

# Structure-Invariant Testing for Machine Translation

PINJIA HE, ETH Zurich, Switzerland

CLARA MEISTER, ETH Zurich, Switzerland

ZHENDONG SU, ETH Zurich, Switzerland

In recent years, machine translation software has increasingly been integrated into our daily lives. People routinely use machine translation for various applications, such as signing lease agreements when studying abroad, describing symptoms to a foreign doctor, and reading political news in a foreign language. However, due to the complexity and intractability of neural machine translation (NMT) models that power modern machine translation systems, these systems are far from being robust. They can return inferior results that lead to misunderstanding, embarrassment, financial loss, medical misdiagnoses, threats to personal safety, or political conflicts. Despite its apparent importance, validating the robustness of machine translation is very difficult and has, therefore, been much under-explored.

To tackle this challenge, we introduce *structure-invariant testing (SIT)*, a novel, widely applicable metamorphic testing methodology for validating machine translation software. Our key insight is that the translation results of similar source sentences should typically exhibit a similar sentence structure. SIT is designed to leverage this insight to test any machine translation system with unlabeled sentences; it specifically targets mistranslations that are difficult-to-find using state-of-the-art translation quality metrics such as BLEU. We have realized a practical implementation of SIT by (1) substituting one word in a given sentence with semantically similar, syntactically equivalent words to generate similar sentences, and (2) using syntax parse trees (obtained via constituency or dependency parsing) to represent sentence structure. To evaluate SIT, we have used it to test Google Translate and Bing Microsoft Translator with 200 unlabeled sentences as input, which led to 64 and 70 buggy translations with 69.5% and 70% top-1 accuracy, respectively. The bugs are diverse, including under-translation, over-translation, incorrect modification, word/phrase mistranslation, and unclear logic, none of which could be detected via existing translation quality metrics.

Beyond testing machine translation, SIT can be adapted to validate sequence-to-text AI software (e.g., figure captioning and speech recognition) in general. This work opens up this exciting, important direction.

Additional Key Words and Phrases: Testing, Machine translation, Structural invariance, Metamorphic testing

## 1 INTRODUCTION

Neural machine translation (NMT) models have driven recent advances in machine translation. As reported by research from Google [Wu et al. 2016] and Microsoft [Hassan et al. 2018], state-of-the-art NMT models are approaching human-level performance in terms of accuracy, i.e., BLEU [Papineni et al. 2002]. Because of these recent breakthroughs, machine translation software (e.g., Google Translate<sup>1</sup> and Bing Microsoft Translator<sup>2</sup>) become indispensable for many in daily life.

In particular, Google Translate, which was launched in 2006, is currently the most widely-used online translation service. In 2016, Google Translate attracted more than 500 million users and translated more than 100 billion words per day [Turovsky 2016]. Machine translation services are also embedded in various software applications, such as Facebook [Facebook 2019] and Twitter [Twitter 2019]. People use machine translation for various applications, such as signing lease agreements

<sup>1</sup><https://translate.google.com/>

<sup>2</sup><https://www.bing.com/translator>

Authors' addresses: Pinjia He, Department of Computer Science, ETH Zurich, Switzerland, pinjia.he@inf.ethz.ch; Clara Meister, Department of Computer Science, ETH Zurich, Switzerland, clara.meister@inf.ethz.ch; Zhendong Su, Department of Computer Science, ETH Zurich, Switzerland, zhendong.su@inf.ethz.ch.

Source sentence	Google Translate result	Target sentence meaning
I live on campus with <a href="#">smart</a> people.	我和聪明的人住在校园里。	I live on campus with smart people. ✓
I live on campus with <a href="#">cute</a> people.	我和可爱的人住在校园里。	I live on campus with cute people. ✓
I live on campus with <a href="#">tall</a> people.	我住在校园里, 身材高大。	I live on campus, I am tall. ✗

Fig. 1. Examples of similar source sentences and Google Translate results.

when studying abroad, describing symptoms to a foreign doctor, and reading political news in a foreign language. Thus, the reliability of machine translation software is of great importance. Unreliable machine translation can return inferior results (i.e., sub-optimal or incorrect translations), leading to misunderstanding, embarrassment, financial loss, medical misdiagnoses, threats to personal safety, or political conflicts [Davies 2017; Macdonald 2015; Okrent 2016; Ong 2017].

However, as the core component of modern machine translation software, NMT models are not as reliable as many may believe. Recently, sub-optimal and incorrect outputs have been found in various software systems with neural networks as their core components [Zhang et al. 2019]. Typical examples include autonomous cars [Evtimov et al. 2018; Pei et al. 2017; Tian et al. 2018], sentiment analysis tools [Alzantot et al. 2018; Iyyer et al. 2018; Li et al. 2019], and speech recognition services [Carlini et al. 2016]. These recent research efforts show that neural networks can easily return inferior results (e.g., wrong class labels) given specially-crafted inputs (i.e., adversarial examples).

NMT models are no exception; they can be fooled by adversarial examples [Ebrahimi et al. 2018] or natural noise (e.g., typos in input sentences) [Belinkov and Bisk 2018]. These examples and noise are found via heuristic rules, such as reordering the characters in a word. Thus, inputs found are mostly "illegal", such as sentences with syntax errors or obvious misspellings that are unlikely given as input. Moreover, as reported by WeChat, a messenger app with over one billion monthly active users, its embedded NMT model sometimes returns inferior results even when the input sentences are syntactically correct [Zheng et al. 2018].

In practice, the typical testing procedure for machine translation software involves three steps [Zheng et al. 2018]: (1) collect bilingual sentence pairs<sup>3</sup> and split them into training, validation, and testing data; (2) calculate translation quality scores (e.g., BLEU [Papineni et al. 2002] and ROUGE [Lin 2004]) of the trained NMT model on the testing data; and (3) compare the scores with predefined thresholds to determine whether the test cases pass. Although widely adopted, this testing procedure has two major limitations. First, while most of the bilingual sentence pairs are used as training data, the testing oracle size is often small, leading to insufficient and ineffective testing. Second, the calculation of translation quality scores (e.g., BLEU) requires bilingual sentence pairs as input, which need to be manually constructed beforehand. However, typically when people use machine translation software, they have no idea what the correct translation should be. To test using real-world user input sentences, which are mostly not in the prepared bilingual dataset, extensive manual effort is needed to translate these sentences. Thus, an effective and efficient testing methodology that can automatically and accurately detect translation bugs<sup>4</sup> in machine translation software is in high demand.

However, testing machine translation software is extremely challenging. First, different from traditional software whose logic is encoded in source code, machine translation software is based

<sup>3</sup>By a sentence pair, we refer to a source sentence and its corresponding target sentence.

<sup>4</sup>By a translation bug, we refer to mistranslation of some parts of a source sentence. The translated sentence (i.e., target sentence) containing translation bug(s) is regarded as a buggy sentence. One buggy sentence could contain multiple bugs. We use "bug in the target sentence" and "bug in the sentence pair" interchangeably in this paper.

Bug type	Source sentence	Target sentence	Target sentence meaning
Under-translation	During a question and answer session, however, Nielsen added that other nation states have adopted a more visible approach <b>to doing the same</b> .	不过，在问答环节，尼尔森补充说，其他民族国家也采取了更明显的做法。 (by Bing)	During a question and answer session, however, Nielsen added that other nation states have adopted a more visible approach.
Over-translation	Entering talks, Brazil hoped to see itself elevated to major non NATO ally status by the Trump administration, a big step that would help it purchase military equipment.	进入谈判，巴西希望看到自己被特朗普政府提升为主要的非北约盟国地位，这是一个帮助其购买军事装备的一大步。 (by Bing)	Entering talks, Brazil hoped to see itself elevated to major non NATO ally status by the Trump administration, <b>one</b> a big step that would help it purchase military equipment.
Incorrect modification	But even so, they remain <b>members of privilege</b> .	但即便如此，他们仍然是特权的成员 (by Google)	But even so, they remain <b>privilege's members</b> .
Word/phrase mistranslation	I am very willing to <b>share</b> my point of view.	我非常愿意同意我的观点。 (by Bing)	I am very willing to <b>agree with</b> my point of view.
Unclear logic	I <b>had a joke to tell and</b> I wanted to finish it, Draper says.	德雷珀说，我开玩笑说，我想完成它。(by Google)	I <b>joked that</b> I want to finish it, Draper says.

Fig. 2. Examples of translation bugs detected by SIT.

on complex neural network structures and enormous amounts of training data. Thus, testing techniques for traditional software, which are mostly based on code, are ineffective. Second, the line of recent research on testing AI (artificial intelligence) software [Alzantot et al. 2018; Carlini et al. 2016; Goodfellow et al. 2015; Iyyer et al. 2018; Jia and Liang 2017; Kim et al. 2019; Li et al. 2019; Ma et al. 2018a; Mudrakarta et al. 2018; Pei et al. 2017; Tian et al. 2018] focuses on much simpler tasks with simple output formats. For example, most existing work aims at testing image classifiers, which output class labels given an image. However, as one of the most difficult natural language processing tasks, the output of machine translation systems (i.e., translated sentences) is significantly more complex. Because of its complexity, there is no effective automated approach to evaluating the correctness of translation. Thus, it is difficult to test machine translation software.

To address this problem, we introduce structure-invariant testing (SIT), a novel, widely-applicable methodology for validating machine translation software. The key insight is that similar source sentences typically have translation results of similar sentence structures. For example, Fig. 1 shows three similar source sentences in English and their target sentences (i.e., translated sentence) in Chinese. The first two translations are correct, while the third is not. We can observe that the structure of the third sentence in Chinese significantly differs from those of the other two. Based on the concept of structural invariance, for each unlabeled sentence, SIT (1) generates a list of its similar sentences by modifying a single word in the source sentence; (2) feeds the source sentence and the generated similar sentences to the machine translation software under test to obtain their translations; (3) uses specialized data structures to represent the syntax structure of each of the translated sentences; and (4) compares the structures of the translated sentences. If a large difference exists between the structures of the translated original and any of the translated modified sentences, we report the modified sentence pair and the original sentence pair because both of them could be buggy.

We apply SIT to test Google Translate and Bing Microsoft Translator with 200 unlabeled sentences crawled from the Web as input. SIT successfully found 64 buggy translations in Google Translate and 70 buggy translations in Bing Microsoft Translator with high accuracy (i.e., 69.5% and 70% top-1 accuracy respectively). The reported bugs are diverse, including under-translation, over-translation,

incorrect modification, word/phrase mistranslation, and unclear logic, none of which could be detected by the widely-used metrics BLEU and ROUGE. Examples of different translation bugs are illustrated in Fig. 2. In addition, all the translation bugs uncovered by SIT will be uploaded online<sup>5</sup> to facilitate the validation and reproduction of our results.

This paper makes the following main contributions:

- It introduces structure-invariant testing (SIT), a novel, widely applicable methodology for validating machine translation software;
- It describes a practical implementation of SIT by adapting BERT [Devlin et al. 2018] to generate similar sentences and leveraging syntax parsers to represent sentence structures;
- It presents the evaluation of SIT using only 200 unlabeled sentences crawled from the Web to successfully find 64 buggy translations in Google Translate and 70 buggy translations in Bing Microsoft Translator with high accuracy; and
- It discusses the diverse bugs found by SIT, including under-translation, over-translation, incorrect modification, word/phrase mistranslation, and unclear logic, of which none could be found by state-of-the-art metrics (i.e., BLEU and ROUGE).

## 2 A REAL-WORLD EXAMPLE

Tom planned to take his son David, who is 14 years old, to the Zurich Zoo. Before their zoo visit, he checked the zoo's website<sup>6</sup> on purchasing tickets and saw the following German sentence:

Kinder bis 15 Jahre erhalten an ihrem Geburtstag gegen Vorweisen eines gültigen Ausweises den Zoeeintritt geschenkt.

Tom is from the United States, and he does not understand German. To figure out the meaning of this sentence, Tom used Google Translate, which is a popular translation service powered by NMT models [Wu et al. 2016]. Google Translate returned the following English sentence:

Children up to the age of 15 are given free admission to the zoo on presentation of a valid ID.

However, David was denied entry by the zoo staff even with a valid ID. They found out that they had misunderstood the zoo's regulation because of the incorrect translation returned by Google Translate. The correct translation should be:

Children up to the age of 15 are given free admission to the zoo *on their birthday* on presentation of a valid ID.

This is a real translation bug that led to a confusing, unpleasant experience. Translation bugs could also cause extremely serious consequences. For example, in 2009, HSBC bank mistranslated its catchphrase "Assume Nothing" as "Do Nothing" in various countries [Macdonald 2015; Tree 2013]. Because of this translation bug, HSBC bank had to launch a \$10 million rebranding campaign to repair the damage. In another case in 2017, a Palestinian man was arrested by Israeli police for a post saying "good morning", which Facebook's machine translation service erroneously translated as "attack them" in Hebrew and "hurt them" in English [Davies 2017; Ong 2017]. To enhance the reliability of machine translation software, this paper introduces a general validation approach called structure-invariant testing, which automatically and accurately detects translation bugs without requiring any labeled data.

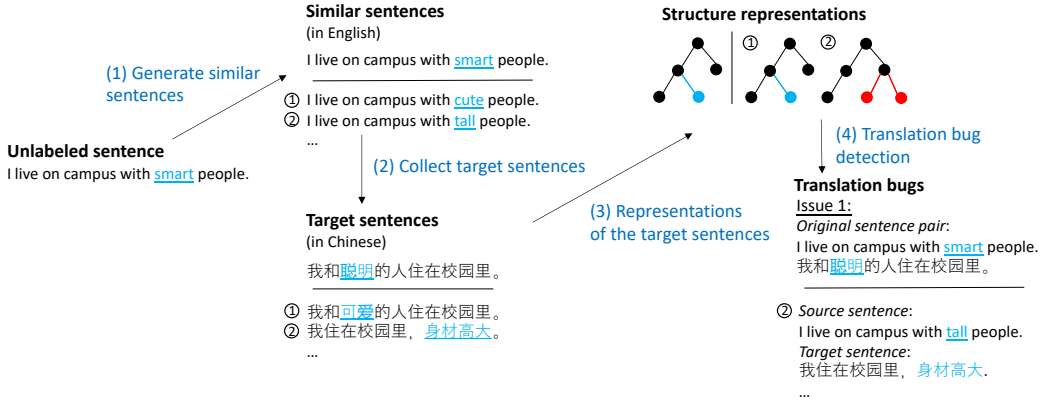


Fig. 3. Overview of our SIT implementation.

### 3 APPROACH AND IMPLEMENTATION

This section introduces structure-invariant testing (SIT) and describes our implementation. The input of SIT is a list of unlabeled, monolingual sentences, while its output is a list of suspicious *issues*. For each original sentence, SIT reports either 0 (i.e., no buggy sentence is found) or 1 *issue* (i.e., at least 1 buggy sentence is found). Each *issue* contains: (1) the original source sentence and its translation; and (2) top- $k$  generated source sentences (i.e. the  $k$  farthest translations from the source sentence translation) and their translations, which SIT thinks are buggy. The source sentence pair is reported for the following reasons: (1) seeing how the original sentence was modified may help the user understand why the translation system made a mistake (2) the bug may actually lie in the translation of the original sentence.

Fig. 3 illustrates the overview of SIT. In this figure, we use one unlabeled sentence as input for simplicity and clarity. The key insight of SIT is that similar source sentences have target sentences of similar syntactic structure. Derived from this insight, SIT carries out the following four steps:

- (1) *Generating similar sentences.* For each unlabeled sentence, we generate a list of its similar sentences by modifying a single word in the sentence.
- (2) *Collecting target sentences.* We feed the original and the generated similar sentences to the machine translation system under test and collect their target sentences (i.e. translations).
- (3) *Representing target sentence structures.* All the target sentences are encoded as data structures specialized for natural language processing.
- (4) *Detecting translation bugs.* The structures of the translated generated sentences are compared to the structure of the translated original sentence. If there is a large difference between the structures of a given translated generated sentence and the translated original sentence, SIT reports the generated sentence pair along with the original sentence pair. For each original sentence, SIT issues a report if at least one of the translated generated sentences demonstrates a large difference.

SIT can be implemented in various ways by choosing different (1) similar sentence generation strategies; (2) representations of target sentence; and (3) distance metrics between two sentence structure representations. We next introduce our realization of SIT using (1) BERT [Devlin et al.

<sup>5</sup><https://github.com/StructureInvariantTesting/StructureInvariantTesting>

<sup>6</sup><https://www.zoo.ch/de/zoobesuch/tickets-preise>

2018] to generate similar sentences; (2) three alternative structures to represent sentences; and (3) three distance metrics. We will also discuss relevant implementation details.

### 3.1 Generating Similar Sentences

In order to test for structural invariance, we must compare two sentences that have the same structure but differ in at least one token. To guarantee this structural similarity, we take an input sentence and modify it to produce a set of different sentences with the same structure. There are a number of approaches that can be taken to produce a list of variations from a single sentence. For example, we could attempt to paraphrase a sentence [Ganitkevitch et al. 2013] and produce some semantically equivalent versions. However, the output of this method is not guaranteed to be structurally similar to the input. We have found that changing one word in the sentence at a time effectively produce structurally identical and semantically similar sentences.

The approach we take modifies a single token, replacing it with another token of the same part of speech, to produce an alternate sentence. For example, we will mask "hairy" in the source sentence in Fig. 4 and replace it with the top- $k$  most similar tokens to generate  $k$  similar sentences. We do this for every candidate token in the sentence; for the sake of simplicity and avoiding strange grammatical phenomena, we only use nouns and adjectives as candidates for replacement. If there are  $m$  words in a sentence that can be replaced, our method will produce a list of  $k \cdot m$  alternate sentences.

Now we discuss the problem of selecting replacement tokens. Perhaps the simplest algorithm for selecting a set of replacement tokens would involve using word embeddings [Mikolov et al. 2013b]. One could choose words that have high vector similarity with and identical part-of-speech (POS)<sup>7</sup> tags to a given token in the original sentence as replacements in the modified sentences. However, since word embeddings have the same value regardless of context, this approach often produced sentences that would not occur in common language. For example, the word "fork" might have high vector similarity with and the same POS tag as the word "plate". However, while the sentence "He came to a fork in the road" makes sense, the sentence "He came to a plate in the road" does not.

Rather, we want a model that considers the surrounding words and comes up with a set of replacements that, when inserted, create realistic sentences. A model that does just this is the masked language model (MLM) [Mikolov et al. 2013a], inspired by the Cloze task [Taylor 1953]. The input to an MLM is a piece of text with a single token masked (i.e., deleted from the sentence and replaced with a special indicator token). The job of the model is then to predict the token in that position given the context. This method forces the model to learn the dependencies between different tokens. Since there are a number of different contexts a single word can fit in, this model, in a sense, allows for a single token to have multiple representations. We therefore get a set of replacement tokens that are context dependent. While the predicted tokens are not guaranteed to have the same meaning as the original token, if the MLM is good, it is highly likely that the sentence with the new, predicted token is both syntactically correct and meaningful.

For our implementation, we use an MLM to generate candidate replacement tokens for a candidate token in our original sentence. There are a number of out-of-the-box options for this task, such as ELMo, the OpenAI transformer, BERT, or training a customized MLM. While training our own MLM would have been quite possible, the above state-of-the-art options would undoubtedly do better as their architectures are incredibly specialized and the models have been trained on massive amounts of data, making them very robust. Since the generation of valid sentences is vital for avoiding false positives, we opt to outsource this task to one of the above options.

<sup>7</sup>[https://en.wikipedia.org/wiki/Part\\_of\\_speech](https://en.wikipedia.org/wiki/Part_of_speech)



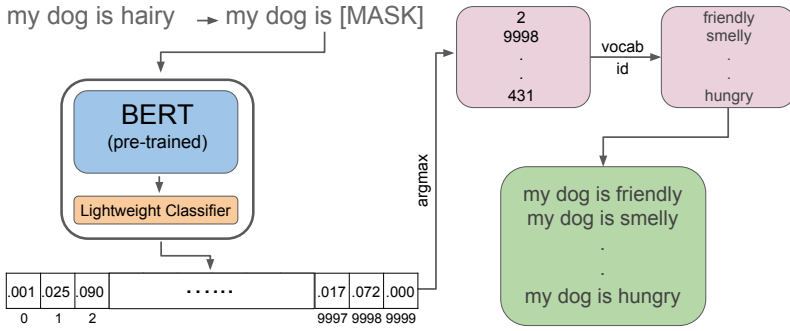


Fig. 4. Similar sentence generation process.

Specifically, we use BERT [Devlin et al. 2018], which is a state-of-the-art language representation model recently proposed and released by Google. BERT provides a robust base onto which a simple classification model can be added to perform NLP tasks. In more technical terms, the out-of-box BERT model provides language representations that can be fine-tuned by adding an additional lightweight softmax classification layer to create models for a wide range of language-related tasks, such as masked language modelling. BERT, which stands for Bidirectional Encoder Representations from Transformers, uses the neural network architecture known as *transformers*, an architecture also used in some state-of-the-art neural translation models [Edunov et al. 2018; Vaswani et al. 2017]. The language representations were trained on a huge amount of data; the corpus used for pre-training was a concatenation of BooksCorpus (800M words) and English Wikipedia (2,500M words), exposing the model to many different domains and writing styles. Additionally, the masked language task was one of two main tasks used to train the base language representation. Thus, we believe that BERT fits this aspect of our approach well. Using BERT as an MLM is incredibly straightforward; example code can be found in Google’s BERT repository.<sup>8</sup> An example of the sentence generation process is demonstrated in Fig. 4.

### 3.2 Collecting Target Sentences

SIT detects translation bugs by comparing the sentence structures of similar sentences after they have been fed through a translation system. Once we have generated a list of similar sentences from our original sentence, the next step is to input all our source sentences to the machine translation software under test and collect the corresponding translation results (i.e., target sentences). We subsequently analyze the results to find bugs.

We use Google’s and Bing’s machine translation systems as test systems for our experiment. To obtain translation results, we invoke the APIs provided by Google Translate<sup>9</sup> and Bing Microsoft Translator<sup>10</sup> which return identical results as their Web interfaces, respectively.

### 3.3 Representations of the Target Sentences

Now we must model the target sentences obtained from the translation system under test since, for SIT, we compare sentence structures post translation in order to detect bugs. Choosing the structure with which to represent our sentences will affect our ability to perform meaningful comparisons.

<sup>8</sup>[https://github.com/google-research/bert/blob/master/run\\_pretraining.py](https://github.com/google-research/bert/blob/master/run_pretraining.py)

<sup>9</sup><https://cloud.google.com/translate/docs/>

<sup>10</sup><https://docs.microsoft.com/en-us/azure/cognitive-services/translator/quickstart-python-translate>

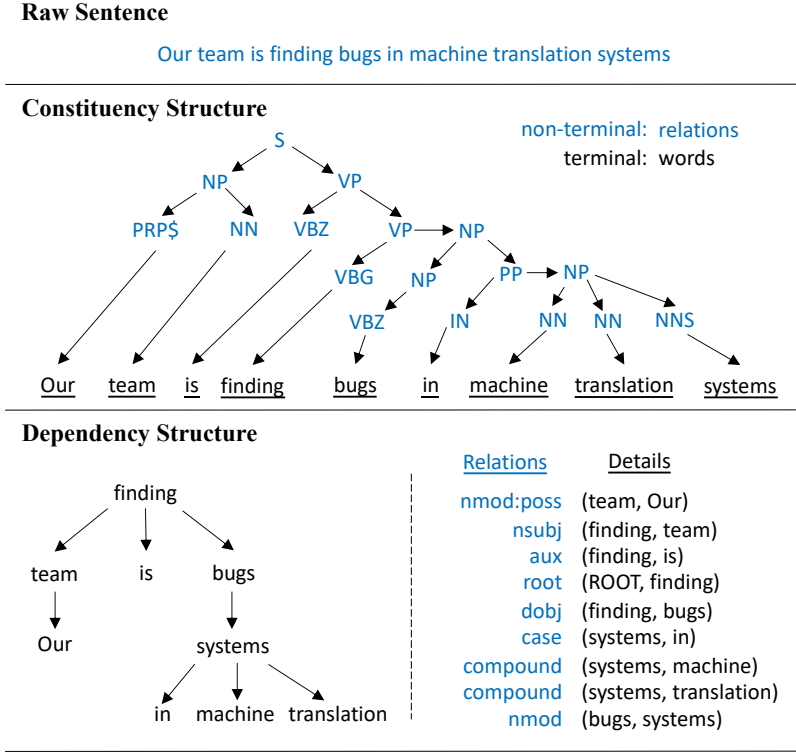


Fig. 5. Representing sentence structures; both dependency & constituency relations can be displayed as trees.

We ultimately want a representation that precisely models the structure of a sentence while offering fast comparison between two values.

The simplest and fastest approach is to compare sentences in their raw form: as strings. Indeed, we test this method and performance is reasonable. However, there are many scenarios in which this method falls short. For example, the prepositional phrase "on Friday" in the sentence "On Friday, we went to the movies" can also be placed on the end of the sentence as follows: "We went to the movies on Friday." The sentences are interchangeable but a metric such as character edit distance (see Levenshtein distance in Section 3.4.1) would indicate a large difference between the strings. We need to devise a method that solves this problem.

One commonly used data structure for representing strings that we considered is word or document embeddings [Dai et al. 2015; Mikolov et al. 2013b]. Embeddings allow us to turn any length string into a  $d$ -dimensional vector. Unfortunately, this approach fails to model the syntax structure of a sentence and is used for comparison primarily when semantic similarity alone is the desired metric, making it ill-suited for our testing.

Syntax parsing overcomes the above issues. Syntax parsing models the syntactic structure of a string and the relationship between words or groups of words. For example, if parsing is done correctly, our two sample sentences above should have identical representations in terms of relation values. There are two main types of syntax parsing: dependency and constituency. In their simplest forms, each method uses a set of context free grammars to process a sentence, deriving the set of relations that describe it. Syntax parsing can be a computationally expensive task. For our



implementation, we test several different representations for sentence comparison, taking into account efficiency vs. bug detection accuracy.

**3.3.1 Raw Target Sentence.** For this method, we leave our target sentence in its original format, i.e., as a string. In most cases, we may expect that editing a single token in a sentence in one language would lead to the change of a single token in the translated sentence. Given the syntactic role of the replacement token is the same, this would ideally happen in all machine translation systems. However, this is not always the case in practice as propositional phrases, modifiers, and other constituents can often be placed in different locations by the translation system and produce a semantically-equivalent, grammatically correct alternate sentence. Nonetheless, this method should provide a good baseline heuristic for our approach.

**3.3.2 Constituency Parse Tree.** Constituency parsing is one method for deriving the syntactic structure of a string. It generates a set of constituency relations, which show how a word or group of words form different units within a sentence. This set of relations is particularly useful for SIT because it will reflect changes to the type of phrases in a sentence. So, for example, while a prepositional phrase can be placed in multiple locations to produce a sentence with the same meaning, the set of constituency relations will remain unchanged.

Formally, in constituency parsing, a sentence is broken down into its constituent parts according to the phrase structure rules [Chomsky 1957] outlined by a given context-free grammar. There are a number of different algorithms for parsing a sentence into a set of constituent relations, but the state-of-the-art model uses shift-reduce parsing [Zhu et al. 2013]. In short, the algorithm parses grammar nonterminals from left to right in a stack-like manner until a complete set of relations is produced. More technically, the shift-reduce parser maintains a state of the current parse relations, with the words of the sentence on a queue and partially completed relation sets on a stack. It applies transitions to the state until the queue is empty and the current stack only contains a finished set. There is ambiguity in this process since often, at any given timestep, there are multiple moves that can be taken. Transitions are therefore determined by featurizing the current state and using a multiclass perceptron, trained on annotated tree banks, to determine the next transition.

Constituency relations can be visualized as a tree, as shown in Fig. 5. A constituency parse tree is an ordered, rooted tree where non-terminal nodes are the constituent relations and terminal nodes are the words. For our experiments, we use the shift-reduce constituency parser by Zhu et al. 2013 and implemented in Stanford's CoreNLP library.<sup>11</sup> It can parse about 100 sentences per second.

**3.3.3 Dependency Parse Tree.** Dependency parsing likewise derives the syntactic structure of a string. However, the set of relations produced describe the direct relationships between words rather than how words constitute a sentence. This set of relations gives us different insights about structure and is intuitively useful for SIT because it will reflect changes between how words interact. Given the way in which we modify a sentence, we should not expect this set of relations to change and such a change could indicate a translation bug.

Dependency parsing also uses context-free grammars, which are composed of dependency grammars (often just called "dependencies") rather than phrase structure grammars. These grammars establish relationships between "head" words and words which modify those heads. Much progress has been made over the past 15 years on dependency parsing. Speed and accuracy increased dramatically with the introduction of neural network based parsers [Chen and Manning 2014]. As with shift-reduce constituency parsers, neural network based dependency parsers use a stack-like system where transitions are chosen using a classifier. The classifier in this case is a neural network, likewise trained on annotated tree banks.

<sup>11</sup><https://stanfordnlp.github.io/CoreNLP/>

For our implementation, we use the most recent neural network based parsers made available by Stanford CoreNLP, which can parse about 1000 sentences per second. We use the Universal Dependencies as our annotation scheme, which has evolved based off the Stanford Dependencies [de Marneffe et al. 2014].

### 3.4 Translation Bug Detection via Structure Comparison

Finally, in order to find translation bugs, we search for structural variation by comparing sentence representations. Whether sentences are modelled as raw strings, word embeddings, or parse trees, there are a number of different metrics for calculating the distance between two values. These metrics tend to be quite domain specific and might have low correlation with each other, making the choice of metric incredibly important.

For example, a metric such as Word Mover's Distance [Kusner et al. 2015] would give us a distance of 0 between the two sentences "He went to the store" and "Store he the went to" while character edit distance would give a distance of 14. Alternatively, "I hate potatoes" would be relatively far from "I love potatoes" with Word Mover's Distance but the character edit distance would only be 3. Both of these metrics present major shortcomings for our task, which additionally places importance on syntactic similarity.

We explore several different metrics for evaluating the distance between sentences: character (Levenshtein) edit distance, dependency set difference, and constituency set difference.

**3.4.1 Levenshtein Distance between Raw Sentences.** The Levenshtein distance [Levenshtein 1966], sometime more generally referred to as the "edit distance," compares two strings and determines how closely they match each other by calculating the minimum number of character edits (deletions, insertions, and substitutions) needed to transform one string into the other. Each of these operations has unit cost (i.e., the cost of inserting a character is the same as substituting a character). For example, Levenshtein distance between "John went out" and "I went home" is 7. While the calculation can intuitively be formulated as a recursive problem, a dynamic programming approach is more commonly used for the sake of computational expense.

We use this metric as a baseline in SIT as it is by far the fastest approach for comparing two sentences. While the method may not demonstrate syntactic similarity between sentences well, it exploits the expectation that editing a single token in a sentence in one language will often lead to the change of only a single token in the translated sentence. Therefore, the Levenshtein distance serves as a good baseline metric.

**3.4.2 Relation Distance between Constituency Parse Trees.** To evaluate the distance between two sets of constituency relations, we must come up with some sort of set comparison. We calculate the distance between two lists of constituency grammars as the sum of absolute difference in the count of each phrasal type, which gives us a basic understanding of how a sentence has changed after modification. The motivation behind this heuristic is that the constituents of a sentence should stay the same between two sentences where only a single token of the same part of speech differs. In a robust machine translation system, this should be reflected in the target sentences as well.

**3.4.3 Relation Distance between Dependency Parse Trees.** Similarly, we must find a metric for calculating the distance between two lists of dependencies. Here, we sum the absolute difference in the number of each type of dependency relations. The motivation here is the same as for constituency relation sets; the relationships between words will ideally remain unchanged when a single token is replaced. We therefore expect the same set of dependency relations between the original target sentence and the modified target sentence. A change in the set is a reasonable indication that structural invariance has been violated and there might be a translation bug.

**3.4.4 Distance Thresholding.** Using one of the above metrics, we calculate the distance between the original target sentence and the generated target sentences. We must then decide whether a given target sentence is far enough from the source to indicate the presence of a translation bug. To do this, we first filter based on a distance threshold, only keeping sentences that are farther from the original sentence than the threshold. For a given original target sentence, we only report the top- $k$  ( $k$  being a chosen parameter) sentences. We leave the distance threshold as a manual parameter since the user may prioritize minimizing false positive reports or minimizing false negative reports depending on their goal. In Section 4.6, we show tradeoffs for different threshold values.

If a translated sentence falls into the "bug" category, it is labeled as a buggy sentence by SIT. For each original sentence, an issue will be reported by SIT if at least one translated generated sentence is considered buggy.

## 4 EVALUATION

In this section, we evaluate our approach by applying it to Google Translate and Bing Microsoft Translator with real-world unlabeled sentences crawled from the Web. Our main research questions are:

- RQ1: How effective is the approach at finding buggy translations in machine translation software?
- RQ2: What kinds of translation bugs can our approach find?
- RQ3: How efficient is the approach?
- RQ4: How do we select the distance threshold in practice?

By answering these research questions with both quantitative results and qualitative analysis, we aim to demonstrate: (1) SIT can detect real-world buggy translations with high accuracy; (2) SIT is a general approach that can detect diverse kinds of translation bugs; (3) SIT is highly efficient; and (4) distance threshold tuning does not require much manual effort in practice.

### 4.1 Experimental Setup

To verify the results of SIT, we manually inspect each issue reported and collectively decide: (1) whether the issue contains buggy sentences; and (2) if yes, what kind of translation bugs it contains. All experiments are run on a Linux workstation with 6 Core Intel Core i7-8700 3.2GHz Processor, 16GB DDR4 2666MHz Memory, and 1TB SATAIII Harddisk Drive 7200rpm. The Linux workstation is running 64-bit Ubuntu 18.04.02 with Linux kernel 4.25.0.

### 4.2 Dataset

Our proposed approach can automatically find translation bugs by mutating unlabeled sentences. Typically, to test a machine translation system, developers can adopt SIT with any unlabeled sentence as input. For example, developers could use sentences extracted from the newest news articles on the Web. Thus, to evaluate the effectiveness of our approach, we collect input sentences from the Web. Specifically, input sentences are extracted from CNN<sup>12</sup> (Cable News Network) articles in two categories: politics and business. The datasets are collected from two categories of articles because we intend to evaluate whether SIT consistently performs well on sentences of different semantic context.

For each category, we crawled the 10 latest articles, extracted their main text contents, and split them into a list of sentences. Then, we randomly select 100 sentences from each sentence list as the experimental datasets (200 in total). In this process, sentences that contain more than 35 words are filtered because we intend to demonstrate that machine translation software can return inferior

<sup>12</sup><https://edition.cnn.com/>

Table 1. Statistics of input sentences for evaluation. Each corpus contains 100 sentences.

Corpus	#Words/ Sentence	Average #Words/Sentence	Words	
			Total	Distinct
Politics	4~32	19.2	1,918	933
Business	4~33	19.5	1,949	944

results even for relatively short, simple sentences. Additionally, before random selection, we remove sentences containing only one word, such as "Wow!", "Why.", "I.", etc. These one-word sentences are either meaningless or very easy to translate (e.g., referring to a dictionary). The details of the collected datasets are illustrated in Table 1.

### 4.3 The Effectiveness of SIT

Our approach aims to automatically find bugs using unlabeled sentences and report them to developers. Thus, the effectiveness of the approach lies in two aspects: (1) how many buggy sentences can SIT find; and (2) How accurate are the reported results? In this section, we evaluate both aspects by applying SIT to test Google Translate and Bing Microsoft Translator using the datasets illustrated in Table 1.

**4.3.1 Evaluation Metric.** The output of SIT is a list of *issues*, each containing (1) an original source sentence and its translation; (2) the top- $k$  reported modified sentences and their translations (i.e. the  $k$  farthest translations from the source sentence translation). In this experiment, we mainly focus on the top-1 reported sentences and top-1 accuracy. For completeness, we also report the top-3 sentences for each issue and collectively label whether a reported sentence is buggy or not.

Here we define top- $k$  accuracy as the percentage of reported issues where at least one reported sentence contains a bug in the set of the top- $k$  reported sentences. We use this as our accuracy metric for SIT. Explicitly, if there is a buggy sentence in the top- $k$  reported sentences of issue  $i$ , we consider the issue to be accurate (i.e., a buggy issue) and set  $buggy(i, k)$  to 1; else we set  $buggy(i, k)$  to 0. Given a list of issues  $I$ , its top- $k$  accuracy is calculated as:

$$Accuracy = \frac{\sum_{i \in I} buggy(i, k)}{|I|}, \quad (1)$$

where  $|I|$  is the number of the issues returned by SIT.

**4.3.2 Results. Top-k accuracy.** The results are summarized in Table 2. SIT (Raw), SIT (Constituency), and SIT (Dependency) are SIT implementations with raw sentence, constituency structure, and dependency structure as sentence structure representation, respectively. Each item in the table presents the top- $k$  accuracy along with the number of buggy issues found. For the top-1 column, the number of buggy issues is the same as the number of buggy sentences since top-1 experiments only contain one sentence – the farthest from an original sentence – per issue. In subsequent discussions, for brevity, we refer SIT (Constituency) and SIT (Dependency) as SIT (Con) and SIT (Dep), respectively.

We observe that SIT (Con) and SIT (Dep) consistently perform better than SIT (Raw), which demonstrates the importance of the structure representation of sentences. SIT (Raw), which is based only on the characters in the sentences, considers too many unnecessary and noisy details. Thus, SIT (raw) may report sentences that are different in word level but similar in sentence structure, leading to false positives. SIT (Con) and SIT (Dep) achieve comparable performance in terms of both top- $k$  accuracy and the number of reported buggy issues. In particular, when testing Google

Table 2. Top-k accuracy of SIT. An issue is buggy if at least one of its top-k reported sentences is buggy.

Google Translate	Top-1 (#buggy issues)	Top-2 (#buggy issues)	Top-3 (#buggy issues)
SIT (Raw)	32.0% (32)	46.0% (46)	55.0% (55)
SIT (Constituency)	42.5% (43)	52.4% (53)	61.3% (62)
SIT (Dependency)	<b>50.0% (46)</b>	<b>57.6% (53)</b>	<b>63.0% (58)</b>
Bing Microsoft Translator	Top-1 (#buggy issues)	Top-2 (#buggy issues)	Top-3 (#buggy issues)
SIT (Raw)	38.2% (39)	53.9% (55)	59.8% (61)
SIT (Constituency)	<b>49.0% (49)</b>	<b>58.0% (58)</b>	65.0% (65)
SIT (Dependency)	48.0% (48)	54.0% (54)	<b>66.0% (66)</b>

Table 3. Number of unique bugs. Top-k unique bugs by SIT are bugs only in translated generated sentences output by SIT (Dependency).

	Original sentences	#Top-1 unique bugs by SIT	#Top-2 unique bugs by SIT	#Top-3 unique bugs by SIT
Google	55	45	64	79
Bing	60	32	43	66

Translate, SIT (Dep) reports 92 suspicious issues, where each issue contains 3 reported sentences and their translations. Among these issues, 46 of them contains translation bugs in the first reported sentence pair, achieving 50% top-1 accuracy. When testing Bing Microsoft Translator, SIT (Con) reports 100 suspicious issues, among which 49 issues have translation bugs within the first reported sentence pair.

**Bugs in original sentence translation.** During the manual inspection on SIT’s reported issues, we noted that a large number of original sentence pairs also contained translation bugs. Specifically, in our experiments, 55 and 60 unique bugs are found in the translation of original sentences by Google Translate and Bing Microsoft Translator respectively.

We believe that in practice, SIT is effective for an input sentence that either is or is not correctly translated by the machine translation software under test. For an original sentence that is correctly translated, SIT generates similar sentences and report the buggy ones, which enhances the reliability of machine translation software by providing meaningful corner cases. For an original sentence that is not correctly translated, SIT can (1) find more unique bugs that the original sentence pair does not contain; (2) provide a list of similar sentences that have similar bugs or are correctly translated, both of which are useful for developers in further analysis and debugging. As illustrated in Table 3, SIT finds 79 and 66 unique translation bugs that are revealed in the generated sentence pairs but not in the original.

**Top-k accuracy considering buggy original sentence translation.** In an issue output by SIT, the original sentence and its translation are also reported. In some cases, the reported top-k similar sentence pairs are bug-free, but the corresponding original sentence is incorrectly translated. Given that the input is an unlabeled sentence, we think successful bug detection involves evaluating both the original sentence pair and the generated sentence pairs. Thus, in Table 4, we evaluate the top-k

Table 4. Top-k accuracy of SIT considering the original sentences.

Google Translate	Top-1 (#buggy issues)	Top-2 (#buggy issues)	Top-3 (#buggy issues)
SIT (Raw)	55.0% (55)	63.0% (63)	66.0% (66)
SIT (Constituency)	61.3% (62)	66.3% (67)	68.3% (69)
SIT (Dependency)	<b>69.5% (64)</b>	<b>71.7% (66)</b>	<b>73.9% (68)</b>
Bing Microsoft Translator	Top-1 (#buggy issues)	Top-2 (#buggy issues)	Top-3 (#buggy issues)
SIT (Raw)	58.8% (60)	69.6% (71)	71.5% (73)
SIT (Constituency)	67.0% (67)	<b>71.0% (71)</b>	74.0% (74)
SIT (Dependency)	<b>70.0% (70)</b>	<b>71.0% (71)</b>	<b>78.0% (78)</b>

Table 5. Number of sentences that have specific bugs in each category for Google Translate via SIT (Dep).

Google Translate	Under translation	Over translation	Incorrect modification	Word/phrase mistranslation	Unclear logic
Top-1	35	9	4	44	27
Top-2	48	12	6	59	44
Top-3	61	15	10	75	53

Table 6. Number of sentences that has specific bugs in each category for Bing Microsoft Translator via SIT (Dep).

Bing Microsoft Translator	Under translation	Over translation	Incorrect modification	Word/phrase mistranslation	Unclear logic
Top-1	17	8	2	54	31
Top-2	23	15	3	60	41
Top-3	35	21	4	93	59

accuracy also considering the original sentence pairs. Intuitively, we count it as a correctly reported buggy top-1 sentence if the original sentence pair contains a translation bug. In particular, for an issue  $i$ , if the reported original sentence pair is buggy, we set *buggy*( $i$ , 1) to 1 even if the generated sentence pair is not buggy.

For this evaluation metric, we can observe that SIT (Con) and SIT (Dep) achieve comparable performance and consistently outperform SIT (Raw). SIT (Dep) has the best performance on Top-1 accuracy for both Google Translate and Bing Microsoft Translator. It successfully finds 64 and 70 buggy issues with 69.5% and 70% top-1 accuracy, respectively. Thus, we think SIT (Dep) achieves the best performance in terms of reporting bugs in original sentence pairs or in the generated similar sentence pairs.



#### 4.4 Translation Bug Reported by SIT

SIT is capable of finding translation bugs of diverse kinds. In our experiments with Google Translate and Bing Microsoft Translator, we mainly find 5 kinds of translation bugs: under-translation, over-translation, incorrect modification, word/phrase mistranslation, and unclear logic. To provide a glimpse of the diversity of the uncovered bugs, this section highlights examples for all the 5 kinds of bugs.

Table 5 and Table 6 present the statistics of the translation bugs SIT found using the unlabeled sentences illustrated in Table 1. Under-translation, word/phrase mistranslation, and unclear logic account for most of the translation bugs found by SIT.

Software	Source sentence	Target sentence	Target sentence meaning
Google Translate	The former employee is now looking for new career opportunities <b>outside their field of expertise</b> after hitting a road block in the job search.	这位前员工在求职时遇到障碍后，正在寻找新的职业机会。	The former employee is now looking for new career opportunities after hitting a road block in the job search.
Bing Microsoft Translator	After pleading guilty in the Manhattan probe, Cohen also later pleaded guilty to lying <b>to Congress</b> in a case brought by Mueller's website.	在曼哈顿调查中认罪后，科恩后来还对穆勒网站提起的一起案件中的撒谎罪供认不讳。	After pleading guilty in the Manhattan probe, Cohen also later pleaded guilty to lying in a case brought by Mueller's website.

Fig. 6. Examples of under-translation bugs detected by SIT.

**4.4.1 Under-Translation.** If some words are mistakenly untranslated (i.e. do not appear in the translation), it is an under-translation bug. Fig. 6 presents two sentence pairs that contain under-translation bugs. In these examples, "outside their field of expertise" and "to Congress" are mistakenly untranslated, which lead to target sentences of largely different semantic meanings. In the second example, "lying to Congress" is illegal while "lying" is just an inappropriate behavior.

Software	Source sentence	Target sentence	Target sentence meaning
Google Translate	The investigators were right that the airplane itself was safe.	调查人员认为飞机本身是安全的。	The investigators <b>thought</b> that the airplane itself was safe.
Bing Microsoft Translator	I am very happy to share my point of view.	我很高兴与大家分享我的观点。	I am very happy to share my point of views <b>with everyone</b> .

Fig. 7. Examples of over-translation bugs detected by SIT.

**4.4.2 Over-Translation.** If some words are unnecessarily translated multiple times or some words in the target sentence are not translated from any words in the source sentence, it is an over-translation bug. Fig. 7 presents two sentence pairs that contain over-translation bugs. In the first example, "thought" in the target sentence is not translated from any words in the source sentence, so it is an over-translation bug. Interestingly, we found that an over-translation bug often happens along with some other kinds of bugs. The first example also contains an under-translation bug because "were right" in the source sentence is mistakenly untranslated. In the second example in Fig. 2, the word "a" is unnecessarily translated twice, which makes it an over-translation bug.

Software	Source sentence	Target sentence	Target sentence meaning
Google Translate	The South has emerged as <b>a hub of new auto manufacturing</b> by foreign makers thanks to lower manufacturing costs and less powerful businesses.	由于制造成本降低和业务不那么强大，南方已成为外国制造商新的汽车制造中心。	The South has emerged as <b>a new hub of auto manufacturing</b> by foreign makers thanks to the reducing manufacturing costs and less powerful businesses.
Bing Microsoft Translator	Anxious gossip about <b>who is and is not mentioned</b> in the latest news reports.	关于最新新闻报道中谁是谁和谁没有被提及的焦虑流言蜚语。	Anxious gossip about <b>who is who and who is not mentioned</b> in the latest news reports.

Fig. 8. Examples of incorrect modification bugs detected by SIT.

**4.4.3 Incorrect Modification.** If some modifiers modify the wrong element in the sentence, it is an incorrect modification bug. Fig. 8 presents two sentence pairs that contain incorrect modification bugs. In the first example, the modifier "new" modifies "auto manufacturing" in the source sentence. However, Google Translate thinks that "new" should modify "hub." In Fig. 2, the third example also shows an interesting incorrect modification bug. In this example ("members of privilege"), "privilege" modifies "members" in the source sentence, while Google Translate thinks "members" should modify "privilege." We think that in the training data of the NMT model, there are some phrases with similar pattern: "A of B", where A modifies B, which leads to an incorrect modification bug in this scenario. Interestingly, the original source sentence that triggers this bug is "But even so, they remain *bastions of privilege*". In the original source sentence, "bastions" modifies "privilege," which fits the supposed archetype. As we might expect, this sentence is correctly translated by Google Translate.

Software	Source sentence	Target sentence	Target sentence meaning
Google Translate	The most elite public universities <b>admit</b> a considerably larger percentage of students from lower income backgrounds than do the elite private schools.	最精英的公立大学承认，与精英私立学校相比，低收入学生的比例要高得多。	The most elite public universities <b>agree unwillingly that</b> considerably larger percentage of students from lower income backgrounds than do the elite private schools.
	It declined to <b>ground</b> the jet.	它拒绝接近喷气式飞机。	It declined to <b>get close to</b> the jet.
Bing Microsoft Translator	But it would keep the new investment in Michigan and focus there on <b>self driving vehicles</b> , according to reports.	但据报道，这将保留对密歇根州的新投资，并将重点放在自动驾驶车辆上。	But it would keep the new investment in Michigan and focus there on <b>vehicles for self driving tours</b> , according to reports.
	The <b>South</b> has emerged as a hub of new auto manufacturing by foreign makers thanks to lower manufacturing costs and less powerful unions.	由于制造成本较低，工会实力较弱，韩国已成为外国制造商新汽车制造业的枢纽。	The <b>South Korea</b> has emerged as a hub of new auto manufacturing by foreign makers thanks to lower manufacturing costs and less powerful unions.

Fig. 9. Examples of word/phrase mistranslation bugs detected by SIT.

**4.4.4 Word/phrase Mistranslation.** If some tokens or phrases are incorrectly translated in the target sentence, it is a word/phrase mistranslation bug. Fig. 9 presents four sentence pairs that

contain word/phrase mistranslation bugs. There are two main sub-categories of this kind of bugs: (1) ambiguity of polysemy and (2) wrong translation.

**Ambiguity of polysemy.** Each token/phrase may have multiple correct translations. For example, "admit" means "allow somebody to join an organization" or "agree with something unwillingly". However, usually in a specific semantic context (e.g., a sentence), a token/phrase only has one correct translation. Modern translation software does not perform well on polysemy. In the first example in Fig. 9, Google Translate thinks the "admit" in the source sentence refers to "agree with something unwillingly," leading to a token/phrase mistranslation bug.

**Wrong translation.** A token/phrase could also be incorrectly translated to another meaning that seems semantically unrelated. For example, in the fourth example in Fig. 9, Bing Microsoft Translator thinks "South" refers to "South Korea", leading to a word/phrase mistranslation bug.

Software	Source sentence	Target sentence	Target sentence meaning
Google Translate	And attacking a dead man who spent five years as a prisoner of war and another three decades <b>serving the country in elected office</b> , is simply wrong.	并且攻击一名死去的人，他在战争中担任战争囚犯五年，另外三十年担任民选职务的国家，这是完全错误的。	And attacking a dead man who spent five years as a prisoner of war and another three decades <b>serving in elected office as a country</b> , is simply wrong.
Bing Microsoft Translator	Since then, <b>the ex employee</b> has been rejected for multiple jobs <b>and is still searching</b> .	此后，前员工因多个工作被拒绝，目前仍在寻找中。	Since then, <b>because the ex employee</b> has been rejected for multiple jobs, <b>he is still searching</b> .

Fig. 10. Examples of unclear logic bugs detected by SIT.

**4.4.5 Unclear Logic.** If all the tokens/phrases are correctly translated but the sentence logic is incorrect, it is an unclear logic bug. Fig. 10 presents two sentence pairs that contain unclear logic bugs. In the first example, Google Translate correctly translates "serving in the elected office" and "country". However, Google Translate generates "serving in the elected office as a country" instead of "serving the country in elected office" because Google Translate does not understand the logical relation between them. Unclear logic bugs exist widely in translations given by NMT models, which is to some extent a sign of whether an NMT model truly understands certain semantic meanings.

Software	Source sentence	Target sentence	Target sentence meaning	Bug type
Google Translate	<b>Sewing</b> may find some room for <b>a merger among European regulators</b> who argue the region has too many banks.	欧洲监管机构认为，该地区拥有太多银行，因此缝制可能会找到合并的空间。	European regulators argue that the region has too many banks, <b>so</b> the <b>sewing skill</b> may find some room for <b>a merger</b> .	<b>word/phrase logic</b>
Bing Microsoft Translator	<b>Covering</b> a memorial service in the nation's capital and then traveling to Texas <b>for</b> another service as well as <b>a funeral train</b> was an honor, he says.	他说，在美国首都举行追悼会，然后前往德州参加另一次礼拜仪式以及葬礼列车，是一种荣誉。	<b>Holding</b> a memorial service in the nation's capital and then traveling to Texas <b>for attending</b> another <b>church service</b> and <b>a funeral train</b> was an honor, he says.	<b>word/phrase logic over</b>

Fig. 11. Examples of sentence with multiple translation bugs detected by SIT.

**4.4.6 Sentences with Multiple Translation Bugs.** A certain percentage of reported sentence pairs contain multiple translation bugs. Fig. 11 presents two sentence pairs that contain multiple bugs.

Table 7. Average running time of SIT on Politics and Business Datasets.

Google	Running time (sec)	Translation time (sec)	#Sentence translated	Time of other SIT steps (sec)
SIT (Raw)	1,469	1,417	2,012	52
SIT (Constituency)	1,524	1,417	2,012	107
SIT (Dependency)	1,488	1,417	2,012	71
Bing	Running time (sec)	Translation time (sec)	#Sentence translated	Time of other SIT steps (sec)
SIT (Raw)	922	870	2,012	52
SIT (Constituency)	981	870	2,012	110
SIT (Dependency)	945	870	2,012	75

In the first example, "Sewing" is a name in the source sentence. However, it is translated to "sewing skill", causing a word/phrase mistranslation bug. Additionally, in the target sentence, it is unclear whether it is a merger among "European regulators" or "banks", so the sentence pair also contains an unclear logic bug. In the second example, "covering" means "reporting news" in the source sentence. However, it is translated to "holding", leading to a word/phrase mistranslation bug. Additionally, "church" in the target sentence is not the translation of any words from the source sentence, so it is an over-translation bug. Bing Microsoft translator also wrongly thinks the subject is "attending a funeral train". But the source sentence actually means the subject is "covering a funeral train", so it is an unclear logic bug.

#### 4.5 The Running Time of SIT

In this section, we evaluate the running time of SIT on the two datasets. Specifically, we apply SIT with 3 different sentence structure representations to test Google Translate and Bing Microsoft Translator. We run each experiment setting 10 times and report their average as the results.

The overall running time of SIT is illustrated in Table 7, and the running time of each step of SIT is presented in Fig. 12 and Fig. 13. We can observe that SIT using raw sentences as structure representation is the fastest for testing both Google Translate (around 25 mins) and Bing Microsoft Translator (around 15 mins). This is because SIT (Raw) does not require any structure representation generation time. Additionally, SIT using a dependency parser achieves comparable running time to SIT (Raw). In particular, SIT (Dependency) uses 19 seconds to parse 2000+ sentences (as opposed to 0 seconds by SIT (Raw)), which we think is efficient and reasonable.

In these experiments, we ran the translation step once per translation system and reused the translation results in all experiment settings since the other settings had no affect on translation time. Thus, in Table 7, the *Translation time* values are the same for different SIT implementations. We can observe that SIT spends most of the time collecting translation results. In this step, for each sentence, we invoked the APIs provided by Google Translator or Bing Microsoft Translator to collect the translated sentence. We did not send all sentences in one API invocation because the APIs have character-level limitations on the size of one translation API call. In practice, if you want to test your own machine translation software with SIT, the running time of this step will be much less. As indicated in a recent study [Zhang et al. 2018a], current NMT model can translate around 20 sentences per second using a single NVIDIA GeForce GTX 1080 GPU. With more powerful computing resource (e.g. TPU [Wu et al. 2016]), modern NMT models can achieve the speed of hundreds of sentences translation per second.

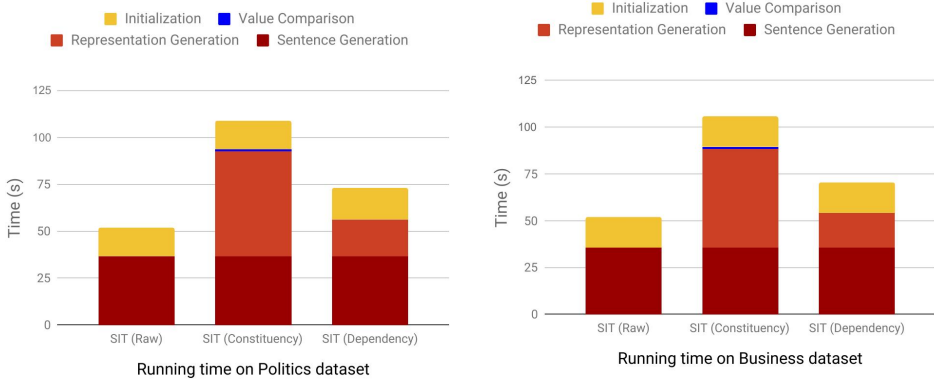


Fig. 12. Running time details of SIT (excluding translation time) in testing Google Translate.

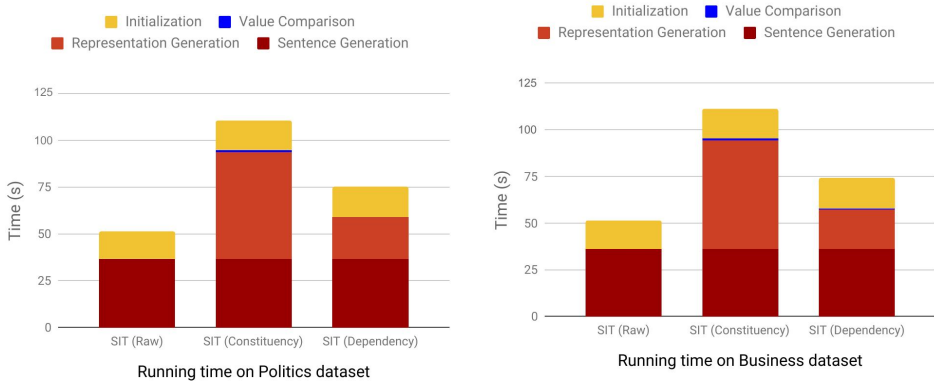


Fig. 13. Running time details of SIT (excluding translation time) in testing Bing Microsoft Translator.

The other steps of SIT are comparatively quite efficient, as indicated in Table 7, Fig. 12, and Fig. 13. Both SIT (Raw) and SIT (Dependency) took around 1 min and SIT (Constituency) took around 2 mins for 2000+ sentences. Compared with SIT (Dep), SIT (Con) is slower because constituency parsing generally takes longer time than dependency parsing. We conclude that as a tool working in an offline manner, SIT is efficient in practice for testing machine translation software.

#### 4.6 The Impact of Distance Threshold

SIT reports the top-k sentence pairs in an issue if the distance between the translated modified sentence and the original target sentence is larger than a distance threshold. Thus, this distance threshold controls (1) the number of issues reported and (2) the top-k accuracy of SIT.

Intuitively, if we lower the threshold, more buggy issues will be reported, while the accuracy will decrease. Fig. 14 demonstrates the impact of the distance threshold on these two factors. In this figure, SIT (Dep) was applied to test the Bing Microsoft Translator on our Politics and Business datasets with different distance thresholds. We can observe that both the number of reported issues and top-1 accuracy remain stable when the threshold is either small or large while the values

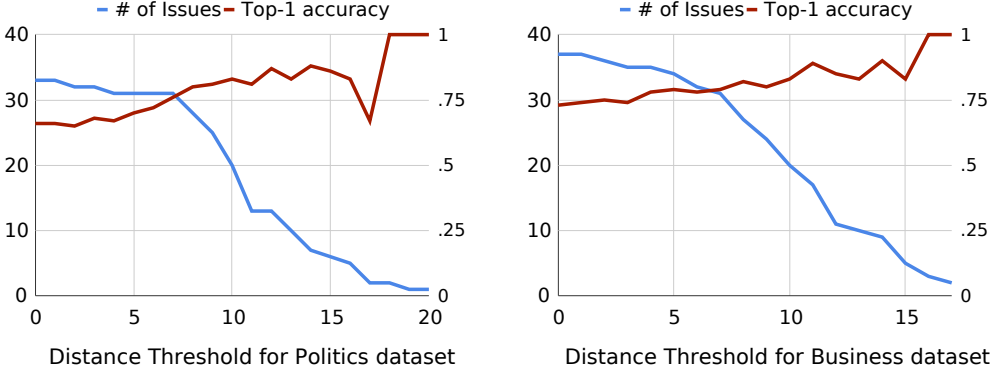


Fig. 14. Impact of distance threshold when testing Bing Microsoft Translator.

fluctuate in the middle. The impact of changing the distance threshold is similar when testing Google Translate.

Based on these results, we present some guidance on using SIT in practice. First, if we intend to uncover as many as translation bugs as possible, we should use a small distance threshold. A small threshold (e.g., 4 for dependency sets) works well on all our experiment settings. In particular, with a small threshold, SIT reports the most issues with decent accuracy (e.g., 70% top-1 accuracy). We adopt this strategy in our accuracy experiments in Section 4.3.2. Developers could use SIT with small distance threshold when they want to intensively test software before a release. Second, if we intend to make SIT as accurate as possible, we could use a large threshold (e.g., 15). With a large threshold, SIT reports fewer issues with very high accuracy (e.g., 86% top-1 accuracy). Given that the number of unlabeled sentences are unlimited on the Web, we could keep running SIT with a large distance threshold in the background and periodically report issues to developers. Thus, we think SIT is useful and easy to use in practice.

## 5 DISCUSSION

### 5.1 False Positives

While SIT can accurately return translation bugs, we believe there is room for further improvement. In particular, the false positives of SIT come from three main sources. First, the generated sentences could contain syntax errors or have strange semantic meanings, leading to changes in the target sentence structure. For example, our SIT implementation may substitute a word in idiomatic phrases. To alleviate the impact of these false positives, we have used BERT, which provides the state-of-the-art masked language model. Second, although the existing syntax parsers are highly accurate, they could produce wrong constituency or dependency structures, leading to false positives. Third, a source sentence could have multiple correct translations of different sentence structures. To lower the impact of these factors, SIT returns the top-k suspicious sentence pairs ranked by their distance to the original target sentence.

To further improve our SIT implementation, we consider three possible approaches: (1) build a dictionary of frequently-used idiomatic phrases and avoid changing any words in them; (2) adopt advanced and phrase-level similar source sentence generation techniques, such as paraphrasing; and (3) develop more accurate syntax parsers by combining neural network-based models with search-based models. We leave such further explorations as future work.



## 5.2 Building Robust Translation Software

The ultimate goal of testing machine translation, similar to testing traditional software, is to build robust machine translation software. Toward this end, SIT's utility is as follow. First, the reported mistranslations typically act as early alarms, and thus developers could hardcode them in advance, which is the quickest bug fixing solution in modern industry. Second, the reported sentences could be used for additional training. Third, for each reported issue, the translations either exhibit similar or different bugs, or only some exhibit bugs. In all cases, developers may find the reported buggy sentences useful for further analysis/debugging since these sentences only differ by one word. This resembles debugging traditional software via input minimization/localization. Additionally, the structural invariance concept could be utilized as inductive bias to design robust NMT models, similar to how [Shen et al. 2019] introduce bias to standard LSTMs. Compared with traditional software, the debugging and bug fixing process of machine translation software is more difficult because the logic of NMT models mainly lies in their model structure and parameters. While this is not the main focus of our paper, we think it is an important research direction and regard it as future work.

## 6 RELATED WORK

This section surveys several strands of relevant related work: robustness of AI software, robustness of NLP algorithms, machine translation, and metamorphic testing.

### 6.1 Robustness of Artificial Intelligence Software

The success of deep neural networks has led to the wide adoption of artificial intelligence (AI) software in our daily lives. Typical examples include virtual assistants, autonomous cars, face recognition systems, fraud detection tools, and machine translation software. However, although the deep models achieved high accuracies evaluated by some automatic metrics, they could generate inferior results that even lead to fatal accidents [Lambert 2016; Levin 2018; Ziegler 2016]. Recently, researchers have design a variety of approaches to attack different DL (deep learning) systems, such as autonomous cars [Athalye et al. 2018; Carlini and Wagner 2017; Goodfellow et al. 2015], 3D object classifiers [Xiong et al. 2019; Yang et al. 2019], and speech recognition services [Carlini et al. 2016; Du et al. 2019]. These work aim at generating inputs that fool the corresponding neural network (i.e., adversarial examples). For example, techniques proposed in [Athalye et al. 2018; Carlini and Wagner 2017; Goodfellow et al. 2015] tried to generate adversarial examples that are context-preserving and lead to incorrect image classification. Compared with these approaches, our paper has two main differences. First, we focus on a novel and important DL system: machine translation system, which has not been studied by these papers. Second, most of these approaches are based on gradients in neural networks, while our approach does not require any internal details of the networks.

To protect DL system against these attacks, a line of research work has been conducted to either train the networks in a robust way [Kannan et al. 2018; Lin et al. 2019; Madry et al. 2018; Papernot et al. 2016] or detect adversarial examples online [Ma et al. 2019; Tao et al. 2018; Wang et al. 2019; Xu et al. 2018]. Different from these papers, our goal is to prevent users suffering from inferior translation results by reporting translation bugs offline.

Recently, testing of image classifiers have also been well studied. Pei et al. 2017 proposed DeepXplore, a differential testing approach for DL image classifiers. They also introduced a novel concept called neuron coverage, which is used to systematically measure the parts of a DNN exercised by test inputs. Inspired by this work, Tian et al. 2018 developed a tool for automated testing of DL self-driving cars. In particular, they designed realistic image transformation rules

and domain-specific metamorphic relations to find erroneous behaviors. Ma et al. 2018a proposed several fine-grained testing metrics for DNNs based on neuron coverage [Pei et al. 2017]. Different from these neuron coverage-based approaches, Kim et al. 2019 developed a DL testing technique guided by surprise adequacy, which is based on the behavior of DL systems with respect to their training data. Ma et al. 2018b proposed MODE, an automated neural network debugging approach that identifies bugs in training data (i.e., misleading images). Different from these papers that focus on testing image classifiers, we aim at testing machine translation software.

## 6.2 Robustness of Natural Language Processing Algorithms

Deep neural networks have boosted the performance of many NLP tasks, such as reading comprehension [Chen 2018; Chen et al. 2016], code analysis [Alon et al. 2019; Iyer et al. 2016; Liu et al. 2018; Pradel and Sen 2018], and machine translation [Hassan et al. 2018; Vaswani et al. 2017; Wu et al. 2016]. However, in recent years, inspired by the attack papers in computer vision field, researchers successfully found flaws in neural networks adopted by various NLP systems.

Most of the existing studies focus on finding adversarial examples (i.e., bugs) for text classification tasks. Typical tasks include sentiment analysis [Alzantot et al. 2018; Iyyer et al. 2018; Li et al. 2019], textual entailment [Iyyer et al. 2018], and toxic content detection [Li et al. 2019]. Compared with text classification, machine translation is a much more difficult task.

Some research studied the robustness of other complex NLP systems. Jia and Liang 2017 proposed an adversarial evaluation scheme for the SQuAD dataset,<sup>13</sup> which is widely used in reading comprehension studies. They found that adversarially inserting sentences to the paragraphs can make reading comprehension systems fail to answer questions about paragraphs correctly. Mudrakarta et al. 2018 further generated adversarial examples for question answering tasks on images, tables, and passages of text. Ribeiro et al. 2018 introduced an approach to generate semantically equivalent adversarial examples and rules for the question answering task. These approaches typically try to perturb normal model inputs (e.g., adding unrelated sentences) and assume that the output should not be affected. If the output changed, the perturbed input will be reported as an adversarial example. However, this assumption does not hold for machine translation software, because one source sentence could have multiple correct target sentences. Thus, we believe testing machine translation software is more difficult and we propose a novel methodology in this paper.

Belinkov and Bisk 2018 found that character-based NMT models fail to translate even moderately noisy texts that humans have no trouble comprehending. Ebrahimi et al. 2018 designed gradient-based methods to attack character-based NMT models. Different from them, our paper aims at testing general NMT models (i.e., character- and token-based NMT models). Besides, this paper aims testing the NMT model by noise-free texts, which is much more common in practice.

Zhou and Sun 2018 proposed a metamorphic testing approach (i.e., MT4MT) for machine translation. They mentioned a concept similar to structural invariance. However, MT4MT is only specialized for very short sentences in *subject-verb-object* pattern (e.g., "Tom likes Nike"). In particular, they change a person name or a brand name in a sentence, and check whether the translation differs more than one place. Thus, MT4MT cannot report any bugs from most real-world sentences, such as the dataset used in our paper. Differently, the approach proposed by us aims at finding bugs in real-world sentences in general. The most relevant work is proposed by Zheng et al. 2018. They introduced two algorithms to detect two specific translation bugs: under-translation and over-translation, respectively. Different from them, our proposed approach is not limited to specific translation bugs. Based on the experimental results, we can find the following translation bugs: under-translation, over-translation, incorrect modification, ambiguity of polysemy, and unclear logic.

<sup>13</sup><https://rajpurkar.github.io/SQuAD-explorer/>

### 6.3 Machine Translation

The past few years have witnessed the rapid growth of neural machine translation (NMT) models. As reported by research from Google [Wu et al. 2016] and Microsoft [Hassan et al. 2018], the state-of-the-art NMT models are approaching human-level performance in terms of accuracy (i.e., BLEU [Papineni et al. 2002]). Typically, an NMT model translates a source sentence into the target language with an encoder-attention-decoder framework [Zhang et al. 2018a]. Under this framework, researchers designed various advanced neural network architectures, ranging from recurrent neural networks (RNN) [Luong et al. 2015; Sutskever et al. 2014], convolutional neural networks (CNN) [Gehring et al. 2017a,b], to full attention networks without recurrence and convolution [Vaswani et al. 2017].

Recently, to further improve model accuracy, researchers have designed different attention mechanisms (e.g., introducing a recurrent cycle on the attention layer for better memorization of translated words [Yang et al. 2017], a GRU-gated attention network [Zhang et al. 2017], modeling intrinsic structures inside attention through graphical models [Kim et al. 2017]). These existing papers aim at improving the accuracy of NMT models. Different from them, this paper focuses on the robustness of NMT models. We believe robustness is as important as accuracy for machine translation in practice. Thus, our proposed approach can complement existing machine translation research.

### 6.4 Metamorphic Testing

Metamorphic testing is a way of generating test cases based on existing ones [Chen et al. 1998, 2018; Segura et al. 2016]. The key idea is to detect violations of domain-specific metamorphic relations across outputs from multiple runs of the program with different inputs. Metamorphic testing has been applied for testing various traditional software, such as compilers [Le et al. 2014; Lidbury et al. 2015], scientific libraries [Zhang et al. 2014], and database systems [Lindvall et al. 2015]. Due to its effectiveness on testing "non-testable" programs, researchers also used it to test AI software, such as statistical classifiers [Murphy et al. 2008; Xie et al. 2009, 2011], search engines [Zhou et al. 2016], and autonomous cars [Tian et al. 2018; Zhang et al. 2018b]. In this paper, we introduce structure-invariant testing for machine translation software. Our approach is inspired by the metamorphic relationship between source sentence and target sentence: if we substitute a token in source sentence with a similar token, the target sentence structure is expected to be consistent with the original one.

## 7 CONCLUSION

We have presented structure-invariant testing (SIT), a new, effective approach to testing machine translation. The distinct benefits of SIT are its simplicity and generality, and thus wide applicability. SIT has been applied to test Google Translate and Bing Microsoft Translators with unlabeled sentences, and successfully found 64 and 70 buggy translations with 69.5% and 70% top-1 accuracy respectively. Moreover, as a general methodology, SIT can uncover diverse kinds of translation bugs that cannot be found by state-of-the-art approaches.

This work focuses on the robustness of machine translation software and complements decades of extensive work on improving the accuracy of NMT models. We believe that this work is the important, first step toward systematic testing of machine translation software. For future work, we will continue refining the general approach and extend it to other AI software (e.g., figure captioning tools and face recognition systems). We will also launch an extensive effort to help continuously test and improve widely-used machine translation systems.

## REFERENCES

- Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. 2019. code2vec: Learning Distributed Representations of Code. *Proceedings of the ACM on Programming Languages* POPL (2019).
- Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. 2018. Generating Natural Language Adversarial Examples. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Anish Athalye, Nicholas Carlini, and David Wagner. 2018. Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*.
- Yonatan Belinkov and Yonatan Bisk. 2018. Synthetic and Natural Noise Both Break Neural Machine Translation. In *Proceedings of the 6th International Conference on Learning Representations (ICLR)*.
- Nicholas Carlini, Pratyush Mishra, Tavish Vaidya, Yuankai Zhang, Micah Sherr, Clay Shields, David Wagner, and Wenchao Zhou. 2016. Hidden Voice Commands. In *Proceedings of the 25th USENIX Security Symposium (USENIX Security)*.
- Nicholas Carlini and David Wagner. 2017. Towards Evaluating the Robustness of Neural Networks. In *IEEE Symposium on Security and Privacy*.
- Danqi Chen. 2018. *Neural Reading Comprehension and Beyond*. Ph.D. Dissertation. Stanford University.
- Danqi Chen, Jason Bolton, and Christopher D. Manning. 2016. A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Danqi Chen and Christopher Manning. 2014. A Fast and Accurate Dependency Parser using Neural Networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Tsong Y. Chen, Shing C. Cheung, and Shiu Ming Yiu. 1998. *Metamorphic testing: a new approach for generating next test cases*. Technical Report. Technical Report HKUST-CS98-01, Department of Computer Science, Hong Kong University of Science and Technology, Hong Kong.
- Tsong Yueh Chen, Fei-Ching Kuo, Huai Liu, Pak-Lok Poon, Dave Towey, T. H. Tse, and Zhi Quan Zhou. 2018. Metamorphic Testing: A Review of Challenges and Opportunities. *ACM Computing Surveys (CSUR)* 51 (2018). Issue 1.
- N. Chomsky. 1957. *Syntactic Structures*. Mouton, The Hague.
- Andrew M. Dai, Christopher Olah, and Quoc V. Le. 2015. Document embedding with paragraph vectors. In *NIPS Deep Learning Workshop*.
- Gareth Davies. 2017. Palestinian man is arrested by police after posting 'Good morning' in Arabic on Facebook which was wrongly translated as 'attack them'. <https://www.dailymail.co.uk/news/article-5005489/Good-morning-Facebook-post-leads-arrest-Palestinian.html>
- Marie-Catherine de Marneffe, Timothy Dozat, Natalia Silveira, Katri Haverinen, Filip Ginter, Joakim Nivre, and Christopher D. Manning. 2014. Universal Stanford dependencies: A cross-linguistic typology. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*. European Language Resources Association (ELRA), Reykjavik, Iceland.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805* (2018).
- Tianyu Du, Shouling Ji, Jinfeng Li, Qinchen Gu, Ting Wang, and Raheem Beyah. 2019. SirenAttack: Generating Adversarial Audio for End-to-End Acoustic Systems. *arXiv preprint arXiv:1901.07846* (2019).
- Javid Ebrahimi, Daniel Lowd, and Dejing Dou. 2018. On Adversarial Examples for Character-Level Neural Machine Translation. In *Proceedings of the 27th International Conference on Computational Linguistics (COLING)*.
- Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding Back-Translation at Scale. *arXiv preprint arXiv:1808.09381* (2018).
- Ivan Evtimov, Kevin Eykholt, Earlene Fernandes, Tadayoshi Kohno, Bo Li, Atul Prakash, Amir Rahmati, and Dawn Song. 2018. Robust Physical-World Attacks on Deep Learning Models. In *Proceedings of the 2018 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Facebook. 2019. How do I translate a post or comment written in another language? [https://www.facebook.com/help/509936952489634?helpref=faq\\_content](https://www.facebook.com/help/509936952489634?helpref=faq_content)
- Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. 2013. PPDB: The Paraphrase Database. In *Proceedings of the 11st Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*.
- Jonas Gehring, Michael Auli, David Grangier, and Yann N. Dauphin. 2017a. A Convolutional Encoder Model for Neural Machine Translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Jonas Gehring, Michael Auli, David Grangier, and Yann N. Dauphin. 2017b. Convolutional Sequence to Sequence Learning. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and Harnessing Adversarial Examples. *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*.

- Hany Hassan, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, et al. 2018. Achieving Human Parity on Automatic Chinese to English News Translation. *arXiv preprint arXiv:1803.05567* (2018).
- Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. 2016. Summarizing Source Code using a Neural Attention Model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. Adversarial Example Generation with Syntactically Controlled Paraphrase Networks. In *Proceedings of the 16th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*.
- Robin Jia and Percy Liang. 2017. Adversarial Examples for Evaluating Reading Comprehension Systems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Harini Kannan, Alexey Kurakin, and Ian Goodfellow. 2018. Adversarial Logit Pairing. *arXiv preprint arXiv:1803.06373* (2018).
- Jinhan Kim, Robert Feldt, and Shin Yoo. 2019. Guiding Deep Learning System Testing using Surprise Adequacy. In *Proceedings of the 41st International Conference on Software Engineering (ICSE)*.
- Yoon Kim, Carl Denton, Luong Hoang, and Alexander M. Rush. 2017. Structured Attention Networks. In *Proceedings of the 5th International Conference on Learning Representations (ICLR)*.
- Matt J. Kusner, Yu Sun, Nicholas I. Kolkin, and Kilian Q. Weinberger. 2015. From Word Embeddings to Document Distances. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37 (ICML'15)*. 957–966.
- Fred. Lambert. 2016. Understanding the fatal Tesla accident on Autopilot and the NHTSA probe. <https://electrek.co/2016/07/01/understanding-fatal-tesla-accident-autopilot-nhtsa-probe/>
- Vu Le, Mehrdad Afshari, and Zhendong Su. 2014. Compiler Validation via Equivalence Modulo Inputs. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*.
- Vladimir I. Levenshtein. 1966. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady* 10 (1966), 707–710. Issue 8.
- Sam Levin. 2018. Tesla fatal crash: 'autopilot' mode sped up car before driver killed, report finds. <https://www.theguardian.com/technology/2018/jun/07/tesla-fatal-crash-silicon-valley-autopilot-mode-report>
- Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2019. TextBugger: Generating Adversarial Text Against Real-world Applications. In *Proceedings of the 26th Annual Network and Distributed System Security Symposium (NDSS)*.
- Christopher Lidbury, Andrei Lascu, Nathan Chong, and Alastair F. Donaldson. 2015. Many-Core Compiler Fuzzing. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*.
- Chin-Yew Lin. 2004. Rouge: A Package for Automatic Evaluation of Summaries. *Text Summarization Branches Out* (2004).
- Ji Lin, Chuang Gan, and Song Han. 2019. Defensive Quantization: When Efficiency Meets Robustness. In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*.
- Mikael Lindvall, Dharmalingam Ganesan, Ragnar Årda, and Robert E. Wiegand. 2015. Metamorphic Model-based Testing Applied on NASA DAT-an experience report. In *Proceedings of the 37th International Conference on Software Engineering (ICSE)*.
- Zhongxin Liu, Xin Xia, Ahmed E. Hassan, David Lo, Zhenchang Xing, and Xinyu Wang. 2018. Neural-machine-translation-based commit message generation: how far are we?. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering (ASE)*.
- Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective Approaches to Attention-based Neural Machine Translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Lei Ma, Felix Juefei-Xu, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Chunyang Chen, Ting Su, Li Li, Yang Liu, et al. 2018a. Deepgauge: Multi-Granularity Testing Criteria for Deep Learning Systems. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering (ASE)*.
- Shiqing Ma, Yingqi Liu, Wen-Chuan Lee, Xiangyu Zhang, and Ananth Grama. 2018b. MODE: Automated Neural Network Model Debugging via State Differential Analysis and Input Selection. In *Proceedings of the 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*.
- Shiqing Ma, Yingqi Liu, Guanhong Tao, Wen-Chuan Lee, and Xiangyu Zhang. 2019. NIC: Detecting Adversarial Samples with Neural Network Invariant Checking. In *Proceedings of the 26th Annual Network and Distributed System Security Symposium (NDSS)*.
- Fiona Macdonald. 2015. The Greatest Mistranslations Ever. <http://www.bbc.com/culture/story/20150202-the-greatest-mistranslations-ever>
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *Proceedings of the 6th International Conference on Learning Representations (ICLR)*.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient Estimation of Word Representations in Vector Space. In *ICLR Workshops*.



- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed Representations of Words and Phrases and their Compositionality. In *Proceedings of the 29th Conference on Neural Information Processing Systems (NeurIPS)*.
- Pramod K. Mudrakarta, Ankur Taly, Mukund Sundararajan, and Kedar Dhamdhere. 2018. Did the Model Understand the Question?. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Christian Murphy, Gail E. Kaiser, Lifeng Hu, and Leon Wu. 2008. Properties of Machine Learning Applications for Use in Metamorphic Testing. In *Proceedings of the 20th International Conference on Software Engineering and Knowledge Engineering (SEKE)*.
- Arika Okrent. 2016. 9 Little Translation Mistakes That Caused Big Problems. <http://mentalfloss.com/article/48795/9-little-translation-mistakes-caused-big-problems>
- Thuy Ong. 2017. Facebook apologizes after wrong translation sees Palestinian man arrested for posting 'good morning'. <https://www.theverge.com/us-world/2017/10/24/16533496/facebook-apology-wrong-translation-palestinian-arrested-post-good-morning>
- Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. 2016. Distillation as a Defense to Adversarial Perturbations against Deep Neural Networks. In *IEEE Symposium on Security and Privacy*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: A Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics (ACL)*.
- Kexin Pei, Yinzhi 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 (SOSP)*.
- Michael Pradel and Koushik Sen. 2018. DeepBugs: A Learning Approach to Name-based Bug Detection. *Proceedings of the ACM on Programming Languages OOPSLA* (2018).
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Semantically Equivalent Adversarial Rules for Debugging NLP Models. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Sergio Segura, Gordon Fraser, Ana B. Sanchez, and Antonio Ruiz-Cort  s. 2016. A Survey on Metamorphic Testing. *IEEE Transactions on Software Engineering (TSE)* 42 (2016). Issue 9.
- Yikang Shen, Shawn Tab, Alessandro Sordoni, and Aaron Courville. 2019. Ordered Neurons: Integrating Tree Structures into Recurrent Neural Networks. In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to Sequence Learning with Neural Networks. In *Proceedings of the 30th Conference on Neural Information Processing Systems (NeurIPS)*.
- Guanhong Tao, Shiqing Ma, Yingqi Liu, and Xiangyu Zhang. 2018. Attacks Meet Interpretability: Attribute-steered Detection of Adversarial Samples. In *Proceedings of the 34th Conference on Neural Information Processing Systems (NeurIPS)*.
- Wilson L. Taylor. 1953. "Cloze Procedure": A New Tool for Measuring Readability. *Journalism Bulletin* 30, 4 (1953), 415–433.
- 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 (ICSE)*.
- Tree. 2013. Adventures in Mistranslation: HSBC's Call to "Do Nothing". <https://contentequalsmoney.com/mistranslation-hsbc-call-to-do-nothing/>
- Barak Turovsky. 2016. Ten years of Google Translate. <https://blog.google/products/translate/ten-years-of-google-translate/>
- Twitter. 2019. About Tweet translation. <https://help.twitter.com/en/using-twitter/translate-tweets>
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, and Illia Kaiser, Lukasz abd Polosukhin. 2017. Attention is All you Need. In *Proceedings of the 33rd Conference on Neural Information Processing Systems (NeurIPS)*.
- Jingyi Wang, Guoliang Dong, Jun Sun, Xinyu Wang, and Peixin Zhang. 2019. Adversarial Sample Detection for Deep Neural Network through Model Mutation Testing. In *Proceedings of the 41st International Conference on Software Engineering (ICSE)*.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's Neural Machine Translation System: Bridging the Gap Between Human and Machine Translation. *arXiv preprint arXiv:1609.08144* (2016).
- Xiaoyuan Xie, Joshua Ho, Christian Murphy, Gail Kaiser, Baowen Xu, and Tsong Yueh Chen. 2009. Application of Metamorphic Testing to Supervised Classifiers. In *Proceedings of the 9th International Conference on Quality Software (QSIC)*.
- Xiaoyuan Xie, Joshua WK Ho, Christian Murphy, Gail Kaiser, Baowen Xu, and Tsong Yueh Chen. 2011. Testing and Validating Machine Learning Classifiers by Metamorphic Testing. *Journal of Systems and Software (JSS)* 84 (2011). Issue 4.
- Chong Xiong, Charles R. Qi, and Bo Li. 2019. Generating 3D Adversarial Point Clouds. In *Proceedings of the 2019 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Weilin Xu, David Evans, and Yanjun Qi. 2018. Feature Squeezing: Detecting Adversarial Examples in Deep Neural Networks. In *Proceedings of the 25th Annual Network and Distributed System Security Symposium (NDSS)*.



- Dawei Yang, Chaowei Xiao, Bo Li, Jia Deng, and Mingyan Liu. 2019. Realistic Adversarial Examples in 3D Meshes. In *Proceedings of the 2019 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zichao Yang, Zhiting Hu, Yuntian Deng, Chris Dyer, and Alex Smola. 2017. Neural Machine Translation with Recurrent Attention Modeling. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*.
- Biao Zhang, Deyi Xiong, and Jinsong Su. 2017. A GRU-gated Attention Model for Neural Machine Translation. *arXiv preprint arXiv:1807.02340* (2017).
- Biao Zhang, Deyi Xiong, and Jinsong Su. 2018a. Accelerating Neural Transformer via an Average Attention Network. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- 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 (ASE)*.
- Jie M. Zhang, Mark Harman, Lei Ma, and Yang Liu. 2019. Machine Learning Testing: Survey, Landscapes and Horizons. *arXiv preprint arXiv:1906.10742* (2019).
- Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. 2018b. Deeproad: Gan-Based Metamorphic Autonomous Driving System Testing. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering (ASE)*.
- Wujie Zheng, Wenyu Wang, Dian Liu, Changrong Zhang, Qinsong Zeng, Yuetang Deng, Wei Yang, Pinjia He, and Tao Xie. 2018. Testing Untestable Neural Machine Translation: An Industrial Case. *arXiv preprint arXiv:1807.02340* (2018).
- Zhi Quan Zhou and Liqun Sun. 2018. Metamorphic Testing for Machine Translations: MT4MT. In *Proceedings of the 25th Australasian Software Engineering Conference (ASWEC)*.
- Zhi Quan Zhou, Shaowen Xiang, and Tsong Yueh Chen. 2016. Metamorphic Testing for Software Quality Assessment: A Study of Search Engines. *IEEE Transactions on Software Engineering (TSE)* 42 (2016). Issue 3.
- Muhua Zhu, Yue Zhang, Wenliang Chen, Min Zhang, and Jingbo Zhu. 2013. Fast and Accurate Shift-Reduce Constituent Parsing. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 434–443.
- Chris. Ziegler. 2016. A Google self-driving car caused a crash for the first time. <https://www.theverge.com/2016/2/29/11134344/google-self-driving-car-crash-report>