a targeted adversarial test input generation method. Its key idea is to construct a special "loss function" based on CTC Loss [30] that takes a desired transcription and an audio file as input, and then minimize the loss by making slight changes to the input through gradient descent. Then, an adversarial input is generated.

3.5 Measurements

To investigate the correlation between test coverage and the ratio of adversarial test inputs or the ratio of error-revealing test inputs, we adopted the widely-used Spearman correlation method [31], which is used to assess how well the relationship between two variables can be described using a *monotonic* function. The sign of the Spearman correlation coefficient represents the positive or negative correlation between two variables, the absolute value of the Spearman correlation coefficient represents the correlation degree between two variables, and the p value reported by the Spearman correlation method represents whether the correlation is statistically significant (at the level of 0.05).

According to the existing work [11], [32], the correlation can be divided into four categories: weak correlation (the absolute value of the correlation coefficient is smaller than 0.4), moderate correlation (the absolute value of the correlation coefficient is between 0.4 and 0.7), strong correlation (the absolute value of the correlation coefficient is between 0.7 and 0.9), and extremely strong correlation (the absolute value of the correlation coefficient is larger than 0.9).

Also, we recorded the time spent on collecting coverage for each test set in order to compare the efficiency of collecting different test coverage.

3.6 Process

We present the process of our study for each subject.

First, we applied each applicable adversarial test input generation method to generate the same number of adversarial test inputs as the number of original test inputs.

Second, for each test-set size (i.e., 100, 500, and 1000) and each adversarial test input generation method, we created

500 test sets that mix adversarial test inputs with natural test inputs. Also, for each test-set size (i.e., 100, 500, and 1000), we created 500 test sets that only contain natural test inputs. This step has been presented in detail in Section 3.3.

Third, for each created test set, we collected its DNC, TKNC, KMNC, NBC, SNAC, LSC, and DSC and recorded the corresponding collection time, respectively.

Finally, for each adversarial input generation method and the corresponding 500 created test sets mixing adversarial test inputs with natural test inputs, we calculated the Spearman correlation between each coverage and the ratio of adversarial test inputs, and the Spearman correlation between each coverage and the ratio of error-revealing test inputs. For the 500 created test sets only containing natural test inputs, we calculated the Spearman correlation between each coverage and the ratio of error-revealing test inputs.

Our study was conducted on a workstation with 32-core Intel Xeon E5-2640, 128GB memory, CentOS 7.6 operating system. In particular, it took more than six months to complete all the experiments through four processes in parallel.

4 Results and Analysis

4.1 RQ1: Is Test Coverage Correlated with the Ratio of Adversarial Test Inputs?

4.1.1 Overall effectiveness of structural coverage.

Table 3 shows the Spearman correlation results between test coverage and the ratio of adversarial test inputs. We used the results of the test-set size of 1,000 as the representative shown in this table, and the influence of different test-set sizes is very slight and will be discussed in Section 4.1.4. From Columns 3-7 in this table, we found that the five structural coverage metrics have moderate to extremely strong positive correlation with the ratio of adversarial test inputs in few cases, but have weak positive correlation, even surprisingly no significant correlation and negative correlation, in most cases. For example, KMNC has weak positive correlation for 20.69% cases, no significant correlation for 31.03% cases, and negative correlation for 44.83% cases, with the ratio of adversarial test inputs, but has moderate to extremely strong positive correlation for only 3.45% cases. That means that the five studied structural coverage metrics do not meet expectations well. That is, when controlling for the test-set size, both natural test inputs and adversarial test inputs have the similar capability to increase these structural coverage, and sometimes the coverage-increasing capability of the former is stronger than that of the latter. The reason is that adversarial test inputs are widely distributed in the regions where natural test inputs are distributed [33], causing that they cannot be distinguished by the studied structural coverage metrics effectively.

Through further observation, we found that SNAC performs better than the other four structural coverage metrics relatively and also performs better on regression models than classification models. More specifically, SNAC has moderate to extremely strong positive correlation for 17.24% cases, weak positive correlation for 31.03% cases, no significant correlation for 41.38%, and negative correlation for only 10.34% cases, with the ratio of adversarial test inputs. Moreover, SNAC has moderate to extremely strong positive correlation for 50% cases on regression models. That

TABLE 3: Spearman correlation between test coverage and the ratio of adversarial test inputs (the test-set size is 1,000)

ID	Adv.	DNC	TKNC	KMNC	NBC	SNAC	LSC	DSC
	FGSM	0.13	0.53	0.14	0.29	0.56	0.91	0.74
1	C&W	-0.10	0.18	0.27	0.03	0.38	0.91	0.46
1	BIM	-0.03	0.18	0.30	0.07	0.42	0.91	0.46
	JSMA	-0.14	0.53	0.58	0.11	0.26	0.92	0.53
	FGSM	0.03	-0.43	0.02	0.26	0.16	0.31	0.26
2	C&W	-0.18	-0.17	0.12	0.10	0.07	0.04	0.04
_	BIM	-0.18	-0.17	0.12	0.10	0.07	0.05	0.04
	JSMA	-0.25	-0.53	-0.04	-0.21	-0.18	-0.01	0.06
	FGSM	-0.34	0.07	-0.33	0.65	0.72	0.71	-0.21
3	C&W	-0.26	-0.20	-0.49	0.01	0.14	0.68	-0.34
9	BIM	-0.25	-0.20	-0.49	0.01	0.14	0.67	-0.34
	JSMA	-0.36	-0.43	-0.62	-0.08	0.06	0.54	-0.40
	FGSM	0.21	-0.12	-0.06	0.29	0.34	-0.19	-0.69
4	C&W	0.06	-0.11	-0.24	-0.02	-0.01	-0.16	-0.74
-	BIM	-0.01	-0.22	-0.24	-0.10	-0.09	-0.18	-0.74
	JSMA	0.01	-0.29	-0.39	-0.12	-0.07	-0.25	-0.74
	FGSM	0.34	-0.04	0.01	-0.20	0.00	0.17	-0.01
5	C&W	0.25	-0.07	0.03	-0.14	0.03	0.16	-0.13
9	BIM	0.25	-0.08	0.03	-0.11	0.04	-0.03	-0.20
	JSMA	0.20	-0.04	-0.11	-0.20	-0.05	0.19	-0.18
	FGSM	0.01	-0.12	-0.08	-0.04	0.04	-0.04	-0.06
6	C&W	-0.26	-0.71	0.24	-0.34	0.12	-0.01	-0.04
O	BIM	-0.02	-0.05	-0.05	-0.03	-0.03	-0.03	-0.04
	JSMA	-0.03	0.20	-0.24	-0.10	-0.08	0.01	0.01
7	Patch	0.19	0.12	-0.03	0.43	0.53	0.99	_
	Light	-0.86	-0.88	-0.43	-0.67	-0.65	0.96	
8	Patch	0.27	-0.98	-0.99	0.79	0.81	0.99	_
	Light	-0.06	-0.98	-0.99	0.17	0.16	0.71	
9	CTCM	-0.79	-0.96	-0.51	-0.04	0.19	0.80	0.78

The **bold** value represents that the corresponding p value is smaller than 0.05, indicating the corresponding correlation has statistical significance. The value marked with the shading represents **moderate positive** correlation. The value marked with the shading represents **strong positive** correlation. The value marked with the shading represents **extremely strong positive** correlation. The cells marked with "—" represent DSC is not applicable to regression models following the existing work [19].

indicates that the upper corner-case regions of neurons is more relevant to the generation of adversarial test inputs, and thus structural coverage may be further improved by better considering this characteristics.

Finding 1: In general, all the five studied structural coverage metrics do not have moderate to extremely strong positive correlation with the ratio of adversarial test inputs when controlling for the size of test sets. Among them, SNAC performs the best relatively, indicating the importance of considering the upper corner-case regions of neurons during the generation of adversarial test inputs.

4.1.2 Overall effectiveness of non-structural coverage.

We then analyzed the results of LSC and DSC from the last two columns in Table 3. We found that the correlation between LSC/DSC and the ratio of adversarial test inputs is very different for different subjects, and LSC performs better than DSC in general. More specifically, LSC has moderate to extremely strong positive correlation for five subjects, weak positive correlation for one subject, no significant correlation for two subjects, and negative correlation for one subject, when controlling for the size of test sets. DSC has moderate to strong positive correlation for two subjects, no significant

correlation for two subjects, negative correlation for three subjects, when controlling for the size of test sets. First of all, the non-structural coverage LSC has been shown to be moderately or strongly correlated with the ratio of adversarial test inputs for more than half of subjects, demonstrating that it outperforms structural coverage. That is, measuring the difference of DNN behaviors between test inputs and training data like LSC is a more promising direction to measure test effectiveness, compared with measuring which structural elements are covered by test inputs.

Further, we analyzed the subjects for which LSC/DSC does not have moderate to extremely strong positive correlation, and found that there are two reasons. First, for some subjects (e.g., ID is 2), the LSA/DSA values of many adversarial inputs are indeed larger than those of natural test inputs, but the LSA/DSA values of these adversarial inputs are very close to each other. Therefore, when calculating LSC/DSC based on LSA/DSA values, many LSA/DSA values are mapped to very few sections, leading to small LSC/DSC values. That is, for these subjects, although LSA/DSA can effectively distinguish adversarial inputs and natural test inputs, LSC/DSC cannot effectively reflect their differences due to the limitation of the calculation method of LSC/DSC. Second, for some other subjects (e.g., ID is 6), the LSA/DSA values of adversarial inputs are close to those of natural test inputs, indicating that both LSA/DSA and LSC/DSC cannot effectively distinguish them. That is, in many cases adversarial inputs do not perform differently from training data in terms of LSA/DSA, indicating the limitation of the capability of LSA/DSA to measure the differences of DNN behaviors. In particular, LSA/DSA suffer from the limitation for more complex datasets, i.e., CIFAR-100 and ImageNet. This is because for more complex datasets, the accuracy of the models is relatively low, but LSA/DSA depend on the predicted class for each test input. Therefore, the performance of LSA/DSA drops when suffering from low accuracy. That also indicates that in the studies of DNN testing, both simple datasets (such as MNIST, Driving, and CIFAR-10) and complex datasets should be considered due to their different characteristics.

Finding 2: The two studied non-structural coverage metrics have different degrees of correlation for different subjects. On the one hand, the non-structural coverage LSC has shown a greater potentiality than structural coverage; On the other hand, the performance of non-structural coverage is still limited by the calculation method of LSC/DSC or the capability of LSA/DSA to measure DNN behaviors' differences caused by low model accuracy.

Finding 3: The datasets MNIST, Driving, and CIFAR-10 are relatively simple, and thus only using them in the studies of DNN testing is not enough. More complex datasets such as ImageNet and CIFAR-100 should be considered since they tend to have different and rich characteristics.

TABLE 4: Spearman correlation between non-structural coverage and the ratio of adversarial test inputs when slightly changing the value of U (the test-set size is 1,000)

	LSC DSC		2		3		4	
Adv.	LSC	DSC	LSC	DSC	LSC	DSC	LSC	DSC
FGSM	0.08	0.07	-0.73	-0.81	-0.86	-0.50	-0.05	-0.76
C&W	-0.21	0.15	-0.72	-0.82	-0.86	-0.57	-0.10	-0.81
BIM	-0.23	0.13	-0.72	-0.81	-0.86	-0.57	-0.10	-0.80
JSMA	-0.15	0.17	-0.75	-0.79	-0.83	-0.58	-0.13	-0.81

Although LSC seems to be promising, the two nonstructural coverage metrics actually suffer from another limitation, i.e., the setting of many parameters. For example, they have to specifically select the layer to calculate LSA/DSA and set the values of U (the upper bound of LSA/DSA) for different subjects. These parameters may influence the performance of LSC/DSC but unfortunately, there is no guide to help set them for different subjects. Here, we conducted a preliminary study to investigate the parameter influence on both LSC and DSC by taking the parameter U as the representative, because this parameter is very important which is used to filter out irrelevant test inputs that have extremely large LSA/DSA values as presented in Section 2.2. In our work, for each subject we set the value of U to be the largest LSA/DSA value of all the test inputs as presented in Section 3.2, indicating that we did not try to filter out irrelevant test inputs. In this preliminary study, for each subject we set U to be a smaller value than the largest LSA/DSA value, i.e., the 10% quantile value of all the ranked LSA/DSA values in their descending order, to try to simulate the filtering process. Please note that, it is very difficult to determine the proper value of U to precisely filter out irrelevant test inputs in practice due to the lack of guide to set the parameter, and thus the setting of *U* is a threat in this study. However, this threat may be not important since the goal of this preliminary study is to investigate whether the results can be largely influenced when the value of *U* is relatively slightly changed.

In this study, we took the subjects whose IDs are $1{\sim}4$ as examples, since the cost spent on collecting non-structural coverage for subjects 5 (VGG-19 based on ImageNet) and 6 (ResNet-50 based on ImageNet) is very huge and DSC is not applicable to regression models. The results are shown in Table 4. From this table, we found that most of the correlation becomes negative after relatively slightly changing U, demonstrating the significant influence of this parameter. Therefore, it could be difficult to apply the two non-structural coverage metrics to the practice due to the significant influence and unclear setting of the parameters.

Finding 4: The two non-structural coverage metrics suffer from the limitations of significant influence and unclear setting of the parameters, leading to the difficulty of applying them to the practice.

4.1.3 Influence of different adversarial test input generation methods.

From Table 3, we found that FGSM has relatively large differences with the other three adversarial test input gen-

TABLE 5: Average time spent on collecting coverage (min)

ID	DNC	TKNC	KMNC	NBC	SNAC	LSC	DSC
1	0.13	0.12	0.25	0.15	0.14	0.12	0.50
2	0.78	0.80	3.42	0.85	0.80	0.60	1.80
3	1.67	2.53	5.98	2.42	1.88	0.17	0.37
4	3.08	2.93	6.30	2.88	3.02	0.80	2.12
5	7.08	7.10	11.90	7.10	7.02	4.03	42.10
6	22.55	22.62	41.08	22.05	22.60	11.67	53.52
7	0.33	0.35	0.69	0.33	0.35	0.37	_
8	1.02	1.32	3.31	1.02	1.01	10.11	_
9	1.35	1.35	2.48	1.32	1.32	0.22	4.02
Avg	4.22	4.35	8.38	4.24	4.24	3.12	14.92

eration methods. In particular, for NBC and SNAC the correlation based on FGSM tends to be obviously stronger than that based on the other three methods. We further analyzed the reason behind this phenomenon, and found that all of C&W, BIM, and JSMA aim to minimize the difference between an original test input and an adversarial test input while FGSM does not consider it. Therefore, FGSM is more likely to generate the adversarial test inputs that have farther distance with the natural test inputs, even irrelevant test inputs (irrelevant test inputs are meaningless in DNN testing). Those more different, even irrelevant test inputs, are much easier to cover corner-case regions of the output values for neurons, leading to stronger positive correlation for NBC and SNAC.

Finding 5: The adversarial test inputs generated by FGSM have different characteristics with those generated by the other three methods, indicating that in the studies of DNN testing, it is not enough to only use FGSM to generate adversarial test inputs.

4.1.4 Influence of different sizes of test sets.

In our study, we studied three different sizes of test sets, i.e., 100, 500, 1,000. By analyzing the results of different test-set sizes, we found that different sizes have no obvious influence. Due to the space limit, we used DeepSpeech based on Speech-Commands as the representative to report the results of different sizes of test sets, whose results are shown in Figure 2. The conclusions from other subjects are the same as them, and their results of different test-set sizes can be found at our project webpage. In Figure 2, the x-axis represents the coverage and the y-axis represents the ratio of adversarial test inputs. We found that the trends of the same coverage metric for different test-set sizes are very similar. The only difference is that the trends become more obvious with the test-set size increasing. This is as expected since larger size tends to bring more statistical significance.

4.1.5 Efficiency

We compared the time spent on collecting each test coverage metric. Table 5 shows the average time for collecting coverage of a test set with the size of 1,000 for each subject, where the last row presents the average time across all subjects.

TABLE 6: Spearman correlation between coverage and the ratio of error-revealing test inputs (the test-set size is 1,000)

FGSM	ID	Adv.	DNC	TKNC	KMNC	NBC	SNAC	LSC	DSC
BIM		FGSM	0.13	0.54	0.13	0.29	0.56	0.91	0.74
JSMA -0.14 0.53 0.58 0.11 0.26 0.92 0.53 Natural -0.07 0.01 0.07 0.02 0.02 0.12 0.13 FGSM 0.03 -0.42 0.02 0.27 0.17 0.12 0.17 C&W -0.18 -0.16 0.11 0.10 0.07 0.09 0.09 JSMA -0.24 -0.53 -0.04 -0.21 -0.17 0.01 0.10 Natural -0.06 -0.06 -0.08 -0.08 -0.03 -0.02 0.10 FGSM -0.34 0.07 -0.34 0.65 0.73 0.71 -0.21 C&W -0.25 -0.20 -0.49 0.01 0.14 0.68 -0.34 JSMA -0.36 -0.43 -0.62 -0.08 0.06 0.54 -0.40 Natural 0.01 -0.05 0.03 -0.03 -0.05 0.15 FGSM 0.21 -0.12 -0.07 0.30 0.34 -0.19 -0.69 C&W 0.07 -0.11 -0.26 -0.01 -0.01 -0.16 -0.75 JSMA 0.02 -0.29 -0.41 -0.12 -0.07 -0.24 -0.74 Natural -0.01 -0.12 -0.17 0.03 0.04 -0.03 -0.05 FGSM 0.26 -0.03 0.05 -0.10 0.04 -0.01 -0.05 C&W 0.20 -0.03 0.06 -0.07 0.05 0.017 0.02 FGSM 0.26 -0.03 0.05 -0.10 0.04 -0.01 -0.05 C&W 0.20 -0.03 0.06 -0.08 0.04 0.08 -0.01 JSMA 0.19 -0.02 -0.04 -0.13 -0.01 0.06 -0.10 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.04 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.05 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.04 Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 FGSM -0.01 0.04 -0.01 0.04 0.05 -0.01 0.03 -0.05 FGSM -0.01 0.02 -0.04 -0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.06 -0.18 0.50 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.06 -0.07 0.05 0.02 -0.06 -0.07 FGSM -0.01 0.02 -0.04 -0.01 0.05 -0.01 -0.03 -0.05 FGSM -0.01 0.04 -0.01 0.06 -0.07 -0.05 -0.01 -0.05 -0.01 -0.05 -0.01		C&W	-0.10	0.18	0.27	0.03	0.38	0.91	0.46
Natural -0.07 0.01 0.07 0.02 0.02 0.12 0.13	1	BIM	-0.03	0.18	0.30	0.07	0.42	0.91	0.46
FGSM		JSMA	-0.14		0.58	0.11	0.26	0.92	0.53
C&W		Natural	-0.07	0.01	0.07	0.02	0.02	0.12	0.13
BIM									
JSMA									
Natural -0.06 -0.06 -0.08 -0.08 -0.03 -0.02 0.10 FGSM -0.34 0.07 -0.34 0.65 0.73 0.71 -0.21 C&W -0.25 -0.20 -0.49 0.01 0.14 0.68 -0.34 BIM -0.25 -0.20 -0.49 0.01 0.14 0.67 -0.35 JSMA -0.36 -0.43 -0.62 -0.08 0.06 0.54 -0.40 Natural 0.01 -0.05 0.03 -0.03 -0.03 -0.05 0.05 FGSM 0.21 -0.12 -0.07 0.30 0.34 -0.19 -0.69 C&W 0.07 -0.11 -0.26 -0.01 -0.01 -0.16 -0.75 BIM -0.01 -0.23 -0.26 -0.09 -0.09 -0.18 -0.74 JSMA 0.02 -0.29 -0.41 -0.12 -0.07 -0.24 -0.74 Natural -0.01 -0.12 -0.17 0.03 0.04 -0.03 -0.05 FGSM 0.26 -0.03 0.05 -0.10 0.04 -0.01 -0.05 C&W 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 SIMA 0.19 -0.02 -0.04 -0.13 -0.01 -0.06 -0.10 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 FGSM -0.01 0.04 -0.01 0.04 0.05 -0.01 -0.03 -0.05 FGSM -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 FGSM -0.01 0.14 -0.12 0.003 -0.05 0.02 -0.05 FGSM -0.01 0.01 -0.00 0.08 0.01 0.09 -0.21 0.03 FGSM -0.01 0.01 -0.00 0.08 0.01 0.09 -0.21 0.03 FGSM -0.01 0.01 -0.00 -0.02 0.05 0.02 -0.05 FGSM -0.01 0.01 -0.00 -0.02 0.05 0.02 -0.05 FGSM -0.01	2								
FGSM									
C&W		Natural	-0.06	-0.06	-0.08	-0.08	-0.03	-0.02	0.10
BIM									
JSMA									
Natural 0.01 -0.05 0.03 -0.03 -0.03 -0.05 0.15 FGSM 0.21 -0.12 -0.07 0.30 0.34 -0.19 -0.69 C&W 0.07 -0.11 -0.26 -0.01 -0.01 -0.16 -0.75 BIM -0.01 -0.23 -0.26 -0.09 -0.09 -0.18 -0.74 JSMA 0.02 -0.29 -0.41 -0.12 -0.07 -0.24 -0.74 Natural -0.01 -0.12 -0.17 0.03 0.04 -0.03 -0.05 FGSM 0.26 -0.03 0.05 -0.10 0.04 -0.01 -0.05 C&W 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 SBIM 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 JSMA 0.19 -0.02 -0.04 -0.13 -0.01 0.06 -0.10 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 SBIM -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA -0.01 0.14 -0.12 0.003 -0.05 0.02 -0.04 Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Fatch 0.20 0.14 -0.12 0.003 -0.05 0.02 -0.04 Natural 0.07 0.03 0.10 0.03 0.04 0.07 -0.08 Patch 0.28 -0.87 -0.43 -0.67 -0.65 0.95 -0.95 Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 -0.01 -0.06 O CTCM -0.77 -0.94 -0.49 -0.03 0.19 0.81 0.78	3								
FGSM									
C&W 0.07 -0.11 -0.26 -0.01 -0.01 -0.16 -0.75 4 BIM -0.01 -0.23 -0.26 -0.09 -0.09 -0.18 -0.74 JSMA 0.02 -0.29 -0.41 -0.12 -0.07 -0.24 -0.74 Natural -0.01 -0.12 -0.17 0.03 0.04 -0.03 -0.05 FGSM 0.26 -0.03 0.06 -0.07 0.05 0.17 0.02 5 BIM 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 5 BIM 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 6 BIM 0.19 -0.02 -0.04 -0.13 -0.01 0.06 -0.19 6 BIM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 6 BIM -0.01 0.02 -0.04 0.05 -0.01 -0.03 -0.05 6		Natural	0.01	-0.05	0.03	-0.03	-0.03	-0.05	0.15
4 BIM JSMA -0.01 0.02 -0.23 -0.29 -0.41 -0.41 -0.09 -0.12 -0.09 -0.41 -0.09 -0.07 -0.18 -0.07 -0.74 -0.24 -0.74 -0.74 Natural -0.01 -0.12 -0.17 0.03 0.04 -0.03 -0.05 FGSM 0.26 -0.03 0.05 -0.10 0.04 -0.01 -0.05 SBIM 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 JSMA 0.19 -0.02 -0.04 -0.03 0.06 -0.08 0.04 0.08 -0.01 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 C&W -0.25 -0.90 -0.66 -0.18 0.50 0.02 -0.05 BIM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 BIM -0.01 0.14									
JSMA									
Natural -0.01 -0.12 -0.17 0.03 0.04 -0.03 -0.05 FGSM 0.26 -0.03 0.05 -0.10 0.04 -0.01 -0.05 C&W 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 BIM 0.20 -0.03 0.06 -0.08 0.04 0.08 -0.01 JSMA 0.19 -0.02 -0.04 -0.13 -0.01 0.06 -0.10 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 JSMA -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA -0.01 0.14 -0.12 0.003 -0.05 0.02 -0.04 Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Patch 0.20 0.14 -0.01 0.44 0.53 0.97 -0.05 Natural 0.07 0.03 0.10 0.03 0.04 0.07 -0.05 Patch 0.28 -0.87 -0.43 -0.67 -0.65 0.95 -0.05 Natural 0.04 -0.95 -0.96 0.18 0.16 0.69 -0.05 Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 -0.06 -0.02 O CTCM -0.77 -0.94 -0.49 -0.03 0.19 0.81 0.78	4								
FGSM 0.26 -0.03 0.05 -0.10 0.04 -0.01 -0.05 C&W 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 BIM 0.20 -0.03 0.06 -0.08 0.04 0.08 -0.01 JSMA 0.19 -0.02 -0.04 -0.13 -0.01 0.06 -0.10 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 0.02 -0.06 -0.07 0.02 -0.06 -0.07 0.05 -0.01 -0.03 -0.05 -0.01 -0.05 -0.01 -0.03 -0.05 -0.01 -0.05 -0.02 -0.01 -0.01 -0.05 -0.02 -0.01 -0.01 -0.05 -0.02 -0.01 -0.01 -0.05 -0.02 -0.01 -0.05									
C&W 0.20 -0.03 0.06 -0.07 0.05 0.17 0.02 BIM 0.20 -0.03 0.06 -0.08 0.04 0.08 -0.01 JSMA 0.19 -0.02 -0.04 -0.13 -0.01 0.06 -0.10 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 BIM -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA -0.01 0.14 -0.12 0.003 -0.05 0.02 -0.04 Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Light -0.86 <td></td> <td>Natural</td> <td>-0.01</td> <td>-0.12</td> <td>-0.17</td> <td>0.03</td> <td>0.04</td> <td>-0.03</td> <td>-0.05</td>		Natural	-0.01	-0.12	-0.17	0.03	0.04	-0.03	-0.05
5 BIM JSMA 0.20 0.19 -0.03 -0.02 0.06 -0.04 -0.08 -0.13 -0.01 -0.01 0.06 -0.10 -0.01 -0.07 0.02 -0.02 -0.01 -0.07 0.02 -0.07 -0.01 -0.07 0.02 FGSM C&W -0.01 -0.25 -0.90 -0.90 -0.66 -0.18 -0.02 -0.05 0.02 -0.05 -0.02 -0.05 -0.02 -0.05 -0.02 -0.05 -0.02 -0.05 -0.02 -0.05 -0.02 -0.05 -0.02 -0.05 -0.02 -0.05 -0.02 -0.05 -0.05 -0.01 -0.05 -0.01 -0.05 -0.01 -0.05 -0.01 -0.01 -0.03 -0.05 -0.02 -0.04 -0.01 -0.01 -0.03 -0.05 -0.02 0.02 -0.04 -0.04 -0.01 -0.04 -0.01 -0.04 -0.03 -0.05 -0.05 0.95 -0.05 -0.21 -0.05 -0.02 -0.01 -0.01 -0.03 -0.05 -0.05 0.95 -0.05 -0.05 -0.05 -0.05 -									
JSMA 0.19 -0.02 -0.04 -0.13 -0.01 0.06 -0.10 Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 0.02 EGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.05 0.02 -0.04 0.05 0.01 0.09 -0.21 0.03 0.05 0.02 -0.04 0.05 0.05 0.02 -0.04 0.05 0.05 0.02 -0.04 0.05 0.05 0.02 -0.04 0.05 0.05 0.05 0.02 -0.04 0.05 0.0									
Natural 0.18 -0.06 0.19 0.12 0.16 -0.07 0.02 FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 BIM -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA -0.01 0.14 -0.12 0.003 -0.05 0.02 -0.04 Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Patch 0.20 0.14 -0.01 0.44 0.53 0.97 - Light -0.86 -0.87 -0.43 -0.67 -0.65 0.95 - Natural 0.07 0.03 0.10 0.03 0.04 0.07 - Patch 0.28 -0.87 -0.89 0.71 0.75 0.89 - Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 -	5								
FGSM -0.01 -0.00 -0.02 0.05 0.02 -0.06 -0.07 C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 BIM -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA -0.01 0.14 -0.12 0.003 -0.05 0.02 -0.04 Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Patch 0.20 0.14 -0.01 0.44 0.53 0.97 - Light -0.86 -0.87 -0.43 -0.67 -0.65 0.95 - Natural 0.07 0.03 0.10 0.03 0.04 0.07 - Patch 0.28 -0.87 -0.89 0.71 0.75 0.89 - Light -0.06 -0.95 -0.96 0.18 0.16 0.69 - Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 - O CTCM -0.77 -0.94 -0.49 -0.03 0.19 0.81 0.78									
C&W -0.25 -0.90 0.66 -0.18 0.50 0.02 -0.05 BIM -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA -0.01 0.14 -0.12 0.003 -0.05 0.02 -0.04 Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Patch 0.20 0.14 -0.01 0.44 0.53 0.97 — Natural 0.07 0.03 0.10 0.03 0.04 0.07 — Patch 0.28 -0.87 -0.89 0.71 0.75 0.89 — B Light -0.06 -0.95 -0.96 0.18 0.16 0.69 — Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 — O CTCM -0.77 -0.94 -0.49 -0.03 0.19 0.81 0.78		Natural	0.18	-0.06	0.19	0.12	0.16	-0.07	0.02
6 BIM JSMA -0.01 0.02 -0.004 0.05 -0.01 -0.03 -0.05 JSMA Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Patch Light -0.86 0.07 0.03 0.10 0.04 0.05 0.95 - Natural 0.07 0.03 0.10 0.44 0.53 0.97 - Natural 0.07 0.03 0.10 0.03 0.04 0.07 - Patch Natural 0.07 0.03 0.10 0.03 0.04 0.07 - Natural 0.04 -0.05 -0.96 0.18 0.16 0.69 - Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 - O CTCM -0.77 -0.94 -0.49 -0.03 0.19 0.81 0.78									
JSMA -0.01 0.14 -0.12 0.003 -0.05 0.02 -0.04 Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Patch 0.20 0.14 -0.01 0.44 0.53 0.97 - Light -0.86 -0.87 -0.43 -0.67 -0.65 0.95 - Natural 0.07 0.03 0.10 0.03 0.04 0.07 - Patch 0.28 -0.87 -0.89 0.71 0.75 0.89 - Light -0.06 -0.95 -0.96 0.18 0.16 0.69 - Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 - O CTCM -0.77 -0.94 -0.49 -0.03 0.19 0.81 0.78									
Natural 0.03 -0.09 0.08 0.01 0.09 -0.21 0.03 Patch 0.20 0.14 -0.01 0.44 0.53 0.97 — Light -0.86 -0.87 -0.43 -0.67 -0.65 0.95 — Natural 0.07 0.03 0.10 0.03 0.04 0.07 — Patch 0.28 -0.87 -0.89 0.71 0.75 0.89 — Natural 0.04 -0.95 -0.96 0.18 0.16 0.69 — Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 —	6								
Patch Light Natural 0.20 0.14 -0.01 0.44 -0.53 0.97 -0.65									
7 Light Natural -0.86 -0.87 -0.87 -0.43 -0.67 -0.65 -0.95 -0.95 -0.03 -0.04 -0.07 -0.03 -0.05 -0.05 -0.95 -0.04 -0.07 -0.07 -0.07 -0.00 Patch Light -0.06 -0.95 -0.96 Natural -0.95 -0.96 -0.96 -0.02 -0.01 -0.01 -0.01 -0.01 -0.06 -0.02 -0.02 -0.01 -0.01 -0.01 -0.01 -0.01 -0.04 -0.03 -0.03 -0.19 -0.81 -0.78		1	ı						0.03
Natural 0.07 0.03 0.10 0.03 0.04 0.07 —	_								—
8 Light -0.06 -0.95 -0.96 -0.96 -0.02 -0.01 -0.06 -0.02 -0.01 -0.06 -0.02 -0.01 -0.06 -0.02 -0.01 -0.01 -0.01 -0.02 -0.01 -0.01 -0.01 -0.01 -0.02 -0.01 -0.0	7	Light							_
8 Light		1	ı						
Natural 0.04 -0.01 -0.06 -0.02 0.01 0.01 —	8								-
O CTCM -0.77 -0.94 -0.49 -0.03 0.19 0.81 0.78									_
		Natural	0.04	-0.01	-0.06	-0.02	0.01	0.01	
Natural 0.08 -0.04 0.20 0.14 0.14 0.07 0.11	Q								
	2	Natural	0.08	-0.04	0.20	0.14	0.14	0.07	0.11

From this table, collecting KMNC is the most costly among the five studied structural coverage, and on average DSC has the largest collection time among all the coverage, i.e., 14.92 minutes. Besides DSC and KMNC, collecting the other coverage is relatively efficient. This is because DSC has to calculate the distance between each test input and almost each training instance. With the vector dimension or the size of training data increasing, the costs could also increase.

4.2 RQ2: Is Test Coverage Correlated with the Ratio of Error-Revealing Test Inputs?

Table 6 shows the Spearman correlation results between test coverage and the ratio of error-revealing test inputs. Similarly, we also used the results of the test-set size of 1,000 as the representative shown in this table, and different test-set sizes also do not influence the conclusions here. Different from Table 3, Table 6 adds the results that the test sets were created using only natural test inputs, denoted as *Natural*. From the rows where the test sets were created by mixing both natural and adversarial test inputs, the conclusions keep consistent with those obtained from Table 3. This is as expected, since the correlation between the ratio of adversarial test inputs and the ratio of error-revealing test inputs is strong. From the rows of *Natural*, test coverage has no significant correlation in most cases (i.e., 73.77%) with

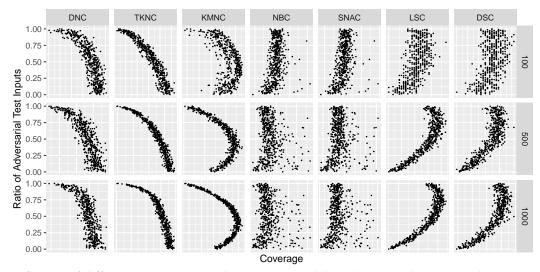


Fig. 2: Influence of different test-set sizes (taking DeepSpeech based on Speech-Commands as an example)

the ratio of error-revealing test inputs, and also has negative correlation in 4.92% cases and weak positive correlation in 21.31% cases. In particular, there is no one case having moderate to extremely strong positive correlation. That is, the correlation obtained from the test sets created based on only natural test inputs is weaker than that obtained from the test sets created based on both natural and adversarial test inputs for each studied test coverage. The reason is that, for many subjects the number of error-revealing natural test input is small, and thus the ratio of error-revealing test inputs for each test set is very close, leading to no statistical significance or weak correlation.

Finding 6: The conclusions keep consistent with those obtained from RQ1 when the test sets are created by mixing both natural and adversarial test inputs, and the correlation between test coverage and the ratio of error-revealing test inputs obtained based on only natural test inputs is weaker than that obtained by mixing both natural and adversarial test inputs due to the small number of error-revealing natural test inputs.

5 FINDINGS SUMMARY AND IMPLICATIONS

We summarize the findings for both structural coverage and non-structural coverage as follows:

- Assumption validation: The assumption fails for structural coverage in general. Among the five studied structural coverage, SNAC performs the best relatively, indicating the importance of the upper corner-case regions of neurons; While the assumption holds for non-structural coverage (especially LSC) on more than half of subjects, indicating that measuring the difference of DNN behaviors between test inputs and training data is more promising than measuring which structural elements are covered by test inputs for measuring test effectiveness.
- Parameters: There are few parameters in the studied structural coverage and the default values of these parameters have been provided; While the studied nonstructural coverage has many parameters, some of which

- are hard to set and have significant influence, leading to the difficulty of applying them to the practice.
- Performance on different subjects: Structural coverage performs stably in general and performs relatively better on simple datasets, i.e., MNIST and Driving; While non-structural coverage performs very differently for different subjects. More specifically, they perform rather well for simple datasets, but their performance drops largely for complex datasets (e.g., CIFAR-100 and ImageNet) since both LSC and DSC rely on the prediction results and these complex datasets tend to have low accuracy.
- Collection cost: The collection cost for structural coverage tends to be small, where collecting KMNC is the most costly among the five structural coverage; While the collection cost for LSC is also small but that for DSC is larger than those for both structural coverage and LSC.
- **Usage restriction**: Structural coverage can be applicable to both classification and regression models, while DSC cannot be applicable to regression models.
- Performance using different adversarial test input generation methods: The adversarial inputs generated by FGSM have different characteristics with those generated by the other three methods (i.e., C&W, BIM, and JSMA), indicating in the studies of DNN testing, it is not enough to only use FGSM to generate adversarial test inputs.

Based on our empirical study, we propose the following implications and directions:

- Among all these studied structural coverage and nonstructural coverage, LSC performs the best relatively by considering both effectiveness and efficiency.
- Although the current non-structural coverage perform better than structural coverage, they still need to be improved from many aspects, such as unfriendly parameters and unstable performance. According to our study, unstable performance may be improved by improving the calculation method of LSC/DSC and proposing more effective strategies to measure the differences of DNN behaviors between test inputs and training data instead of LSA/DSA. In particular, the strategy should not depend on the predicted class for each test input, in order to make it perform stably for both simple and complex datasets.

- From the findings obtained from structural coverage, the upper corner-case regions of neurons is important, and thus it may be promising to integrate this strength of structural coverage to non-structural coverage. That is, combining the strengths of both structural coverage and non-structural coverage may achieve better performance.
- As there are different characteristics between simple datasets and complex datasets and different characteristics between FGSM and the other three adversarial input generation methods (i.e., C&W, BIM, and JSMA), we suggest to consider such diversity when evaluating a new DNN testing metric or a new DNN testing technique.
- Our study controlling for the test-set size obtained different conclusions with those existing studies that do not control for the test-set size, and thus we highlight the necessity of controlling for the test-set size when evaluating the effectiveness of a new metric for DNN testing.

To sum up, although many test coverage have been proposed over the years, there is still a lot of room for improvement of measuring test effectiveness in DNN testing, and our study has pointed out some promising directions.

6 Discussion

6.1 Generality of Our Conclusions to State-of-the-art Structural Coverage

As presented in Section 2.1, MCDC-inspired neuron coverage has been proposed recently, which is the state-of-theart structural coverage. Due to its huge costs described in Section 2.1, we cannot study it in our systematical study. Definitely, it is interesting to investigate whether our conclusions still hold for the state-of-the-art structural coverage, and thus we conducted a small experiment by taking VGG-16 based on CIFAR-10, C&W, and the MCDC-inspired coverage metric – Sign-Sign Coverage, as the representative. We repeated the process in RQ1 by constructing 200 test sets with the size of 1,000. By calculating the Spearman correlation between the coverage and the ratio of adversarial test inputs, we found that the correlation coefficient is -0.54 and the p value is 5.92e-18 (< 0.05). The results demonstrate, there is also no moderate to extremely strong positive correlation between the state-of-the-art Sign-Sign Coverage and the ratio of adversarial inputs, which is consistent with the conclusions from the five structural coverage in our study, indicating the generality of our conclusions to some degree.

6.2 Threats to Validity

6.2.1 Internal threat to validity.

We adopted the released tools to collect DNC, LSC, and DSC, and re-implemented tools to collect TKNC, KMNC, NBC, and SNAC according to the descriptions in their paper [9] since their tools are unavailable. Therefore, the internal threat mainly lies in our re-implementation and our experimental scripts. To reduce it, we have carefully checked our code. In particular, we checked the correctness of our re-implementation by applying it to collect the coverage of the subjects used in the existing work [9], and found that the results are indeed consistent within the margin of error.

6.2.2 External threat to validity.

Our study has five external threats to validity.

Subjects: Although we used the most large-scale and comprehensive subjects to investigate DNN test coverage, the currently used subjects may not represent other subjects. To further reduce this threat, we will use more diverse subjects in the future.

Created test sets: In our study, we created 500 test sets for each subject under each given test-set size, which may not represent real-world test sets. However, it is difficult to collect many real-world test sets since the labeling process is very costly. Moreover, the created test-set sizes involve three sizes, i.e., 100, 500, and 1,000, which may not represent other sizes. In our study, the conclusions obtained from the three sizes are consistent, and thus this threat may not be serious.

Studied test coverage: We considered both structural and non-structural coverage. More specifically, we studied seven typical test coverage, which have been explained in Section 2, and also discussed the state-of-the-art structural coverage (i.e., MCDC-inspired coverage) in Section 6.1. In the future, we will study more coverage metrics, such as t-way combination coverage [34] and state coverage [35].

Adversarial test input generation methods: We adopted the most widely-used FGSM, BIM, JSMA, C&W for general image-domain DNN models, the widely-used Patch and Light for Driving following the existing work [19], [25], and CTCM for Speech-Commands, which may not represent other adversarial test input generation methods. However, this threat may not be serious since all the widely-used methods have been used in our study and the conclusions obtained from BIM, JSMA, and C&W are almost consistent. In the future, we will use more advanced methods, such as DeepRoad [36] for Driving.

Used errors: In our study, we regard a test input predicted wrongly as an error-revealing test input, but we do not distinguish whether these error-revealing test inputs trigger the same model bugs. This threat exists in all existing studies to evaluate the corresponding proposed test coverage since it is challenging to distinguish them due to involving many manual efforts. Actually, the number of error-revealing test inputs is also very important in DNN testing since identifying them can be helpful to improve the accuracy of the DNN model [37].

6.2.3 Construct threat to validity.

Our study has three construct threats to validity.

Parameters of studied test coverage: The specific parameter settings have been presented in Section 3.2. Also, our study has investigated the influence of some parameter (i.e., the upper bound U for LSC and DSC). In the future, we will investigate the influence of other parameters.

Method of test set creation: As presented in Section 3.3, we randomly selected test inputs with replacement. Actually, both the sampling method with replacement and that without replacement have been used in the existing studies [11], [12], and the influence is slight. To reduce this threat, we repeated our study by using the method without replacement by taking CIFAR-10 and ResNet-20 using JSMA as an example. Here, we created 50 test sets with the size of 400 without replacement. The correlation coefficients between

the five structural coverage (DNC, TKNC, KMNC, NBC, and SNAC) and the ratio of adversarial test inputs are -0.42, -0.72, -0.74, -0.19, and 0.01 respectively. The conclusion from the results is consistent with that of using the method with replacement. Therefore, this threat may not be serious.

Measurements: We calculated Spearman correlation to measure the correlation between test coverage and the ratio of adversarial test inputs/the ratio of the error-revealing test inputs. In the future, we can use more measurements (e.g., Kendall τ [38]) to reduce this threat.

7 RELATED WORK

7.1 DNN Testing

There are a lot of research on DNN testing [23], [39], [40], [41], [42], [43], [44], [45]. We classify them into: measuring test effectiveness, generating test inputs, and improving test efficiency.

Regarding measuring test effectiveness, we have introduced five structural coverage and two non-structural coverage in Section 2. Besides, Ma et al. [34] proposed tway combination coverage for DNN testing by borrowing the idea of traditional combinatorial testing. Du et al. [35] proposed state coverage for RNN-based stateful DNN testing. Ma et al. [46] proposed a mutation-based metric by designing a set of mutation operators for DNN source programs and models. In particular, Li et al. [33] conducted a preliminary study on structural coverage using only two datasets, and found that adversary-oriented search may make more contributions to the error-revealing capability than structural coverage. Different from this preliminary study, our study is the first one to systematically investigate the correlation between test coverage and the generation of adversarial test inputs/the error-revealing capability of test inputs. More specifically, we designed our study more systematically, i.e., considering both structural and nonstructural coverage, both the ratio of adversarial test inputs and the ratio of error-revealing test inputs, more comprehensive benchmark, and more formal measurements. Furthermore, there are some metrics to measure DNN robustness [47], [48], [49], [50]. Since they are different from the metrics measuring test effectiveness, which are the target of our work, we do not discuss them in detail.

Regarding generating test inputs, many approaches have been proposed in recent years. For example, Guo et al. [51] proposed *DLFuzz*, the first differential testing framework for DNN by maximizing neuron coverage. Xie et al. [20] proposed *Deephunter*, a coverage-guided fuzz testing framework for DNN. Sun et al. [52] proposed to generate test inputs for DNN through concolic testing.

Regarding improving test efficiency, some test input selection and prioritization approaches have been proposed for DNN [25], [26], [53], [54]. For example, Li et al. [25] proposed to select test inputs by minimizing the cross entropy between the selected test inputs and the whole test set, so as to save labeling costs. Feng et al. [54] proposed *DeepGini* to prioritize test inputs by measuring the purity of test inputs, so as to label error-revealing test inputs earlier.

Different from them, our work aims to investigate the correlation between test coverage and the generation of adversarial test inputs/the error-revealing capability of test inputs through *a systematical study*.

7.2 Empirical Studies on Traditional Test Coverage

In traditional software testing, there are some studies to investigate the correlation between traditional test coverage and test effectiveness. For example, Namin and Andrews [55] conducted a study to explore the relationship among the test-suite size, structural coverage, and test effectiveness based on C/C++ programs, and found that coverage is sometimes correlated with test effectiveness when controlling for the test-suite size. Inozemtseva and Holmes [11] conducted a study to investigate the above relationship based on Java programs, and found that the correlation between test coverage and test effectiveness is low to moderate when controlling for the test-suite size. Zhang and Mesbah [12] conducted a study to investigate the correlation between assertions and test effectiveness based on Java programs, and found that both assertion numbers and assertion coverage are strongly correlated with test effectiveness. Zhang et al. [56] explored the influence of pseudo test suites on these studies, and found that the correlation obtained from pseudo test suites is stronger than that obtained from original test suites.

Different from these studies, our study is to investigate the correlation between test coverage and test effectiveness in *DNN testing*. Traditional software testing and DNN testing have different test coverage and the corresponding studied subjects have totally different characteristics.

8 Conclusion

In this work, we systematically conducted the first extensive study to validate the assumption of test coverage metrics in DNN testing, i.e., they are correlated with the generation of adversarial test inputs or the error-revealing capability of test inputs. In the study, we studied five structural coverage and two non-structural coverage based on 9 pairs of datasets and DNN models with great diversity. The results demonstrate that in general the assumption fails for structural coverage while it holds for non-structural coverage (especially LSC) on more than half of subjects, when controlling for the test-set size, indicating that non-structural coverage is more promising than structural coverage. In the future, we need to further improve the current non-structural coverage from several aspects (such as unfriendly parameters and unstable performance) or propose new non-structural coverage to more effectively measure test effectiveness in DNN testing.

REFERENCES

- Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao.
 Deepdriving: Learning affordance for direct perception in autonomous driving. In ICCV, pages 2722–2730, 2015.
- [2] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al. Deep speech 2: End-to-end speech recognition in english and mandarin. In ICML, pages 173–182, 2016.
- [3] Yi Sun, Yuheng Chen, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation by joint identification-verification. In NeurIPS, pages 1988–1996, 2014.
- [4] Yong Cheng. Semi-supervised learning for neural machine translation. In *Joint Training for Neural Machine Translation*, pages 25–40. 2019.
- [5] Ziad Obermeyer and Ezekiel J Emanuel. Predicting the future—big data, machine learning, and clinical medicine. N. Engl. J. Med., 375(13):1216, 2016.

- [6] Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. Deep code search. In ICSE, pages 933–944, 2018.
- [7] Xiaodong Gu, Hongyu Zhang, Dongmei Zhang, and Sunghun Kim. Deep API learning. In FSE, pages 631–642, 2016.
- [8] Kexin Pei, Yinzhi Cao, Junfeng Yang, and Suman Jana. Deepxplore: Automated whitebox testing of deep learning systems. In SOSP, pages 1–18, 2017.
- [9] Lei Ma, Felix Juefei-Xu, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Chunyang Chen, Ting Su, Li Li, Yang Liu, et al. Deepgauge: Multi-granularity testing criteria for deep learning systems. In ASE, pages 120–131, 2018.
- [10] Augustus Odena, Catherine Olsson, David Andersen, and Ian Goodfellow. Tensorfuzz: Debugging neural networks with coverage-guided fuzzing. In *ICML*, pages 4901–4911, 2019.
- [11] Laura Inozemtseva and Reid Holmes. Coverage is not strongly correlated with test suite effectiveness. In ICSE, pages 435–445, 2014
- [12] Yucheng Zhang and Ali Mesbah. Assertions are strongly correlated with test suite effectiveness. In FSE, pages 214–224, 2015.
- [13] Thierry Titcheu Chekam, Mike Papadakis, Yves Le Traon, and Mark Harman. An empirical study on mutation, statement and branch coverage fault revelation that avoids the unreliable clean program assumption. In ICSE, pages 597–608, 2017.
- [14] G. C. Morrison, Cornelia P. Inggs, and W. C. Visser. Automated coverage calculation and test case generation. In SAICSIT, pages 84–93. ACM, 2012.
- [15] Michael Hilton, Jonathan Bell, and Darko Marinov. A large-scale study of test coverage evolution. In ASE, pages 53–63, 2018.
- [16] Milos Gligoric, Alex Groce, Chaoqiang Zhang, Rohan Sharma, Mohammad Amin Alipour, and Darko Marinov. Comparing nonadequate test suites using coverage criteria. In ISSTA, pages 302– 313, 2013.
- [17] Milos Gligoric, Alex Groce, Chaoqiang Zhang, Rohan Sharma, Mohammad Amin Alipour, and Darko Marinov. Guidelines for coverage-based comparisons of non-adequate test suites. TOSEM, 24(4):22:1–22:33, 2015.
- [18] Gregory Gay, Ajitha Rajan, Matt Staats, Michael W. Whalen, and Mats Per Erik Heimdahl. The effect of program and model structure on the effectiveness of MC/DC test adequacy coverage. TOSEM, 25(3):25:1–25:34, 2016.
- [19] Jinhan Kim, Robert Feldt, and Shin Yoo. Guiding deep learning system testing using surprise adequacy. In ICSE, pages 1039– 1049, 2019.
- [20] Xiaofei Xie, Lei Ma, Felix Juefei-Xu, Minhui Xue, Hongxu Chen, Yang Liu, Jianjun Zhao, Bo Li, Jianxiong Yin, and Simon See. Deephunter: a coverage-guided fuzz testing framework for deep neural networks. In ISSTA, pages 146–157, 2019.
- [21] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *S&P*, pages 39–57, 2017.
- [22] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In ICLR, 2015.
- [23] Youcheng Sun, Xiaowei Huang, Daniel Kroening, James Sharp, Matthew Hill, and Rob Ashmore. Structural test coverage criteria for deep neural networks. TECS, 18(5s):94:1–94:23, 2019.
- [24] Matt P Wand and M Chris Jones. *Kernel smoothing*. Chapman and Hall/CRC, 1994.
- [25] Zenan Li, Xiaoxing Ma, Chang Xu, Chun Cao, Jingwei Xu, and Jian Lü. Boosting operational dnn testing efficiency through conditioning. In FSE, pages 499–509, 2019.
- [26] Junjie Chen, Zhuo Wu, Zan Wang, Hanmo You, Lingming Zhang, and Ming Yan. Practical accuracy estimation for efficient deep neural network testing. TOSEM, 2020. to appear.
- [27] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In ICLR, 2017.
- [28] Nicolas Papernot, Patrick D. McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In S&P, pages 372–387, 2016.
- [29] Nicholas Carlini and David A. Wagner. Audio adversarial examples: Targeted attacks on speech-to-text. In S&P Workshops, pages 1–7, 2018.
- [30] Alex Graves, Santiago Fernández, Faustino J. Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *ICML*, pages 369–376, 2006.
 [31] Leann Myers and Maria J Sirois. Spearman correlation coeffi-
- [31] Leann Myers and Maria J Sirois. Spearman correlation coefficients, differences between. Encyclopedia of statistical sciences, 12, 2004

- [32] Junjie Chen, Yanwei Bai, Dan Hao, Lingming Zhang, Lu Zhang, and Bing Xie. How do assertions impact coverage-based test-suite reduction? In *ICST*, pages 418–423, 2017.
- [33] Zenan Li, Xiaoxing Ma, Chang Xu, and Chun Cao. Structural coverage criteria for neural networks could be misleading. In *ICSE (NIER)*, pages 89–92, 2019.
- [34] Lei Ma, Felix Juefei-Xu, Minhui Xue, Bo Li, Li Li, Yang Liu, and Jianjun Zhao. Deepct: Tomographic combinatorial testing for deep learning systems. In *SANER*, pages 614–618, 2019.
- [35] Xiaoning Du, Xiaofei Xie, Yi Li, Lei Ma, Yang Liu, and Jianjun Zhao. Deepstellar: model-based quantitative analysis of stateful deep learning systems. In ESEC/SIGSOFT FSE, pages 477–487, 2019
- [36] Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. Deeproad: Gan-based metamorphic testing and input validation framework for autonomous driving systems. In *ASE*, pages 132–142, 2018.
- [37] Shiqing Ma, Yingqi Liu, Wen-Chuan Lee, Xiangyu Zhang, and Ananth Grama. MODE: automated neural network model debugging via state differential analysis and input selection. In FSE, pages 175–186, 2018.
- [38] Hervé Abdi. The kendall rank correlation coefficient. *Encyclopedia of Measurement and Statistics*. *Sage, Thousand Oaks, CA*, pages 508–510, 2007.
- [39] Youcheng Sun, Xiaowei Huang, and Daniel Kroening. Testing deep neural networks. *arXiv preprint arXiv:1803.04792*, 2018.
- [40] Jie M Zhang, Mark Harman, Lei Ma, and Yang Liu. Machine learning testing: Survey, landscapes and horizons. *TSE*, 2020. to appear.
- [41] Fuyuan Zhang, Sankalan Pal Chowdhury, and Maria Christakis. Deepsearch: Simple and effective blackbox fuzzing of deep neural networks. *CoRR*, abs/1910.06296, 2019.
- [42] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. Deeptest: automated testing of deep-neural-network-driven autonomous cars. In *ICSE*, pages 303–314, 2018.
- [43] Youcheng Sun, Xiaowei Huang, Daniel Kroening, James Sharp, Matthew Hill, and Rob Ashmore. Deepconcolic: testing and debugging deep neural networks. In ICSE, pages 111–114, 2019.
- [44] Simos Gerasimou, Hasan Ferit Eniser, Alper Sen, and Alper Cakan. Importance-driven deep learning system testing. In *ICSE*, 2020. to appear.
- [45] Seokhyun Lee, Sooyoung Cha, Dain Lee, and Hakjoo Oh. Effective white-box testing of deep neural networks with adaptive neuron-selection strategy. In ISSTA, pages 165–176, 2020.
- [46] Lei Ma, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Felix Juefei-Xu, Chao Xie, Li Li, Yang Liu, Jianjun Zhao, et al. Deepmutation: Mutation testing of deep learning systems. In ISSRE, pages 100–111, 2018.
- [47] Osbert Bastani, Yani Ioannou, Leonidas Lampropoulos, Dimitrios Vytiniotis, Aditya Nori, and Antonio Criminisi. Measuring neural net robustness with constraints. In *NeurIPS*, pages 2613–2621, 2016.
- [48] Susmit Jha, Sunny Raj, Steven Lawrence Fernandes, Sumit Kumar Jha, Somesh Jha, Brian Jalaian, Gunjan Verma, and Ananthram Swami. Attribution-based confidence metric for deep neural networks. In *NeurIPS*, pages 11826–11837, 2019.
- [49] Divya Gopinath, Guy Katz, Corina S. Pasareanu, and Clark W. Barrett. Deepsafe: A data-driven approach for assessing robustness of neural networks. In ATVA, volume 11138, pages 3–19, 2018.
- [50] Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. Towards proving the adversarial robustness of deep neural networks. In FVAV@iFM, volume 257 of EPTCS, pages 19–26, 2017.
- [51] Jianmin Guo, Yu Jiang, Yue Zhao, Quan Chen, and Jiaguang Sun. Dlfuzz: Differential fuzzing testing of deep learning systems. In FSE, pages 739–743, 2018.
- [52] Youcheng Sun, Min Wu, Wenjie Ruan, Xiaowei Huang, Marta Kwiatkowska, and Daniel Kroening. Concolic testing for deep neural networks. In ASE, pages 109–119, 2018.
- [53] Wei Ma, Mike Papadakis, Anestis Tsakmalis, Maxime Cordy, and Yves Le Traon. Test selection for deep learning systems. CoRR, abs/1904.13195, 2019.
- [54] Yang Feng, Qingkai Shi, Xinyu Gao, Jun Wan, Chunrong Fang, and Zhenyu Chen. Deepgini: prioritizing massive tests to enhance the robustness of deep neural networks. In ISSTA, pages 177–188, 2020.

- [55] Akbar Siami Namin and James H Andrews. The influence of size and coverage on test suite effectiveness. In ISSTA, pages 57–68, 2009.
- [5] Jie M Zhang, Lingming Zhang, Dan Hao, Meng Wang, and Lu Zhang. Do pseudo test suites lead to inflated correlation in measuring test effectiveness? In ICST, pages 252–263, 2019.