

RNN-Test: Towards Adversarial Testing for Recurrent Neural Network Systems

Jianmin Guo, Yue Zhao, Quan Zhang, and Yu Jiang

Abstract—While massive efforts have been investigated in adversarial testing of convolutional neural networks (CNN), testing for recurrent neural networks (RNN) is still limited and leaves threats for vast sequential application domains. In this paper, we propose an adversarial testing framework RNN-Test for RNN systems, focusing on the main sequential domains, not only classification tasks. First, we design a novel search methodology customized for RNN models by maximizing the inconsistency of RNN states to produce adversarial inputs. Next, we introduce two state-based coverage metrics according to the distinctive structure of RNNs to explore more inference logics. Finally, RNN-Test solves the joint optimization problem to maximize state inconsistency and state coverage, and crafts adversarial inputs for various tasks of different kinds of inputs.

For evaluations, we apply RNN-Test on three sequential models of common RNN structures. On the tested models, the RNN-Test approach is demonstrated to be competitive in generating adversarial inputs, outperforming FGSM-based and DLfuzz-based methods to reduce the model performance more sharply with 2.78% to 32.5% higher success (or generation) rate. RNN-Test could also achieve 52.65% to 66.45% higher adversary rate on MNIST-LSTM model than relevant work testRNN. Compared with the neuron coverage, the proposed state coverage metrics as guidance excel with 4.17% to 97.22% higher success (or generation) rate.

Index Terms—Adversarial testing, Recurrent neural networks, Coverage metrics

1 INTRODUCTION

As the core part of the current artificial intelligence applications, deep learning has made great breakthroughs in computer vision [1], natural language processing (NLP) [2], and automatic speech recognition (ASR) [3]. With the increasing deployments of deep neural network (DNN) systems in the safety- and security-critical domains, such as autonomous driving [4] and medical diagnose [5], ensuring the robustness of DNNs becomes an essential concern in both academic research and security communities.

However, it is demonstrated that state-of-the-art DNN systems [6] are easy to suffer attacks and produce completely wrong predictions, when fed with adversarial inputs which are nearly indistinguishable from original test inputs. This inspired numerous adversarial testing works devoted to generating adversarial inputs for DNNs, aiming to provide rich sources to train the DNNs to be more robust. The majority of these works [7], [8], [9] try to fool popular image classifiers by applying minute perturbations to the inputs. They exhibit high efficiency but achieve low testing completeness [10]. Recently, researchers [10], [11], [12] try to apply traditional software testing techniques over DNNs, with various neuron-based coverage criteria proposed to measure the testing completeness of DNN logics. These works could reach high testing coverage and produce numbers of adversarial inputs.

In spite of the efficiency of these works, they are largely limited to CNNs and image classification tasks. Overall, there are two main types of DNN, the convolutional neural networks (CNN) [13] and recurrent neural networks (RNN) [14]. They are of distinct structures and preferred for different kinds of tasks. CNN exhibits excellent competence in dealing with image processing tasks [15], [16], with thousands of neurons good at extracting image features. RNN is known for the iterative structure over cells and specific components dealing with context semantics, hence expert in handling tasks with sequential data, like natural language processing [17] and speech recognition [18]. Owing to the huge gap, the testing techniques and coverage metrics for the two types of DNNs are hard to fit the other.

So far, adversarial testing for RNN systems has received limited attention, especially those of sequential tasks (outputs without explicit class labels). Existing works [19], [20], [21], [22] still concentrate more on classification domains, thanks to advances in classification tasks of CNNs. They perform well over tasks such as sentiment analysis [19], [20], [22], image classification [21], [22], and lipophilicity prediction [22], etc. But the core sequential tasks of RNNs leave tested insufficiently, threatening their large-scale applications. Moreover, existing coverage criteria [10], [11], [12] are mostly designed for CNNs and neurons, with a large gap to fit for RNNs. Lastly, recent works [21], [22] also proposed coverage guided testing methods with specific coverage criteria for RNN systems. They could generate lots of adversarial inputs for tested models (mainly of classification tasks), by mutating inputs directly (e.g. random noise) and employing coverage values as constraints to terminate testing. However, the specific inference logics of RNNs are not utilized in searching for adversarial inputs.

- J. Guo, Q. Zhang and Y. Jiang are with School of Software, Tsinghua University, Beijing National Research Center for Information Science and Technology, and Key Laboratory for Information System Security, Ministry of Education, Beijing, 100084, China.
- Y. Zhao is with Huawei Technologies Co., Ltd.
- J. Guo is corresponding author. Email: guojm17@mails.tsinghua.edu.cn.

Manuscript submitted XX XX, 2021.

Therefore, challenges for RNN testing are summarized as threefold. First, adversarial testing methods for RNNs with the main sequential domains are rather inadequate, leaving threats for majority application scenarios. For tasks with sequential outputs, there is no standard to recognize the inputs as adversarial inputs without obvious class labels. Second, neuron-based coverage metrics fail to consider characteristics of RNN structures and could not be adopted directly. Third, existing testing methods are limited in making use of distinct logics of RNN models.

Approach: In this paper, we propose an adversarial testing framework RNN-Test for RNN systems, especially those with sequential outputs. RNN-Test concentrates on the kernel structure for sequential processing and rids of remain structures for particular applications. According to the unique features of RNNs, we put forward a specific search methodology, which maximizes the inconsistency of RNN states to obtain adversarial inputs. Meanwhile, we also design two state-based coverage metrics for different RNN models to explore more logics and guide to discover adversarial inputs in irregular space. They are then combined as a joint optimization problem, which is to maximize the state inconsistency and state coverage. Finally, it will be solved to acquire perturbations in a gradient-based way.

When obtained the perturbations, adversarial inputs will be crafted by applying perturbations to original test inputs in different ways for various kinds of inputs. In the end, we adopt model performance metrics to identify the adversarial inputs and assess their qualities, which are generally available for RNN variants. In this way, RNN-Test provides a scalable and extensible solution for RNN testing.

Evaluation: We evaluate the RNN-Test approach over three RNN models of sequential tasks, including two NLP models with discrete inputs (PTB language model [23], spell checker model [24]) and one ASR model with continuous inputs (DeepSpeech ASR model [25]). The RNN-Test approach could efficiently acquire adversarial inputs of high quality, which reduce the model performance sharply while nearly imperceptible to original inputs. Compared with baselines which are two popular techniques (FGSM [7] and DLFuzz [26]) adapted here for RNN testing, our approach achieves more performance reduction with higher success (or generation) rate. Taking DeepSpeech ASR model as an example, RNN-Test could decline the model performance by 17.29% higher WER, 3.61% lower BLEU with 10% higher success rate, in contrast with the FGSM-based method. As for the most relevant work testRNN [22], RNN-Test could achieve 52.65% to 66.45% higher adversary rate on MNIST-LSTM model [27].

Furthermore, coverage guidance as pure optimization is first demonstrated with diverse searching capability for adversarial inputs compared with other methods. The proposed state coverage guidance achieved 4.17% to 97.22% higher success (or generation) rate than neuron coverage guidance, and even best performance on the spell checker model. Finally, adversarial inputs obtained by RNN-Test are of high quality. Besides reducing the model performance sharply with minute perturbations and high time efficiency, they could also improve the model by retraining, such as

12.582% improvement of test perplexity on PTB language model.

Contribution: Our work has the following contributions:

- We design a novel search methodology based on the inner logics of RNNs, which maximizes the inconsistency of RNN states to produce adversarial inputs efficiently.
- We propose two state-based coverage metrics customized for RNNs, mainly as guidance for adversarial testing. We first demonstrate that coverage guidance has a diverse searching capability for adversarial inputs compared with other methods. Our state coverage guidance is also superior to neuron coverage guidance in RNN testing.
- We design and implement the adversarial testing framework RNN-Test. To our best knowledge, it is the first step towards systematically testing the majority sequential domains for RNN systems, referred to as tasks of sequential outputs without class labels. It is effective and scalable for variants of RNNs, outperforming FGSM- and DLFuzz-based methods as well as testRNN (on MNIST-LSTM model).

2 BACKGROUND

2.1 Deep Neural Network

In the following, we will describe the two main kinds of DNN, convolutional neural networks (CNN) [13] and recurrent neural networks (RNN) [14].

Convolutional neural network and neuron. Fig. 1a shows the simplified structure of a typical CNN. CNN keeps the fundamental feed-forward structure, where each neuron is connected with neurons of adjacent layers while no connections with those of the same layer. They are broadly used in image processing tasks [15], [16], with specific convolution layers good at extracting image features. Besides, CNN adopts classical DNN neurons, as shown in Fig. 1b. The neuron output is a single value transformed by a non-linear activation function, which is usually RELU (Rectified Linear Unit).

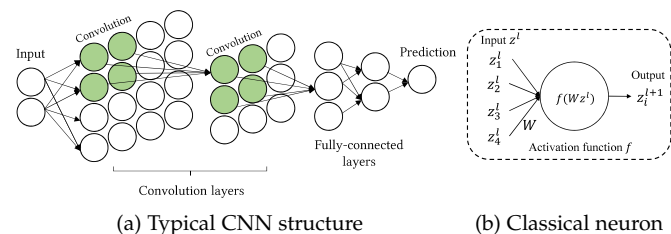


Fig. 1: CNN structure and neuron

Recurrent neural network and cell. Fig. 2 illustrates the basic structure of RNN, where an elementary RNN cell (noted as the square) iteratively makes predictions \hat{y} based on inputs x and intermediate outputs h , which are referred to as hidden states. When it is unfolded, it can be seen that the input sequence x is fed to the RNN cell as a series of time steps, where x could be a sentence and x_t is the t -th word. Moreover, each prediction \hat{y}_t could be the predicted word right after x_t based upon the received sequence.

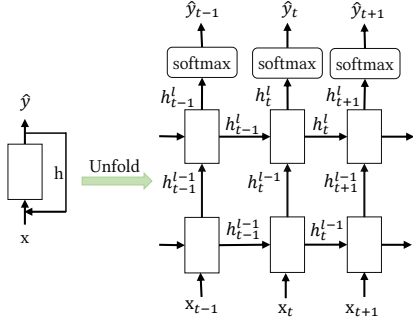


Fig. 2: RNN structure

In contrast to CNN neuron, the hidden state output h_t^l of the cell at time step t of layer l is decided by current input x_t as well as h_{t-1}^l from the previous layer as well as h_{t-1}^l from the prior step in the same layer, and then passed forward to compute the softmax predictions. Due to this key design, RNN excels in making use of the interior contextual semantics of sequential inputs.

Nevertheless, the basic RNN is unable to learn the semantic dependency within longer time steps. LSTM (Long short-term memory) [28], [29] and GRU (Gated recurrent unit) [30] networks are generally deployed solutions, bringing gate mechanisms to RNN cell. Fig. 3 provides the general RNN cell and LSTM cell as the example, in which f, i, n, o stand for various gates¹, and σ for sigmoid function.

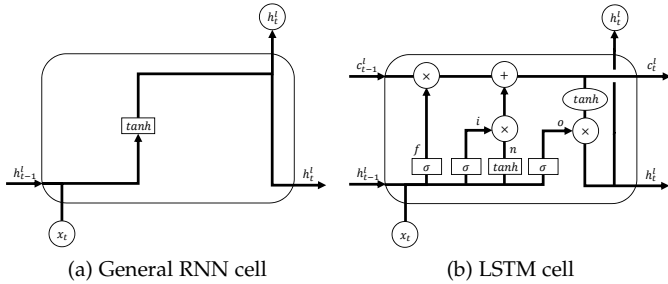


Fig. 3: RNN cell

2.2 Limitations of existing coverage metrics

While numerous coverage metrics [10], [11], [12] are proposed for DNN testing, they are mostly based on CNN neurons and hard to satisfy RNN testing. First, an RNN usually comprises two or three layers each with several cells when unfolded, much fewer than a CNN usually of ten more layers each with hundreds of neurons. Second, sigmoid and tanh are conventionally used for an RNN cell whereas RELU is the most choice for a CNN neuron. This is critical because that the value range of RELU is $[0, \infty)$ while $[0, 1]$ for sigmoid and $[-1, 1]$ for tanh, leading to infinite values of CNN neurons whereas narrowed value ranges of RNN states.

1. Gates of GRU are different but similar in design, not listed here.

Unfortunately, neuron-based coverage metrics fail to consider these critical characteristics. We evaluate neuron coverage [10] on our tested models, which treats the hidden states of each RNN cell as the equivalent output of a CNN neuron. As for the PTB language model comprising two layers each with 10 time steps, there will be only 20 neurons in total. The neuron coverage reaches 100% with at most 4 inputs even taking a higher threshold 0.5. When for adversarial testing on this model, it fails to find any adversarial inputs, shown in Table 4 of § 5. As for coverage criteria [12], [31] defined over multi-section (e.g. 1000) neuron value ranges discriminated by training and testing inputs, the narrow value ranges of RNN states (hidden state value not exceeding ± 1) limit their applications to RNN testing.

Recently, DeepStellar [21] abstracts RNN models as Discrete-Time Markov Chain (DTMC) models and then adapts coverage metrics of [12] for testing. Another work testRNN [22] designs novel coverage metrics for LSTM models to quantify temporal relations in RNNs. As to these coverage criteria, they primarily measure over abnormal values in testing, which are recognized by thresholds set according to training data. In this paper, we define coverage metrics trying to capture key inference logics inspired by distinctive roles of core RNN components, with no aid of training data.

3 STATE COVERAGE METRICS

In this section, we propose two state coverage metrics based on unique features of RNNs. Hidden state coverage (HS_C) is designed to capture RNN prediction logics, which could be universally applied to RNN models. Due to the widespread deployments of LSTM networks, cell state coverage (CS_C) is specially designed for LSTM models.

3.1 Hidden State Coverage

As discussed in § 2.2, RNN states cannot be regarded to be identical to CNN neurons. Fig. 4 provides a simple illustration of the inner logics of hidden states and cell states. In Fig. 4a, h_t^l represents the output of each RNN cell, which is a vector containing hundreds or thousands of hidden states. Here each rectangle represents each hidden state, where a darker color means a higher value. When to predict the next word following “a”, these hidden states will be mapped to a list of candidates. If a hidden state is the maximum, its mapped candidate will probably be the prediction result.

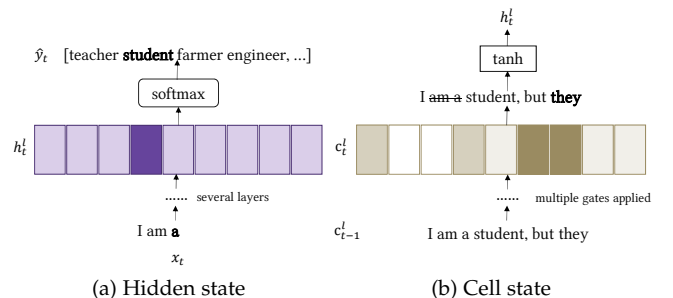


Fig. 4: Illustrations of RNN states.

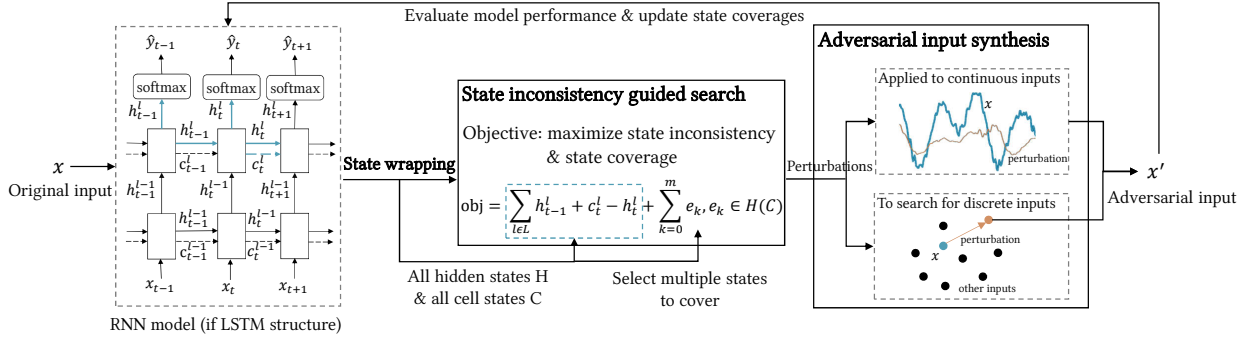


Fig. 5: Architecture of RNN-Test

Therefore, hidden states of each vector h lead to varied prediction results for each step. To explore the inference logics thoroughly, it is meaningful for each hidden state to be the maximum in the vector and perform the key inference. We define hidden state coverage as the ratio of such hidden states of all the hidden states during testing. Formally, the definition is given in formula (1).

$$HS_C = \frac{|\{e \mid \forall e \in h, \forall h \in H, e = \max(h)\}|}{|H|} \quad (1)$$

Assume H denotes all the hidden states of an RNN model of given inputs, which is a four-dimensional matrix of shape (T, L, B, E) , where T, L, B are the number of the time steps², layers and batch, respectively. E is the state size. Though H varies among RNN models, T, L, B, E are always the necessary components, where batch is to accelerate computations by feeding multiple inputs simultaneously. Thus, $|H|$ will be $T \times L \times B \times E$. Here, a specific hidden state vector $h \in H$ contains E hidden states and is denoted as its index of H , which is $[t, l, b]$, where $t \in \{1, 2, \dots, T\}, l \in \{1, 2, \dots, L\}, b \in \{1, 2, \dots, B\}$. That means, h is the output of the RNN cell of the t -th step in the l -th layer for the b -th input. For a state $e \in h$, e is covered if its value $e = \max(h)$.

3.2 Cell State Coverage

As in Fig. 3b, the cell states and gates are activated by functions sigmoid and tanh. The sigmoid function of three gates outputs values between 0 and 1, determining how much of each cell state to keep. The tanh function pushes the states between -1 and 1, for gate n to add information to cell states, and for cell states c to compute the hidden states. Thus, each cell state value protects the contexts. Fig. 4b illustrates the cell states, where c_t^l is the output of each LSTM cell which is also a vector of cell states. Similarly, a rectangle is also a cell state and a darker color for a higher value. When to predict the predicate after ‘‘they’’, c_t^l will receive contexts from c_{t-1}^l to keep key semantics and remove those invalid, and then to compute predictions.

In this paper, we design cell state coverage over different value ranges standing for degrees to keep contexts. In the experiments, cell state values mostly fall into the central range while few be the boundary value. We suppose that

2. For models could be fed with inputs of non-equal steps, T will be adjusted according to the length of each input.

covering more of each section (5 sections in this paper), especially boundary sections, could explore more context space. The formal definition is given in formula (2).

$$CS_{C_{sec_i}} = \frac{|\{e \mid \forall e \in c, \forall c \in C, \tanh(e) \in sec_i\}|}{|C|} \quad (2)$$

Here, all the cell states of an RNN model fed with given inputs are denoted as C , which is also a matrix of shape (T, L, B, E) . The value range of function tanh is split to Sec sections and each section is $sec_i = [v_{i-1}, v_i]$, where $-1 \leq v_i \leq 1$. For a specific cell state vector $c \in C$, it is also denoted by the index of C as $[t, l, b]$. If a cell state $e \in c$ and its activation value $\tanh(e) \in sec_i, i \in \{1, 2, \dots, Sec\}$, then e is covered in sec_i .

Thereby, we could measure how extensively the test inputs exercise RNN logics and benefit adversarial testing with state coverage metrics as guidance, without additional resources extracted from training data.

4 RNN-TEST DESIGN

In this section, we present a technical description of RNN-Test in detail. Fig. 5 depicts the overall workflow. Given the tested model, RNN-Test will focus on the kernel sequential part of the RNN structure without other components for particular tasks. For original test input x , we first extract the hidden states and cell states of each RNN cell by state wrapping, without affecting its inherent process. These states are crucial for the subsequent state inconsistency guided search, maximizing state inconsistency and state coverage to generate adversarial inputs. Unlike the usual idea of increasing the model cost [7], [32] or the probabilities of targeted classes [9], [26], RNN-Test boosts the inconsistency of RNN states opposite to their inner dependencies (The part of obj rounded with blue dashed frame violates data dependencies marked with blue lines in the model). In this way, RNN-Test could search for adversarial inputs in a lightweight and scalable means. Meanwhile, RNN-Test also tries to cover more states and explore more inference logics during testing, guided with specific state coverage information for different models. Then, the joint optimization problem will be solved in a gradient-based manner and acquire minute perturbations.

Once obtained the perturbations, adversarial inputs are easy to acquire for models with continuous inputs like speech, by applying perturbations directly to original inputs. For models with discrete inputs like NLP tasks, the

perturbation applied to the test input probably will not lead to a legal input. Here, we adopt the nearest one as the adversarial input after iteratively scaling the perturbation, thus avoiding the invalid input. Finally, these adversarial inputs will be assessed concerning the tested model for the performance and coverage, to improve subsequent testing efficiency. The detailed descriptions for each step are given below.

4.1 State wrapping

In the inherent implementation of an RNN model, there are two data structures accessible in the inference: all the hidden states of the last layer, and all the hidden states and cell states (if LSTM underlying) of the last time step. For RNN testing, exploiting all the states should be a better choice for thoroughly searching for adversarial inputs. Therefore, we wrap the RNN cell implementation and keep all the hidden states and cell states of RNN cells in every layer and time step. With straightforward configurations, state wrapping will not interfere inner computation of the tested models. Note that this step is not expensive and will not affect the time efficiency, with an open-source library (20 lines of Python code). It is based on fundamental ‘‘RNNCell’’, the parent class of various cell implementation, making it possible to be generalized to most RNN models.

4.2 State inconsistency guided search

The state inconsistency guided search is the core portion of RNN-Test, an optimization problem composed of two parts. It is formulated in equation (3), where the first part (obj_1) is referred to as adversary search, and the second part (obj_2) is as coverage guidance. In addition, the two parts are free to be united together or alone, or even substituted with those of other methodologies, thus offering multiple possibilities of discovering adversarial inputs.

$$\begin{cases} obj = obj_1 + obj_2 \\ obj_1 = \sum_{l \in L} h_{t-1}^l + c_t^l - h_t^l \\ obj_2 = \sum_{k=0}^m e_k, e_k \in H(C) \end{cases} \quad (3)$$

Adversary search. In RNN-Test, a novel methodology is designed to craft adversarial inputs specially for RNN models. As illustrated in Fig. 3, hidden state vector h_t^l has a positive correlation with the inputs h_{t-1}^l, h_t^{l-1} and intermediate outputs c_t^l , if c implemented. Here, RNN-Test tries to increase h_{t-1}^l and c_t^l while decrease h_t^l simultaneously, which intentionally violates their inner dependencies to lead the model to exhibit unusual behaviors. Then the violated dependencies will spread across the whole model. Thus, RNN-Test is able to search for adversarial inputs distributed outside of the regular inference space.

As for the time step t in the objective, one step selected randomly out of each input will be adequate to achieve considerable performance. For the model with inputs always of hundreds of time steps, several more steps can be employed to increase the state inconsistency. Moreover, states of multiple layers $l \in L$ (L for all the layers) with respect to the same time step t will be leveraged to accelerate the search efficiency.

Coverage guidance. This part aims to cover the uncovered states, exercising more decision logics to produce adversarial inputs. RNN-Test leverages the proposed HS_C and CS_C metrics to guide adversarial testing, where CS_C guidance for LSTM models and HS_C for general RNN models. To boost the specific coverage, RNN-Test selects m hidden states or cell states to compose the optimization objective, as in formula (3). Rather than merely selecting uncovered states randomly, RNN-Test mainly chooses the states with values near to be covered so as to reach a higher coverage value at an earlier stage. Since CS_C is defined over a series of sections and the boundary sections are hardly covered, the states with values near the boundary section endpoints will be the targets to cover, thus leading RNN-Test to search in more sensitive space.

Then, the joint optimization objective will be maximized by mutating the test inputs, unlike the training course minimizing the prediction error by tuning the parameters. Given the objective, its derivative for the input x will be the perturbation, which is the gradient direction along which it increases or decreases most. Afterwards, the perturbations will be exploited to generate adversarial inputs.

Algorithm 1 Adversarial input synthesis for discrete inputs

Input: $x \leftarrow$ original test input
 $t \leftarrow$ one time step selected to modify
 $grad \leftarrow$ perturbations obtained
 $embs \leftarrow$ embeddings of the vocabulary
 $MAX_SCALE \leftarrow$ maximum degree of scaling the gradient
.....
1: */*generate adversarial inputs for NLP tasks.*/**
2: **procedure** GEN_ADV($x, t, grad, embs$)
3: $x' \leftarrow x$
4: $dist_vec \leftarrow \emptyset$
5: **for** $scale \in [1, MAX_SCALE]$ **do** *//search along the gradient*
6: $pert = grad_t \times scale$ *//perturbation for the time step*
7: $t_emb = x_t + pert$ *//get invalid embedding by gradient ascent*
8: **for** $emb \in embs$ **do**
9: $dist = norm(t_emb - emb)$ *//distance of t_emb to emb*
10: $dist_vec = dist_vec \cup \{dist\}$
11: $nearest_emb = argmin(dist_vec)$ *//the nearest embedding*
12: **if** $nearest_emb \neq x_t$ **then**
13: $x'_t = nearest_emb$ *//modify the time step*
14: **break**
15: **return** x' *//acquire the adversarial input*

4.3 Adversarial input synthesis

For continuous inputs like speech, the perturbations could be applied directly to acquire the adversarial input. For NLP tasks whose inputs are words or characters scattered in discrete embedding space, procedure GEN_ADV is presented in Algorithm 1. In the procedure, we iteratively scale the gradient to be applied as the perturbation and then search for the nearest word/character in the embedding space to mutate the input step (lines 8 to 14 in Algorithm 1). This is a straightforward way to obtain valid adversarial inputs, ridding of embeddings which are not equivalent to legal words/characters. Besides, the embedding representations of words or characters in each NLP task are acquired after enough training, which could unveil their semantic properties. Therefore, searching along the gradient for the nearest embedding could get desired adversarial inputs with existing semantic information.

TABLE 1: Summary of RNN models to evaluate RNN-Test. The first three models with sequential outputs are for major evaluation on sequential domains. The last model with classification outputs is constructed for comparison to testRNN.

Model	Description	Architecture	Performance		
			Metric	Reported	Ours
PTB language model	General language model	Two-layer LSTM in its small configuration(i.e. fewer steps)	Train perplexity ¹ Test perplexity	37.99 115.91	43.316 117.122
Spell checker model	Simple seq2seq model	Two-layer bi-direction LSTM for the encoder	Sequence loss	15%	10%
DeepSpeech ASR model	State-of-the-art ASR model	One-layer bi-direction LSTM with CNN layers around	WER ²	16%	16%
MNIST-LSTM model	Handwritten digit recognition of LSTM network	One-layer LSTM	Test accuracy	98.3%	96.88%

¹ Perplexity, the universal metric for language models, where lower perplexity corresponds to a better model.

² Word error rate, a common performance metric for seq2seq and ASR models, where higher WER means worse predictions.

In the literature, the adversarial input is identified for the imperceptibility from the original input but with the distinct class label. In sequential domains with no classifications, it is hard to recognize a generated sequence as the adversarial input avoiding false-positives, which has no standards yet [33], [34], [35]. Fortunately, model performance metrics are a good choice to exhibit qualities of adversarial inputs, which are supposed to be accessible in all the tasks. Consequently, adversarial inputs obtained will be fed into the model assessing whether to decay the performance and updating the coverage, where coverage information will be exploited to guide subsequent testing.

5 EXPERIMENT

5.1 Experiment Setup

Implementation. We developed the framework RNN-Test on the widely deployed framework Tensorflow 1.3.0, and evaluated RNN-Test on a computer having Ubuntu 16.04 as the host OS, with an Intel i7-7700HQ@3.6GHz processor of 8 cores, 16GB of memory and an NVIDIA GTX 1070 GPU.

As for hyperparameters in RNN-Test algorithms, such as m and MAX_SCALE, they are tuned for each tested model and not listed here for simplicity. We will release our code and datasets upon publication for further discussions.

Tested models. A summary of four tested models is presented in Table 1. We mainly evaluated RNN-Test on the first three RNN models dealing with different sequential tasks. The last model of classification task is particularly constructed for comparison to testRNN. These models of common structures and various tasks provide more confidence for the generalization of RNN-Test to other RNN models.

PTB language model [23] is a well-known RNN model, basically to generate subsequent texts taking previous texts as input. It is the implementation of the fundamental LSTM [29] without particular adaptations for specific applications. The training and testing data are provided by the Penn Tree Bank dataset [36], and we extracted the first 25 sentences of the testing data for evaluation. We trained this model to achieve comparable performance to that reported using the same training course.

Spell checker model [24] is one of the widespread seq2seq models in NLP tasks, which receives a sentence with spelling mistakes as input and outputs the sentence with mistakes corrected. The training data are twenty popular

books from project Gutenberg [37]. For testing, we constructed 160 sentences with spelling mistakes like examples of developers, thanks to rich sources from Tatoeba [38]. Since their pre-trained model is unavailable, we trained this model in the same way.

DeepSpeech ASR model [25] is a state-of-the-art speech-to-text RNN model employed in lots of security-critical scenarios. Its pre-trained model DeepSpeech-0.1.1 (Mozilla’s implementation) could be deployed conveniently, and our testing data are the first 20 samples extracted from the Common Voice corpus [39].

Baselines. For comparison, we customize adversarial testing methodologies FGSM [7] and DLFuzz [26] to work for RNN models. They both generate adversarial inputs by solving optimization problems in a gradient-based manner. We implement their optimization objectives and coverage metrics on RNNs, while other procedures are the same as RNN-Test. For FGSM-based method, its optimization objective only contains the adversary search part without coverage guidance, whereas DLFuzz-based methodology also makes use of coverage guidance to obtain adversarial inputs, where neuron coverage (NC) is the underlying metric. Here, the NC definition of RNN models is the same as DeepTest [40]. Note that the customization will not degrade their performance since our optimized searching procedures are also used for them.

For relevant works on RNN testing, there is a significant gap to conduct comparisons due to framework incompatibility. Tested models of testRNN [22] and DeepStellar [21] are built on Keras and corresponding models on TensorFlow are mostly unavailable. Ultimately, we developed RNN-Test on an LSTM network [27] of MNIST dataset, which is an image classifier both evaluated in the two works. This MNIST-LSTM network is constructed on TensorFlow and achieves comparable test accuracy over the default MNIST dataset.

Research questions: We constructed experiments to answer the following research questions.

- **RQ1.** How is the effectiveness of the RNN-Test? (§ 5.2)
- **RQ2.** How is the effectiveness of coverage guidance for adversarial testing? (§ 5.3)
- **RQ3.** How is the quality of adversarial inputs obtained by RNN-Test? (§ 5.4)

5.2 Effectiveness of RNN-Test (RQ1)

To conduct a thorough evaluation, we compare our RNN-Test approach with other methodologies over the tested

TABLE 2: Effectiveness of RNN-Test and other methods in their default settings, measured over adversarial inputs obtained by each method. Note that worse performance values (e.g. WER) indicate stronger test capability of methods (The best result across each row is denoted bold). The coverage guidance used by RNN-Test is given following w. (with), HS_C is hidden state coverage and CS_C is cell state coverage. The same symbols are used in below tables.

Model	Performance	Methodology					
		Original	Random testing	FGSM-based	DLFuzz-based	RNN-Test (w. HS_C)	RNN-Test (w. CS_C)
PTB language model	Perplexity	150.46	229.97	240.07	233.35	285.13	277.44
	Generation Rate ¹	-	100.00%	95.78%	93.59%	100.00%	100.00%
Spell checker model	WER	5.63	7.10	7.19	7.07	7.40	7.49
	BLEU ²	0.870	0.830	0.829	0.826	0.827	0.822
	Success Rate ³	-	64.58%	73.61%	73.61%	73.61%	76.39%
DeepSpeech ASR model	WER	5.50	5.35	6.65	6.18	8.10	7.80
	BLEU	0.796	0.800	0.747	0.786	0.703	0.720
	Success Rate	-	40.67%	90.00%	67.50%	100.00%	100.00%

¹ Generation rate. Ratio of the test set the methodology has managed to produce the adversarial input.

² BLEU (Bilingual evaluation understudy). Correspondence of prediction and the ground truth, where higher BLEU means better predictions.

³ Success Rate. Ratio of the generated adversarial inputs to successfully reduce the model performance, not used for the first model as its performance is recorded over all inputs.

TABLE 3: Effectiveness of RNN-Test compared to testRNN on MNIST-LSTM model over 500 original inputs. The results are listed in a similar way of testRNN [22], where those for testRNN are exactly that they reported.

MNIST-LSTM model	Methodology			
	testRNN (RM)	testRNN (TM)	RNN-Test (w. HS_C)	RNN-Test (w. CS_C)
Test Cases Generated	2000	2000	500	500
# Adv. Inputs	26	63	348	279
Avg. Perturb. (L2 norm)	1.051	4.028	1.74	1.69
Adversary Rate	1.3%	3.15%	69.6%	55.8%

models, measuring the model performance fed with adversarial inputs obtained by each methodology. Besides, we also provide results of original test inputs, and those of random testing that randomly replaces a word/character of text input or applies Gaussian noise to speech input. We run them on each tested model over the same original test set three times, to alleviate the uncertainty each time. The same settings are adopted for below evaluations of other methods.

Overall results. Table 2 summarizes the overall results, from which we could derive the following inferences. Firstly, adversarial inputs can decline the model performance, since tested models all achieve worse performance over adversarial input sets than the original test sets.

Secondly, random testing methods can also obtain adversarial inputs, but they are far from satisfactory. For the 100% generation rate on PTB language model, random replacement could always get mutated inputs while FGSM- and DLFuzz-based methods may fail to find adversarial inputs for some inputs.

Thirdly, the RNN-Test approach outperforms FGSM-based and DLFuzz-based approaches, with more performance reduction and higher success (or generation) rate. For instance, in comparison with FGSM-based method, RNN-Test (w. CS_C) achieves 15.57% higher perplexity and 4.22% higher generation rate on PTB language model, 4.17% higher WER, 0.84% lower BLEU and 2.78% higher success rate on spell checker model, and 17.29% higher WER, 3.61%

lower BLEU and 10% higher success rate on DeepSpeech ASR model. As for the slight improvement on the spell checker model, the sparse embedding space may be the primary cause, since it largely limits the searching capability.

How to choose the appropriate coverage guidance. As shown in Table 2, RNN-Test guided with HS_C and CS_C both gain better effectiveness than other methodologies, with no one always superior to the other. When applied for realistic RNN tasks, it is straightforward to choose the appropriate coverage guidance. For LSTM models, both are good alternatives. For common RNN and GRU models, HS_C is the choice, since hidden states are universal across these structures.

Comparison to relevant RNN testing works. As described in § 5.1, we conduct comparisons to other RNN testing methods over one MNIST-LSTM model due to framework incompatibility. Table 3 presents the comparison of RNN-Test to testRNN on MNIST-LSTM model, both over 500 original inputs as the evaluation setting of testRNN. TestRNN employed random mutation (RM) and targeted mutation (TM) to generate 2000 test cases for evaluation, respectively. In RNN-Test, once one adversarial input is obtained for the corresponding test input, the testing procedure starts for another test input, and thus acquired 500 test cases. In Table 3, RNN-Test could generate much more adversarial inputs out of fewer test cases, obtaining 52.65% to 66.45% higher adversary rate than TM algorithm of testRNN, a refined means of RM. Although testRNN (RM) obtains smallest perturbations, testRNN (TM) with largest perturbations still reached limited adversary rate. From another point of view, RNN-Test is also competitive to be applied to classification domains.

As to DeepStellar, it reported that thousands of adversarial inputs were obtained for 100 test inputs after 6 hours on a high-performance server. Despite our design, RNN-Test is also able to test each input continuously and generate much more adversarial inputs. If with a same high-performance server (with 28-core CPU, 196 GB RAM, and 4 NVIDIA Tesla V100 16G GPUs), RNN-Test is supposed to achieve comparable results.

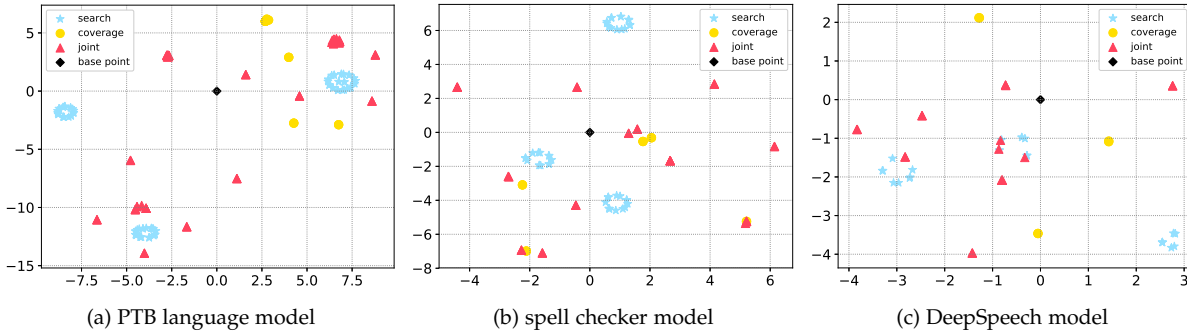


Fig. 6: TSNE transformations of perturbations of the above approaches with different optimization objectives for one same test input. Search is adversary search, coverage is coverage guidance, joint is joint objective, same in below figures. The divergent distribution represents various perturbations and thus adversarial inputs.

The answer to RQ1: The RNN-Test approach is effective in generating adversarial inputs, with the ability to reduce the model performance sharply with high success (or generation) rate.

5.3 State coverage guidance contributes to adversarial testing (RQ2)

Divergent perturbations of coverage guidance. The previous work [41] suggests that perturbations obtained by neuron coverage guidance are similar to adversary-based search methods (e.g. FGSM) and so the coverage guidance does not add too much, which is concluded based upon analyses over popular coverage guided testing methodologies [10], [12], [31]. But the conclusion may not work for the proposed state coverage metrics, since those criteria assessed are all over CNN neurons.

Here, we recorded perturbation vectors obtained by approaches we evaluated over the same inputs of each RNN model. Besides their default settings, we run each approach with either pure adversary search or coverage guidance, as well as the joint way. To visualize, we leverage the state-of-the-art high-dimensional reduction technique TSNE [42] to transform multi-dimensional perturbation vectors into two dimensions. As Fig. 6 shows, there is no evident similarity of perturbations of the adversary search, coverage guidance or joint objectives. In contrast, the divergent distribution implies that coverage guidance is capable to offer alternative perturbations, whether utilized independently or jointly, thus providing varied adversarial inputs. Therefore, the coverage guidance is worthy to be applied to adversarial testing.

Effectiveness of pure coverage guidance for adversarial testing. State coverage guidance can also be adopted to discover adversarial inputs independently, due to the unique perturbations of coverage guidance. Table 4 presents the results of HS_C , CS_C , and NC as guidance to be the only optimization respectively. As shown, state coverage metrics as guidance can acquire adversarial inputs on the tested models, while NC guidance fails to obtain any on PTB language model. Overall, both CS_C and HS_C outperform NC as guidance, especially on PTB language model and DeepSpeech ASR model. Therefore, neuron-based coverage metrics will not be appropriate for RNN models, as

TABLE 4: Effectiveness of state coverage guidance, compared to neuron coverage guidance.

Model	Performance	Methodology		
		NC	HS_C	CS_C
PTB language model	Perplexity	150.46	238.07	236.96
	Generation Rate	0%	94.66%	97.22%
Spell checker model	WER	7.03	7.48	8.00
	BLEU	0.830	0.827	0.823
	Success Rate	73.61%	75.00%	77.78%
DeepSpeech ASR model	WER	5.45	5.65	5.35
	BLEU	0.801	0.791	0.794
	Success Rate	10.00%	40.00%	55.00%

discussed in § 2.2. Surprisingly, when cross-referenced with Table 2 for the spell checker model, CS_C guidance exhibits the best performance with the highest WER and success rate.

Enhancement of coverage guidance to other methods. We also demonstrate that both FGSM-based and DLFuzz-based approaches with state coverage guidance could gain higher effectiveness than themselves and those jointed with NC guidance. For instance, Table 5 provides results of DLFuzz-based methodology jointed with HS_C and CS_C guidance, as well as its NC guidance. Compared with NC guidance, HS_C and CS_C guidance improve DLFuzz-based technique over the tested models in varying degrees. Additionally, similar results are attained for FGSM-based approach, where state coverage guidance improves more than NC guidance, as presented in Table 6. Thus, state coverage guidance is proven to be able to enhance other adversarial testing methodologies. Though these two methods could be improved with state coverage guidance, the most powerful means is still the RNN-Test, as cross-referenced with Table 2.

TABLE 5: Effectiveness of DLFuzz-based methodology with state coverage guidance, compared to its NC guidance.

Model	Performance	Methodology (DLFuzz-based)		
		w. NC	w. HS_C	w. CS_C
PTB language model	Perplexity	233.35	243.19	238.71
	Generation Rate	95.42%	97.44%	99.15%
Spell checker model	WER	7.07	7.44	7.10
	BLEU	0.826	0.825	0.825
	Success Rate	73.61%	75.00%	76.39%
DeepSpeech ASR model	WER	6.18	6.15	6.10
	BLEU	0.786	0.785	0.778
	Success Rate	67.50%	75.00%	70.00%

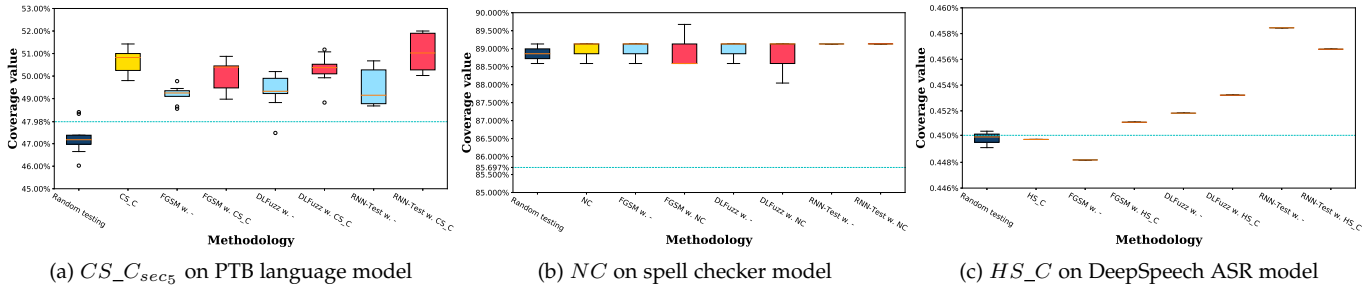


Fig. 7: Value ranges of coverage metrics among different approaches (w.- denotes no coverage guidance) over the same amount of adversarial inputs. The blue dashed lines denote the corresponding coverage value of original test input sets.

TABLE 6: Effectiveness of FGSM-based methodology jointed with the coverage metrics, compared to its default setting with no coverage guidance.

		Methodology (FGSM-based)			
Model	Performance	w. -	w. <i>NC</i>	w. <i>HS_C</i>	w. <i>CS_C</i>
PTB language model	Perplexity	240.07	241.23	256.91	256.91
	Generation Rate	95.78%	98.69%	97.65%	100.00%
Spell checker model	WER	7.19	6.99	7.46	7.14
	BLEU	0.829	0.832	0.828	0.831
	Success Rate	73.61%	73.61%	73.61%	70.83%
DeepSpeech ASR model	WER	6.65	6.70	6.75	6.80
	BLEU	0.747	0.747	0.748	0.746
	Success Rate	90.00%	90.00%	90.00%	90.00%

Coverage value may not be a strong indicator of methodology effectiveness. Numerous works [10], [12], [22], [31] adopt the coverage value as an indicator of effectiveness for adversarial testing methodologies. Meanwhile, researchers [41], [43] raised doubts that there may be limited correlations between coverage and robustness of DNNs.

Here, we have analyzed correlations between the model performance and values of coverage metrics on the first three models, but found out weak positive or even negative correlations (not given for brevity). Therefore, we could not draw the conclusion that obtaining higher coverage definitely results in higher effectiveness. For example, higher coverage value does not indicate higher success rate or WER on DeepSpeech ASR model. Though testRNN proved that their adversarial input set is with higher coverage rate than the normal input set, whether a higher coverage rate leads to a higher adversary rate also remains unsettled.

From the current point of view, we suggest that more efforts should be put into coverage guidance for adversarial testing, but not just to boost the coverage value. Indicators related to the model performance should be considered with more weight to assess an adversarial testing method.

Simple illustration of value ranges of coverage metrics. Fig. 7 presents the value ranges of *HS_C*, *CS_C* and *NC* achieved by different methods on each tested model respectively. We also provide results of each methodology with and without the corresponding coverage guidance, since the coverage guidance still tends to improve the value. For each box, it represents a set of coverage values of the methodology at different times, marked with bounds and the median. Note that boxes in Fig. 7c resemble lines be-

cause of the limited stochasticity and equal coverage values obtained (except random testing) on that model.

It must be claimed that coverage values strongly depend on the number of test inputs, and the same amount of inputs are supposed to be with similar value ranges. As presented, the value ranges of these coverage metrics among methodologies vary not much (within 5%), especially *NC* ranges are almost the same. It is the same case for figures not given here. Furthermore, *HS_C* values are always very low since it is inherently hard to boost *HS_C* with a few inputs, similar to boundary sections of *CS_C* over larger models. Meanwhile, it also supplies evidence that methodology effectiveness may be affected little by coverage values. However, coverage guidance is still worthy of more research investment. In summary, we could get the following answer.

The answer to RQ2: State coverage metrics as guidance are able to acquire adversarial inputs, superior to neuron coverage guidance whether independently or jointed with adversary search. The coverage guidance has the potential to be more effective, since the divergent perturbations and best performance on the spell checker model.

5.4 Quality of adversarial inputs of RNN-Test (RQ3)

Samples of adversarial inputs. Table 7 lists samples of adversarial inputs on the two NLP models, with each approach to modify the same word. For both models, RNN-Test tends to generate different words with other methods, offering diverse adversarial inputs. For PTB language model, our adversarial inputs could result in the model sampling words farther from semantics and generating higher perplexity texts, where results of RNN-Test are with totally wrong semantics. Meanwhile, adversarial inputs for the spell checker model could result in the corrected mistakes in original inputs appearing again in the predictions of adversarial inputs. For DeepSpeech ASR model, they could result in the model making wrong predictions, as depicted in Fig. 8. Such adversarial inputs of misleading semantics may harm the security requirements of tested models.

Time efficiency, tiny perturbations and reality concerns. RNN-Test also has high time efficiency, producing each adversarial input costs 3s, 11s, and 24s on average on PTB language model, spell checker model, and DeepSpeech ASR model, respectively. That means RNN-Test has the potential to be applied in industrial practice. In terms of

TABLE 7: Samples of adversarial inputs on the tested models, the targeted words to modify are in red and underlined. The affected results for the spell checker model are also underlined.

Methodology	PTB language model	spell checker model
Original	Input: no it was n't black <u>Monday</u> ... Perplexity: 259.67 Generate: economy goes forward on behalf of ...	Input: I would swim through <u>theocean</u> just to see your smile again. Predict: I would swim through <u>the ocean</u> just to see your smile again. Input: The sound of yur voice <u>islike</u> siren's songto me. Predict: The sound of yur voice <u>is like</u> siren' song to me.
FGSM-based	Input: no it was n't black <u>co.</u> ... Perplexity: 376.46 Generate: economy goes forward on behalf of ...	Input: I would swim through <u>theootean</u> just to see your smile again. Predict: I would swim through <u>the otean</u> just to see your smile again. Input: The sound of yur voice <u>isliee</u> siren's songto me. Predict: The sound of yur voice <u>is liee</u> siren' song to me.
DLFuzz-based	Input: no it was n't black <u>due</u> ... Perplexity: 357.38 Generate: soviets appear reluctant between france 's ...	Input: I would swim through <u>theootean</u> just to see your smile again. Predict: I would swim through <u>the otean</u> just to see your smile again. Input: The sound of yur voice <u>islske</u> siren's songto me. Predict: The sound of yur voice <u>issle</u> siren' song to me.
RNN-Test	Input: no it was n't black <u>\$</u> ... Perplexity: 513.91 Generate: soviets appear reluctant toward nov. a.m. ...	Input: I would swim through <u>theoKcean</u> just to see your smile again. Predict: I would swim through <u>the cocean</u> just to see your smile again. Input: The sound of yur voice <u>isltke</u> siren's songto me. Predict: The sound of yur voice <u>istle</u> siren' song to me.

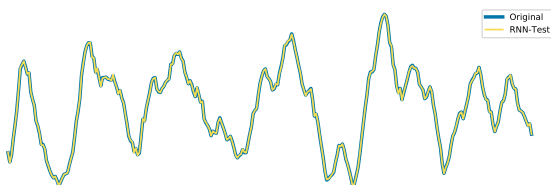


Fig. 8: An example adversarial input of RNN-Test for DeepSpeech ASR model. The waveform of a test input (blue, thick line) is overlapped with the waveform of the adversarial input (yellow, thin line). Each waveform is 500 samples long and was chosen randomly from the corresponding inputs. The original prediction is “the shop as closed on mondays” while the prediction for adversarial input is “the shop as close tan monas”.

the perturbations, only one word or character is modified for the former two models of text inputs. For DeepSpeech ASR model, the averaged perturbations are smaller than 0.04 in L2 norm. Therefore, adversarial inputs of RNN-Test appear benign to humans and ensure imperceptibility. As for realistic concerns about adversarial images of tasks like autonomous driving, adversarial cases will probably not encounter in real-world circumstances and weaken the security significance. In RNN testing, users are likely to mistype text inputs and fail NLP models in reality. Nevertheless, this is an urgent issue for ASR models at present, since adversarial audios [35] always become invalid when played over-the-air, which we will be devoted to in future works.

Improve the model by retraining. Last but not the least, adversarial inputs obtained by RNN-Test are also capable to improve the model performance by retraining. We tried on PTB language model and incorporated adversarial inputs (82.5 KB) to the training set (5.1 MB).

Table 8 presents the perplexity of PTB language model before and after retraining, where train perplexity indicates the performance on the training set while valid perplexity for the valid set. Here the data are averaged over 5 times of the same retraining process with 12 epochs, to mitigate affects due to the intrinsic indeterminism of neural networks. From columns 4 and 7, results show that the train perplexity

TABLE 8: The perplexity before and after retraining on PTB language model. Columns 3 and 5 are for the augmented training set. Columns 4 and 7 are for the improvement of retraining results w.r.t original results.

epoch	train perplexity			valid perplexity		
	original	w. adv.	increment	original	w. adv.	decrement
0	290.584	288.579	-0.690%	190.004	192.096	-1.101%
2	113.216	113.712	0.439%	140.328	140.339	-0.008%
4	86.290	87.195	1.049%	132.589	132.969	-0.287%
6	56.282	56.961	1.207%	121.410	120.566	0.695%
8	46.549	47.082	1.146%	122.981	121.611	1.114%
10	43.991	44.474	1.096%	123.065	121.385	1.365%
12	43.227	43.695	1.082%	122.440	121.020	1.159%

of the model after retraining increases by 1.082% whereas the valid perplexity decreases by 1.159% in end. Moreover, the test perplexity after retraining is 102.75, which is also declined by 12.582% compared to the original test perplexity 117.53. Notice that even by incorporating fewer adversarial inputs (1.6KB), the valid perplexity still declines by 0.058%. Therefore, adversarial inputs could alleviate the over-fitting issue in training by reducing little train performance, but improving the valid and test performance and thus the robustness of RNN models.

The answer to RQ3: RNN-Test could efficiently produce adversarial inputs of high quality, declining the model performance sharply and improving the model by retraining.

6 THREATS TO VALIDITY

Though RNN-Test exhibits appreciable effectiveness with the default setting in evaluations, its performance is inevitably influenced by the parameters, including the number of states selected to boost, the weights applied to joint objectives and the scaling degree of perturbations, especially the ways of sections splitting of CS_C . They are worthy to be well explored in the future work. Furthermore, the uncertainty running each time still exists, owing to stochastic word/character to modify, which could be diminished by fixing the target.

In addition, RNN-Test is devoted to being general and scalable for variants of RNNs, but we could not exhaustively evaluate all the variants and applications. In this paper, the structures of tested models are general to some extent, but training the spell checker model still costs hard work, due to its bad reproducibility of the training results given. Moreover, the state wrapping is designed to avoid interfering with model logics, but the adaptation efforts may be necessary for some variants with complex structures.

7 RELATED WORK

Adversarial deep learning. The concept of adversarial attacks was first introduced in [6]. It discovered that state-of-the-art DNNs would misclassify the input images by applying imperceptible perturbations, where these mutated inputs are called adversarial examples/inputs. Their work FGSM [7] and numerous following works [8], [9], [44] generate adversarial inputs by maximizing the prediction error in a gradient-based manner. They provide rich input resources for CNNs to improve their robustness.

Afterwards, [33] explains adversarial inputs for RNNs, but presents rough qualitative descriptions for those of sequential outputs. As mentioned before, existing works mostly focus on tasks of categorical outputs like sentiment analysis [20]. For those of text inputs, some works [19], [20] add, delete or substitute a word/character to construct adversarial inputs, leading models to give wrong classifications.

Unfortunately, few works are evaluating the majority RNN models processing sequential outputs. Due to no explicit class labels and no standards for such adversarial inputs, present works attack these models to perform abnormally in various ways. TensorFuzz [34] crafts adversarial inputs to lead the language model to sample words from the blacklist. Several works [35] fool well-known ASR models to produce targeted phrases given by the attacker. In this paper, we propose RNN-Test as an effective and scalable methodology for diverse models of sequential outputs.

Coverage guided testing. Based upon the exposed threats of DNNs, traditional software testing techniques are subsequently applied to test DNN systems, where coverage guided testing is of a popular trend. DeepXplore [10] first introduces neuron coverage which is defined over CNN neurons with pre-defined thresholds. Then, DeepGauge [12] defines a set of coverage metrics with finer-grained granularity, where neuron value ranges are split as thousands of sections according to training data. DeepCT [31] is even fine-grained to measure over combinations of neuron outputs. As stated in § 2.2, these neuron-based coverage metrics can not be applied to RNN states directly.

As for works also among the first attempts of adversarial testing for RNN systems, DeepStellar [21] adapts coverage metrics of DeepGauge to test RNN models, which need to be abstracted as a Markov Chain first. Despite its effectiveness, it is inevitable to miss key features and introduce computation overhead owing to intrinsic properties of abstraction. Another work testRNN [22] designs novel coverage metrics according to structures of LSTM models, some of which are special to quantify temporal relations. Besides of similar classification tasks, these two works both mutate inputs

directly (e.g. random noise) and use saturated coverage values to terminate testing. Our RNN-Test approach mutates inputs based on RNN logics and adopts coverage guidance to help search for adversarial inputs, which is effective for both sequential and classification domains.

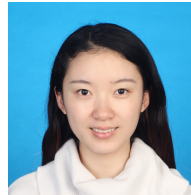
8 CONCLUSIONS

We design and implement the adversarial testing framework RNN-Test for recurrent neural networks. RNN-Test focuses on testing the main sequential structure without limit to tasks, aggregating advantages of both the proposed search method and novel state coverage metrics as guidance. It is superior to existing methodologies for DNN testing and could effectively produce adversarial inputs over RNN models of various applications, reducing model performance evidently with high success (or generation) rate. We also first demonstrate that coverage guidance has a diverse searching capability for adversarial inputs compared with other methods and our state coverage guidance outperforms neuron coverage guidance in RNN testing.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *International Conference on Neural Information Processing Systems*, 2012.
- [2] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proceedings of the 25th International Conference on Machine Learning*, ser. ICML '08. New York, NY, USA: ACM, 2008, pp. 160–167. [Online]. Available: <http://doi.acm.org/10.1145/1390156.1390177>
- [3] G. Hinton, G. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, B. Kingsbury, and T. Sainath, "Deep neural networks for acoustic modeling in speech recognition," *IEEE Signal Processing Magazine*, vol. 29, pp. 82–97, November 2012. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/deep-neural-networks-for-acoustic-modeling-in-speech-recognition/>
- [4] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.
- [5] D. Shen, G. Wu, and H.-I. Suk, "Deep learning in medical image analysis," *Annual review of biomedical engineering*, vol. 19, pp. 221–248, 2017.
- [6] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.
- [7] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *Computer Science*, 2015.
- [8] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *Security and Privacy (EuroS&P), 2016 IEEE European Symposium on*. IEEE, 2016, pp. 372–387.
- [9] S. M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," in *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, no. EPFL-CONF-218057, 2016.
- [10] K. Pei, Y. Cao, J. Yang, and S. Jana, "Deepxplore: Automated whitebox testing of deep learning systems," in *Proceedings of the 26th Symposium on Operating Systems Principles*. ACM, 2017, pp. 1–18.
- [11] Y. Sun, X. Huang, and D. Kroening, "Testing deep neural networks," *arXiv preprint arXiv:1803.04792*, 2018.
- [12] L. Ma, F. Juefei-Xu, F. Zhang, J. Sun, M. Xue, B. Li, C. Chen, T. Su, L. Li, Y. Liu *et al.*, "Deepgauge: Multi-granularity testing criteria for deep learning systems," in *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. ACM, 2018, pp. 120–131.

- [13] Y. LeCun, Y. Bengio *et al.*, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.
- [14] P. Rodriguez, J. Wiles, and J. L. Elman, "A recurrent neural network that learns to count," *Connection Science*, vol. 11, no. 1, pp. 5–40, 1999.
- [15] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [17] T. Mikolov, S. Kombrink, L. Burget, J. Černocký, and S. Khudanpur, "Extensions of recurrent neural network language model," in *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2011, pp. 5528–5531.
- [18] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013, pp. 6645–6649.
- [19] M. Sato, J. Suzuki, H. Shindo, and Y. Matsumoto, "Interpretable adversarial perturbation in input embedding space for text," *arXiv preprint arXiv:1805.02917*, 2018.
- [20] S. Samanta and S. Mehta, "Towards crafting text adversarial samples," *arXiv preprint arXiv:1707.02812*, 2017.
- [21] X. Du, X. Xie, Y. Li, L. Ma, Y. Liu, and J. Zhao, "Deepstellar: Model-based quantitative analysis of stateful deep learning systems," in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. ACM, 2019, pp. 477–487.
- [22] W. Huang, Y. Sun, J. Sharp, W. Ruan, J. Meng, and X. Huang, "Coverage guided testing for recurrent neural networks," *arXiv preprint arXiv:1911.01952*, 2019.
- [23] (2016, Sep.) Variant of ptb word lm that could save and restore the model and display the predictions. <https://github.com/nelken/tf>.
- [24] (2017, June) A seq2seq model that can correct spelling mistakes. <https://github.com/Currie32/Spell-Checker/>.
- [25] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Satheesh, S. Sengupta, A. Coates *et al.*, "Deep speech: Scaling up end-to-end speech recognition," *arXiv preprint arXiv:1412.5567*, 2014.
- [26] J. Guo, Y. Jiang, Y. Zhao, Q. Chen, and J. Sun, "Dlfuzz: differential fuzzing testing of deep learning systems," in *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. ACM, 2018, pp. 739–743.
- [27] U. ZCC, "Understanding inputs and training process of lstm," <https://www.cnblogs.com/USTC-ZCC/p/11171209.html>, 2019.
- [28] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [29] W. Zaremba, I. Sutskever, and O. Vinyals, "Recurrent neural network regularization," *arXiv preprint arXiv:1409.2329*, 2014.
- [30] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.
- [31] L. Ma, F. Zhang, M. Xue, B. Li, Y. Liu, J. Zhao, and Y. Wang, "Combinatorial testing for deep learning systems," *arXiv preprint arXiv:1806.07723*, 2018.
- [32] Y. Gong and C. Poellabauer, "Crafting adversarial examples for speech paralinguistics applications," *arXiv preprint arXiv:1711.03280*, 2017.
- [33] N. Papernot, P. McDaniel, A. Swami, and R. Harang, "Crafting adversarial input sequences for recurrent neural networks," in *Military Communications Conference, MILCOM 2016-2016 IEEE*. IEEE, 2016, pp. 49–54.
- [34] A. Odena and I. Goodfellow, "Tensorfuzz: Debugging neural networks with coverage-guided fuzzing," *arXiv preprint arXiv:1807.10875*, 2018.
- [35] N. Carlini and D. Wagner, "Audio adversarial examples: Targeted attacks on speech-to-text," *arXiv preprint arXiv:1801.01944*, 2018.
- [36] A. Taylor, M. Marcus, and B. Santorini, "The penn treebank: an overview," in *Trebanks*. Springer, 2003, pp. 5–22.
- [37] (2016, Apr.) Free ebooks - project gutenber. http://www.gutenberg.org/ebooks/search/?sort_order=downloads.
- [38] (2006) Tatoeba is a collection of sentences and translations. <https://tatoeba.org/cmn/downloads>.
- [39] (2017, June) Common voice is a project to help make voice recognition open to everyone. <https://voice.mozilla.org/>.
- [40] Y. Tian, K. Pei, S. Jana, and B. Ray, "Deeptest: Automated testing of deep-neural-network-driven autonomous cars," in *Proceedings of the 40th International Conference on Software Engineering*. ACM, 2018, pp. 303–314.
- [41] Z. Li, X. Ma, C. Xu, and C. Cao, "Structural coverage criteria for neural networks could be misleading," 2019.
- [42] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.
- [43] Y. Dong, P. Zhang, J. Wang, S. Liu, J. Sun, J. Hao, X. Wang, L. Wang, J. S. Dong, and D. Ting, "There is limited correlation between coverage and robustness for deep neural networks," *arXiv preprint arXiv:1911.05904*, 2019.
- [44] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 38th IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 39–57.



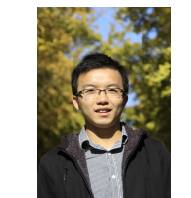
Jianmin Guo is a Ph.D. candidate at School of Software Engineering, Tsinghua University, Beijing, China. She received the BS degree in software engineering from Beijing University of Posts and Telecommunications, Beijing, China, in 2017. Her research interests are software testing, mainly focusing on deep learning testing and adversarial testing of recurrent neural networks.



Yue Zhao is a software engineer in Huawei Technologies Co., Ltd. He received MS degree at School of Software Engineering, Tsinghua University, Beijing, China, in 2020. He received the BS degree in software engineering from Beijing University of Posts and Telecommunications, Beijing, China, in 2017. His research interests are deep learning testing and backdoor detection of deep learning systems.



Quan Zhang is a Ph.D. student at School of Software Engineering, Tsinghua University, Beijing, China. He received the BS degree in computer science from Beijing University of Posts and Telecommunications, Beijing, China, in 2020. His research interests are backdoor detection of deep learning systems.



Yu Jiang received the BS degree in software engineering from Beijing University of Posts and Telecommunications, Beijing, China, in 2010, and the PhD degree in computer science from Tsinghua University, Beijing, China, in 2015. He was a Postdoc researcher in the department of computer science of University of Illinois at Urbana-Champaign, IL, USA, in 2016. He is now an associate professor in Tsinghua University, Beijing, China. His current research interests include domain specific modeling, formal computation model, formal verification and their applications in embedded systems.

putation model, formal verification and their applications in embedded systems.