# Testing Untestable Neural Machine Translation: An Industrial Case

Wujie Zheng*, Wenyu Wang†, Dian Liu*, Changrong Zhang*, Qinsong Zeng*,
Yuetang Deng*, Wei Yang‡, Pinjia He§, Tao Xie†

\* Tencent Inc., China

{wujiezheng, dianeliu, chrongzhang, qinzzeng, yuetangdeng}@tencent.com

† University of Illinois at Urbana-Champaign, USA

{wenyu2, taoxie}@illinois.edu

‡ University of Texas at Dallas, USA

wei.yang@utdallas.edu

§ ETH Zurich, Switzerland

pinjiahe@gmail.com

*Abstract*—**Neural Machine Translation (NMT) has been widely adopted recently due to its advantages compared with the traditional Statistical Machine Translation (SMT). However, an NMT system still often produces translation failures due to the complexity of natural language and sophistication in designing neural networks. While in-house black-box system testing based on reference translations (*i.e.*, examples of valid translations) has been a common practice for NMT quality assurance, an increasingly critical industrial practice, named *in-vivo testing*, exposes unseen types or instances of translation failures when real users are using a deployed industrial NMT system. To fill the gap of lacking test oracle for in-vivo testing of an NMT system, in this paper, we propose a new approach for automatically identifying translation failures, without requiring reference translations for a translation task; our approach can directly serve as a test oracle for in-vivo testing. Our approach focuses on properties of natural language translation that can be checked systematically and uses information from both the test inputs (*i.e.*, the texts to be translated) and the test outputs (*i.e.*, the translations under inspection) of the NMT system. Our evaluation conducted on real-world datasets shows that our approach can effectively detect targeted property violations as translation failures. Our experiences on deploying our approach in both production and development environments of WeChat (a messenger app with over one billion monthly active users) demonstrate high effectiveness of our approach along with high industry impact.**

## I. INTRODUCTION

Neural Machine Translation (NMT) [1], [2], [3] has been widely adopted over recent years due to its simple models and satisfactory effectiveness on various machine translation tasks compared with the traditional Statistical Machine Translation (SMT) [4], [5]. Although NMT has shown great advantage and is becoming increasingly popular, in practice it often produces unexpected translation failures in its translations. For example, Google Translate was recently reported to have various translation failures [6] such as producing the same translation for "A is worse than B" and "B is worse than A" [7]. In our own empirical investigation, we also find various translation failures in translations from other highly popular providers of translation services. Past various incidents [8]

show that translation failures could lead to serious consequences. For example, in 2009, a translation failure (mistranslating a catchphrase "Assume Nothing" to "Do Nothing" in various countries) caused the HSBC bank $10 million to repair the caused damage [8].

Aiming to eliminate in-field translation failures for an NMT system, in-house quality assurance techniques such as black-box system testing commonly help NMT developers discover translation failures, subsequently enabling the developers to address these translation failures before the NMT system is deployed to serve the users. In particular, the developers construct test tests from parallel corpora, which are large collections of bilingual text pairs. Each test case corresponds to a bilingual text pair, which consists of (1) the original text (*i.e.*, the text to be translated), being the test input, and (2) one or more reference translations, being the expected test output. Note that reference translations are usually provided by human and are considered as examples of valid translations (*i.e.*, samples of expected outputs from the NMT system). Therefore, the developers cannot simply adopt the traditional test oracle here: checking whether the NMT system's actual output (*i.e.*, the translation under inspection) is equal to one of the reference translations. Instead, the developers typically adopt a special test oracle that (1) calculates the translation quality score (*e.g.*, BLEU score [9]) for measuring multi-granular natural-language similarity between the actual output and the reference translations, and (2) checks whether the calculated translation quality score is equal to or greater than a threshold (either a fixed or varying value), in order to determine whether the test case passes.

While black-box system testing has been shown useful for in-house quality assurance of an industrial NMT system, *in-vivo testing* [10], which executes tests in the production environment, is also becoming increasingly critical in industrial cases. In particular, by leveraging translation requests from real users in the production environment as test inputs, in-vivo testing can expose unseen types or instances of translation failures not revealed by in-house black-box system test cases,

whose quantity and variety are usually much limited. In addition, in-vivo testing enables the developers to monitor and handle translation failures instantly (*e.g.*, by switching to apply a backup model or multiple models of higher quality but with higher cost) in the production environment, improving the overall translation quality on translation requests from real users.

However, during in-vivo testing for an industrial NMT system, existing test-oracle techniques cannot be applied in the production environment. First, due to the lack of samples of expected outputs (*i.e.*, reference translations), we cannot simply apply the special test oracle that relies on the translation quality score (as widely used in black-box system testing) to in-vivo testing. Second, although there have been heuristics [11], [12], [13] on speculating the expected outputs for unlabeled data on classification-oriented or prediction-oriented machine learning systems, these heuristics cannot be directly applied to testing an NMT system due to the intractability of natural languages.

To fill the gap of lacking test oracle for in-vivo testing of an industrial NMT system, in this paper, we propose a new approach for automatically identifying translation failures; our approach can directly serve as a test oracle for in-vivo testing. Different from using the translation quality score, our approach conducts statistical analysis on both the inputs (*i.e.*, the original texts) and outputs (*i.e.*, the translations under inspection) of the NMT system, aiming to directly spot out translation failures without samples of the expected outputs (*i.e.*, the reference translations). Our key insight underlying our approach is that there are some properties of natural language translation that can be checked systematically. If there is any violation of these properties in the translation under inspection, there are likely translation failures. In addition, we can leverage the information from the inputs of the NMT system (*i.e.*, the original texts) to provide hints for finding translation failures, whereas such information is overlooked when the translation quality score is calculated. In this paper, we focus on one property that generally holds for translations: the original text and the translation generally have one-to-one mappings in terms of their constituents, *e.g.*, words/phrases.

In particular, our approach includes two algorithms for detecting two specific types of violations of this property, respectively: (1) *under-translation*, where some words/phrases from the original text are missing in the translation, and (2) *over-translation*, where some words/phrases from the original text are unnecessarily translated multiple times. Based on our experience in deploying NMT models in large-scale industrial settings, many translation failures can be reflected by these two types of violations. For under-translation, our detection algorithm leverages the recommendation algorithm of Item-based Collaborative Filtering [14] to build a word/phrase translation dictionary from massive parallel corpora, and uses such algorithm for examining translations. Our detection algorithm additionally addresses two challenges: efficiently identifying and translating phrases from texts, and recognizing implicit word/phrase translations. For over-translation, our detection

algorithm is based on word frequencies and also considers the cases where multiple occurrences of words/phrases exist in the text to be translated. We implement and evaluate our two algorithms as well as baseline algorithms on manually labeled datasets, which contain sets of original texts, corresponding translations, and corresponding marks of the existence for two types of violations. Our algorithms achieve better overall performance compared with the baseline algorithms.

Our experiences on deploying our approach in both production and development environments of WeChat (a messenger app with over **one billion** monthly active users [15]) demonstrate that our approach is both practical and scalable. By using our approach for in-vivo testing in the production environment, the developers are able to collect translation tasks with unseen types or instances of translation failures and observe the performance of deployed models in real-world usages. The developers are also able to handle the translation failures instantly through switching to backup models (*e.g.*, SMT models), improving the overall translation quality nearly effortlessly. On top of that, our approach is able to process about **12 million** unique translation tasks every day.

In addition, not limited to only in-vivo testing, our approach also helps enhance in-house black-box system testing during the NMT-system development. By using our approach, the developers manage to find outputs from the NMT system (*i.e.*, the translations) that contain translation failures and are previously unable to be processed by the special test oracle based on the translation quality score (due to missing reference translations). The developers are also able to quickly locate and diagnose translation failures when test cases fail, based on the information provided by our approach.

Finally, our approach helps build in-house testing sets containing **130,000** English and **180,000** Chinese words/phrases. With such valuable test sets, we test multiple other highly popular translation services and find various translation failures. The detected translation failures indicate potential defects of the design, implementation, or training data in machine translation systems used by these service providers and help these providers improve their services. All these preceding results demonstrate high effectiveness of our approach along with high industry impact. We also build an online demonstration[1] for our approach on English and Chinese translation tasks.

## II. BACKGROUND

In this section, we introduce background of NMT quality assurance. We first explain why testing neural network (NN) software is a challenging task. Then, we discuss recent work on testing NN software. After that, we introduce two common types of translation-property violations that we focus on (*i.e.*, under-translation and over-translation) with examples. Finally, we present BLEU score [9], one of the most widely used translation quality scores.

---

[1] http://bit.ly/2P4hEB4

## A. Why NN Software is Hard to Test

Although there are various studies on automatically testing traditional software, testing NN software still largely relies on manual and non-systematic strategies. Typically, NN software developers often conduct testing with a randomly selected testing dataset. Because the size of the testing dataset is limited and the testing dataset is constructed in a random manner, these existing strategies are limited at finding and locating potential issues. In addition, it is inappropriate for NN software developers to directly adopt automated testing techniques proposed for traditional software for the following two main reasons.

(1) *Learning instead of implementing program logics* [16]. For traditional software, improving the coverage of program control-flow structures (*e.g.*, path coverage) is the main target of automated testing techniques, because for traditional software, the program logics lie in its control-flow statements. However, the core logics of NN software are embedded in the network structure and parameters, which cannot be revealed by existing coverage metrics based on control-flow structures. In particular, network structure is often defined by only a few lines of code, whose coverage can easily reach 100% without finding any potential issues.

(2) *Non-linearity.* Constraint solving is a crucial component in automated testing techniques for traditional software. Solvers using the Satisfiability Modulo Theory (SMT) [17] have been quite successful for different demands such as generating high-coverage test inputs. However, a neural network might consist of thousands or even millions of neurons with non-linear activation functions. These neurons might also be composed in a variety of ways according to different network structures. Thus, it is significantly more challenging to find the connections between inputs and outputs (or any other intermediate values), making constraint solving nearly infeasible there.

## B. Existing Testing Techniques for NN Software

To address the aforementioned difficulties, a line of research work has been recently proposed on testing NN software, including whitebox testing [11], [18], [19], metamorphic testing [20], [21], mutation testing [22], and empirical studies [23]. Specifically, inspired by path coverage for traditional software, Pei at al. [11] propose *Neuron Coverage*, which evaluates how many neurons in the neural network under test have been activated (*i.e.*, covered) in testing. This novel whitebox testing technique is further applied by Tian et al. [13] to testing NN-based autonomous cars. Based on the core idea of neuron coverage, Ma et al. [18] further propose five fined-grained coverage criteria for whitebox testing of NN software. Sun et al. [19] propose the first concolic testing technique for NN software. Zhang et al. [20] develop a GAN-based technique using the idea of metamorphic testing to synthesize test images with the same label under different weather conditions to test NN-based autonomous cars. Dwarakanath et al. [21] design metamorphic rules to detect implementation faults for image classifiers (*e.g.*, SVM and ResNet [24]). Ma et al. [22] study mutation testing strategies for both network structure and training data.

Although researchers have designed various testing techniques specialized for NN software, these techniques consider only the neural networks whose outputs can be judged using simple equivalence oracles. For example, most of these techniques focus on image classification, where the correctness of NN outputs can be judged by trivially referring to image labels or voting among multiple models. However, for NMT systems, there lacks a rigorous definition on the correctness of its output (*i.e.*, whether the translation is valid) even when some of the expected outputs (*i.e.*, reference translations) are provided. As the result, it is difficult to adapt the existing testing techniques (for classification/prediction-oriented NN systems) to NMT quality assurance.

## C. Translation-Property Violations

As introduced in Section I, we focus on two specific types of violating the one-to-one mapping property for natural language translation: under-translation and over-translation.

*Under-translation.* For a specific translation task, if the translation misses one or more words/phrases from the text to be translated, then the translation is considered under-translated. Table III shows two examples of under-translation. In the first example, the underlined Chinese word corresponding to 'mother' is missing in the English translation, leading to the wrong interpretation of what the person 'she' refers to at the beginning of the second sentence. In the second example, the underlined Chinese word corresponding to 'desk' is missing in the English translation, making readers unaware of the object *desk*.

*Over-translation.* For a specific translation task, if any word/phrase is unnecessarily translated multiple times, then the translation is considered over-translated. Table V shows an example of over-translation, where the underlined Chinese phrase representing 'can never be changed' appears four times in the translation, while there is only one occurrence (as indicated by italicized words) in the original sentence. Such repetition is redundant since it does not change the meaning of the translation.

As can be seen, these two specific types of violations can directly affect the user experience and even lead to misunderstanding of the translation. Also according to our observation, these two specific types of violations can indicate more types of translation failures, including wrong names and incorrect interpretation of numbers. In this paper, we propose two algorithms that can automatically and effectively detect these violations.

## D. Translation Quality Scores

There are various translation quality scores that can be used to measure the overall quality of translations. BLEU (BiLingual Evaluation Understudy) score [9] is one of the most widely-used quality scores for machine translation. BLEU score measures the similarity between the translation under inspection and the reference translation(s). In particular, a

higher BLEU score indicates that translation under inspection is closer to reference translation(s), being considered of higher translation quality. The range of a BLEU score is $[0,1]$, which is often presented as a percentage value (*i.e.*, $[0,100]$). For example, if none of the tokens in the translation appears in any reference translation, the BLEU score is 0. On the contrary, if the translation is exactly the same as any of the reference translations, the BLEU score is 100.

## III. Detection Algorithms for Under- and Over-translation

In this section, we introduce our detection algorithms for under-translation and over-translation. Our algorithms are based on the one-to-one constituent mapping property of natural language translation. Specifically, we first build the mappings between both words and phrases from the source and destination languages on massive training translation task sets (*i.e.*, parallel corpora). Then, the two algorithms can leverage such mappings to examine the property on the translations under inspection.

### A. Building Mappings Between Bilingual Words and Phrases

For the very beginning, based on massive training translation task sets (*i.e.*, parallel corpora), we aim to build the mappings from each word/phrase in the source language to a list of words/phrases in the destination language, where each word/phrase in the destination language is a valid translation for the corresponding word/phrase in the source language. There are two steps to achieve this goal: phrase identification and mapping learning.

**Phrase identification.** The purpose of identifying phrases from texts is to handle the situations where the phrase meanings cannot be conveyed by the combination of word meanings. The identified phrases from texts of both source and destination languages are essential to build an accurate mapping. As a naive approach to identify phrases, we can assume that each phrase contains up to $l_{ph}$ words, extract all 2-grams, 3-grams, ..., $l_{ph}$-grams (note that here an $n$-gram is a contiguous sequence of $n$ words) from the texts, and regard them all as candidate phrases. However, this naive approach would probably require an enormous amount of computation due to a huge number of candidate phrases, whose number would be approximately $\sum_{i=2}^{l_{ph}} |\mathcal{W}_s|^i$ for only the source language, where $\mathcal{W}_s$ is the word vocabulary of the source language. Besides, a phrase might have different forms, adding difficulties to statistical analysis when the whole word sequences are considered.

In order to improve the efficiency make phrase identification most robust, we adopt a heuristic by only considering the representative word pairs (*i.e.*, keyword pairs) of any phrase. In other words, we abstract a phrase $w_1 w_2 \cdots w_n$ with an ordered (potentially not unique) word pair $\langle w_a, w_b \rangle$, where $1 \leqslant a < b \leqslant n$ and the word pair appears in different forms of the phrase. By using the heuristic, we reduce the problem of enumerating all potential candidate phrases from enumerating all short word sequences to enumerating all word pairs in

a sentence. Specifically, assume we want to identify phrases whose keywords are at most $d_{ph}$ words away from each other, then we only need to extract each word pair $\langle w_a, w_b \rangle$ ($a, b$ denotes the word order in the same sentence) of both source and destination languages from the training dataset, where $a < b$ and $b - a \leqslant d_{ph} + 1$. For example, if we set $l_{ph} = 1$, we would obtain at most 5 unique word pairs from a 4-word sentence $w_1 w_2 w_3 w_4$: $\langle w_1, w_2 \rangle, \langle w_1, w_3 \rangle, \langle w_2, w_3 \rangle, \langle w_2, w_4 \rangle, \langle w_3, w_4 \rangle$. Using such heuristic, we reduce the number of candidate phrases to at most $|\mathcal{W}_s|^2$ for only the source language. We further discard any word pair with less than $c_{ph}$ occurrences to reduce false positives (*i.e.*, the word pairs that are not keyword pairs of any phrase). Additionally, *mapping learning* (as introduced later) also helps identify real phrases.

**Mapping learning.** We employ the recommendation algorithm of Item-based Collaborative Filtering [14] to build the word/phrase translation mappings from massive training translation task sets (*i.e.*, parallel corpora). The recommendation algorithm of Item-based Collaborative Filtering was originally used in scenarios such as product and movie recommendation [25], [26], [27]. The algorithm predicts an user's prospective interested items (*e.g.*, movies) by looking at the user's rating history and searching the items that are similar to the user's highly-rated items. The algorithm's key insight of finding similar items is that similar items should have similar ratings from a lot of users. Thus, it determines the similarity between two items by looking at the similarity between their ratings. To reduce our mapping problem to the recommendation problem, we can view each training translation task as a user, each word/phrase appearing the original text (of the source language) or reference translation (of the destination language) as an item with a positive rating, and all other words/phrases as items with negative ratings. In this way, we adapt the recommendation algorithm to uncover the connections between words/phrases in the source and destination languages, *i.e.*, the word/phrase translations. Specifically, we define the user rating matrix $M$ as follows:

$$M_{k,w} = \begin{cases} 1 & \text{if } w \text{ appears in } P_s^k \text{ or } P_d^k \\ 0 & \text{otherwise} \end{cases}$$

where $w$ is a word/phrase, $P_s^k$ and $P_d^k$ denote the original text and the reference translation in the training translation task numbered $k$ correspondingly. Then, we calculate the relevance/similarity between $w_s$ (a word/phrase in the source language) and $w_d$ (a word/phrase in the destination language) using the Cosine similarity [28]:

$$Rel(w_s, w_d) = \frac{\overrightarrow{M_{\cdot, w_s}} \cdot \overrightarrow{M_{\cdot, w_d}}}{||\overrightarrow{M_{\cdot, w_s}}||_2 \cdot ||\overrightarrow{M_{\cdot, w_d}}||_2}$$
$$= \frac{\sum_k M_{k,w_s} M_{k,w_d}}{\sqrt{\sum_k M_{k,w_s}^2} \sqrt{\sum_k M_{k,w_d}^2}}$$

Tables II shows an instance of $M$ for the translation task in Table I. We can calculate $Rel(\text{I}, 我) = ([1,1,1,0] \cdot [1,1,1,0])/(||[1,1,1,0]|| \cdot ||[1,1,1,0]||) = 1$, indicating that we speculate '我' to be the translation of 'I'. In fact, such translation is exactly the correct Chinese translation of 'I'.

TABLE I: Example training translation tasks

| | English (original) | Chinese (translated) |
|---|---|---|
| 1 | I love my family. | 我爱我的家人。 |
| 2 | I have a lovely son. | 我有一个可爱的儿子。 |
| 3 | I have worked for seven years. | 我工作七年了。 |
| 4 | They have a big house. | 他们有一个大房子。 |

TABLE II: User rating matrix $M$ corresponding to Table I

| | I | love | have | $\cdots$ | 我 | 爱 | 有 | $\cdots$ | 。 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | $\cdots$ | 1 | 1 | 0 | $\cdots$ | 1 |
| 2 | 1 | 0 | 1 | $\cdots$ | 1 | 0 | 1 | $\cdots$ | 1 |
| 3 | 1 | 0 | 1 | $\cdots$ | 1 | 0 | 0 | $\cdots$ | 1 |
| 4 | 0 | 0 | 1 | $\cdots$ | 0 | 0 | 1 | $\cdots$ | 1 |

Based on such relevance between words/phrases in the source and destination languages, we build the mappings directly. Specifically, for each word $w_s$ in the source language, we consider a total of $c_{tr}$ words in the destination language with the highest relevance to $w_s$ to be valid translations of $w_s$, where $c_{tr}$ is a predefined threshold value. Such technique has two advantages compared with using generic translation dictionaries. First, the words in an existing generic translation dictionary may be limited. If a word is not in the existing dictionary, then the mapping cannot be constructed. Second, the word translations included in the generic dictionary might also be too limited, too generic, or outdated with respect to various translation tasks, causing the algorithms to miss correct translations.

### B. Detection Algorithm for Under-translation

With the word/phrase translation mappings derived from the previous steps, under-translation detection can now be achieved by checking the existence of word/phrase translations in the whole translation text. Specifically, for each word/phrase in the original text, we obtain the list of its corresponding word/phrase translations, and check if any of these word/phrase translations appear in the whole translation text. We use two real-world examples to illustrate this process. Table III shows two example translations with under-translation violations. The underlined Chinese words (corresponding to 'mother' and 'desk' in English) are missing in the English translation. When our algorithm detects potential violations related to these two words, two translation lists are available for each of them. Specifically, our algorithm produces the contents in Table IV, which is part of the mappings constructed by the previous steps. In the table, 'origin' indicates the words to be translated and 'trans $k$' denotes the $k$-th most relevant translation. Our algorithm goes through each translation for both Chinese words and checks whether it appears in the translation being examined. Then, the algorithm finds that for both Chinese words, none of translations in Table IV appears in the translation being examined (as shown in Table III). The algorithm subsequently marks both translations with potential under-translation violations.

However, there might be some words/phrases of the source language whose translations usually do not show up in the translations of the whole sentences, i.e., some words/phrases

might incur implicit translations. Examples include *the* and *to* in English. The techniques specified earlier can likely produce numerous false positives due to such phenomenon. We introduce *error rate filtering* to address such limitation. Specifically, we calculate the error rate $e_{w_s}$ for each word $w_s$ of the source language (similar for any phrase) on the training dataset using $e_{w_s} = N_{w_s}^e / N_{w_s}$, where $N_{w_s}^e$ denotes the number of sentence pairs that are considered under-translated on $w_s$, and $N_{w_s}$ is the number of translation sentence pairs involving $w_s$. Then, we disregard any under-translation caused by missing $w_s$ on real translation results if we find that $e_{w_s} > e_{th}$, where $e_{th}$ is a predefined threshold value that controls the tolerance of such error-prone $w_s$. Table III shows the technique being used in two example translations, where 'error_rate' shows the error rates for two words in the source language. The algorithm also checks whether the error rates are beyond $e_{th}$ (usually set to be $0.2$ in our production environment) for both Chinese words before marking them with potential under-translation violations. In the examples, the algorithm finds that both error rates are acceptable and confirm the marks. On the contrary, if the algorithm finds that the error rate on one word/phrase is too high, it will skip the word/phrase and continue checking other words/phrases.

### C. Detection Algorithm for Over-translation

Our over-translation detection algorithm is based on frequency of words appearing in the translation. Specifically, we count the occurrences for each word appearing in the translation under inspection, and mark the word as over-translated if more than one occurrence is found and the occurrences are near to each other. In addition, our algorithm includes techniques to address two main challenges being faced.

First, particles such as *have*, *the*, and *to* in English tend to have multiple occurrences and no actual meaning in many well-formed expressions. These words would confuse over-translation detection and cause false positives. To alleviate this challenge, we remove all such common words (*i.e.*, stop words) from the translation under inspection before conducting over-translation detection.

Second, some words/phrases might have multiple occurrences in the original text, and they are probably also supposed to appear multiple times in the translation under inspection. Aiming to differentiate such situations from those with real over-translation violations, our algorithm includes a technique to estimate the number of words/phrases (in the original text) that are translated to each word $w_d$ in the translation under inspection, and compare such number with the number of occurrences of $w_d$. Specifically, for each word $w_d$ in the translation under inspection, we use the word/phrase translation mappings introduced by the previous steps in the other direction to find out a set of words/phrases (of the source language) that can be all translated to $w_d$ (*i.e.*, for each word $w_s$ in the set, $w_d$ is its translation according to the dictionary), and count the number of words/phrases in the original text (denoted as $count_s(w_d)$) that fall in the set. Let

TABLE III: Example translations with under-translation violations

| Chinese (original) | English (translated) | English (reference) |
|---|---|---|
| 三姑给你的红包给你妈妈了，她见了你会给你的。 | Third Aunt gave you a red envelope. She'll give it to you when she sees you. | Third Aunt gave your red envelope to your *mother*. She'll give it to you when she sees you. |
| 放寒假一个月宿舍地上桌上床上全是灰...都不知道该从哪里开始收拾.. | Winter vacation a month dormitory on the floor of the bed is all gray ... I don't know where to start .. | After a month's winter vacation, dust is everywhere on the floor, the *desk*, and the bed in the dormitory... I don't know where to start cleaning... |

TABLE IV: Example translation lists for two under-translated words in Table III

| origin | 妈妈 | 桌 |
|---|---|---|
| error_rate | 0.04116 | 0.07293 |
| trans 1 | mother | table |
| trans 2 | mom | desk |
| trans 3 | mum | papers |
| trans 4 | mama | tables |
| trans 5 | mommy | coffee |
| trans 6 | moms | placed |
| trans 7 | mothers | piled |
| trans 8 | mummy | put |
| trans 9 | my | breakfast |
| trans 10 | her | laid |

TABLE V: Example translation with over-translation violation

| English (original) | Chinese (translated) |
|---|---|
| U have to admit that something *can never be changed*, live with it or break with it! | 你必须承认，有些东西是永远无法改变的，无法改变的，无法改变的，无法改变的！ |

us denote the number of occurrences of $w_d$ in the translation under inspection as $count(w_d)$. Finally, we can consider a $w_d$ to be over-translated if we find $count_s(w_d) < count(w_d)$, indicating that we find more occurrences of $w_d$ than there should be, *i.e.*, redundant translations.

Table V shows an example with an over-translation violation to illustrate our algorithm. The Chinese translation contains 3 more repetitions for the translation of "can never be changed". Our algorithm first discovers that there are 4 instances of "无法" and "改变". Then, the algorithm looks up the word/phrase translation mappings to find whether there is any word/phrase in the original text with such translations. The algorithm finds that "can never" and "changed" have such translations. Then it also finds that these two English words/phrases appear only once in the original text, fewer than the occurrences in the translation under inspection. Finally our algorithm marks the translation with an over-translation violation.

## IV. EVALUATION

In this section, we present our evaluation on assessing our proposed approach. Specifically, we try to answer the following two research questions (RQs):

*RQ1*: How accurate are the two proposed violation detection algorithms (included in our approach) in testing benchmark datasets consisting of real-world translation tasks?

*RQ2*: How useful is our approach in assisting our NMT system improvement during both in-vivo testing and in-house testing?

*RQ3*: How effective is our approach on testing NMT systems from other public service providers?

We next discuss the evaluation setup (Section IV-A) and present the accuracy comparison and analysis results (Section IV-B) *w.r.t.* baseline algorithms, for RQ1. We then share our experience of deploying our proposed approach in both the production and development environment of WeChat, a messenger app with over one billion of monthly active users (Section IV-C), for RQ2. We also demonstrate the effectiveness of our approach on testing NMT systems from other public service providers by showing the instances of detected translation failures in the translations from those systems (Section IV-D), for RQ3.

### A. Evaluation Setup

We implement our proposed detection algorithms (denoted as 'proposed') from scratch, and evaluate our proposed algorithms. We also evaluate baseline algorithms on both under- and over-translation detection for comparison with our proposed algorithms:

*Algorithms based on primitive dictionary looking-up* (denoted as 'std-dict'). For under-translation, we look up each word/phrase of the text to be translated in an existing generic translation dictionary and check whether any translation of the word appears in the translation. For over-translation, we replace the learned word/phrase translation mappings in our algorithm with the existing generic translation dictionary. We use the software *StarDict* [29] to provide existing generic translation dictionaries and implement such baseline algorithm by ourselves.

*Algorithms based on word alignment* [30], [31] (named as 'word-align') from the traditional SMT models. For under-translation, we let the word-alignment model provide a list of candidate translations for each word of the text to be translated, and check whether any translation of the word appears in the translation. For over-translation, we replace the learned word/phrase translation mappings in our algorithm with the candidate translation lists provided by the word-alignment model. Note that the candidate translation list is also derived by choosing the translations with top $c_{tr}$ alignment probabilities. We use the software *fast_align* [32] with default settings for the word-alignment model implementation.

We build the word/phrase translation mappings (for our algorithms) and train the word-alignment model (for the baseline algorithm) with 16 million sentence pairs crawled from various sources (*e.g.*, example word/phrase usages in dictionaries) on the Internet. We evaluate all the algorithms

TABLE VI: Overview of evaluation datasets

| Name | Lang | Type | $\#_{all}$ | $\#_{ws}$ | $\#_U$ | $\#_O$ |
|---|---|---|---|---|---|---|
| ench_news | en-cn | News | 200 | 7497 | 54 | 4 |
| chen_news | cn-en | News | 200 | 7418 | 31 | 8 |
| ench_oral | en-cn | Oral | 300 | 3237 | 42 | 1 |
| chen_oral | cn-en | Oral | 300 | 2918 | 37 | 5 |

on datasets that are randomly sampled from larger benchmark datasets (crawled online and independent of training datasets) and are manually labeled by us. Each dataset contains a number of sentence pairs, each of which consists of the sentence to be translated (*i.e.*, the input to our NMT system) and the translation under inspection (*i.e.*, the corresponding output of our NMT system). Each sentence pair is also marked with the existence of two types of violations in the translation identified by manual inspection. Table VI shows the overview of the evaluation datasets, where $\#_{all}$ denotes the total number of sentence pairs, $\#_{ws}$ denotes the total number of words of the source language, and $\#_U$ and $\#_O$ denote the numbers of sentence pairs flagged with under- and over-translation violations, respectively. 'en-cn' and 'cn-en' indicate translating from English to Chinese and from Chinese to English, respectively.

For any algorithm on any dataset, let $\mathcal{S}$ be the manually labeled sentence pair sets containing a specific violation type (*i.e.*, ground truth), and let $\mathcal{S}'$ be the sentence pair sets identified by the algorithm detecting the specific violation type, then *precision*, *recall*, and *F-measure* can be calculated as:

$$\text{precision} = \frac{|\mathcal{S} \cap \mathcal{S}'|}{|\mathcal{S}'|} \quad \text{recall} = \frac{|\mathcal{S} \cap \mathcal{S}'|}{|\mathcal{S}|}$$

$$\text{F-measure} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

### B. Effectiveness Comparison and Analysis

**General comparison.** For the algorithms under comparison, Table VII shows the result summary of under-translation detection, and Table VIII shows the result summary of over-translation detection. The precision, recall, and F-measure are abbreviated as 'P', 'R', and 'F', respectively, in both tables. In Table VII, 'Count' indicates the number of under-translated words identified by the three algorithms under comparison. We set $c_{tr} = 10$ and $c_{ph} = 10$ for all the experiments, $e_{th} = 0.15$ for experiments on the 'ench_news' and 'chen_news' datasets, and $e_{th} = 0.25$ for experiments on the 'ench_oral' and 'chen_oral' datasets. We turn on phrase identification for 'en-cn' tasks while keeping it off for 'cn-en' tasks due to the fact that most Chinese words are already phrases after word segmentation. As shown in Tables VII and VIII, our proposed algorithm achieves the highest F-measures on all the datasets, compared with the two baseline algorithms ('std-dict' and 'word-align'). Such result indicates the effectiveness of our proposed algorithm.

For under-translation detection, as shown in Table VII, the higher effectiveness of our proposed algorithm mainly comes from higher precisions. Our algorithm's precisions are up to about $2.6\times$ of those by the baseline algorithms
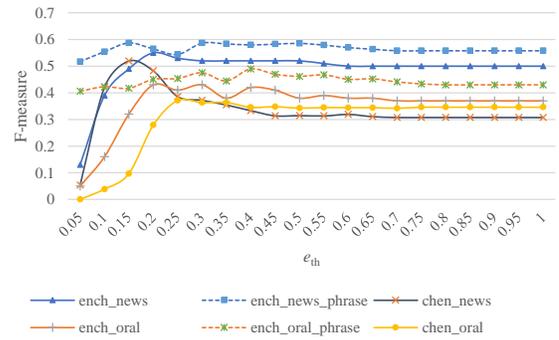


Fig. 1: Trends of F-measures of the detection results on all datasets with different $e_{th}$ values

('std-dict' and 'word-align'). Such result indicates the high accuracy of our combined techniques for under-translation detection. We also notice that precisions achieved by 'word-align' are higher than those by 'std-dict'. Such result shows the necessity of word/phrase translation mapping learning, which helps avoid too generic or limited word/phrase translations. Meanwhile, we find that our algorithm has relatively lower recalls, potentially as the cost of achieving higher accuracy. On the contrary, 'std-dict' achieves the highest recalls on all the datasets. However, such scores are accompanied by a large number of false positives, making it difficult to leverage the detection results.

For over-translation detection, as shown in Table VIII, the higher effectiveness of our proposed algorithm also mainly comes from higher precisions. However, different from the results of under-translation detection, rankings of recalls vary across the datasets. Such result might be due to a relatively small number of over-translation instances in the datasets.

**Analysis of Error Rate Thresholds** $e_{th}$**.** Changing the value of $e_{th}$ might influence the detection results. Thus, we further investigate how different threshold values would affect the overall effectiveness of our algorithm. As stated in Section III-B, $e_{th}$ controls the tolerance of error-prone words in the proposed under-translation detection algorithm. We also know from the definition that $e_{th} \in (0, 1]$. Aiming to understand the influence of $e_{th}$ on the algorithm effectiveness, we run the under-translation algorithm on all the datasets with $e_{th}$ set to $\{0.05, 0.1, \cdots, 0.95, 1\}$. Figure 1 shows the trends of F-measure of the detection results under different $e_{th}$ values. Note that we run the algorithm on 'en-ch' datasets both when phrase identification is enabled and disabled (corresponding to 'ench_news', 'ench_news_phrase', 'ench_oral', and 'ench_oral_phrase' in Figure 1), while on 'ch-en' datasets the feature is always turned off (corresponding to 'chen_news' and 'chen_oral' in Figure 1).

As shown in Figure 1, the trends of F-measures on the detection results of different datasets are very similar under the same setting. Such result indicates the commonality of under-translation detection in various scenarios. However, differences still exist. As shown in Figure 1, F-measures on the detection results of 'ench_news', 'ench_news_phrase', and 'chen_news' peak when $e_{th}$ is around $0.15 \sim 0.2$. In contrast,

TABLE VII: Results of under-translation detection on the evaluation datasets

| | proposed | | | | std-dict | | | | word-align | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P | R | F | Count | P | R | F | Count | P | R | F | Count |
| ench_news | **0.51** | 0.69 | **0.59** | 113 | 0.28 | **1.00** | 0.43 | 1853 | 0.35 | 0.85 | 0.49 | 250 |
| chen_news | **0.43** | 0.65 | **0.52** | 50 | 0.16 | **1.00** | 0.27 | 1827 | 0.18 | **1.00** | 0.30 | 523 |
| ench_oral | **0.52** | 0.40 | **0.45** | 32 | 0.15 | **1.00** | 0.26 | 689 | 0.20 | 0.79 | 0.32 | 309 |
| chen_oral | **0.30** | 0.49 | **0.37** | 70 | 0.12 | **0.97** | 0.22 | 895 | 0.14 | 0.73 | 0.24 | 259 |

TABLE VIII: Results of over-translation detection on the evaluation datasets

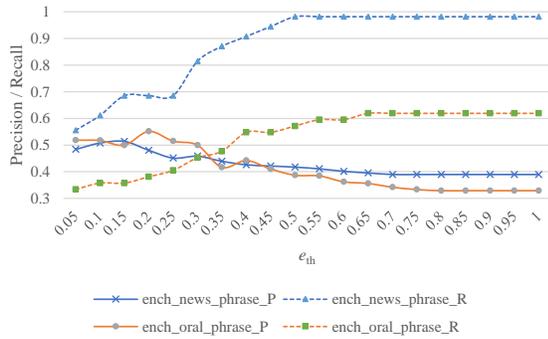| | proposed | | | std-dict | | | word-align | | |
|---|---|---|---|---|---|---|---|---|---|
| | P | R | F | P | R | F | P | R | F |
| ench_news | **0.38** | 0.75 | **0.50** | 0.13 | 0.50 | 0.20 | 0.24 | **1.00** | 0.38 |
| chen_news | **0.73** | **1.00** | **0.84** | 0.13 | **1.00** | 0.23 | 0.40 | **1.00** | 0.57 |
| ench_oral | **0.33** | **1.00** | **0.50** | 0.00 | 0.00 | 0.00 | 0.14 | **1.00** | 0.25 |
| chen_oral | **0.80** | 0.80 | **0.80** | 0.38 | **1.00** | 0.56 | 0.28 | **1.00** | 0.43 |



Fig. 2: Trends of precisions and recalls of the detection results on 'en-ch' datasets with different $e_{th}$ values



Fig. 3: Statistics of unique English-Chinese translations with translation-property violations identified by our algorithms
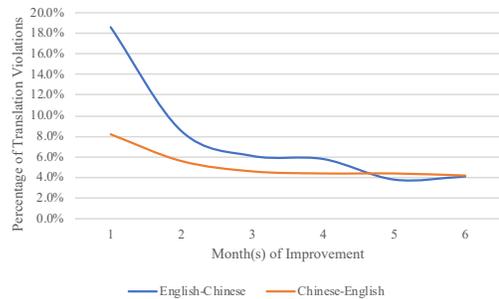


Fig. 4: Combined statistics of under- and over-translation violations in real-world English and Chinese translations by our NMT model

for 'ench_oral' and 'chen_oral', F-measures peak when $e_{th}$ is around $0.2 \sim 0.25$, and for 'ench_oral_phrase' the peak value is even higher. This finding is expected: expressions in news reports tend to be more rigorous than those in casual conversations, limiting the scope of word/phrase meanings and thus making it easier to achieve precise translations. Such finding also shows that our algorithm can be easily tuned on different scenarios. In addition, the algorithm effectiveness is worse at $e_{th} = 1$ for all curves, where error rate filtering is not in use. Such finding indicates that error rate filtering contributes to the algorithm effectiveness.

However, as shown in Figure 1, the influence on F-measures of the detection results is much smaller when phrase identification is on, although the overall F-measures are higher than those in the situation where phrase identification is off. A potential explanation is that disabling phrase identification causes more false positives on words that actually belong to phrases, given that words can have totally different meanings in phrases. Filtering the detection results based on error rates helps reduce these false positives and consequently improve the F-measures. Such finding indicates that phrase identification could actually contribute to the algorithm effectiveness. To understand why the F-measures change insignificantly, in Figure 2, we show the trends of precisions and recalls with different $e_{th}$ values for 'en-ch' datasets when phrase identification is enabled. As can be seen from Figure 2, as $e_{th}$ goes up, precisions gradually become lower while recalls steadily increase. The trend is intuitive: less strict filters allow

more true positives that contribute to higher recalls and miss more false positives, which lead to lower precisions. Such finding still shows that error-rate filtering is useful: one of our goals for detecting violations in translations is to build in-house test suites for regression testing and manual inspection. Keeping test suites small and precise saves time and manual efforts from running or inspecting false positives.

### C. Experience of Deployment

By using our approach for in-vivo testing in the production environment, the developers are able to collect translation tasks with unseen types or instances of translation failures and observe the performance of deployed models in real-world

usages. The developers are also able to handle the translation failures instantly through switching to backup models (*e.g.*, SMT models), improving the overall translation quality nearly effortlessly. Figure 3 shows an example 2-day statistics of unique English-Chinese translations in which our algorithms identify any translation-property violation, where Y-axis corresponds to the statistics within each period of 5 minutes. As can be seen, our approach is able to find over 8 translations with translation-property violations each second during the busiest periods. On top of that, our algorithms are able to process about 12 million unique translation tasks every day, with over 200 translations processed each second during the busiest periods. Such result indicates the good performance and applicability of our algorithms and approaches on real-world tasks. The statistics also show that these translation-property violations are non-trivial in the real world, and in-vivo testing on NMT systems is of great significance.

In addition, not limited to only in-vivo testing, our approach also helps enhance in-house black-box system testing during the NMT-system development. By using our approach, the developers manage to find outputs from the NMT system (*i.e.*, the translations) that contain translation failures and are previously unable to be processed by the special test oracle based on the translation quality score (due to missing reference translations). Our result also shows the violations decrease significantly after the deployment of our tool. Such result suggests that the developers are able to quickly locate and diagnose translation failures when test cases fail, based on the information provided by our approach. Figure 4 shows the combined statistics of under- and over-translation violations in real-world English and Chinese translations by our NMT model within 6 months after the deployment of our approach. As can be seen, the percentage of translations with under- or over-translation violations decreases significantly over time. More specifically, for English-Chinese translations, the percentage drops from 18.6% to 4.1%, while for Chinese-English translations, the percentage drops from 8.2% to 4.2%. Such statistics reflect the effectiveness of our approach in helping development and improvement of a machine translation system.

### D. Large-scale Test Suite Generation and Test Results

Our algorithms also help collect 130,000 English and 180,000 Chinese meaningful words/phrases (*i.e.* words/phrases with low error rates by learning from training data). These words are used as in-house test cases for testing and improving WeChat's continuously-improved machine translation model. Using our approach not only helps diagnose issues in WeChat's NMT system but also helps diagnose issues in other competing machine translation systems released by other providers. We show some examples of translation-property violations in English and Chinese translations provided by systems from various providers in Table IX. By analyzing the translations, we gain some insights on potential causes of those violations, showing that our approach is able to detect potential defects lying in the design, implementation, or training data in machine translation systems, with examples as below.

TABLE IX: Example issues found in machine translation systems by various providers

| Provider Name | Original Text | Given Translation | Expected Translation |
|---|---|---|---|
| Prvd. A | 成人 | mature people | adult |
| Prvd. A | 太好了 | what fun | great |
| Prvd. B | large-scale | large-scale | 大规模 |
| Prvd. B | long-term | long-term | 长期 |
| Prvd. B | U.S. | U.S. | 美国 |
| Prvd. C | 蛋糕 | Runeberg torte | cake |
| Prvd. C | 酸奶 | Viili | yoghurt |
| Prvd. D | 疟原虫 | p. | plasmodium |
| Prvd. D | 酶原 | The original enzyme | zymogen |

For *Provider A*, translations in both examples seem to be too "straightforward", *i.e.*, two Chinese words are translated character by character instead of being as a whole. There might be an issue in the design of the model, causing the model to disregard some context information.

For *Provider B*, three English phrases with dots or hyphens remain untranslated. It might also be the problem of model implementation that causes the model to disregard such phrases.

For *Providers C and D*, translations for common food names and terminologies are too specific (*e.g.*, *Runeberg torte* is one type of cakes) or inaccurate (*e.g.*, *p.* is the abbreviation of *plasmodium*, but there are too many words beginning with the letter p). Both providers might need to investigate into their training data.

## V. RELATED WORK

Tu et al. [33] propose coverage-based NMT models, which keep track of the attention history by maintaining coverage vectors. The coverage mechanism assists attention adjustment to provide more chances for untranslated source words. Both under- and over-translation issues are addressed by the coverage mechanism. The coverage mechanism is also adopted by Google NMT models [34]. However, such coverage mechanism only alleviates (but cannot eliminate or detect) under-translation and over-translation issues from the perspective of NMT model design. It does not reveal defects lying in the implementation or the training data of the model. Our approach is applicable on general machine translations, aiming to provide common-issue detection regardless of specific translation models, and such common-issue detection could be useful for model development and improvement. In other words, even when coverage-based NMT models are used, our approach is still needed and applicable to detect remaining issues. In addition, our approach can be used even when other NMT models (not being coverage-based ones) are used.

Our detection algorithms make use of the Item-based Collaborative Filtering [14] algorithm to construct the word/phrase translation dictionary. Linden et al. [25] describe how to use the algorithm to recommend products to Amazon customers, and Davidson et al. [26] use a similar approach to recommend videos to YouTube users. We show that such algorithm is also suitable for recommending words/phrases in the scenario of translation quality assurance. It is also possible to apply

more advanced recommendation algorithms, such as those summarized by Adomavicius and Tuzhilin [27].

## VI. CONCLUSION

In this paper, we present a novel approach for in-vivo testing of an NMT system. Our approach automatically identifies translation failures without requiring reference translations. Our approach focuses on properties of natural language translation that can be checked systematically and uses information from both the test inputs (*i.e.*, the texts to be translated) and the test outputs (*i.e.*, the translations under inspection) of the NMT system. Our evaluation conducted on real-world datasets shows that our approach can effectively detect targeted property violations as translation failures. Our experiences on deploying our approach in both production and development environments of WeChat, a messenger app with over one billion monthly active users, demonstrate high effectiveness of our approach along with high industry impact.

## REFERENCES

[1] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin, "A neural probabilistic language model," *Journal of machine learning research*, vol. 3, no. Feb, pp. 1137–1155, 2003.

[2] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.

[3] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in neural information processing systems*, 2014, pp. 3104–3112.

[4] P. F. Brown, V. J. D. Pietra, S. A. D. Pietra, and R. L. Mercer, "The mathematics of statistical machine translation: Parameter estimation," *Computational linguistics*, vol. 19, no. 2, pp. 263–311, 1993.

[5] P. Koehn, F. J. Och, and D. Marcu, "Statistical phrase-based translation," in *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1*. Association for Computational Linguistics, 2003, pp. 48–54.

[6] D. Hofstadter. (2018) The shallowness of Google Translate. [Online]. Available: https://www.theatlantic.com/technology/archive/2018/01/the-shallowness-of-google-translate/551570/

[7] DuanJiShu. (2018) Fudan and SJTU, PKU and THU, which one is better? Google Translate has the answer. (Chinese version). [Online]. Available: https://baijiahao.baidu.com/s?id=1597259364260803730

[8] A. Okrent. (2016) 9 little translation mistakes that caused big problems. [Online]. Available: http://mentalfloss.com/article/48795/9-little-translation-mistakes-caused-big-problems

[9] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: A method for automatic evaluation of machine translation," in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ser. ACL '02. Stroudsburg, PA, USA: Association for Computational Linguistics, 2002, pp. 311–318. [Online]. Available: https://doi.org/10.3115/1073083.1073135

[10] C. Murphy, G. Kaiser, I. Vo, and M. Chu, "Quality assurance of software applications using the in vivo testing approach," in *2009 International Conference on Software Testing Verification and Validation*. IEEE, 2009, pp. 111–120.

[11] K. Pei, Y. Cao, J. Yang, and S. Jana, "Deepxplore: Automated whitebox testing of deep learning systems," in *Proceedings of the 26th Symposium on Operating Systems Principles*. ACM, 2017, pp. 1–18.

[12] S. Srisakaokul, Z. Wu, A. Astorga, O. Alebiosu, and T. Xie, "Multiple-implementation testing of supervised learning software," in *Proc. AAAI-18 Workshop on Engineering Dependable and Secure Machine Learning Systems (EDSMLS)*, 2018.

[13] Y. Tian, K. Pei, S. Jana, and B. Ray, "Deeptest: Automated testing of deep-neural-network-driven autonomous cars," in *Proceedings of the 40th International Conference on Software Engineering*. ACM, 2018, pp. 303–314.

[14] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Item-based collaborative filtering recommendation algorithms," in *Proceedings of the 10th international conference on World Wide Web*. ACM, 2001, pp. 285–295.

[15] R. Hollander. (2018) WeChat has hit 1 billion monthly active users. [Online]. Available: https://www.businessinsider.com/wechat-has-hit-1-billion-monthly-active-users-2018-3

[16] W. Yang and T. Xie, "Telemade: A testing framework for learning-based malware detection systems," in *Proc. AAAI-18 Workshop on Engineering Dependable and Secure Machine Learning Systems (EDSMLS)*, 2018.

[17] L. De Moura and N. Bjørner, "Satisfiability modulo theories: Introduction and applications," *Commun. ACM*, vol. 54, no. 9, pp. 69–77, Sep. 2011. [Online]. Available: http://doi.acm.org/10.1145/1995376.1995394

[18] L. Ma, F. Juefei-Xu, F. Zhang, J. Sun, M. Xue, B. Li, C. Chen, T. Su, L. Li, Y. Liu, J. Zhao, and Y. Wang, "Deepgauge: multi-granularity testing criteria for deep learning systems," in *ASE'18: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, 2018.

[19] Y. Sun, M. Wu, W. Ruan, X. Huang, M. Kwiatkowska, and D. Kroening, "Concolic testing for deep neural networks," in *ASE'18: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, 2018.

[20] M. Zhang, Y. Zhang, L. Zhang, C. Liu, and S. Khurshid, "Deeproad: Gan-based metamorphic testing and input validation framework for autonomous driving systems," in *ASE'18: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, 2018.

[21] A. Dwarakanath, M. Ahuja, S. Sikand, R. M. Rao, R. P. J. C. Bose, N. Dubash, and S. Podder, "Identifying implementation bugs in machine learning based image classifiers using metamorphic testing," in *ISSTA'18: Proceedings of 27th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2018.

[22] L. Ma, F. Zhang, J. Sun, M. Xue, B. Li, F. Juefei-Xu, C. Xie, L. Li, Y. Liu, J. Zhao, and Y. Wang, "Deepmutation: Mutation testing of deep learning systems," in *ISSRE'18: Proceedings of 29th IEEE International Symposium on Software Reliability Engineering*, 2018.

[23] Y. Zhang, Y. Chen, S.-C. Cheung, Y. Xiong, and L. Zhang, "An empirical study on tensorflow program bugs," in *ISSTA'18: Proceedings of 27th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2018.

[24] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR'16: Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.

[25] G. Linden, B. Smith, and J. York, "Amazon.com recommendations: Item-to-item collaborative filtering," *IEEE Internet computing*, vol. 7, no. 1, pp. 76–80, 2003.

[26] J. Davidson, B. Liebald, J. Liu, P. Nandy, T. Van Vleet, U. Gargi, S. Gupta, Y. He, M. Lambert, B. Livingston *et al.*, "The youtube video recommendation system," in *Proceedings of the fourth ACM conference on Recommender systems*. ACM, 2010, pp. 293–296.

[27] G. Adomavicius and A. Tuzhilin, "Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions," *IEEE transactions on knowledge and data engineering*, vol. 17, no. 6, pp. 734–749, 2005.

[28] P. Dangeti, *Statistics for Machine Learning*. Packt Publishing Ltd, 2017.

[29] H. Zheng. (2015) StarDict - The best dictionary program in linux and windows. [Online]. Available: http://stardict-4.sourceforge.net/

[30] P. Koehn, *Statistical machine translation*. Cambridge University Press, 2009.

[31] S. Vogel, H. Ney, and C. Tillmann, "Hmm-based word alignment in statistical translation," in *Proceedings of the 16th conference on Computational linguistics-Volume 2*. Association for Computational Linguistics, 1996, pp. 836–841.

[32] C. Dyer, V. Chahuneau, and N. A. Smith, "A simple, fast, and effective reparameterization of ibm model 2." Association for Computational Linguistics, 2013.

[33] Z. Tu, Z. Lu, Y. Liu, X. Liu, and H. Li, "Modeling coverage for neural machine translation," *arXiv preprint arXiv:1601.04811*, 2016.

[34] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey *et al.*, "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.