



# Fuzzing Deep-Learning Libraries via Automated Relational API Inference

Yinlin Deng\*

University of Illinois at Urbana-Champaign  
USA  
yinlind2@illinois.edu

Anjiang Wei

Stanford University  
USA  
anjiang@stanford.edu

Chenyuan Yang\*

University of Illinois at Urbana-Champaign  
USA  
cy54@illinois.edu

Lingming Zhang

University of Illinois at Urbana-Champaign  
USA  
lingming@illinois.edu

## ABSTRACT

Deep Learning (DL) has gained wide attention in recent years. Meanwhile, bugs in DL systems can lead to serious consequences, and may even threaten human lives. As a result, a growing body of research has been dedicated to DL model testing. However, there is still limited work on testing DL libraries, e.g., PyTorch and TensorFlow, which serve as the foundations for building, training, and running DL models. Prior work on fuzzing DL libraries can only generate tests for APIs which have been invoked by documentation examples, developer tests, or DL models, leaving a large number of APIs untested. In this paper, we propose DeepREL, the first approach to automatically inferring relational APIs for more effective DL library fuzzing. Our basic hypothesis is that for a DL library under test, there may exist a number of APIs sharing similar input parameters and outputs; in this way, we can easily “borrow” test inputs from invoked APIs to test other relational APIs. Furthermore, we formalize the notion of value equivalence and status equivalence for relational APIs to serve as the oracle for effective bug finding. We have implemented DeepREL as a fully automated end-to-end relational API inference and fuzzing technique for DL libraries, which 1) automatically infers potential API relations based on API syntactic/semantic information, 2) synthesizes concrete test programs for invoking relational APIs, 3) validates the inferred relational APIs via representative test inputs, and finally 4) performs fuzzing on the verified relational APIs to find potential inconsistencies. Our evaluation on two of the most popular DL libraries, PyTorch and TensorFlow, demonstrates that DeepREL can cover 157% more APIs than state-of-the-art FreeFuzz. To date, DeepREL has detected 162 bugs in total, with 106 already confirmed by the developers as previously unknown bugs. Surprisingly, DeepREL has detected 13.5% of the high-priority bugs for the entire PyTorch issue-tracking system

\*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ESEC/FSE '22, November 14–18, 2022, Singapore, Singapore

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9413-0/22/11...\$15.00

<https://doi.org/10.1145/3540250.3549085>

in a three-month period. Also, besides the 162 code bugs, we have also detected 14 documentation bugs (all confirmed).

## CCS CONCEPTS

- Software and its engineering → Software testing and debugging; Software reliability.

## KEYWORDS

Fuzz Testing, Oracle Inference, Deep Learning, Differential Testing

### ACM Reference Format:

Yinlin Deng, Chenyuan Yang, Anjiang Wei, and Lingming Zhang. 2022. Fuzzing Deep-Learning Libraries via Automated Relational API Inference. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '22), November 14–18, 2022, Singapore, Singapore*. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3540250.3549085>

## 1 INTRODUCTION

Recent years have witnessed the surge of deep learning (DL) in a variety of applications, including computer vision [32, 60], natural language processing [24, 28], robotics [23, 35], bioinformatics [46, 57], and software engineering [25, 29, 36, 37, 45, 68, 78, 79, 81, 82]. Meanwhile, similar to traditional software systems, DL systems can also have bugs, which can lead to serious consequences and may even threaten human lives [4].

To date, most prior work on DL testing focused on testing/verifying DL models, with an emphasis on adversarial attacks [14, 19, 27, 42, 47, 51], metrics for model testing [31, 34, 41, 53, 77], application-specific model testing [62, 85, 88], and verifying certain properties of models [15, 38]. Meanwhile, there is limited work targeting the reliability of DL libraries, which serve as the central infrastructures for building DL models, and are the foundation for training and deploying DL models. CRADLE [54] is one of the trailblazing work on testing DL libraries, which resolves the oracle challenge with differential testing of various DL models on multiple backends of Keras [1]. AUDEE [30] and LEMON [65] further augment CRADLE by leveraging search-based mutation strategies to generate more diverse DL models/inputs for testing library code. Different from the above model-level DL library testing techniques, more recently, FreeFuzz [66] has been proposed to mine example inputs from open source (including code snippets from the library documentation, developer tests, and DL models in the wild) to directly test each

DL library API in isolation. FreeFuzz has been evaluated on PyTorch [52] and TensorFlow [13], currently the two most popular DL libraries (with 54K/162K stars on Github). The experimental results show that FreeFuzz can cover 9× more APIs than state-of-the-art LEMON and detect various previously unknown bugs.

Despite the promising results, existing techniques for fuzzing DL libraries still suffer from the following limitations. First, the input generation is still far from optimal. CRADLE and AUDEE can only test APIs that are covered in the original models, and LEMON can cover slightly more APIs with layer mutations; furthermore, although FreeFuzz can cover up to 1158 APIs for PyTorch and TensorFlow (which is already a huge improvement over other work), it is still unable to test an API if there is no code snippet directly invoking the API. Second, there is still a lack of powerful test oracles. Existing techniques typically perform differential testing across different DL libraries or hardware backends (e.g., GPU/CPU) to address the test oracle issue. However, differential testing across DL libraries is typically applied at the model level and suffers from the limited effectiveness of model-level testing (e.g., limited API coverage and accumulated floating-point precision loss) [30, 54, 65], while different backends often share common code logic/design (and thus may also share similar bugs) [66]. Thus, it is also crucial to investigate novel test oracles for effective DL library fuzzing.

To address the aforementioned limitations, in this work, we open a new dimension for testing DL libraries via automated relational API inference. The inspiration stems from the fact that prior work [16, 26, 43, 64] has discovered a number of equivalent APIs in traditional software systems (e.g., Java projects)<sup>1</sup>. We envision such relational API inference also to be an inspiring direction for fuzzing DL libraries. In this way, given the same inputs generated via fuzzing, APIs that are equivalent in functionality should produce the same numerical results (i.e., *value equivalence*). Moreover, besides the previously studied equivalent APIs, we further leverage the fact that DL APIs with similar functionality should behave similarly in terms of program status (i.e., *status equivalence*) for more effective fuzzing. For example, although AdaptiveAvgPool3d and AdaptiveMaxPool3d in PyTorch are not equivalent, they are functionally similar APIs; thus, we can feed any valid input of the first API to the second API and expect its invocation to also be successful. Based on this intuition, we can easily “borrow” test inputs generated for one API to test other relational APIs. Also, API relations can directly serve as test oracles for differential testing. Therefore, we can easily overcome the aforementioned limitations.

We have built a fully-automated technique, DeepREL, which infers such API relations without human intervention for fuzzing DL libraries. One key challenge is how to obtain the API relations automatically and accurately. Existing work [16, 26, 43, 64] on equivalent API inference for traditional software systems can hardly be applied for DL library testing, e.g., the most recent MeMo work [16] heavily relies on well-documented API relations, which are rare in DL libraries. To this end, DeepREL first automatically infers all possible candidate API pairs that are matched based on API syntactic and semantic information. Then, DeepREL synthesizes concrete test programs for those potentially relational APIs. After that, DeepREL

<sup>1</sup>Some of such existing work([16, 26, 43]) treated the entire software systems under test as the test objects, and thus viewed this as *metamorphic testing* [20]. In this paper, we treat each API as a test object and view this as *differential testing* [44].

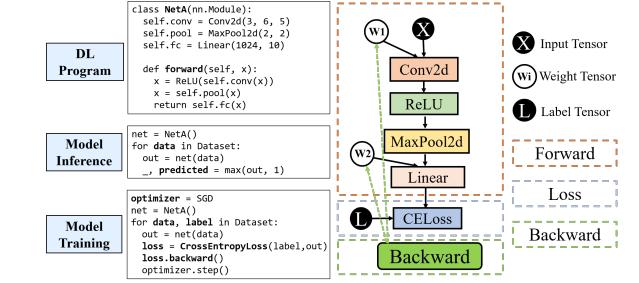


Figure 1: Background knowledge on DL Models and APIs

leverages a set of representative valid inputs (automatically traced during prior normal API executions) to check whether the inferred API relations hold or not. Lastly, DeepREL takes the validated API pairs and leverages mutation-based fuzzing to generate a much more diverse and extensive set of test inputs for detecting potential inconsistencies among relational APIs. Our study has shown for the first time that there can be a surprising number of equivalent or similar APIs within popular DL libraries (e.g., 4290/8808 verified relational API pairs by DeepREL for PyTorch/TensorFlow), which can substantially help with fuzzing DL libraries (and beyond). In summary, our paper makes the following contributions:

- **Dimension.** This paper opens a new dimension for fully-automated DL library fuzzing via relational API inference.
- **Technique.** We build DeepREL, a *fully-automated* end-to-end framework for DL library testing. DeepREL automatically infers all possible candidate relational APIs based on both API syntactic and semantic information, and then dynamically verifies them via test program synthesis. While this work focuses on DL libraries, the basic idea of DeepREL is general and can also be applied to other software systems.
- **Evaluation and Impact.** DeepREL covers 1815 more APIs than prior work (i.e., 157% improvement), and has detected 162 bugs in total, with 106 already confirmed by the developers as previously unknown bugs. Surprisingly, DeepREL was able to detect 13.5% of the high-priority bugs for the entire PyTorch issue-tracking system in a three-month period. Also, besides the 162 code bugs, we were able to detect 14 documentation bugs (all confirmed) as a by-product of our experimentation.

## 2 BACKGROUND

### 2.1 Basics about DL Models and APIs

For DL models, *inference* is the process of using a fixed DL model to complete a specific task, while *training* is the process of updating the weights of a neural network to better perform a certain task given the labeled data (under the scenario of supervised learning [72]). We next shed light on how this is achieved by APIs from DL libraries.

**DL Models.** To build and run a DL model, developers first need to define the model by writing a DL program in DL libraries (e.g., PyTorch [52] and TensorFlow [13]). Take the DL program NetA written in PyTorch (shown on the left side of Figure 1) as an example, it includes a convolution layer (Conv2d), a max pooling layer (MaxPool2d), and a linear layer (Linear). The function forward defines how the input tensor (x) should flow in the defined layers (and other related APIs). Besides input tensors, there are also *weight*

*tensors* (e.g.,  $w1$  and  $w2$  shown on the right side of Figure 1). Their values will be updated during training, the process of which is called *back-propagation* [69], a procedure natively supported by DL libraries. Training the model is achieved by first running the forward part (i.e., inference) of the neural network ( $out = net(data)$ ), computing the loss ( $loss = \text{CrossEntropyLoss}(label, out)$ ), computing the gradient ( $loss.backward()$ ), and invoking the optimizer ( $\text{optimizer.step}()$ ) for back-propagation.

**DL APIs.** When running a DL model, the APIs involved in building the model are also executed. Essentially, writing a DL program can be viewed as defining a *computation graph*. It is a directed acyclic graph (DAG) whose nodes are DL APIs while edges are the flow of the tensors. Figure 1 also shows the computation graph, composed of three parts: forward part (taking input tensors and weight tensors as input), loss computation part (requiring label tensors), and backward part (for updating weight tensors). Actually, the backward part needs to construct a very complex graph [87], but we omit it in Figure 1 for simplicity. Essentially, running a whole DL model can be broken down into invoking a series of DL APIs based on the topological sorting [74] of the computation graph.

## 2.2 Fuzzing DL Libraries

To our knowledge, there are mainly two categories of work for fuzzing DL libraries, model-level testing, and API-level testing.

**Model-level Testing.** CRADLE [54] is one of the first to apply differential testing for DL libraries. Since Keras [1] is a library featuring high-level APIs for building DL models, APIs in Keras may have multiple implementations in its supported lower-level libraries. Thus, CRADLE takes 30 pre-trained DL models as input and runs differential testing to find inconsistencies between different low-level libraries for Keras. More recently, AUDEE [30] and LEMON [65] have proposed to use search-based mutation strategies to generate mutated DL models for differential testing on different backends. While AUDEE focuses on mutating parameters of layers, weight tensors, and input tensors, LEMON applies mutation rules by adding/deleting layers and changing the values of weight tensors. In this way, LEMON’s mutation rules are more general and can cover more APIs than the original models. However, even LEMON’s mutation can only be applied to a limited number of APIs. For instance, its intact-layer mutation rule [65] requires that the output tensor shape of the API to be added/deleted should be identical to its input tensor shape, making a large number of APIs inapplicable for mutation. Researchers have recently shown that LEMON can hardly invoke additional library code/APIs with its mutation rules, e.g., covering only 35 APIs in total for TensorFlow [66]. Very recently, Liu et al. [39] proposed NNSmith, which leverages lightweight specifications for more diverse model generation. Meanwhile, it also mainly focuses on commonly used APIs/operators.

**API-level Testing.** Different from prior work on DL library fuzzing, the recent FreeFuzz work [66] proposes to directly mine test inputs from open source for API-level fuzzing, a much finer-grained level than model-level testing. One challenge is that Python is a dynamically-typed language, making it hard to determine the types of API parameters for fuzzing Python APIs. Prior work has to manually set up the API arguments, and thus can only test a small

```
1 result1 = torch.broadcast_shapes(*shapes)
2 result2 = torch.broadcast_tensors(*map(torch.empty, shapes))[0].shape
```

Figure 2: API pair with the same output

number of APIs, e.g., Predoo [86] can only test 7 APIs for TensorFlow. A very recent work DocTer [76] constructs rules to extract DL-specific input constraints from API documentation and uses them to generate valid/invalid inputs for testing DL libraries; meanwhile, it requires manual annotation for 30% of API parameters. In contrast, FreeFuzz resolves this challenge fully automatically via dynamically tracing API executions in code snippets from the documentation, developer tests, and 202 DL models. More specifically, FreeFuzz records the traced argument values in a database and further performs mutation-based fuzzing to mutate those traced values to generate even more inputs for fuzzing DL library APIs. Lastly, FreeFuzz applies differential testing on different hardware backends (i.e., CPU/GPU) for detecting potential consistency bugs. Despite its big improvement over prior work, FreeFuzz can only test 1158 APIs for PyTorch and TensorFlow, which are the ones covered in its input mining stage, leaving a total of 6815 APIs uncovered. Also, different hardware backends may still share code logic/design, causing the differential testing oracle used by FreeFuzz to miss various bugs. In this work, we propose to test relational APIs to further overcome such limitations. We build our technique (DeepREL) upon FreeFuzz to automatically infer relational APIs and leverage them to fuzz DL libraries. Note, however, that our DeepREL idea is general and can be built on any API-level fuzzer for DL libraries (e.g., DocTer [76]). We choose FreeFuzz since it is a recent state-of-the-art technique that is both publicly available and fully automated.

## 3 PRELIMINARIES

We first introduce the preliminaries for our fuzzing technique in this section. Given the set of all possible APIs,  $\mathcal{A}$ , for a DL library under test, we aim to define the relational property between the invocation results of a source API  $S \in \mathcal{A}$  and a target API  $T \in \mathcal{A}$ .

Intuitively, we can directly check whether  $S$  and  $T$  produce equivalent outputs. For example, Figure 2 shows an API pair which, according to the PyTorch documentation [3], should always produce the same results. The `broadcast_shapes` API applies broadcasting on a list of compatible shapes to align them, while the `broadcast_tensors` API applies broadcasting on a list of shape-compatible tensors to align their shapes. In fact, the first API can be rewritten as 1) creating intermediate empty tensors from tensor shapes with `map` and `torch.empty`, 2) calling `broadcast_tensors` with these tensors, and 3) getting the shape of the output tensor. Since the source and target APIs can achieve the same functionality with totally different implementations given the same input (shapes in Figure 2), they provide a great opportunity for differential testing [44]. Therefore, we have the following formal definition:

*Definition 3.1. Equivalence<sub>value</sub>.* Given a set of inputs  $\mathcal{D}$ , source API  $S \in \mathcal{A}$  and target API  $T \in \mathcal{A}$  satisfy Equivalence<sub>value</sub> (modulo  $\mathcal{D}$ ) iff their invocations always output the same results given any input in  $\mathcal{D}$ . Formally,

$$S \equiv T \pmod{\mathcal{D}} \iff \forall x \in \mathcal{D}. S(x) = T(x) \quad (1)$$

While this can be effective in detecting potential consistency bugs, the checking is too strict and may not apply to a large number

```

1 layer1 = torch.nn.AdaptiveAvgPool3d(output_size)
2 result1 = layer1(input)
3 layer2 = torch.nn.AdaptiveMaxPool3d(output_size)
4 result2 = layer2(input)

```

**Figure 3: API pair with different outputs but same status**

of APIs. In fact, it could be possible that  $S$  and  $T$  produce totally different results, but tend to behave similarly given similar inputs. For example, the API pair shown in Figure 3 does not satisfy the Equivalence<sub>value</sub> property since the output of AdaptiveAvgPool3d is different from AdaptiveMaxPool3d. The first API applies a 3D adaptive *average* pooling over an input but the latter applies a 3D adaptive *maximum* pooling. However, these two APIs do share certain functionality in common since both of them apply a pooling operation, which is also valuable for testing. Therefore, we further abstract the invocation result of a program into a set of coarse-grained *statuses*: Success, Exception, and Crash. Success denotes that program executions terminate normally, while Exception means that program executions throw known exceptions. Lastly, Crash represents the cases where the program executions crash with unexpected errors, e.g., segmentation faults or INTERNAL ASSERT FAILED errors (which are “*never acceptable*” as commented by PyTorch developers). We then further introduce the notation of  $\llbracket \cdot \rrbracket \in \{\text{Success}, \text{Exception}, \text{Crash}\}$  to return the execution status of an input program. For example,  $\llbracket S(x) \rrbracket = \text{Success}$  indicates that  $S$  terminates normally with the input  $x$ . In this way, we can define another property for checking potential consistency:

**Definition 3.2. Equivalence<sub>status</sub>.** Given a set of inputs  $\mathcal{D}$ , source API  $S \in \mathcal{A}$  and target API  $T \in \mathcal{A}$  satisfy Equivalence<sub>status</sub> (modulo  $\mathcal{D}$ ) iff their invocation always output the same statuses given any input in  $\mathcal{D}$ . Formally,

$$S \sim T(\text{mod } \mathcal{D}) \iff \forall x \in \mathcal{D}. \llbracket S(x) \rrbracket = \llbracket T(x) \rrbracket \quad (2)$$

To conclude, the Equivalence<sub>status</sub> relation is a relaxed notation for the Equivalence<sub>value</sub> relation, which is in turn a relaxed notation of semantic equivalence (denoted as  $S \equiv T$ ). Formally,

$$S \equiv T \implies S \equiv T(\text{mod } \mathcal{D}) \implies S \sim T(\text{mod } \mathcal{D}) \quad (3)$$

One crucial component of these two definitions is the domain  $\mathcal{D}$  on which the properties are constrained. Aiming for more accurate API relations, it would be beneficial to cover more representative test inputs within the intersection of the valid input space of the source and target APIs as the domain.

## 4 FUZZING RELATIONAL APIs

Figure 4 shows the overview of our DeepREL technique for fuzzing relational APIs of DL libraries. DeepREL takes as input the targeted DL library, its API documentation, and a database of valid historical API invocations (e.g., automatically collected via running documentation examples, library tests, and DL models [66]). Each entry of the database contains the concrete argument values passed into an API during an invocation and is obtained through dynamic tracing. Overall, DeepREL performs the following four phases iteratively:

**API Matcher** (Section 4.1). In order to test a DL library which typically has hundreds or even thousands of APIs, the first challenge is to identify the API pairs that are likely to satisfy the desired properties Equivalence<sub>value</sub> or Equivalence<sub>status</sub>. API Matcher maps

each API into embeddings based on API documentation, and uses embedding similarity to identify candidates of API pairs.

**Invocation Synthesizer** (Section 4.2). Given a collection of potential API pairs, Invocation Synthesizer decides how to invoke them. To construct valid invocations for later verification, we impose a constraint on the source API: it must have at least one (valid) invocation in the database. In this way, given an invocation of the source API, Invocation Synthesizer aims to synthesize the invocation code for the target API.

**API Match Verifier** (Section 4.3). Given the invocation code of matched APIs, this phase would check whether each API pair satisfies property Equivalence<sub>value</sub> or Equivalence<sub>status</sub> with a set of representative inputs as the verifying test inputs. If the result values (resp. execution statuses) are consistent for all tests, then API Match Verifier accepts the API pair as Equivalence<sub>value</sub> (resp. Equivalence<sub>status</sub>). If API Match Verifier detects any inconsistency in this phase, it then rejects the API pair.

**API Fuzzer** (Section 4.4). The last step is to leverage the verified API pairs to detect potential consistency bugs. API Fuzzer uses mutation-based fuzzing to generate a large number of test inputs for source APIs, and tests the verified API pairs with oracles (Section 3).

Lastly, recall that in order to generate valid inputs, the source API must have at least one (valid) invocation in the database. DeepREL further adopts an iterative process to cover more API pairs (Section 4.5). The newly generated valid target API invocations can be added to the database to serve as the source APIs for the next iterations to detect more potential API pairs. The following sections would explain each phase in detail.

### 4.1 API Matcher

In this phase, DeepREL identifies potential matched API pairs from the documentation. DeepREL uses API Matcher to infer similar API pairs as candidate pairs (which will be further verified later). API Matcher would map each API into its embeddings, and compute the cosine similarity of embeddings to be the similarity of each API pair. We consider *signature similarity* and *document similarity* that cover both API syntactic and semantic information for similarity computation. Overall, for each API pair, the similarity is defined as the maximum of the two:

$$\text{Sim}_{\text{API}}(S, T) = \text{Max}(\text{Sim}_{\text{sig}}(S, T), \text{Sim}_{\text{doc}}(S, T)) \quad (4)$$

We compute the pair-wise similarity for every API pair, and pair each API with its  $K$ -closest neighbors as the candidate pairs.  $K$  is a hyper-parameter and is set to 10 in the default setting of DeepREL. Notably, we also analyze the impact of different values of  $K$  in our experimental study (Section 6.3).

**Signature Similarity.** The signature of an API contains the API name and an ordered list of argument names. The APIs that will be paired typically have signatures that follow a similar syntactic pattern. For example, the API pair `tf.maximum` and `tf.minimum` satisfies Equivalence<sub>status</sub>, and their signatures are very similar: `tf.maximum(x, y)` and `tf.minimum(x, y)`. We map an API signature into its TF-IDF (term frequency-inverse document frequency) embedding [73], and compute the embedding similarity.

TF-IDF has been widely adopted in the field of information retrieval, and it reflects the importance of each word in a document. Some common words like `tf` and `torch` in the API signature are

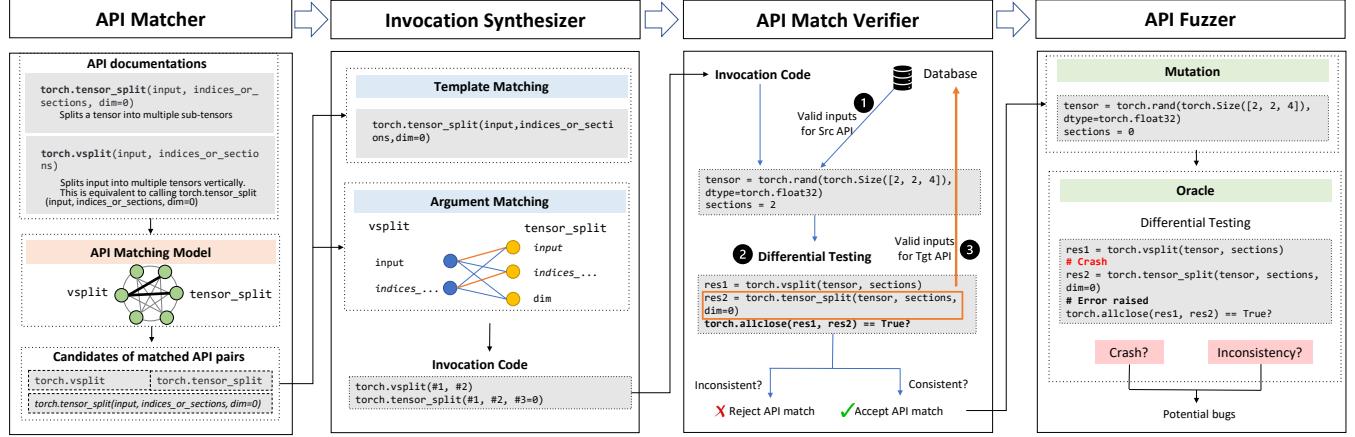


Figure 4: DeepREL overview

less informative, so their TF-IDF weights tend to be smaller. To obtain the TF-IDF embedding for each API, we first break the API signatures into subwords (also called tokens) and then standardize them. Let  $n$  denote the size of the vocabulary from all the tokenized API signatures, an API  $S$  can be represented as an unnormalized term frequency embedding  $[c_1^S, c_2^S, \dots, c_n^S]$ , where  $c_j^S$  is the number of occurrences of word  $j$  in the API signature of  $S$ . We further normalize it with the inverse document frequency for each word to get the TF-IDF embedding:

$$\text{Rep}_{\text{sig}}(S) = \left[ \frac{c_1^S}{\sum_{S' \in \mathcal{A}} c_1^{S'}}, \frac{c_2^S}{\sum_{S' \in \mathcal{A}} c_2^{S'}}, \dots, \frac{c_n^S}{\sum_{S' \in \mathcal{A}} c_n^{S'}} \right] \quad (5)$$

The cosine similarity of two vectors is the cosine of the angle between them, and thus always belongs to the interval  $[-1, 1]$ . The cosine similarity of two arbitrary vectors  $x, y$  is defined as follows:

$$\text{Cos}(x, y) = \frac{x \cdot y}{\|x\| \|y\|} \quad (6)$$

We then compute the cosine similarity between the TF-IDF embeddings to be the signature similarity of two APIs  $(S, T)$ :

$$\text{Sim}_{\text{sig}}(S, T) = \text{Cos}(\text{Rep}_{\text{sig}}(S), \text{Rep}_{\text{sig}}(T)) \quad (7)$$

**Document Similarity.** To complement the signature similarity, we further model the semantic similarity between API documents. First, we extract all API descriptions from the documentation, each of which is a one-sentence summary of an API given at the beginning of the document. For example, API `torch.vsplit` is described as “*Splits input, a tensor with two or more dimensions, into multiple tensors vertically according to indices\_or\_sections.*” The description succinctly and surgically states the expected input, the transformation applied, and the expected output. Next, we use Sentence-BERT [58] to encode these informative description sentences into semantically meaningful sentence embeddings. The Sentence-BERT encoder takes a natural language sentence as input and outputs a vector in high-dimensional space; it targets specifically at generating embeddings whose cosine similarity reflects the semantic textual similarity. We use  $S.\text{description}$  to denote the description for the API  $S$ , and use  $\text{SBEncoder}$  to denote the Sentence-BERT

encoder. For each API  $S$ , we obtain its document embedding as:

$$\text{Rep}_{\text{doc}}(S) = \text{SBEncoder}(S.\text{description}) \quad (8)$$

For each API pair  $(S, T)$ , we compute the cosine similarity between their document embeddings as their document similarity:

$$\text{Sim}_{\text{doc}}(S, T) = \text{Cos}(\text{Rep}_{\text{doc}}(S), \text{Rep}_{\text{doc}}(T)) \quad (9)$$

## 4.2 Invocation Synthesizer

In this phase, we leverage argument matching and template matching to synthesize the invocation code for each matched API pair. Note that the invocation code of source API is simply the code snippet that directly invokes the source API. Therefore, we will next focus on generating the target API invocation code.

**Argument Matching.** For each candidate API pair, DeepREL first synthesizes the invocation code based on API definitions. It maps the arguments of the source API to the arguments of the target API to synthesize the invocation code of the target API (with the arguments from the source API).

We transform the argument matching problem into a *maximum weighted bipartite matching* problem [71]. DeepREL generates the invocation code based on the best argument match. More formally, given the source API  $S$ , the target API  $T$ , and their argument lists  $S.\text{args}$  and  $T.\text{args}$ , the corresponding bipartite graph is  $G = (L, R, E)$ , where  $L = \{a \mid a \in S.\text{args}\}$ ,  $R = \{b \mid b \in T.\text{args}\}$  and  $E = \{(a, b) \mid a \in L, b \in R\}$ . The weight of each edge  $(a, b) \in E$  is the similarity  $\text{Sim}_{\text{arg}}(a, b)$  of  $a$  and  $b$ , which is defined as:

$$\text{Sim}_{\text{arg}}(a, b) = \text{Sim}_{\text{name}}(a, b) + \text{Sim}_{\text{type}}(a, b) + \text{Sim}_{\text{pos}}(a, b) \quad (10)$$

The similarity  $\text{Sim}_{\text{arg}}(a, b)$  is determined by the names, potential types, and positions of the arguments. First, the similarity of the names is computed based on the following formula:

$$\text{Sim}_{\text{name}}(a, b) = 1 - \frac{\text{Levenshtein}(a.\text{name}, b.\text{name})}{\text{Max}(\text{Len}(a.\text{name}), \text{Len}(b.\text{name}))} \quad (11)$$

where  $a.\text{name}$  is the name of argument  $a$ . This similarity is based on the Levenshtein Distance [70] between the two names. Next, we compute the similarity of two type sets as:

$$\text{Sim}_{\text{type}}(a, b) = \frac{|a.\text{type} \cap b.\text{type}|}{|a.\text{type}|} \quad (12)$$

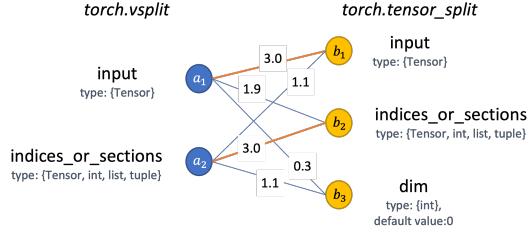


Figure 5: Example weighted bipartite graph

```
torch.vsplit(#1, #2)
torch.tensor_split(#1, #2, dim=0)
```

Figure 6: Invocation synthesis via argument matching

where  $a_{type}$  is the set of possible types of  $a$ . If the set of types that argument  $b$  can take contains all possible types of  $a$ ,  $a$  is more likely to be mapped to  $b$  since all types of  $a$  are legal for  $b$ .

We also compute the positional similarity of two arguments as:

$$Sim_{pos}(a, b) = 1 - \frac{|a_{idx} - b_{idx}|}{\text{Max}(\text{Len}(S.\text{args}), \text{Len}(T.\text{args}))} \quad (13)$$

where  $a_{idx}$  is the index of  $a$  in  $S$ 's argument list. For example, if  $a$  and  $b$  are both the first argument for  $S$  and  $T$ , then  $|a_{idx} - b_{idx}|$  equals to 0 and thus their positional similarity is 1.

After constructing the graph, DeepREL leverages Kuhn–Munkres algorithm [48] to find the best argument match and synthesizes the invocation code based on it. When the source and target APIs have the same number of arguments, DeepREL will synthesize the invocation code based on the best argument match directly. Otherwise, if any non-optional argument is unmatched, DeepREL will abort for the current API pair since the search space for determining the values of unmatched non-optional arguments is huge. That said, DeepREL only considers the case where the optional arguments of the source or target API are unmatched. For the unmatched optional arguments, DeepREL just uses their default values (Python optional arguments always have default values).

Take an API pair `vsplit` and `tensor_split` for example, whose weighted bipartite graph is shown in Figure 5. For each vertex, its name and type information are marked next to it, such as vertex  $a_1$ , whose argument name is `input`, possible type set is composed of `Tensor`, and index is 1. The weight on each edge is the similarity between the two arguments. For vertices  $a_1$  and  $b_1$ , since they have exactly the same name, possible type, and index, their similarity is 3. For  $a_1$  and  $b_2$ , their names are different but the type set of  $a_1$  is the subset of  $b_2$ 's, so they have a relatively high similarity (1.9). However, for  $a_2$  and  $b_1$ , only one type of  $a_2$  is legal for  $b_1$ , causing their similarity to be low (1.1). Thus, the best match in the graph is  $\{(a_1, b_1), (a_2, b_2)\}$ , leaving  $b_3$ , the optional argument `dim`, unmatched. Then DeepREL sets `dim` as its default value 0 to generate the invocation code, as shown in Figure 6 (where placeholders `#1`, `#2` indicate the argument mapping between source and target APIs).

Note that for each argument of each API, we gather its possible types from open source. To be precise, we extract the argument type information from all the traced invocations of the API to form the possible type set. The type set could be incomplete since it relies on the traces. When no traces cover a particular argument  $b$  of a target API, the type set  $b_{type}$  will be empty; thus, the type similarity between  $b$  and any other argument  $a$  will be  $Sim_{type}(a, b) = 0$ . Note

```
tf.scatter_nd(indices, updates, shape, name=None)
Calling tf.scatter_nd(indices, updates, shape) is identical to calling
tf.tensor_scatter_nd_add(tf.zeros(shape, updates.dtype), indices, updates).
```

#### (a) Documentation of `tf.scatter_nd`

```
tf.scatter_nd(#1, #2, #3)
tf.tensor_scatter_nd_add(tf.zeros(#3, #2.dtype), #1, #2)
```

#### (b) Invocation code from template

Figure 7: Invocation synthesis via template matching

that the argument matching algorithm still works in this case, with only the name and positional similarities being considered.

**Template Matching.** For some complex API pairs, DeepREL cannot leverage argument matching to generate the correct invocation code of the target API. Figure 7b shows the right invocation code between `scatter_nd` and `tensor_scatter_nd_add`. Obviously, argument matching fails to synthesize the complex invocation code of `tensor_scatter_nd_add`, in which case a matching template can be useful. A *matching template* is a code snippet elaborating a matched API pair. It presents an invocation of the target API, whose inputs are obtained from the arguments of the invocation of the source API. Figure 7a shows an example matching template (highlighted with underline) from the documentation of `tf.scatter_nd`. It suggests that invoking `scatter_nd` is equivalent to invoking `tensor_scatter_nd_add` with proper argument mappings as shown in Figure 7b.

To automatically detect such templates, DeepREL examines each code block in the documentation. Whenever a code snippet contains the invocation of another API, DeepREL extracts it as a potential candidate. Note that not every API pair has such templates. If DeepREL failed to extract any, the matching template will simply be `None`. For API pairs with a matching template, regardless of whether it is in top- $K$  similar pairs, DeepREL synthesizes additional invocation code using the matching template as the target API invocation code.

### 4.3 API Match Verifier

In this phase, DeepREL runs the invocation code synthesized for each matched API pair over a set of *verifying inputs* to validate properties `Equivalence_value` and `Equivalence_status` (defined in Section 3). The verifying inputs are a collection of valid inputs for the source API. These inputs can be collected from the documentation, library tests, and existing DL models, so that they are representative of the valid input space of the source API. For example, in our implementation, we leverage state-of-the-art FreeFuzz [66], which can collect all such information fully automatically.

**Equivalence<sub>value</sub>.** DeepREL first checks whether a candidate API pair satisfies the `Equivalence_value` property given the verifying inputs. If so, DeepREL accepts this matching pattern and marks it as `Equivalence_value`. Take the API pair in Figure 7b for example, `scatter_nd` always has the same output as `tensor_scatter_nd_add` for all verifying inputs. Hence, it is labeled as `Equivalence_value`.

**Equivalence<sub>status</sub>.** If the API pair violates `Equivalence_value`, DeepREL further checks whether the source and target API have the same status given verifying inputs. If so, DeepREL accepts this API relation as `Equivalence_status`. For instance, the `Equivalence_value` property does not hold for the pattern shown in Figure 3 since the output of `AdaptiveAvgPool3d` is different from `AdaptiveMaxPool3d`.

over the input set. However, the Equivalence<sub>status</sub> property holds as they always have the same status over the verifying inputs.

If an API pair is verified as Equivalence<sub>value</sub> or Equivalence<sub>status</sub>, it will be accepted by the API Match Verifier and further tested in the next fuzzing phase; otherwise, DeepREL rejects this pair. Please note that it is crucial to have a set of representative verifying inputs. If the verifying inputs are not representative, the API Match Verifier could mistakenly accept a wrong API relation when the verifying inputs do not cover certain important regions of the possible input space. We will study such false positives in detail in our experimental study (Section 6.2). On the other hand, the API Match Verifier may also reject some true matched API pairs if the verifying inputs directly trigger real consistency bugs in this phase. Although it is hard to avoid such false negatives, our experimental results show that DeepREL detects 162 bugs for popular DL libraries fully automatically, demonstrating the effectiveness of this design.

#### 4.4 API Fuzzer

In the last phase, DeepREL further leverages the verified API pairs to detect bugs for DL libraries. Specifically, DeepREL applies off-the-shelf mutation-based fuzzing techniques [66] to mutate the source API inputs for generating diverse inputs for differential testing. For Equivalence<sub>value</sub>, the source and target APIs are expected to have the same output. We can detect consistency bugs by comparing the results of the source and target APIs. For Equivalence<sub>status</sub>, the source and target APIs are expected to have the same status. Thus, we compare their statuses to detect bugs.

#### 4.5 The Iterative Process

In order to cover more APIs by matched API pairs, DeepREL performs the above four phases iteratively until the fixed point or a given number of iterations  $I$  (by default 10). In the API Match Verifier phase, if the target API is not covered in the current iteration and its invocation generated by the synthesizer has Success status, we will add this invocation into the database (even when the corresponding API relation got falsified) and label the target API as “newly covered API”. After this iteration, if there is any newly covered API, DeepREL will re-run the framework *with these newly covered APIs as the source APIs*; otherwise, the fixed point has been reached, and DeepREL will terminate. It is worth mentioning that the entire iterative DeepREL approach is fully automated. For the newly covered APIs, the verifying inputs are also automatically “borrowed” from the source APIs’ valid inputs.

Figure 4 presents one example of this iterative process. In the first iteration, the API Match Verifier takes API pair (`torch.vsplit`, `torch.tensor_split`) as input, and queries the invocation database ① to get a record for `torch.vsplit`. Invoking `torch.tensor_split` ② results in Success. Assuming that `torch.tensor_split` is not covered in the database at the beginning of the iteration, this successful invocation is then inserted into the database ③. In the next iteration, this invocation record will be retrieved ① to verify the matched API pairs with `torch.tensor_split` as the source API.

### 5 EXPERIMENTAL SETUP

In the experiments, we address the following research questions:

- RQ1: How effective is DeepREL in terms of API coverage?

- RQ2: What is the false positive rate of DeepREL?
- RQ3: How do different configurations affect DeepREL?
- RQ4: Can DeepREL detect real-world bugs?

Our experiments are performed on PyTorch 1.10 [52] and TensorFlow 2.7 [13], the latest stable release versions for the two most popular DL libraries (with 54K/162K stars on GitHub), which are the most widely studied DL libraries in prior DL library testing work [30, 54, 65, 66]. We run all experiments on a machine with an 8-core 2.20GHz Intel CPU, 16GB RAM, and Ubuntu 20.04.

#### 5.1 Implementation

**API Matcher.** To find high-quality matched API pairs, we first use the bs4 Python package [2] to parse the documentation of all 7973 APIs from TensorFlow and PyTorch. We collect both API signatures and descriptions from the documentation. Before computing the TF-IDF embedding, we use the Snowball stemmer [55] to convert tokens into word stems. As for the document embedding, we use the SentenceTransformer Python package [12] and adopt the pre-trained model all-MiniLM-L6-v2 as our *SBEncoder*.

**Invocation Synthesizer.** For argument matching, we use the munkres Python package [10] (implementing the Kuhn–Munkres algorithm) to solve the maximum weighted bipartite matching problem. For template matching, we automatically search and extract the code snippets for matching templates from the documentation.

**API Match Verifier.** We verify the relation of candidate API pairs with a set of representative valid inputs. We obtain the valid inputs traced from various input sources used by state-of-the-art FreeFuzz [66], which include the documentation, developer tests, and 202 DL models, and are representative to verify the function of APIs. We feed a maximum of 100 valid invocations of the source API from the FreeFuzz database into both the source API and target API as the verifying inputs and check if they have consistent behaviors.

**API Fuzzer.** We leverage the fuzzing strategies of FreeFuzz to mutate all the valid inputs traced for each source API (within the FreeFuzz database) and run all the generated inputs (1000 for each source API following the default setting of FreeFuzz [66]) on both the source API and target API to detect consistency bugs.

#### 5.2 Metrics

**# of Covered API.** Following prior work in DL library testing [66], we report the number of covered APIs. An API is covered by DeepREL if it is successfully invoked by API Match Verifier either as a source API or a target API (i.e., invocations with the Success status). Since DL libraries contain a large number of APIs, API coverage is an important metric of test adequacy.

**False Positive Rate.** If an API pair satisfies that 1) at least one of its invocation code is accepted by the API Match Verifier, and 2) at least one of its accepted invocation code is against its labeled property during the fuzzing phase, it is named an inconsistent API pair. The false positive rate for inconsistent API pairs is the proportion of detected inconsistent API pairs which are false alarms. It is commonly used in prior work on fuzzing or testing [22, 61, 80].

**# of Detected Bugs.** Bug finding is the ultimate goal for fuzzing, and thus we also report the number of distinct bugs DeepREL finds.

**Table 1: Comparison with FreeFuzz on API coverage**

	#Total	#FreeFuzz	#DeepREL	Improvement(%)
PyTorch	1592	470	1071	601 (128%)
TensorFlow	6381	688	1902	1214 (176%)
Total	7973	1158	2973	1815 (157%)

**Table 2: Verified API pairs**

	Equivalence <sub>value</sub>	Equivalence <sub>status</sub>	Total
PyTorch	1357	2933	4290
TensorFlow	5132	3676	8808
Total	6489	6609	13098

## 6 RESULT ANALYSIS

### 6.1 RQ1: Effectiveness in API Coverage

In this RQ, we aim to study the effectiveness of DeepREL in terms of covering more APIs with API relations. Table 1 shows the number of DL library APIs covered by DeepREL and state-of-the-art FreeFuzz. Column “#Total” presents the total number of APIs in DL libraries, while Column “Improvement” presents the improvement of DeepREL over FreeFuzz. In the fuzzing stage, DeepREL covers 2973 APIs, which is a huge improvement (157%) over FreeFuzz that covers only 1158 APIs. For example, there are totally 1592 PyTorch APIs, and DeepREL can cover 1071 APIs, 128% more than FreeFuzz. A large number of APIs are not covered by FreeFuzz because they are less frequently used and not covered by any of the three sources of FreeFuzz. Leveraging API relations, DeepREL can successfully invoke these APIs with their relational APIs’ inputs. The huge API coverage improvement demonstrates the potential of DeepREL.

**Table 3: Source distribution of inferred API pairs**

		Tgt	Seed	New
		Src	Seed	New
PyTorch	Seed	902	544	
	New	598	2246	
TensorFlow	Seed	777	1633	
	New	830	5568	

Table 2 shows the number of verified API pairs detected by DeepREL. Columns “Equivalence<sub>value</sub>” and “Equivalence<sub>status</sub>” present the number of value-equivalent and status-equivalent API pairs accepted by the API Match Verifier respectively. DeepREL accepts 13098 API pairs in total, showing that such API relations are common in DL libraries. On PyTorch, DeepREL accepts more status-equivalent API pairs (2933) than value-equivalent (1357). The reason is that the latter relation is stricter than the former, and status-equivalent API pairs are more common. For example, in term of splitting a tensor, PyTorch provides a set of APIs: `torch.split`, `torch.tensor_split`, `torch.vsplit` (splits the tensor vertically), and `torch.dsplit` (splits the tensor depthwise). It is worth noting that TensorFlow has much more APIs than PyTorch, and DeepREL detects more value-equivalent API pairs than status-equivalent ones on TensorFlow. This is because TensorFlow contains lots of APIs for compatibility and low-level access operations and thus has higher functional overlap: (1) The `tf.compat` module [8] contains 2579 redundant APIs to support forward and backward compatibility across versions. For example, `tf.compat.v1.layers.conv2d` is an alias for `tf.layers.conv2d`, and it allows user to use the `conv2d` layer with TensorFlow v1 behavior in TensorFlow v2; (2) The `tf.raw_ops`

module [9] contains 1339 low-level APIs to provide direct access to all TensorFlow ops. For example, `tf.raw_ops.Pad` adds padding to tensors and is a low-level API compared to the high-level API `tf.pad` with the same functionality.

Table 3 further presents a detailed distribution of the API pairs inferred by DeepREL based on whether the source/target APIs involve newly covered APIs. Column “Src” and Row “Tgt” present the categorization of the source and target APIs, respectively; Columns/Rows “Seed” and “New” indicate whether an API is from “seed APIs” covered by FreeFuzz or is newly covered by DeepREL. Out of all the 13098 API pairs verified by DeepREL, 1679 (902 + 777) only involve APIs covered by the original FreeFuzz, and all the remaining 11419 pairs involve newly covered APIs, which shows the importance of leveraging API relations to cover more APIs.

We also conduct a manual study to investigate why there are so many value-equivalent API pairs. Note that we do not look into status-equivalent API pairs because they are more intuitive (e.g., many APIs may share similar input parameter types and/or output behaviors). Since the number of verified value-equivalent API Pairs is huge, we select the set of equivalent API pairs explicitly specified in the documentation for our study. We mine the documentation for all 7973 PyTorch and TensorFlow APIs to extract API pairs when one API explicitly references another API in the documentation. In this way, we automatically extract 2942 API pairs and we manually categorize them in terms of why such relational APIs exist. 2692 out of 2942 API pairs are value-equivalent pairs. Note that 1828 of them are backward compatibility pairs, and are unique to TensorFlow. Thus, we group the remaining pairs into 5 main categories as shown in Table 4. Ease of programming is the main reason for Equivalence<sub>value</sub> API pairs in both DL libraries. For example, `torch.det` is an alias for `torch.linalg.det`, and users can use the two symbols interchangeably. Meanwhile, Deprecation is one of the minor reasons. It is worth noting that we include the deprecated APIs in our study since they are still in the code base and can help cover more new APIs as well as find potential consistency bugs. Table 4 also provides examples to demonstrate each reason, where Column “Example (S, T)” refers to the (source API, target API) pair. The last column explains the difference between the relational APIs. For the Performance example, `tf.stack` and `tf.parallel_stack` are equivalent APIs which pile a list of tensor up. `parallel_stack` is more efficient than `stack` as the documentation of `tf.parallel_stack` says “*parallel\_stack will copy pieces of the input into the output as they become available, in some situations this can provide a performance benefit.*” [7].

### 6.2 RQ2: False Positives

False Positive Rate (FPR) is a common metric to evaluate the effectiveness of fuzz testing. Table 5 shows the FPR of DeepREL on the DL libraries. We analyze all inconsistencies reported by API Fuzzer. Column “All” presents the total number of inconsistencies detected by API Fuzzer. “TP” (True Positive) means the number of true inconsistencies, and “FP” (False Positive) means the number of false alarms (e.g., inconsistencies due to incorrectly inferred API relations). We separately report the statistics of the Equivalence<sub>value</sub> and Equivalence<sub>status</sub> oracles, while the “Overall” statistics report the merged results.

**Table 4: Classification of reasons for Equivalence<sub>value</sub> API pairs**

Reason	# TF	% TF	# PT	% PT	Example ( $S, T$ )	Diff. between $S$ and $T$
Ease of programming	460	95.44%	349	91.36%	(torch.det, torch.linalg.det)	$S$ is an alias for $T$
Performance	10	2.07%	9	2.36%	(tf.parallel_stack, tf.stack)	$S$ uses parallelism for efficiency.
Special cases	10	2.07%	10	2.62%	(tf.boolean_mask, tf.ragged.boolean_mask)	$T$ extends $S$ to ragged tensors.
Numerical stability	1	0.21%	4	1.05%	(torch.linalg.inv, torch.linalg.solve)	$T$ is faster and more numerically stable.
Deprecation	1	0.21%	10	2.62%	(torch.qr, torch.linalg.qr)	$S$ is deprecated.

**Table 5: False positive rate of DeepREL**

	Oracle	# All	TP	FP	FPR
PyTorch	Equivalence <sub>value</sub>	412	364	48	11.65%
	Equivalence <sub>status</sub>	853	554	299	35.05%
	Overall	1265	918	347	27.43%
TensorFlow	Equivalence <sub>value</sub>	97	68	29	29.90%
	Equivalence <sub>status</sub>	363	209	154	42.42%
	Overall	460	277	183	39.78%
Total	Equivalence <sub>value</sub>	509	432	77	15.13%
	Equivalence <sub>status</sub>	1216	763	453	37.25%
	Overall	1725	1195	530	30.72%

The FPR of DeepREL is only 30.72%, which implies that our API relation detection and verification techniques are effective. The false positives mainly originate from incorrect API relation verification. API Match Verifier leverages valid inputs of the source API to decide whether an API relation candidate is correct or not. However, the valid inputs come from the FreeFuzz database and are not complete. Hence, API Match Verifier can mistakenly accept a wrong API relation when the inputs cannot cover a certain part of the possible input space and thus is insufficient to distinguish the different behaviors of the API pair. We can also observe that the FPR of Equivalence<sub>value</sub> is lower than Equivalence<sub>status</sub>. This is because that Equivalence<sub>value</sub> is stricter than Equivalence<sub>status</sub>, making it easier to reject wrong API relations over verifying inputs.

### 6.3 RQ3: Impacts of Configuration

In this RQ we analyze how different configurations affect the performance of DeepREL, including the API coverage, FPR, and running time. We focus on two hyper-parameters,  $K$  (the number of matched pairs for each source API, discussed in Section 4.1) and  $I$  (the number of iterations, discussed in Section 4.5). The default values for  $K$  and  $I$  are both 10 in DeepREL. To figure out the impact of  $K$ , we run our experiments with different  $K$  values of 5, 10, 15, 20. We also run DeepREL for up to 10 iterations and show the impact of different  $I$  values from 1 to 10. Figures 8 and 9 show the results, where the  $x$  axis is the iterations, and the  $y$  axis is the number of covered APIs, FPR, and the running time respectively. The results for different  $K$  values are shown in different lines. We can observe that DeepREL terminates (i.e., reaching fixed points) in at most 8 iterations on PyTorch, and on TensorFlow it either terminates in at most 8 iterations, or only covers 2 new APIs in the last iteration.

The impact of  $K$  is similar for PyTorch and TensorFlow. First, the number of covered APIs of  $K = 10, 15, 20$  are close, all significantly higher than  $K = 5$ . Second, the FPR increases as  $K$  increases. The reason is that less similar API pairs can incur more false positives. Third, the running time is increasing approximately linearly with  $K$ . Therefore, it is reasonable to set  $K$  as 10 in our default setting, as setting  $K$  to be higher than 10 would bring marginal benefit in terms of API coverage, but degrade FPR and time cost.

For the impact of  $I$ , we have the following observations. First, as expected, the number of covered APIs increases at a lower and lower pace with the increase of  $I$ , and converges within 10 iterations. Second, the FPR generally increases with iteration. The reason is that source APIs in later iterations are typically target APIs from earlier iterations, and may have fewer and fewer valid inputs (since they may fail on some inputs from the original source APIs) for verifying inferred API pairs. The only exception is that FPR drops when  $I$  increases to 2 with  $K = 5$  or 10 on PyTorch. We look into it and observe that DeepREL can accidentally detect a large number of true positives in the 2<sup>nd</sup> iteration (e.g., torch.Tensor.\* APIs and their value-equivalent APIs torch.\*). Third, the running time grows more rapidly at early iterations and more slowly later, which is consistent with the growth of the number of covered APIs.

We can also observe that the total running time of the default DeepREL is 12.8h for PyTorch and 26.3h for TensorFlow. Such cost is actually quite common for fuzzing techniques, e.g., various fuzzing techniques have been applied for 24h or even more [17, 40, 59, 67, 80] (including the recent LEMON work [65]).

### 6.4 RQ4: Bugs Detected

Table 6 presents the summary of real-world bugs detected by DeepREL for the studied DL libraries. Column “Total” shows the total number of bugs detected by DeepREL, and Columns “Value” and “Status” show the number of bugs detected with the Equivalence<sub>value</sub> and Equivalence<sub>status</sub> oracle respectively. We also present the number of bugs rejected by developers (Column “Rejected”), confirmed as previously unknown (Column “Confirmed”), and the number of previously unknown bugs that have already been fixed (Column “Fixed”). We can observe that DeepREL is able to detect 162 bugs in total, with only 7 rejected by developers, 106 confirmed by developers as previously unknown bugs (48 already fixed), and all others pending. Among those 106 confirmed bugs, only 9 can be found by FreeFuzz, and none of them can be detected by CRADLE, AUDEE, or LEMON. Furthermore, we have also found 10 documentation bugs for PyTorch and 4 for TensorFlow during the experiment (these document bugs are not included in Table 6).

The bugs detected by DeepREL can also be categorized into single API bugs<sup>2</sup> and consistency bugs between relational API pairs. Table 7 shows the breakdown of the 106 confirmed bugs based on the two cases. For the 19 single API bugs detected by DeepREL, 9 are in the “seed APIs” (Column “Seed”), while 10 are bugs in the newly covered APIs (Column “New”), emphasizing the importance of leveraging API relations to cover more APIs. Meanwhile, all the remaining 87 ones are consistency bugs, demonstrating the effectiveness of using API relations as the oracle for DL library

<sup>2</sup>The single API bugs are unexpected crashes caused by single buggy APIs. In Table 6, we categorize such bugs as “Value”/“Status” if they were detected when testing API pairs with the Equivalence<sub>value</sub>/Equivalence<sub>status</sub> oracle (although they can be detected without such relational API oracles).

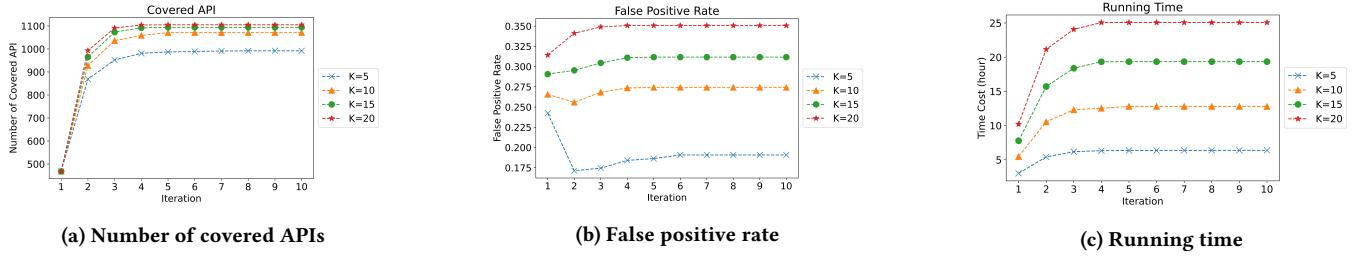
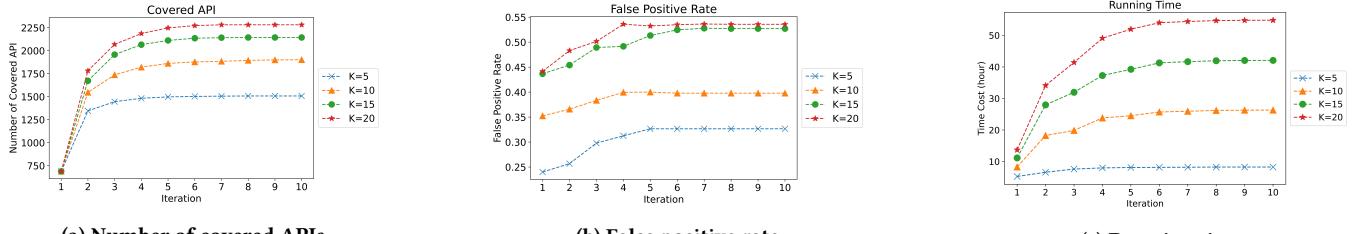
Figure 8: Analysis of hyper-parameter top- $K$  and iteration  $I$  on PyTorchFigure 9: Analysis of hyper-parameter top- $K$  and iteration  $I$  on TensorFlow

Table 6: Summary of detected bugs

	Total			Rejected			Confirmed			Fixed		
	Total	Value	Status	Total	Value	Status	Total	Value	Status	Total	Value	Status
PyTorch	121	76	45	6	3	3	71	33	38	29	16	13
TensorFlow	41	22	19	1	0	1	35	22	13	19	6	13
Total	162	98	64	7	3	4	106	55	51	48	22	26

Table 7: Source distribution of confirmed bugs

	Single API bugs			Relational API bugs			
	Total	Total	Seed	New	Total	Seed-only	Others
PT	71	18	9	9	53	26	27
TF	35	1	0	1	34	9	25
Total	106	19	9	10	87	35	52

testing. More specifically, out of the 87 consistency bugs, 35 are inconsistencies between “seed APIs” (Column “Seed-only”) while 52 involve newly covered APIs (Column “Others”), further indicating that covering new APIs contributes to consistency bug detection.

Notably, DeepREL is able to detect 23 high-priority bugs for PyTorch (note that TensorFlow is not discussed here as it does not have such labels). These bugs are marked as “high priority” because they are critical and require urgent resolution. The reported “high priority” bugs are highly valued by PyTorch developers and have raised heated discussions. During the three months from our first bug issue (November 23, 2021) to March 1, 2022, there have been a total of 170 high-priority issues in the entire PyTorch issue-tracking system [11], of which we contributed 23 (13.5%), emphasizing the effectiveness of DeepREL.

We next present some example bugs detected by DeepREL:

**Out-of-Bounds Read (Equivalence<sub>value</sub>, ✓)** The bug shown in Figure 10 is detected in API pair `torch.kthvalue(tensor, k)` and `tensor.kthvalue(k)` via Equivalence<sub>value</sub>. The latter API is a method of `torch.Tensor` [5], which is the fundamental class in PyTorch. Given the same tensor (`input`) and argument (`k`), the returned results are not equal (Line 4). After debugging, we find that the returned results can actually be different across runs, indicating that it reads values from memory locations outside of user-controlled

```

1 input = torch.tensor([0,1,2,3,4])
2 result1 = torch.kthvalue(input,k=6)
3 result2 = input.kthvalue(k=6)
4 torch.testing.assert_close(result1, result2) # fail

```

Figure 10: Out-of-Bound Read in `torch.kthvalue`

```

1 output_size = [-36, 0, 0]
2 tensor = torch.rand([4, 4, 128, 2048, 4])
3 result1 = torch.nn.AdaptiveAvgPool3d(output_size)(tensor) # no exception
4 result2 = torch.nn.AdaptiveMaxPool3d(output_size)(tensor) # RuntimeError

```

Figure 11: Inconsistent check for `torch.nn.AdaptiveAvgPool3d`

data! This bug is a silent error, and has a severe security implication: without proper range checking of `k`, users may be able to read data outside of the allocated memory bounds (i.e., out-of-bound read). This bug is labeled as “high priority” and fixed immediately.

**Inconsistent Check (Equivalence<sub>status</sub>, ✓)** Figure 11 shows a bug for `AdaptiveAvgPool3d` and `AdaptiveMaxPool3d` in PyTorch (found via Equivalence<sub>status</sub>). With exactly the same input tensor (`tensor`) and argument (`output_size`), `AdaptiveAvgPool3d` runs without exception (Line 3) but `AdaptiveMaxPool3d` throws `RuntimeError`: Trying to create tensor with negative dimension (Line 4). After inspection, we find that `AdaptiveAvgPool3d` lacks checking for dimensions with negative integers. This can be disastrous since users may inadvertently introduce a bug into their model, but no warning/exception is raised. This bug has also been fixed.

**Wrong Computation (Equivalence<sub>value</sub>, ✗)** The bug shown in Figure 12 is detected by Equivalence<sub>value</sub> oracle for API pair `torch.std_mean` and `torch.mean`. Both of them could compute the mean of all elements in the input tensor. However, DeepREL has found that for certain input tensors, `torch.std_mean` has a different returned value from `torch.mean`. However, the developers said these two APIs are not expected to output the same mean and rejected

```

1 input = torch.tensor([[0.5786, 0.1719, 0.3760, 0.2939, 0.3984],
2 [0.5361, 0.7104, 0.8765, 0.0903, 0.0483]], dtype=float16)
3 result1 = torch.std_mean(input) # 0.4080
4 result2 = torch.mean(input) # 0.4082

```

**Figure 12: Inconsistency for `torch.std_mean` and `torch.mean`**

this bug report. Note that the other 6 rejected bugs are similar (there are indeed inconsistencies, but developers feel unnecessary to fix).

## 6.5 Threats to Validity

The threats to internal validity mainly lie in the correctness of the DeepREL implementation. To reduce such threats, the first three authors have performed extensive testing and code review of DeepREL. Moreover, we have released our code/scripts on our project website for public review [6]. The threats to external validity mainly lie in the evaluation benchmarks used. To demonstrate the generalizability of DeepREL, we have evaluated DeepREL on two widely used DL libraries, TensorFlow and PyTorch. Last but not least, the threats to construct validity mainly lie in the metrics used. To reduce such threats, we adopt the number of covered APIs and detected bugs, following prior work on DL library testing [54, 65, 66]. Moreover, we have also included false positive analysis.

## 7 RELATED WORK

We have talked about related work on DL library testing in Section 2. Thus, in this section, we mainly focus on existing work targeting API relations for other software systems. API mappings refer to the process of finding equivalent APIs between different libraries or programming languages. Prior work [49, 50] proposed to train neural networks to mine API mappings by learning transformations between the vector spaces of APIs from different languages (e.g., Java and C#). Another work [18] leveraged an unsupervised domain adaptation approach to automatically construct and align vector spaces for identifying API mappings with much less human effort. Our work is different from all such work in that DeepREL focuses on fuzzing while the aforementioned work targets code migration. Also, DeepREL utilizes not only equivalent APIs, but also other relational APIs for fuzzing. Moreover, DeepREL obtains relational APIs purely from documentation, and then verifies the API relations dynamically; thus, DeepREL does not need extensive API usages or training data of existing API pairs (which can be hard to obtain).

Function synonyms [21] are functions that play a similar (*not* necessarily semantically equivalent) role in code. UC-KLEE [56] leverages symbolic execution to check two different implementations/versions of the same function in C open source libraries. Func2vec [21] is a technique for finding function synonyms in Linux file systems and drivers via learning function embeddings to improve the quality of error handling specifications for Linux code. More recently, FPDiff [64] aims to automatically identify synonymous functions across multiple numerical libraries, and performs differential testing to detect discrepancies in these functions. While FPDiff is closely related, it mainly considers synonymous APIs with the same functionalities and argument lists across *different* libraries and is applied when there exist libraries with close design; in contrast, DeepREL is more general – it considers arbitrary relational APIs with value/status equivalence within *any given* library. Also, FPDiff can only target numerical library APIs taking double or int variables as input, whereas the input of APIs in DL libraries can be

much more complex (e.g., tensors of different dimensions/types, objects, strings, etc.). Moreover, FPDiff is a purely dynamic technique that finds synonymous functions via running them on an *elementary* set of double/int values, while it is impossible to find such a universal elementary input set for DL library APIs. Lastly, our study has found that DL library documentations contain various valuable information and shown for the first time that there can be plenty of equivalent/similar APIs within a given DL library, which can substantially help with fuzzing DL libraries (and beyond).

Throughout this paper, we have viewed each API as a test object, and thus our technique of fuzzing relational API pairs falls into the differential testing category. Meanwhile, if we treat the entire systems under test (e.g., TensorFlow or PyTorch) as test objects, this work can also be viewed as an instance/extension of the traditional metamorphic oracle generation work, which aims to automatically infer metamorphic relations as test oracle [16, 26, 33, 43, 63, 75, 83, 84]. For example, SBES [26, 43] synthesizes sequences of method invocations that are equivalent to a target method based on the observed dynamic program behaviors. The recent state-of-the-art technique MeMo [16] mainly targets mature open-source Java projects and automatically derives metamorphic relations from natural language specifications. Meanwhile, MeMo can only infer metamorphic relations that are well documented and lacks an invocation synthesizer or API relation verifier; thus the effectiveness of MeMo strictly depends on documentation completeness, correctness, and preciseness. In fact, very few metamorphic relations in DL libraries are explicitly described with informative sentences (with code templates) and can be directly translated into valid oracles. We substituted the first three phases of DeepREL with MeMo, and found that it can infer only 17/28 API relations for TensorFlow/PyTorch (compared to 8808/4290 relations found by DeepREL), and can at most detect 8 (including 6 found by MeMo and 2 found by FreeFuzz on the APIs involved in the MeMo inferred relations) of the 106 confirmed bugs detected by DeepREL.

## 8 CONCLUSION

We have introduced DeepREL, the first fully automated end-to-end approach to fuzzing DL libraries via inferring relational APIs. DeepREL can “borrow” test inputs from any API to test its relational APIs, and can leverage relational APIs as reference implementations for performing differential testing. The extensive study of DeepREL on PyTorch and TensorFlow shows that DeepREL is able to detect 162 bugs in total, with 106 already confirmed by the developers as previously unknown bugs. Notably, DeepREL has detected 13.5% of the high-priority bugs for the entire PyTorch issue-tracking system in a three-month period. Also, besides the 162 code bugs, we were also able to detect 14 documentation bugs (all confirmed).

## ACKNOWLEDGMENTS

We appreciate the insightful comments from the anonymous reviewers. This work was partially supported by National Science Foundation under Grant Nos. CCF-2131943 and CCF-2141474, as well as the Google Research Scholar Program.

## REFERENCES

- [1] 2015. Keras. <https://keras.io>.
- [2] 2021. bs4. <https://pypi.org/project/beautifulsoup4/>.
- [3] 2021. Definition of torch.broadcast\_shapes from Pytorch official documentation. [https://pytorch.org/docs/stable/generated/torch.broadcast\\_shapes.html](https://pytorch.org/docs/stable/generated/torch.broadcast_shapes.html).
- [4] 2021. News. [https://www.vice.com/en\\_us/article/9kga85/uber-is-giving-up-on-self-driving-cars-in-california-after-deadly-crash](https://www.vice.com/en_us/article/9kga85/uber-is-giving-up-on-self-driving-cars-in-california-after-deadly-crash).
- [5] 2022. Class torch.Tensor from PyTorch official documentation. <https://pytorch.org/docs/stable/tensors.html#torch.Tensor>.
- [6] 2022. DeepREL Repository. <https://github.com/ise-uic/DeepREL>.
- [7] 2022. Definition of tf.parallel\_stack from TensorFlow official documentation. [https://www.tensorflow.org/api\\_docs/python/tf\\_parallel\\_stack](https://www.tensorflow.org/api_docs/python/tf_parallel_stack).
- [8] 2022. Module tf.compat from TensorFlow official documentation. [https://www.tensorflow.org/api\\_docs/python/tf\\_compat](https://www.tensorflow.org/api_docs/python/tf_compat).
- [9] 2022. Module tf.raw\_ops from TensorFlow official documentation. [https://www.tensorflow.org/api\\_docs/python/tf\\_raw\\_ops](https://www.tensorflow.org/api_docs/python/tf_raw_ops).
- [10] 2022. munkres. <https://pypi.org/project/munkres/>.
- [11] 2022. PyTorch Repository. <https://github.com/pytorch/pytorch>.
- [12] 2022. SentenceTransformer. <https://www.sbert.net/>.
- [13] Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*. 265–283.
- [14] Naveed Akhtar and Ajmal Mian. 2018. Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access* 6 (2018), 14410–14430.
- [15] Guy Amir, Haoze Wu, Clark Barrett, and Guy Katz. 2020. An SMT-Based Approach for Verifying Binarized Neural Networks. *arXiv preprint arXiv:2011.02948* (2020).
- [16] Arianna Blasi, Alessandra Gorla, Michael D Ernst, Mauro Pezzè, and Antonio Carzaniga. 2021. MeMo: Automatically identifying metamorphic relations in Javadoc comments for test automation. *Journal of Systems and Software* 181 (2021), 111041. <https://doi.org/10.1016/j.jss.2021.111041>
- [17] Marcel Böhme, Van-Thuan Pham, and Abhik Roychoudhury. 2017. Coverage-based greybox fuzzing as markov chain. *IEEE Transactions on Software Engineering* 45, 5 (2017), 489–506. <https://doi.org/10.1145/2976749.2978428>
- [18] Nghi Bui. 2019. Towards zero knowledge learning for cross language API mappings. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*. IEEE, 123–125. <https://doi.org/10.1109/ICSE-Companion.2019.00054>
- [19] Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. 2019. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705* (2019).
- [20] Tsong Y Chen, Shing C Cheung, and Shiu Ming Yiu. 2020. Metamorphic testing: a new approach for generating next test cases. *arXiv preprint arXiv:2002.12543* (2020).
- [21] Daniel DeFreez, Aditya V Thakur, and Cindy Rubio-González. 2018. Path-based function embedding and its application to error-handling specification mining. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 423–433. <https://doi.org/10.1145/3236024.3236059>
- [22] Alastair F Donaldson, Hugues Evrard, Andrei Lascu, and Paul Thomson. 2017. Automated testing of graphics shader compilers. *Proceedings of the ACM on Programming Languages* 1, OOPSLA (2017), 1–29. <https://doi.org/10.1145/3133917>
- [23] Andreas Eitel, Jost Tobias Springenberg, Luciano Spinello, Martin Riedmiller, and Wolfram Burgard. 2015. Multimodal deep learning for robust RGB-D object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 681–687. <https://doi.org/10.1109/IROS.2015.7353446>
- [24] Felix A Gers, Jürgen Schmidhuber, and Fred Cummins. 2000. Learning to forget: Continual prediction with LSTM. *Neural computation* 12, 10 (2000), 2451–2471.
- [25] Alberto Goffi, Alessandra Gorla, Michael D Ernst, and Mauro Pezzè. 2016. Automatic generation of oracles for exceptional behaviors. In *Proceedings of the 25th international symposium on software testing and analysis*. 213–224. <https://doi.org/10.1145/2931037.2931061>
- [26] Alberto Goffi, Alessandra Gorla, Andrea Mattavelli, Mauro Pezzè, and Paolo Tonella. 2014. Search-based synthesis of equivalent method sequences. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*. 366–376. <https://doi.org/10.1145/2635868.2635888>
- [27] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014).
- [28] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*. Ieee, 6645–6649.
- [29] Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. 2018. Deep code search. In *2018 IEEE/ACM 40th International Conference on Software Engineering (ICSE)*. IEEE, 933–944. <https://doi.org/10.1145/3180155.3180167>
- [30] Qianyu Guo, Xiaofei Xie, Yi Li, Xiaoyu Zhang, Yang Liu, Xiaohong Li, and Chao Shen. 2020. Audee: Automated testing for deep learning frameworks. In *2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 486–498. <https://doi.org/10.1145/3324884.3416571>
- [31] Fabrice Harel-Canada, Lingxiao Wang, Muhammad Ali Gulzar, Quanquan Gu, and Miryung Kim. 2020. Is neuron coverage a meaningful measure for testing deep neural networks? In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 851–862. <https://doi.org/10.1145/3368089.3409754>
- [32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR.2016.90>
- [33] Upalée Kanewala. 2014. Techniques for automatic detection of metamorphic relations. In *2014 IEEE Seventh International Conference on Software Testing, Verification and Validation Workshops*. IEEE, 237–238. <https://doi.org/10.1109/ICSTW.2014.62>
- [34] Jinhan Kim, Robert Feldt, and Shin Yoo. 2019. Guiding deep learning system testing using surprise adequacy. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 1039–1049. <https://doi.org/10.1109/ICSE.2019.00108>
- [35] Ian Lenz, Honglak Lee, and Ashutosh Saxena. 2015. Deep learning for detecting robotic grasps. *The International Journal of Robotics Research* 34, 4–5 (2015), 705–724. <https://doi.org/10.1177/0278364914549607>
- [36] Xia Li, Wei Li, Yuqun Zhang, and Lingming Zhang. 2019. Deepfl: Integrating multiple fault diagnosis dimensions for deep fault localization. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 169–180. <https://doi.org/10.1145/3293882.3330574>
- [37] Yi Li, Shaohua Wang, and Tien N Nguyen. 2020. Dlfix: Context-based code transformation learning for automated program repair. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 602–614. <https://doi.org/10.1145/3377811.3380345>
- [38] Changliu Liu, Tomer Arnon, Christopher Lazarus, Christopher Strong, Clark Barrett, and Mykel J Kochenderfer. 2019. Algorithms for verifying deep neural networks. *arXiv preprint arXiv:1903.06758* (2019).
- [39] Jiawei Liu, Jinkun Lin, Fabian Ruffy, Cheng Tan, Jinyang Li, Aurojit Panda, and Lingming Zhang. 2022. Finding Deep-Learning Compilation Bugs with NNSmith. *arXiv preprint arXiv:2207.13066* (2022).
- [40] Jiawei Liu, Yuxiang Wei, Sen Yang, Yinlin Deng, and Lingming Zhang. 2022. Coverage-Guided Tensor Compiler Fuzzing with Joint IR-Pass Mutation. *Proc. ACM Program. Lang.* 6, OOPSLA1, Article 73 (apr 2022), 26 pages. <https://doi.org/10.1145/3527317>
- [41] Lei Ma, Felix Juefei-Xu, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Chunyang Chen, Ting Su, Li Li, Yang Liu, et al. 2018. Deepgauge: Multi-granularity testing criteria for deep learning systems. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. 120–131. <https://doi.org/10.1145/3238147.3238202>
- [42] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083* (2017).
- [43] Andrea Mattavelli, Alberto Goffi, and Alessandra Gorla. 2015. Synthesis of equivalent method calls in Guava. In *International Symposium on Search Based Software Engineering*. Springer, 248–254.
- [44] William M McKeeman. 1998. Differential testing for software. *Digital Technical Journal* 10, 1 (1998), 100–107.
- [45] Ali Mesbah, Andrew Rice, Emily Johnston, Nick Glorioso, and Edward Aftandilian. 2019. Deepdelta: learning to repair compilation errors. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 925–936. <https://doi.org/10.1145/3338906.3340455>
- [46] Seonwoo Min, Byunghan Lee, and Sungroh Yoon. 2017. Deep learning in bioinformatics. *Briefings in bioinformatics* 18, 5 (2017), 851–869.
- [47] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. 2016. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2574–2582. <https://doi.org/10.1109/CVPR.2016.282>
- [48] James Munkres. 1957. Algorithms for the assignment and transportation problems. *Journal of the society for industrial and applied mathematics* 5, 1 (1957), 32–38.
- [49] Trong Duc Nguyen, Anh Tuan Nguyen, and Tien N Nguyen. 2016. Mapping API elements for code migration with vector representations. In *2016 IEEE/ACM 38th International Conference on Software Engineering Companion (ICSE-C)*. IEEE, 756–758. <https://doi.org/10.1145/2889160.2892661>
- [50] Trong Duc Nguyen, Anh Tuan Nguyen, Hung Dang Phan, and Tien N Nguyen. 2017. Exploring API embedding for API usages and applications. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*. IEEE, 438–449. <https://doi.org/10.1109/ICSE.2017.47>
- [51] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. 2016. The limitations of deep learning in adversarial

- settings. In *2016 IEEE European symposium on security and privacy (EuroS&P)*. IEEE, 372–387. <https://doi.org/10.1109/EuroSP.2016.36>
- [52] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019), 8026–8037.
- [53] Kexin Pei, Yinzheng Cao, Junfeng Yang, and Suman Jana. 2017. Deepxplore: Automated whitebox testing of deep learning systems. In *proceedings of the 26th Symposium on Operating Systems Principles*, 1–18. <https://doi.org/10.1145/3361566>
- [54] Hung Viet Pham, Thibaud Lutellier, Weizhen Qi, and Lin Tan. 2019. CRADLE: Cross-Backend Validation to Detect and Localize Bugs in Deep Learning Libraries. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*, 1027–1038. <https://doi.org/10.1109/ICSE.2019.00017>
- [55] Martin F Porter. 2001. Snowball: A language for stemming algorithms.
- [56] David A Ramos and Dawson R Engler. 2011. Practical, low-effort equivalence verification of real code. In *International Conference on Computer Aided Verification*. Springer, 669–685.
- [57] Daniele Ravi, Charence Wong, Fani Deligianni, Melissa Berthelot, Javier Andreu-Perez, Benny Lo, and Guang-Zhong Yang. 2016. Deep learning for health informatics. *IEEE journal of biomedical and health informatics* 21, 1 (2016), 4–21. <https://doi.org/10.1109/JBHI.2016.2636665>
- [58] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 3982–3992. <https://doi.org/10.18653/v1/D19-1410>
- [59] Dongdong She, Rahul Krishna, Lu Yan, Suman Jana, and Baishakhi Ray. 2020. MTFuzz: fuzzing with a multi-task neural network. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 737–749. <https://doi.org/10.1145/3368089.3409723>
- [60] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [61] Ting Su, Yichen Yan, Jue Wang, Jingling Sun, Yiheng Xiong, Geguang Pu, Ke Wang, and Zhendong Su. 2021. Fully automated functional fuzzing of Android apps for detecting non-crashing logic bugs. *Proceedings of the ACM on Programming Languages* 5, OOPSLA (2021), 1–31. <https://doi.org/10.1145/3485533>
- [62] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. 2018. DeepTest: Automated testing of deep-neural-network-driven autonomous cars. In *Proceedings of the 40th international conference on software engineering*, 303–314. <https://doi.org/10.1145/3180155.3180220>
- [63] Javier Troya, Sergio Segura, and Antonio Ruiz-Cortés. 2018. Automated inference of likely metamorphic relations for model transformations. *Journal of Systems and Software* 136 (2018), 188–208. <https://doi.org/10.1016/j.jss.2017.05.043>
- [64] Jackson Vanover, Xuan Deng, and Cindy Rubio-González. 2020. Discovering discrepancies in numerical libraries. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 488–501. <https://doi.org/10.1145/3395363.3397380>
- [65] Zan Wang, Ming Yan, Junjie Chen, Shuang Liu, and Dongdi Zhang. 2020. Deep learning library testing via effective model generation. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 788–799. <https://doi.org/10.1145/3368089.3409761>
- [66] Anjiang Wei, Yinlin Deng, Chenyuan Yang, and Lingming Zhang. 2022. Free Lunch for Testing: Fuzzing Deep-Learning Libraries from Open Source. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*, 995–1007. <https://doi.org/10.1145/3510003.3510041>
- [67] Cheng Wen, Haijun Wang, Yuekang Li, Shengchao Qin, Yang Liu, Zhiwu Xu, Hongxu Chen, Xiaofei Xie, Geguang Pu, and Ting Liu. 2020. Memlock: Memory usage guided fuzzing. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*, 765–777. <https://doi.org/10.1145/3377811.3380396>
- [68] Martin White, Michele Tufano, Christopher Vandome, and Denys Poshyvanyk. 2016. Deep learning code fragments for code clone detection. In *2016 31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 87–98. <https://doi.org/10.1145/2970276.2970326>
- [69] Wikipedia contributors. 2022. Backpropagation. <https://en.wikipedia.org/wiki/Backpropagation>.
- [70] Wikipedia contributors. 2022. Levenshtein distance. [https://en.wikipedia.org/wiki/Levenshtein\\_distance](https://en.wikipedia.org/wiki/Levenshtein_distance).
- [71] Wikipedia contributors. 2022. Maximum weight matching. [https://en.wikipedia.org/wiki/Maximum\\_weight\\_matching](https://en.wikipedia.org/wiki/Maximum_weight_matching).
- [72] Wikipedia contributors. 2022. Supervised Learning. [https://en.wikipedia.org/wiki/Supervised\\_learning](https://en.wikipedia.org/wiki/Supervised_learning).
- [73] Wikipedia contributors. 2022. Tf-idf. <https://en.wikipedia.org/wiki/Tf%E2%80%93idf>.
- [74] Wikipedia contributors. 2022. Topological Sorting. [https://en.wikipedia.org/wiki/Topological\\_sorting](https://en.wikipedia.org/wiki/Topological_sorting).
- [75] Zhenglong Xiang, Hongrun Wu, and Fei Yu. 2019. A genetic algorithm-based approach for composite metamorphic relations construction. *Information* 10, 12 (2019), 392. <https://doi.org/10.3390/info10120392>
- [76] Danning Xie, Yitong Li, Mijung Kim, Hung Viet Pham, Lin Tan, Xiangyu Zhang, and Michael W Godfrey. 2022. DocTer: documentation-guided fuzzing for testing deep learning API functions. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis (Virtual, South Korea) (ISSTA 2022)*. Association for Computing Machinery, New York, NY, USA, 176–188. <https://doi.org/10.1145/3533767.3534220>
- [77] Shenao Yan, Guanhong Tao, Xuwei Liu, Juan Zhai, Shiqing Ma, Lei Xu, and Xiangyu Zhang. 2020. Correlations between deep neural network model coverage criteria and model quality. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 775–787. <https://doi.org/10.1145/3368089.3409671>
- [78] Yanming Yang, Xin Xia, David Lo, and John Grundy. 2020. A survey on deep learning for software engineering. *arXiv preprint arXiv:2011.14597* (2020).
- [79] Guixin Ye, Zhanyong Tang, Shin Hwei Tan, Songfang Huang, Dingyi Fang, Xiaoyang Sun, Lizhong Bian, Haibo Wang, and Zheng Wang. 2021. Automated conformance testing for JavaScript engines via deep compiler fuzzing. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, 435–450. <https://doi.org/10.1145/3453483.3454054>
- [80] Wei You, Xueqiang Wang, Shiqing Ma, Jianjun Huang, Xiangyu Zhang, XiaoFeng Wang, and Bin Liang. 2019. Profuzzer: On-the-fly input type probing for better zero-day vulnerability discovery. In *2019 IEEE symposium on security and privacy (SP)*. IEEE, 769–786. <https://doi.org/10.1109/SP.2019.00057>
- [81] Zhengran Zeng, Han Zhuo Tan, Haotian Zhang, Jing Li, Yuqun Zhang, and Lingming Zhang. 2022. An Extensive Study on Pre-trained Models for Program Understanding and Generation. In *Proceedings of the ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. To appear.
- [82] Zhengran Zeng, Yuqun Zhang, Haotian Zhang, and Lingming Zhang. 2021. Deep just-in-time defect prediction: how far are we?. In *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 427–438.
- [83] Bo Zhang, Hongyu Zhang, Junjie Chen, Dan Hao, and Pablo Moscato. 2019. Automatic discovery and cleansing of numerical metamorphic relations. In *2019 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 235–245. <https://doi.org/10.1109/ICSME.2019.00035>
- [84] Jie Zhang, Junjie Chen, Dan Hao, Yingfei Xiong, Bing Xie, Lu Zhang, and Hong Mei. 2014. Search-based inference of polynomial metamorphic relations. In *Proceedings of the 29th ACM/IEEE international conference on Automated software engineering*, 701–712. <https://doi.org/10.1145/2642937.2642994>
- [85] Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. 2018. DeepRoad: GAN-based metamorphic testing and input validation framework for autonomous driving systems. In *2018 33rd IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 132–142. <https://doi.org/10.1145/3238147.3238187>
- [86] Xufan Zhang, Ning Sun, Chunrong Fang, Jiawei Liu, Jia Liu, Dong Chai, Jiang Wang, and Zhenyu Chen. 2021. Predoo: precision testing of deep learning operators. In *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 400–412. <https://doi.org/10.1145/3460319.3464843>
- [87] Sizhe Zheng, Renze Chen, Yicheng Jin, Anjiang Wei, Binyang Wu, Xiuhong Li, Shengen Yan, and Yun Liang. 2021. NeoFlow: A Flexible Framework for Enabling Efficient Compilation for High Performance DNN Training. *IEEE Transactions on Parallel and Distributed Systems* (2021). <https://doi.org/10.1109/TPDS.2021.3138862>
- [88] Husheng Zhou, Wei Li, Zelun Kong, Junfeng Guo, Yuqun Zhang, Bei Yu, Lingming Zhang, and Cong Liu. 2020. Deepbillboard: Systematic physical-world testing of autonomous driving systems. In *2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE)*. IEEE, 347–358. <https://doi.org/10.1145/3377811.3380422>