

S²LCT: Specification- and Syntax-based Automated Testing Linguistic Capabilities of NLP Models

Anonymous Author(s)

ABSTRACT

Natural language processing has been widely used for developing a variety of text analytic applications with high reliability. For any NLP-based applications to be trustworthy, assessment through measuring the accuracy of the holdout datasets is needed. Recent research introduced methods that assess an NLP model on a set of linguistic capabilities. However, we observe that state-of-the-art methods rely on manually created test cases, which may be limited in terms of scalability and diversity. In this work, we introduce S²LCT, an automated linguistic capability-based testing technique. Given a set of linguistic capabilities that assess an NLP model, S²LCT first searches for suitable seed inputs from existing datasets, and then generates a sufficient number of new test inputs by expanding the seed inputs based on their context-free grammar (CFG). S²LCT guarantees the correctness of the generated test input labels and that these test inputs belong to the correct linguistic capabilities. Evaluation results show that S²LCT can generate test cases that are 100% more diverse than state-of-the-art while ensuring the correctness of labels and linguistic capabilities of the test inputs. We also show that S²LCT facilitates the identification of critical failures and its origins in the NLP models.

ACM Reference Format:

Anonymous Author(s). 2023. S²LCT: Specification- and Syntax-based Automated Testing Linguistic Capabilities of NLP Models. In *Proceedings of The 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2023)*. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Natural language processing (NLP) applications are growing rapidly. As a result, trustworthiness in how the quality of NLP applications is assessed becomes critical for its practical use in the real world. Traditionally, the quality of NLP models is assessed using aggregated metrics. In particular, accuracy (i.e., the fraction of outputs that the model correctly predicts) is the most widely used metric for assessing the quality of classification models: higher accuracy suggests better quality of a model. However, all NLP models have their strengths and weaknesses; using a single, aggregated metric (i.e., accuracy) makes it difficult for the users to assess the capabilities of NLP models. Such a metric fails to validate how well a model

supports linguistic capabilities [8, 11] (e.g., when the model always fails on certain type of inputs).

JL: To address this limitation, increasing number of evaluation methodologies has been conducted to measure and evaluate the a NLP model in different aspects. First, the NLP model can be fragile to adversarial attacks, and it decreases trustworthiness of the NLP model. Thus, robustness on adversarial examples are evaluated by generating the examples using perturbations to an input text. [1, 21, 26, 41, 47]. Second, fairness, specific bias in an NLP model toward certain subpopulation groups such as gender and races, has been measured [38, 40]. Additionally, methodologies to assess an NLP model on a set of linguistic capabilities have been introduced [32, 33]. Typically, it describes certain type of inputs and outputs observed in real world for the target NLP task. A linguistic capability includes comprehensive functionalities of the input and output behavior for an NLP model, and various task-specific model behaviors are evaluated compared with the aforementioned evaluation metrics. Therefore, Testing the linguistic capabilities allows the model developers to better understand the capabilities and potential issues of the NLP models. For example, CHECKLIST [32], a behavioral testing framework for evaluating NLP models, defines a linguistic capability called “Negated neutral should still be neutral” to measure how accurately a sentiment analysis model understands that the negated neutral input has neutral sentiment [32]. It requires the sentiment analysis model to output neutral sentiment on the negated neutral inputs. Such evaluation methodology avoids the overestimation of the model performance as it measures the model performance on each functionality. In the end, testing through linguistic capabilities provides not only the overall model performance, but also the malfunction facets of the model.

However, existing linguistic capability-based testing approaches generate inputs only relying on word substitution, and need to be preset before generating tests. As a result, these templates are limited in their sentence structures. This restricts existing approaches’ ability to comprehensively test linguistic capabilities.

To address this issue, we present S²LCT, an automated linguistic capability-based testing framework for NLP models. The goal of S²LCT is to generate a diverse linguistic capability-based test suite for a given model, both in terms of covering diverse linguistic structures and diverse model behaviors. To achieve the goal, S²LCT needs to address two challenges.

(i) *Capability-based text generation*. Each test case should be automatically categorized into a linguistic capability; and (ii) *Automated Test Oracle*. The label of each test case should be automatically and accurately defined.

Capability-based text generation. The suitability of test sentences for evaluating NLP models on a linguistic capability is determined by its relevancy to this linguistic capability. It is challenging to maintain relevancy between a test sentence and its linguistic capability when transforming and expanding the test sentence. This

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

ESEC/FSE 2023, 11 - 17 November, 2023, San Francisco, USA

© 2023 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

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

is because the linguistic capability is defined on a specific mixture of syntax and semantics of the sentence. Due to the inherent ambiguity of natural language sentences, there exists no automatic way to check the consistency of each sentence with the semantics and structures specified in the corresponding linguistic capability. S²LCT defines specifications (search rules and transformation templates) from linguistic capabilities, and generates seed test cases by searching and transforming test cases that conform to the specifications. In addition, S²LCT analyzes the parse tree of the seed sentence to identify possible expansion and validates coherence of linguistic capability relevancy between the seed and its expansion. **Automated Test Oracle.** For NLP model testing, test oracle is usually determined by understanding the semantics of texts and requires domain knowledge of different NLP tasks. Thus, the current testing practice requires manual efforts to create the test oracles of the holdout data. The manual work is costly in terms of time consumption. Therefore, it necessitates automated test oracle generation for improving the testing process of NLP models. However, automatically generating the test oracles remains one of the main challenges in NLP software testing [15]. A few metamorphic testing approaches have been proposed in image recognition domain, but they require understanding the characteristics of metamorphic relations between inputs and outputs, which require domain expertise and non-trivial manual efforts [36]. In addition, it is even more challenging to design metamorphic relations in textual data because the semantics of nature language sentences can be greatly changed even by a slight perturbation to the sentences. Motivated by software metamorphic testing approach [36], we address the challenge by only accepting validated test oracles by the linguistic capability specification and linguistic rules determined by input-output relations.

In this work, as a first step, we consider sentiment analysis and hate speech detection as the NLP tasks for the models under test. S²LCT obtains the appropriate test oracle by implementing domain-specific knowledge on the word sentiment dataset and word suggestion model for validating the generated sentence and its oracle.

We demonstrate the generality of S²LCT and its utility as a NLP model evaluation tool by evaluating three well-known sentiment analysis models: BERT-base [7], RoBERTa-base [22] and DistilBERT-base [35].

We made the following contributions in this work:

- We design and implement an automated testing tool, S²LCT, and compared it with CHECKLIST. We find that the test cases generated by S²LCT achieve 100% higher coverage than CHECKLIST when used for testing sentiment analysis models.
- We perform a manual study to measure the correctness of the sentiment labels and their linguistic capabilities produced by S²LCT. We find that S²LCT generates test cases that consistently label their sentiment correctly as human.
- We analyze the root causes of misclassification in the sentiment analysis models, guided by S²LCT testing result. We find that seeds and expanded test cases produced by S²LCT are rather useful in helping developers understand bugs in the model.

```

1  air_noun = [
2      'flight', 'seat', 'pilot', 'staff', ...
3  ]
4  pos_adj = [
5      'good', 'great', 'excellent', 'amazing', ...
6  ]
7  neg_adj = [
8      'awful', 'bad', 'horrible', 'weird', ...
9  ]
10 t = editor.template('{it} {air_noun} {be} {pos_adj}.',
11                     it=['The', 'This', 'That'], be=['is', 'was'], labels
12                     =2, save=True)
13 t += editor.template('{it} {be} {a:pos_adj} {air_noun}.',
14                     it=['It', 'This', 'That'], be=['is', 'was'], labels=2,
15                     save=True)
16 t += ...
17 t += editor.template('{it} {be} {a:neg_adj} {air_noun}.',
18                     it=['It', 'This', 'That'], be=['is', 'was'], labels=0,
19                     save=True)

```

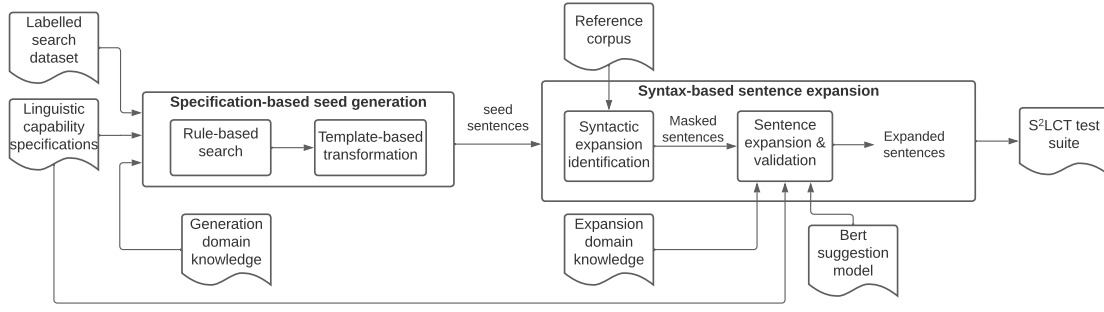
Figure 1: Example of CHECKLIST templates on the linguistic capability of “Sentiment change over time, present should prevail”.

2 BACKGROUND

In this section, we provide a brief background of CHECKLIST’s [32] approach of test case generation via an example. CHECKLIST introduces task-dependent linguistic capabilities for monitoring model performance on each linguistic capability. It assumes that the linguistic phenomena can be represented in the behavior of the model under test, because each linguistic capability specifies the desired behavior between model inputs and outputs. For each linguistic capability, CHECKLIST manually creates test case templates and generates sentences by filling in the values of each placeholder.

We show an example of CHECKLIST templates in Figure 1. These templates are used for evaluating the sentiment analysis model on the linguistic capability of “Sentiment change over time, present should prevail”. For the templates defined at lines 10 to 16, they have placeholders such as *it*, *air_noun*, *pos_adj*. Values for the placeholders are defined at lines 1 to 9. For each template, CHECKLIST fills in all the combinations of the values of placeholders to generate sentences under this linguistic capability. For example, sentences such as “The flight is good” and “That airline was happy” are generated using the template at line 10.

Despite the simple approach used for test case generation, CHECKLIST has limitations. First, it relies on manual work for defining template structure and its placeholder values, making the test case generation a costly process. Moreover, such manual work generates a limited number of templates, produces similar test cases, and induces bias in the test cases. These limitations motivated the design of our approach.

Figure 2: Overview of S²LCT.

3 SPECIFICATION- AND SYNTAX-BASED LINGUISTIC CAPABILITY TESTING

We design and implement a new NLP model testing method, *Specification- and Syntax-based Linguistic Capability Testing (S²LCT)*, which automatically generates test cases with oracles to assess the quality of NLP models.

Figure 2 depicts an overview of S²LCT, which consists of two phases. The *specification-based seed generation* phase performs rule-based searches from a real-world dataset and template-based transformation to obtain the initial seed sentences. The search rules (e.g., search for neutral sentences that do not include any positive or negative words) and transformation templates (e.g., negating a sentence) are defined in the *linguistic capability specifications* (Table 1). Using the specification will result in seed sentences that are most likely to conform to a specific linguistic capability and is labelled correctly.

The *syntax-based sentence expansion* phase expands the seed sentences with additional syntactic elements (i.e., words) to cover many more real-world syntactic structures. It first performs a syntax analysis to identify the part-of-speech (PoS) tags that can be inserted to each seed, by comparing the PoS parse trees between the seed sentence and many other sentences from a large reference dataset. Each identified tag is inserted into the seed as a *mask*. It then uses an NLP recommendation model (i.e., BERT [7]) to suggest possible words. Finally, the resulting sentence is validated to be consistent with the specification which additionally defines the rules for expansion (e.g., the expanded word should be neutral).

We now describe each phase of S²LCT in detail.

3.1 Specification-based Seed Generation

The seed generation phase of S²LCT starts by searching sentences in a real-world dataset, which match the rules defined in the linguistic capability specification, and then transforming the matched sentences using templates to generate seed sentences that conform to individual linguistic capabilities. The reasons for this design choice are twofold. First, while it is generally infeasible to judge which linguistic capability any sentence falls into and which label it should have, there exist simple rules and templates that allow for classifying the resulting sentences into individual linguistic capabilities and with the correct labels, with high confidence. This enables us to test each linguistic capability individually. Second, searching from a real-world dataset ensures that the sentences used as test

cases for testing linguistic capabilities are realistic and diverse. The diverse test cases are more likely to achieve a high coverage of the target model’s functionality in each linguistic capability, thus detecting more errors.

Specifically, S²LCT first searches and selects sentences applicable to the linguistic capability in a given real-world dataset with search rules. In case the search rules only fulfill part of the linguistic capability specifications (i.e., LC3-LC10 in Table 1), we then transform the selected sentences into seed sentences using the heuristic templates.

Table 1 shows the search rules and the transformation templates of all 10 linguistic capabilities we implemented in S²LCT. The first column shows the linguistic capability type and its description, and the second column shows the search rule and transformation template used in each linguistic capability. For LC1 and LC2, the NLP models are evaluated in the scope of short sentences with selective sentiment words. It does not require any transformation because the search rules alone are sufficient to find sentences that conform to the linguistic capabilities. On the other hand, search rules of LC3 to LC10 are not enough to match their linguistic capability specification; thus, S²LCT additionally uses templates to transform the searched sentences to match the corresponding linguistic capability. For example, in LC3’s transformation template, the searched sentences are appended as part of a sentence to generate a seed. In LC4’s transformation template, the searched demonstrative sentences are negated.

3.2 Syntax-based Sentence Expansion

The simple search rules and transformation templates used to generate the seed sentences may limit the syntactic structures these seeds may cover. To address this limitation, the syntax-based sentence expansion phase extends the seed sentences to cover syntactic structures commonly used in real-life sentences. Our idea is to differentiate the parse trees between the seed sentences and the reference sentences from a large real-world dataset. The extra PoS tags in the reference parse trees are identified as potential syntactic elements for expansion and inserted into the seed sentences as masks. We then use masked language model to suggest the fill-ins. If the resulting sentences still conform to the linguistic capability specification, they are added to S²LCT’s test suite.

3.2.1 Syntax Expansion Identification. Algorithm 1 shows how masks are identified for each seed sentence. It takes the parse trees

Table 1: Search rules and transformation templates for linguistic capabilities of sentiment analysis task.

Linguistic capability	Search rule and transformation template
LC1: Short sentences with neutral adjectives and nouns	Search seed={length: < 10; include: neutral adjs & neutral nouns; exclude: pos adjs & neg adjs & pos nouns & neg nouns; label: neutral} Transform N/A
LC2: Short sentences with sentiment-laden adjectives	Search seed={length: < 10; include: pos adjs; exclude: neg adjs & neg verbs & neg nouns; label: pos} {length: < 10; include: neg adjs; exclude: pos adjs & pos verbs & pos nouns & neg verbs & neg nouns; label: neg} Transform N/A
LC3: Sentiment change over time, present should prevail	Search pos_sent={label: pos, length: < 20}, neg_sent={label: neg, length: < 20} Transform seed=[['Previously, I used to like it saying that','Last time, I agreed with saying that','I liked it much as to say that']+[pos_sent neg_sent]+['but', 'although', 'on the other hand']+['now I don't like it', 'now I hate it.']] [['I used to disagree with saying that','Last time, I didn't like it saying that','I hated it much as to say that']+[neg_sent, pos_sent]+['but', 'although', 'on the other hand']+['now I like it.']]
LC4: Negated negative should be positive or neutral	Search demonstrative_sent={start: [This, That, These, Those] + [is, are]; label: neg} Transform seed=negation of demonstrative_sent (['is'] → ['is not', 'isn't'], ['are'] → ['are not', 'aren't'])
LC5: Negated neutral should still be neutral	Search demonstrative_sent={start: [This, That, These, Those] + [is, are]; label: neutral} Transform negation of demonstrative_sent
LC6: Negation of negative at the end, should be positive or neutral	Search neg_sent={label: neg} Transform seed=[['I agreed that', 'I thought that']+[neg_sent]+['but it wasn't', 'but I didn't']]
LC7: Negated positive with neutral content in the middle	Search pos_sent={length: < 20; label: pos}, neutral_sent={length: < 20; label: neutral} Transform seed=[['I wouldn't say', 'I do not think', 'I don't agree with,']+[neutral_sent] +['']+[pos_sent]]
LC8: Author sentiment is more important than of others	Search pos_sent={label: pos}, neg_sent={label: neg} Transform seed=[{temp1}+[pos_sent]+{temp2}+[neg_sent]] [{temp1}+[neg_sent]+{temp2}+[pos_sent]] where temp1=[['Some people think that', 'Many people agree with that', 'They think that', 'You agree with that'], temp2=[['but I think that']]
LC9: Parsing sentiment in (question, yes) form	Search pos_sent={label: pos}, neg_sent={label: neg} Transform seed=[['Do I think that', 'Do I agree that']+[pos_sent neg_sent]+['? yes']]
LC10: Parsing sentiment in (question, no) form	Search pos_sent={label: pos}, neg_sent={label: neg} Transform seed=[['Do I think that', 'Do I agree that']+[pos_sent neg_sent]+['? no']]

of the seeds, generated by the Berkeley Neural Parser [18, 19], and a reference context-free grammar (CFG) from the Penn Treebank corpus dataset [25] as inputs. This reference CFG was learned from a large dataset [39] that is representative of the distribution of real-world language usage. The algorithm identifies the discrepancy between the seed syntax and the reference grammar to decide how a seed can be expanded.

Algorithm 1 Syntax expansion identification algorithm.

```

1: Input: Parse trees of seed sentences  $S$ , reference context-free grammar  $R$ 
2: Output: Set of masked sentences  $M$ 
3: for each part tree  $s$  from  $S$  do
4:   for each production  $s\_prod$  from  $s$  do
5:      $s\_lhs = s\_prod.lhs$ 
6:      $s\_rhs = s\_prod.rhs$ 
7:     for each  $r\_rhs$  from  $R[s\_lhs]$  do
8:       if  $s\_rhs \subset r\_rhs$  then
9:          $M = M \cup insertMask(r\_rhs - s\_rhs, s)$ 
10:      end if
11:    end for
12:  end for
13: end for
14:
15: return  $random(M, k)$ 

```

For each production in each seed's parse tree (lines 3 and 4), we extract its non-terminal at the left-hand-side (line 5), s_lhs , and the grammar symbols at the right-hand-side (line 6), s_rhs . In line 7, the algorithm iterates through all productions in the reference context-free grammar and matches these that have the same non-terminal at the left-hand-side as s_lhs . The right-hand-side of each matched production is called r_rhs . If s_rhs consists of a subset of the grammar symbols in r_rhs (line 8), the additional symbols in the r_rhs are inserted as masks in the parse tree of the seed sentence, in their respective positions in the expanded production. The left to right traversal of the leaves of an expanded parse tree forms a masked sentence. Due to a potential large number of masked sentences created by this algorithm, we randomly select k masked sentences for the next sentence expansion and validation phase (line 14).

Running example. Figure 3 shows an example using Algorithm 1 to generate a masked sentence. The sentence "Or both." is a seed of the linguistic capability of "Short sentences with neutral adjectives and nouns". The tree on the left shows the parse tree of this seed; it consists of two productions: " $FRAG \rightarrow [CC, NP, .]$ " and " $NP \rightarrow [DT]$ ". When matching the left-hand-side non-terminal of the second production (i.e., " NP ") in the reference CFG, we found that it includes a production " $NP \rightarrow [DT, NNS]$ " which has an additional symbol "NNS" on the right-hand-side. The algorithm thus expands the parse tree with this symbol, shown in the second

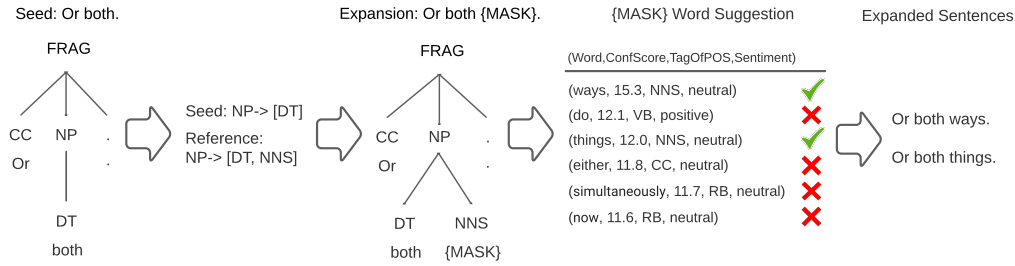


Figure 3: Example of sentence expansion.

tree. The masked sentence “Or both {MASK}.” is the result of the left-to-right traversal of this expanded parse tree.

3.2.2 Sentence Expansion and Validation. In this phase, the words to fill in the masks in the masked sentences are suggested by the BERT model [7]. BERT is a transformer-based natural language model that is pre-trained on two tasks: masked token prediction and next sentence prediction. BERT model is capable of suggesting words for the masked token according to its surrounding context in a sentence. For each masked token, multiple words may be suggested, ranked by their confidence scores. Because BERT model is not aware of the linguistic capability specification and the grammar symbol in the expanded parse tree, an expanded sentence using the suggested words may no longer satisfy the linguistic capability specification. Therefore, we perform validation on the suggested words and only accept them if the following three criteria are met.

First, the PoS tag of the suggested word must match the PoS tag of the expanded symbol in the parse tree. For the example in Figure 3, the masked symbol is a “NNS” (i.e., plural noun); thus, the suggested word must also be a “NNS”. We use SpaCy [14], a free open-source library for natural language processing, to extract the PoS tag for each suggested word. Second, the sentiment of the expanded sentence must be the same as the seed sentence. To ensure this, we require the suggested words be neutral. Third, we additionally verify that the expanded sentences satisfy the same search rules for the seed sentences in LC1 and LC2. This criteria cannot be applied to other linguistic capabilities because they have additional transformation templates.

Running example. The third step in Figure 3 shows the words suggested by BERT. For this masked sentence, BERT suggested six words. Each word is associated with the confidence score provided by BERT, the PoS tag, and the sentiment. Among the six words, only “ways” and “things” are validated by S²LCT because they have the Pos tag “NNS” and are neutral. In addition, both sentences still meet the search rule of the associated linguistic capability of “Short sentences with neutral adjectives and nouns”. In the end, the syntax-based sentence expansion results in two sentences, “Or both ways.” and “Or both things.”, from the seed “Or both.”.

4 EXPERIMENTAL SETUP

In this section, we present the setup of our experiments to evaluate the effectiveness of S²LCT. We answer the following research questions (RQs):

RQ1 : Test Efficiency. Can S²LCT generated test cases evaluate an NLP model more efficiently?

RQ2 : Test Oracle. Can S²LCT generate test sentences with correct sentiment labels?

RQ3 : Capability Testing. Are S²LCT generated test sentences correctly categorized into a linguistic capability?

4.1 Experimental Subjects

NLP Models & Dataset. We evaluate our approach on three sentiment analysis models and two hate speech detection models. All are learning-based models released from the Hugging Face centralized model hub:¹ BERT-base (textattack/bert-base-uncased-SST-2), RoBERTa-base (textattack/bert-base-uncased-SST-2), and DistilBERT-base (distilbert-base-uncased [?]) for sentiment analysis task, and dehate-BERT (Hate-speech-CNERG/dehatebert-mono-english [?]) and tweeter-RoBERTa (cardiffnlp/twitter-roberta-base-hate) for hate speech detection.

These models were pre-trained on English language using a masked language modeling (MLM) objective, and were fine-tuned on the sentiment analysis and hate speech detection task. In this experiment, we used the SST [37] and HateXplain [?] datasets, for sentiment analysis and hate speech detection respectively, for searching the seeds in S²LCT. SST is a corpus of movie review, consisting of 11,855 sentences with sentiment scores. We split the scores into ranges [0, 0.4], (0.4, 0.6] and (0.6, 1.0] to assign them negative, neutral and positive labels, respectively. HateXplain is a hate speech dataset collected from Twitter and Gab. It consists of 20,148 sentences (9,055 from Twitter and 11,093 from Gab). Each data has sentence, class of hate speech, target community, and rationale for the labelling decision. In addition, SentiWordNet [3] is a publicly available dataset for English sentiment lexicons. We used SentiWordNet as both the generation and the expansion domain knowledge as shown in Figure 2.

Comparison Baseline. We compared S²LCT with CHECKLIST,² a manual template and linguistic capability based approach to generate test cases. In this evaluation, we used CHECKLIST’s sentiment analysis test cases that are generated from its publicly available Jupyter Notebook implementation. JL: *Do we need to explain why we used checklist only as a baseline?*

¹<https://huggingface.co/models>

²<https://github.com/marcotcr/checklist>

4.2 Experimental Process

RQ1. Recall that CHECKLIST relies on significant manual efforts and may not generate comprehensive test cases in a linguistic capability. S²LCT, instead, automatically generates test cases based on a search dataset and the syntax in a large reference corpus. We expect S²LCT can generate more efficient test suite than CHECKLIST for finding errors with regards to quantity and diversity of test suite.

We first evaluate the three sentiment analysis models and two hate speech detection models by testing them on the S²LCT test cases, reporting number of test cases and each model’s failure rate in each linguistic capability. Specifically, for each linguistic capability, we expand test cases from all generated seeds. We test the five models using both seed and expanded test cases.

We use two metrics to compare the diversity between the S²LCT generated seeds and the CHECKLIST test cases.

Self-BLEU. We reuse the input diversity metric, called Self-BLEU, that Zhu et al. introduced [48]. BLEU evaluates the token-level similarity. The Self-BLEU is defined as the average BLEU scores over all reference sentences. A higher Self-BLEU score indicates less diversity in the test suite. In the experiment, we collected 50, 100 and 200 randomly selected S²LCT seeds and reported the median Self-BLEU score over three trials for each group of seeds.

Production rule coverage. We propose a new metric to evaluate the syntactic diversity of the generated test suite. It is defined as the number of production rules used in a set of test sentences. In our experiments, we used the Berkeley Neural Parser [18, 19] to parse and collect all the production covered in a set of test sentences. We compared the production rule coverage between 50, 100 and 200 randomly selected S²LCT seeds and the CHECKLIST test cases.

In addition, we follow the approach presented by Ma et al. [23], where the authors measure the coverage of NLP model intermediate states as corner-case neurons. Because the matrix computation of intermediate states impacts NLP model decision-making, a test suite that covers a greater number of intermediate states can represent more NLP model decision-making, making it more diverse. Specifically, we used two coverage metrics by Ma et al. [23], *boundary coverage* (BoundCov) and *strong activation coverage* (SActCov), to evaluate the test suite diversity. It is worth noting that a test sample with a statistical distribution similar to the training data is rarely found in the corner case region. Thus, covering a larger corner case region indicates that the test suite is more likely to be buggy.

$$\begin{aligned} \text{UpperCorner}(\mathcal{X}) &= \{n \in N | \exists x \in \mathcal{X} : f_n(x) \in (\text{high}_n, +\infty)\}; \\ \text{LowerCorner}(\mathcal{X}) &= \{n \in N | \exists x \in \mathcal{X} : f_n(x) \in (-\infty, \text{low}_n)\}; \end{aligned} \quad (1)$$

Equation 1 defines the corner-case neuron of the NLP model $f(\cdot)$, where \mathcal{X} is the given test suite, N is the number of neurons in model $f(\cdot)$, $f_n(\cdot)$ is the n^{th} neuron’s output, and high_n and low_n are the n^{th} neurons’ output bounds on the model training dataset. Equation 1 can be interpreted as the collection of neurons that emit outputs beyond the model’s numerical boundary.

$$\begin{aligned} \text{BoundCov}(\mathcal{X}) &= \frac{|\text{UpperCorner}(\mathcal{X})| + |\text{LowerCorner}(\mathcal{X})|}{2 \times |N|} \\ \text{SActCov}(\mathcal{X}) &= \frac{|\text{UpperCorner}(\mathcal{X})|}{|N|} \end{aligned} \quad (2)$$

The definition of our coverage metrics is shown in Equation 2, where BoundCov measures the coverage of neurons that produces outputs exceeding the upper or lower bounds, and SActCov measures the coverage of neurons that creates outputs exceeding the lower bound. Higher coverage indicates the test suite is better for triggering the corner-case neurons, thus better test suite diversity.

To answer RQ1, for each NLP model under test, we first feed its training dataset to compute each neuron’s lower and upper bounds. After that, we select the same number of test cases from S²LCT and CHECKLIST as the test suite and compute the corresponding coverage metrics.

RQ2 and RQ3. JL: need to do the manual study for the hate speech detection task also? As described in Section 3, S²LCT generates test cases in two steps: specification-based seed generation and syntax-based sentence expansion. These automated steps may generate seed/expanded sentences marked with incorrect sentiment labels or categorized into wrong linguistic capabilities. To answer RQ2 and RQ3, we performed a manual study to measure the correctness of the sentiment labels and linguistic capabilities associated with the seed/expanded sentences, produced by S²LCT.

In the manual study, we randomly sample 100 S²LCT seed sentences, each of which has at least one expanded sentence, and divide these seeds to two sets (i.e., 50 in each set). For each sampled seed sentence, we randomly obtain one of its expanded sentences. This forms the two sets of sentences (200 sentences in total) we use for this study, each with 50 seeds and 50 corresponding expanded sentences. We recruited three participants for this study; all of them are graduate students with no knowledge about this work. 2 of them were assigned one set of sentences and the third was assigned the other set. Each participant was asked to provide two scores for each sentence. (1) *Relevancy score between sentence and its associated linguistic capability*: this score measures the correctness of S²LCT linguistic capability categorization. The scores are discrete, ranging from 1 (“strongly not relevant”) to 5 (“strongly relevant”). (2) *sentiment score of the sentence*: this score measures the sentiment level of the sentence. It is also discrete, ranging from 1 (“strongly negative”) to 5 (“strongly positive”). We measure the following:

$$\text{Label_consistency} = \frac{1}{\#Sample} \cdot \sum_i \delta(\text{label}_{S^2LCT} = \text{label}_{human}) \quad (3)$$

$$\text{LC_relevancy}_{AVG} = \frac{1}{\#Sample} \cdot \sum_i \text{Norm}(\text{LC_relevancy}_i) \quad (4)$$

Equation 3 represents the percentage of the test cases that S²LCT and the participants produce the same sentiment labels. High value of this metric indicates S²LCT generates test cases with correct labels. Equation 4 represents the average of the normalized relevancy

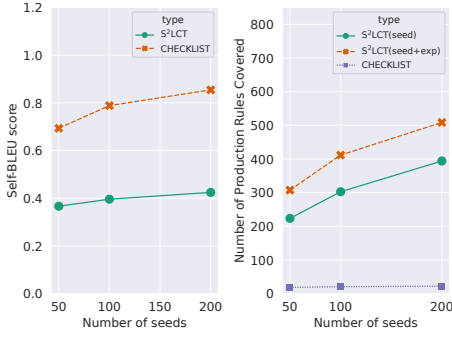


Figure 4: Results of Self-BLEU (left) and Production Rule Coverage (right) of S²LCT and CHECKLIST test cases.

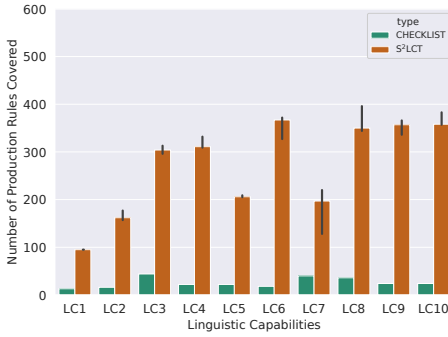


Figure 5: Results of production rule coverage between S²LCT with 50 seeds and CHECKLIST test cases.

score between a sentence and its associated linguistic capability. The relevancy score is to evaluate the relevancy of the sentence to be used for testing model under test on the corresponding linguistic capability. Higher average score means indicate the linguistic capability categorization by S²LCT is correct. We answer RQ2 and RQ3 using the metrics defined by Equation 3 and Equation 4, respectively.

5 EXPERIMENTAL RESULTS

This section presents the experimental results and the answers to the RQs. More results are available at <https://github.com/csresearcher27/s2lct-project>.

5.1 RQ1: Test Case Diversity

Our results show that *S²LCT generated many test cases that the sentiment analysis models failed to predict the correct labels, and it produced significantly more diverse test cases than CHECKLIST did.*

Self-BLEU. Left in Figure 4 compares the Self-BLEU scores between S²LCT seed sentences and CHECKLIST test cases. The x-axis shows sizes of random samples of S²LCT seeds and CHECKLIST test cases, and y-axis shows the Self-BLEU scores. Median of the Self-BLEU scores over all linguistic capabilities is shown in left in figure 4. We observe that Self-BLEU scores of CHECKLIST test cases is significantly higher than those of S²LCT seeds, for all numbers of seeds. This result indicates CHECKLIST test cases are less

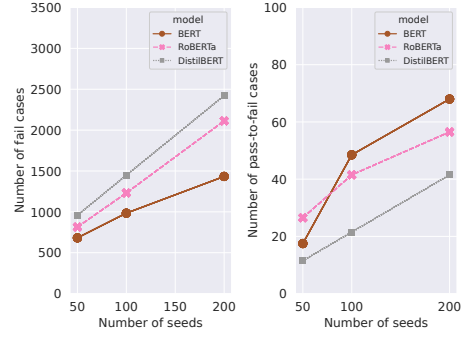


Figure 6: Number of failure (left) and Pass-to-Fail (right) results of NLP models on S²LCT test cases over numbers of seeds.

diverse than S²LCT seeds, which demonstrates the benefit of searching from a real-world search dataset over creating test cases from limited number of preset templates.

Production Rule Coverage. Right in figure 4 compares production rule coverage between the S²LCT seed sentences, S²LCT seed and expanded sentences, and the CHECKLIST test cases. The x-axis shows the sizes of random samples of S²LCT seed sentences, S²LCT seed and expanded sentences, and CHECKLIST test cases, and y-axis shows the production rule coverage scores. Median of the production rule coverage scores over all linguistic capabilities is shown in the right in the figure 4. The figure 4 shows that S²LCT seed and/or expanded sentences produce significantly higher production rule coverage scores than CHECKLIST test cases did. Also, more S²LCT seeds cover more production rules.

Figure 5 compares the production rule coverage between 50 S²LCT seeds and all CHECKLIST test cases for each linguistic capability. The x-axis shows each one linguistic capability, and the y-axis is the production rule coverage for these test cases. We observed that even with only 50 seeds in each linguistic capability, S²LCT seeds cover significantly higher number of production rules than all CHECKLIST test cases (ranging from 95 for LC1 to 367 for LC6 for S²LCT seeds and from 13 for LC1 to 44 for LC3 for CHECKLIST seeds) in each linguistic capability.

Overall, the above results show that S²LCT test cases are significantly more diverse in terms of syntactic structure than CHECKLIST.

Model Test Results. Table ?? shows the testing results of the three NLP models on S²LCT test cases using 50 seeds. First column lists linguistic capabilities for the sentiment analysis task, and Column 2 shows the numbers of seed test cases, and Column 3 shows the median numbers of expanded test cases over 3 trials. Columns 4 and 5 show the median numbers of failed test cases and the failure rates (i.e., percentage of test cases that a model predicts incorrect labels) of each NLP models on the seed and expanded test cases over the 3 trials, respectively. Column 6 shows the median number of expanded test cases that failed, but their corresponding seed test cases passed over the 3 trials(Pass-to-Fail). We observe that in all linguistic capabilities, S²LCT produces hundreds of test cases,

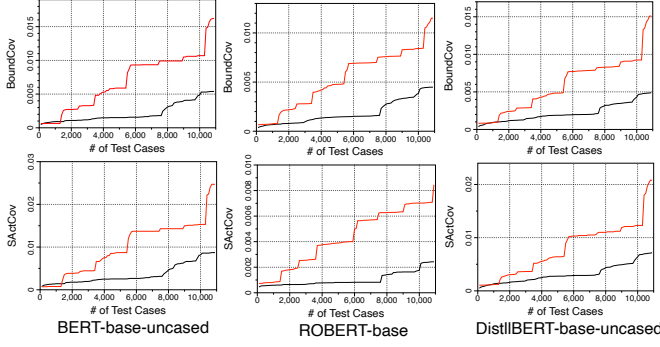


Figure 7: Coverage results of S^2LCT and CHECKLIST test cases.

expanding at least an order of magnitude more test cases than the seeds. LC1 and LC5 have 19 and 26 seeds, respectively. This is because S^2LCT 's search rules and transformation templates of these two linguistic capabilities produced few seeds. Nevertheless, the syntax-based sentence expansion phase generated of 224 and 1090 test cases, respectively. In terms of model performance, all three models achieve low failure rates in LC2 and LC9 while the failure rates are high in all other linguistic capabilities (13.52%-99.95%). We also observe that there are many expanded test cases that failed, but their corresponding seeds did not (last column). This shows that the syntax-based sentence expansion phase indeed introduces more diverse sentence structures in the test cases that cause the models to fail.

In addition, Figure 6 shows the numbers of fail and pass-to-fail cases with different number of seeds used. We observed that more S^2LCT seeds introduce more failed cases and more Pass-to-Fail cases, demonstrating the benefit of using more seeds when resource permits.

Model Coverage. Figure 7 shows the coverage results of S^2LCT and CHECKLIST test cases. The red line represents S^2LCT coverage and the black line represents CHECKLIST coverage. Each column in Figure 7 represents the results for one NLP model. The first row is the *BoundCov* results and the second row is the *SActCov* results.

We made three observations. First, for *all* experimental settings (i.e., NLP model and coverage metric), S^2LCT achieves higher coverage than CHECKLIST. Recall that a higher coverage implies the test cases are more diverse and do not have a similar statistical distribution to the model training data. As a result, a test suite with greater coverage complements the model training data distribution (i.e. holdout testing data) better. For example, for the first NLP model under test, S^2LCT can achieve a higher coverage than CHECKLIST with only half the number of test cases. This result confirms that S^2LCT can generate more diverse test cases to complement the holdout dataset for testing NLP models. Second, as the number of test cases increases, the test suite can achieve better coverage. Such observation is intuitive. However, generating a more extensive test suite is not easy, particularly for CHECKLIST, which is a manually template-based approach. Third, for each NLP model, there is no fixed relationship between *BoundCov* and *SActCov*. In other words, while a test suite may produce higher *BoundCov* for some models,

the same test suite may get higher *SActCov* for other NLP models. Recall that *BoundCov* measures both the upper and lower corner neurons and *SActCov* measures only the upper corner neurons. Such observation implies that the upper and lower corner neurons are distributed unevenly, and measuring only one of them is not enough.

5.2 RQ2: Correctness of Sentiment Labels

Table 5 shows results of our manual study. The first column distinguishes the seed and expanded sentences. The number of test cases used for the study is represented in the second column. The label consistency score defined in Equation 3 is shown in column 3.

We observe that S^2LCT generates test cases that consistently label their sentiment correctly. Column 3 shows that the label consistency scores are 0.83 and 0.84 for the seed and expanded sentences, respectively. This means that S^2LCT generates test oracles consistent with human understanding most of the time. Also, there is little difference of the scores between the seed and expanded sentences. This implies that the syntax-based sentence expansion in S^2LCT preserves the sentiment as its seed.

5.3 RQ3: Correctness of Linguistic Capability Categorization

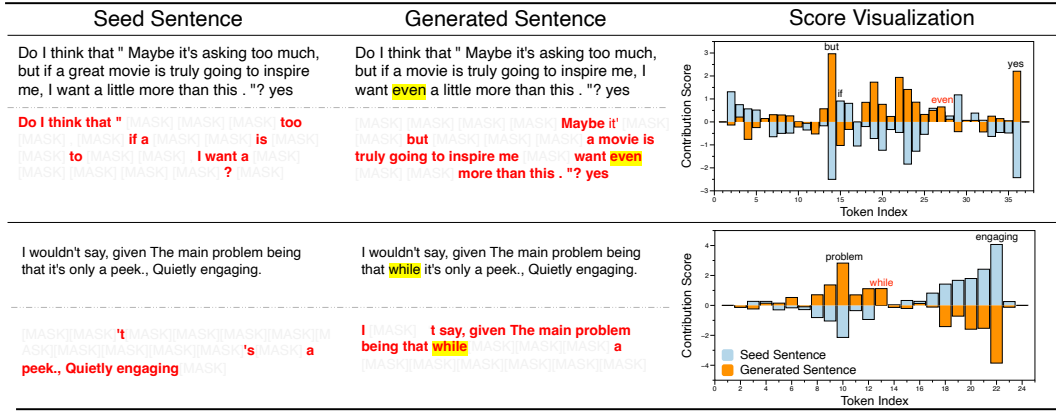
The linguistic capability relevancy score defined in Equation 4 is shown in Column 4 of Table 5. The result shows that S^2LCT generates test cases that are correctly categorized to the corresponding linguistic capabilities most of the time. The linguistic capability relevancy scores for the seed and expanded sentences are both 0.9, achieving high order of agreement with human assessment. The fact that the expanded sentences generated by S^2LCT also have the same level of linguistic capability relevancy as the seed sentences shows that the syntax-based sentence expansion retains the linguistic capabilities.

6 APPLICATION OF S^2LCT

In this section, we use one case study to show S^2LCT can be useful to help developers to find root causes of bugs in the sentiment analysis models.

Experimental Process. we conduct experiments to demonstrate that S^2LCT can help developers understand the bugs in the NLP models. Recall that S^2LCT generates test cases by mutating seed sentences (e.g. by expanding one token in the seed input). Still, it is unclear why mutating one token will cause the model to produce misclassified results. We seek to help developers understand why such mutation will result in the misclassification. Existing work [5, 10, 29] has demonstrated that the ML model prediction is dominated by a minimal set of input features (i.e. tokens in input sentences). Motivated by such intuition, we identify a minimal set of input tokens that dominate the model prediction.

Formally, given an input sentence $x = [tk_1, tk_2, \dots, tk_n]$, and the NLP model under test $f(\cdot)$, our goal is to find a masking template $T = [t_1, t_2, \dots, t_n]$, where t_i is 0 or 1, representing masking the i^{th} token (i.e. tk_i) in x or not. The template T can mask some tokens in x with attribute tokens, and the masked input has a high probability

Figure 8: Visualization of the buggy reason of two S²LCT generated test cases.

of retaining the original prediction x , denoted as

$$P(f(T(x)) = f(x)) \geq P_{\text{thresh}} \quad (5)$$

To create such a template T , we first compute the contribution score of each input token using an existing explainable ML technique [10]. We then begin with the full mask template (i.e., all tokens are masked); such full mask template definitely does not satisfy Equation 5. We then iteratively shift one position from mask to non-mask based on the order of each token’s contribution score, until the template T satisfies Equation 5. Because we iterate the size of the mask, the generated template T will keep the minimum number of tokens in x . Moreover, since the input x is an incorrect prediction, the generated template T is likely to produce misclassification (i.e., the probability to be misclassified is larger than P_{thresh}).

Results. We generate a template that dominates the NLP model prediction to assist developers in understanding the false predictions. Figure 8 shows the generated templates for two randomly selected seeds and their corresponding expanded test inputs. The first example tests “sentiments in (question, yes) form” (LC9), and the second example tests “negative positive with neutral content in the middle” (LC7). The first column shows the seed sentence, the second shows the expanded sentence, and the third shows each token’s contribution score. The blue bar indicates the score for seed inputs, whereas the orange bar reflects the score for the expanded sentences. We highlight the mutated token with yellow background and generated templates with red text.

The results show that after mutating the seed sentence with one token, the token set that dominates the NLP model prediction has changed. We can trace the root causes to the bias of the training data on the linguistic capability under test; as a result, for the linguistic capability under test, the model has a bias towards positive/negative for certain token sequence patterns. For example, LC9 has a bias toward the token sequence pattern that includes “maybe ... but ... even... yes”. Thus, adding the token “even” to the seed sentence will match one of those biased sequence patterns. Sentences with such pattern in the training dataset are dominantly positive; thus, the models make the wrong decision on the sentence with “even” as positive. The visualization of each token’s contribution score in the third column confirms our observation. Once “even” is added, scores of other tokens such as “but” and “yes” all change from

negative to positive. To fix the issue for LC9, we need to add more negative training samples with the format of “maybe ... but ... even ... yes”.

7 THREATS TO VALIDITY

There are two potential threats to validity. First, the dataset used in our study might not be representative of all English grammatical structures and word sentiments. We mitigate this threat by using widely used dataset in the NLP domain [16]. Second, using ranges of the sentiment scores in SST-2 to determine the sentiment labels of search dataset may introduce incorrect labels. We performed a manual study to confirm that this is rare. [25, 39].

8 RELATED WORK

NLP Testing. Apart from the accuracy-based testing scheme, recent works have considered model robustness as an aspect for evaluation. Belinkov and Bisk *et al.* [4] aimed to fail Neural Machine Translation (NMT) model by intentionally introducing noise in the input text. Pinjia *et al.* [12, 13] measured the robustness of NMT model by assuming syntactic and semantic relation between input and output. Ribeiro *et al.* [31] proposed an approach to generalize semantically equivalent adversarial rules. Rychalska *et al.* [34] measured drops in BLEU by corruption operation, and compared model robustness based on the amount of the drops. Iyyer *et al.* [17] introduced learning-based model for adversarial data augmentation.

In addition to the robustness on adversarial set, various other aspects of the NLP model have been considered for the robustness evaluation. Prabhakaran *et al.* [27] developed an evaluation framework to detect unintended societal bias in NLP models. Rottger *et al.* [33] introduced a functional test suite for hate speech detection in the NLP model. Ribeiro *et al.* [28] measured logical consistency of NLP model. However, in this work, we focused on the model evaluation over multiple perspectives and produced debugging information by comparing seed and expanded test cases.

Text Generation. Different data augmentation approaches have been proposed for text generation. One of methods in this category is easy data augmentation (EDA) [44]. It uses simple operations such as synonym replacement and random swap. Coulombe [6] implements rule-based transformations using regular expressions for text

transformations such as insertion of spelling mistakes, data alterations, entity names, and abbreviations. Mixup interpolation, which was introduced in computer vision category, was implemented in text domain [9]. Learning-based language models has been used to generate text data for multiple tasks [2, 46]. While these approaches focus on obtaining synthetic training data, not testing NLP model over linguistic capabilities, our work generates text data for testing NLP model over linguistic capabilities. In addition, text generation methods has been used for attacking NLP model [1, 21, 26, 47]. Morris *et al.* [26] transform input text by leveraging existing methods such as insertion, deletion and swap of word/character within the similarity-based constraints. Zang *et al.* [47] finds sememe-based synonyms and generates adversarial examples by the word substitution. Alzantot *et al.* [1] substitutes word in an input text with its neighboring words in embedding space that fools an NLP model. Linyang *et al.* [21] generates word substitutes with BERT predicted words for attacking an NLP model. Also, Grammar-based text generation has been also introduced [24, 38, 41]. Ogma [41] and Astraea [38] leverages predefined grammar, and search adversarial examples by replacing word in an input text with another one that conforms to same production rule in the grammar until the replacement leads to error in the model. Ma *et al.* [24] applies a simple expansion that inserts an adjective word before noun besides to replacement of analogy word found in embedding space. Our work generates text data does not only rely on the word substitution, but also changes diverse structures from the input text. In addition, we investigate input semantic and syntax to generate text conforming to the same context as input text and its corresponding linguistic capability.

Linguistic Capability Evaluation. Wang *et al.* [42, 43] proposed multiple diagnostic datasets to evaluate NLP models via natural language inference problems. More recently, CHECKLIST proposed an evaluation method of input-output behavior defined as linguistic capabilities. CHECKLIST generates behavior-guided inputs for validating the behaviors. [32]. Unlike prior work that relies on manual test case generation, we used structural information in text to generate test cases automatically.

NLP Model Debugging. Researchers have been explaining NLP model prediction and analyzing it for debugging the model. Ribeiro *et al.* [30] evaluated model prediction guided by human feedback providing relevance scores for words on the model prediction. Interactive error analysis [45] also has been proposed to evaluate model robustness. Zylberajch [49] used influence functions to generate model explanation, and it enables interactive debugging incorporating humans feedback explanation. Lertvittaya [20] proposed an approach to understand behavior of text classifier model and improve the model by disabling irrelevant hidden features. In this work, S²LCT is useful for identifying the sources of model failure as shown in RQ4.

9 CONCLUSIONS

In this paper, we present S²LCT that automatically searches and generates seed test cases for each linguistic capability. S²LCT expands the seeds into a sufficient number of test cases based on context-free grammar, covering diverse syntactic structures. We evaluated sentiment analysis models on the generated test cases. It is shown that S²LCT ensures quality of test cases by showing

consistency of test oracle with its input sentence and its linguistic capability. In addition, we showed that diversity in S²LCT test cases leads to cover more region in sentiment analysis models, resulting in higher reliability on the testing. We also examined failure-inducing cases to understand the root causes of the bugs. These results demonstrate the correctness and usefulness S²LCT for model evaluation.

Table 2: Search rules and transformation templates for linguistic capabilities of hate speech detection task.

Linguistic capability	Search rule and transformation template
LC1: Expression of strong negative emotions (explicit)	Search hatecheck Transform N/A
LC2: Description using very negative attributes (explicit)	Search hatecheck Transform N/A
LC3: Dehumanisation (explicit)	Search hatecheck Transform N/A
LC4: Implicit derogation	Search hatecheck Transform N/A
LC5: Direct threat	Search hatecheck Transform N/A
LC6: Threat as normative statement	Search hatecheck Transform N/A
LC7: Hate expressed using slur	Search seed={include: slur; label: toxic} Transform N/A
LC8: Non-hateful use of slur	Search seed={include: slur; label: non-toxic} Transform N/A
LC9: Hate expressed using profanity	Search seed={include: profanity; label: toxic} Transform N/A
LC10: Non-Hateful use of profanity	Search seed={include: profanity; label: non-toxic} Transform N/A
LC11: Hate expressed through reference in subsequent clauses	Search toxic_sent={include: identity; label: toxic}, non_toxic_sent={include: identity; label: non-toxic} Transform seed={{['people say', 'I heard that']+[toxic_sent]+'and'+['i hate them too', 'i also hate them']} [['people say', 'I heard that']+[toxic_sent]+'but'+['i do not hate them', 'i like them']] [['people say', 'I heard that']+[non_toxic_sent]+'but'+['i hate them', 'i dislike them']] [['people say', 'I heard that']+[non_toxic_sent]+'and'+['i also do not hate them', 'i do not hate them too']]
LC12: Hate expressed through reference in subsequent sentences	Search toxic_sent={include: identity; label: toxic}, non_toxic_sent={include: identity; label: non-toxic} Transform seed={{['people say', 'I heard that']+[toxic_sent]+'.'+['i hate them too', 'i also hate them']] [['people say', 'I heard that']+[toxic_sent]+'.'+['however i do not hate them', 'regardless i like them']] [['people say', 'I heard that']+[non_toxic_sent]+'.'+['however i hate them', 'i dislike them though']] [['people say', 'I heard that']+[non_toxic_sent]+'.'+['i also do not hate them', 'i do not hate them too']]
LC13: Hate expressed using negated positive statement	Search non_toxic_sent={label: non-toxic} Transform seed={{['I do not think that', 'I do not agree with that', 'I do not believe that']+[non_toxic_sent]}}
LC14: Non-hate expressed using negated hateful statement	Search toxic_sent={label: toxic} Transform seed={{['I do not think that', 'I do not agree with that', 'I do not believe that']+[toxic_sent]}}
LC15: Hate phrased as a question	Search toxic_sent={label: toxic} Transform seed={{['why do you agree with that', 'why do you like to say that']+[toxic_sent]+'?'}}
LC16: Hate phrased as an opinion	Search toxic_sent={label: toxic} Transform seed={{['i think that', 'i believe that']+[toxic_sent]}}
LC17: Neutral statements using protected group identifiers	Search seed={include: identity; exclude: pos adjs & neg adjs & pos nouns & neg nouns & pos verbs & neg verbs; label: non-toxic} Transform N/A
LC18: Positive statements using protected group identifiers	Search seed={include: identity; exclude: neg adjs & neg nouns & neg verbs; label: non-toxic} Transform N/A
LC19: Denouncements of hate that quote it	Search toxic_sent={label: toxic} Transform seed={{['statements like']+[toxic_sent]+'are hurtful', 'are harmful', 'are deleterious']} [['you should stop saying like']+[toxic_sent]}}

Table 3: Results of BERT-base, RoBERTa-base and DistilBERT-base sentiment analysis models on S^2 LCT test cases using all seeds. BERT-base, RoBERTa-base and DistilBERT-base models are denoted as BERT, RoBERTa and dstBERT, respectively.

Linguistic capability	Ck1st#TCs	S^2 LCT #Seeds	S^2 LCT #Exps	S^2 LCT/Ck1st #Fail	S^2 LCT/Ck1st Fail rate[%]	S^2 LCT #PassToFail
LC1: Short sentences with neutral adjectives and nouns	1,716	19	51	BERT: 60/1,330 RoBERTa: 55/1,391 dstBERT: 68/1,661	BERT: 85.71/77.51 RoBERTa: 78.57/81.06 dstBERT: 97.14/96.79	BERT: 9 RoBERTa: 2 dstBERT: 0
LC2: Short sentences with sentiment-laden adjectives	8,658	160	262	BERT: 25/26 RoBERTa: 39/139 dstBERT: 18/125	BERT: 5.92/0.30 RoBERTa: 9.24/1.61 dstBERT: 4.27/1.44	BERT: 5 RoBERTa: 14 dstBERT: 10
LC3: Sentiment change over time, present should prevail	8,000	75,159	507,091	BERT: 134,641/1,680 RoBERTa: 283,614/829 dstBERT: 363,619/2,532	BERT: 23.12/21.00 RoBERTa: 48.71/10.36 dstBERT: 62.45/31.65	BERT: 18,510 RoBERTa: 17,142 dstBERT: 16,527
LC4: Negated negative should be positive or neutral	6,786	67	503	BERT: 523/799 RoBERTa: 498/218 dstBERT: 494/734	BERT: 91.75/11.77 RoBERTa: 87.37/3.21 dstBERT: 86.67/10.82	BERT: 20 RoBERTa: 9 dstBERT: 6
LC5: Negated neutral should still be neutral	2,496	26	194	BERT: 207/2,427 RoBERTa: 204/2,304 dstBERT: 213/2,450	BERT: 94.09/97.24 RoBERTa: 92.73/92.31 dstBERT: 96.82/98.16	BERT: 11 RoBERTa: 6 dstBERT: 10
LC6: Negation of negative at the end, should be positive or neutral	2,124	18,576	102,304	BERT: 120,451/1,871 RoBERTa: 120,054/445 dstBERT: 118,910/2,124	BERT: 99.65/88.09 RoBERTa: 99.32/20.95 dstBERT: 98.37/100.00	BERT: 71 RoBERTa: 93 dstBERT: 334
LC7: Negated positive with neutral content in the middle	1,000	24,328	190,723	BERT: 195,525/860 RoBERTa: 157,635/416 dstBERT: 180,865/865	BERT: 90.92/86.00 RoBERTa: 73.30/41.60 dstBERT: 84.10/86.50	BERT: 2,014 RoBERTa: 7,289 dstBERT: 5,130
LC8: Author sentiment is more important than of others	8,528	68,284	492,817	BERT: 158,854/3,741 RoBERTa: 112,390/2,692 dstBERT: 170,198/3,535	BERT: 28.31/43.87 RoBERTa: 20.03/31.57 dstBERT: 30.33/41.45	BERT: 9,748 RoBERTa: 9,453 dstBERT: 13,870
LC9: Parsing sentiment in (question, yes) form	7,644	15,465	106,907	BERT: 7,361/253 RoBERTa: 6,510/32 dstBERT: 5,702/52	BERT: 6.02/3.31 RoBERTa: 5.32/0.42 dstBERT: 4.66/0.68	BERT: 1,646 RoBERTa: 1,519 dstBERT: 1,193
LC10: Parsing sentiment in (question, no) form	7,644	15,483	108,412	BERT: 93,329/4,056 RoBERTa: 105,319/4,576 dstBERT: 117,664/6,440	BERT: 75.33/53.06 RoBERTa: 85.01/59.86 dstBERT: 94.97/84.25	BERT: 1,891 RoBERTa: 1,764 dstBERT: 594

Table 4: Results of dehate-BERT and tweeter-RoBERTa hate speech detection models on S²LCT test cases using all seeds. dehate-BERT and tweeter-RoBERTa are denoted as BERT, RoBERTa respectively.

Linguistic capability	Htck #TCs	S ² LCT #Seeds	S ² LCT #Exps	S ² LCT/Htck #Fail	S ² LCT/Htck Fail rate[%]	S ² LCT #PassToFail
LC1: Expression of strong negative emotions (explicit)	140	140	799	alex-BERT: 823/123 CNERG-BERT: 829/124	alex-BERT: 87.65/87.86 CNERG-BERT: 88.29/88.57	alex-BERT: 12 CNERG-BERT: 38
LC2: Description using very negative attributes (explicit)	140	140	2,959	alex-BERT: 2,731/123 CNERG-BERT: 2,465/108	alex-BERT: 88.13/87.86 CNERG-BERT: 79.54/77.14	alex-BERT: 31 CNERG-BERT: 140
LC3: Dehumanisation (explicit)	140	140	3,124	alex-BERT: 2,646/119 CNERG-BERT: 2,854/124	alex-BERT: 81.07/85.00 CNERG-BERT: 87.44/88.57	alex-BERT: 15 CNERG-BERT: 96
LC4: Implicit derogation	140	140	5,664	alex-BERT: 5,029/118 CNERG-BERT: 5,481/133	alex-BERT: 86.65/84.29 CNERG-BERT: 94.43/95.00	alex-BERT: 77 CNERG-BERT: 19
LC5: Direct threat	133	133	2,689	alex-BERT: 2,056/90 CNERG-BERT: 2,455/121	alex-BERT: 72.86/67.67 CNERG-BERT: 87.00/90.98	alex-BERT: 49 CNERG-BERT: 25
LC6: Threat as normative statement	140	140	4,163	alex-BERT: 3,207/108 CNERG-BERT: 3,591/118	alex-BERT: 74.53/77.14 CNERG-BERT: 83.45/84.29	alex-BERT: 31 CNERG-BERT: 50
LC7: Hate expressed using slur	144	805	17,318	alex-BERT: 13,383/129 CNERG-BERT: 15,710/135	alex-BERT: 73.85/89.58 CNERG-BERT: 86.69/93.75	alex-BERT: 92 CNERG-BERT: 68
LC8: Non-hateful use of slur	111	395	7,881	alex-BERT: 5,600/22 CNERG-BERT: 6,815/76	alex-BERT: 67.67/19.82 CNERG-BERT: 82.35/68.47	alex-BERT: 11 CNERG-BERT: 26
LC9: Hate expressed using profanity	140	1,842	49,743	alex-BERT: 38,673/119 CNERG-BERT: 46,325/114	alex-BERT: 74.97/85.00 CNERG-BERT: 89.80/81.43	alex-BERT: 177 CNERG-BERT: 145
LC10: Non-Hateful use of profanity	100	701	15,761	alex-BERT: 3,583/6 CNERG-BERT: 7,267/3	alex-BERT: 21.77/6.00 CNERG-BERT: 44.14/3.00	alex-BERT: 44 CNERG-BERT: 60
LC11: Hate expressed through reference in subsequent clauses	140	11,968	44,890	alex-BERT: 33,230/114 CNERG-BERT: 44,034/122	alex-BERT: 58.44/81.43 CNERG-BERT: 77.45/87.14	alex-BERT: 1,434 CNERG-BERT: 1,441
LC12: Hate expressed through reference in subsequent sentences	133	11,968	42,840	alex-BERT: 32,167/104 CNERG-BERT: 43,256/102	alex-BERT: 58.69/78.20 CNERG-BERT: 78.92/76.69	alex-BERT: 1,410 CNERG-BERT: 1,348
LC13: Hate expressed using negated positive statement	140	23,379	126,742	alex-BERT: 130,886/123 CNERG-BERT: 126,635/124	alex-BERT: 87.19/87.86 CNERG-BERT: 84.36/88.57	alex-BERT: 1,888 CNERG-BERT: 4,361
LC14: Non-hate expressed using negated hateful statement	133	34,212	243,739	alex-BERT: 86,576/25 CNERG-BERT: 203,241/33	alex-BERT: 31.15/18.80 CNERG-BERT: 73.12/24.81	alex-BERT: 3,134 CNERG-BERT: 3,274
LC15: Hate phrased as a question	140	22,808	185,267	alex-BERT: 163,340/117 CNERG-BERT: 182,688/131	alex-BERT: 78.50/83.57 CNERG-BERT: 87.80/93.57	alex-BERT: 3,435 CNERG-BERT: 5,531
LC16: Hate phrased as an opinion	133	22,808	160,927	alex-BERT: 142,434/108 CNERG-BERT: 158,127/121	alex-BERT: 77.52/81.20 CNERG-BERT: 86.06/90.98	alex-BERT: 3,544 CNERG-BERT: 4,584
LC17: Neutral statements using protected group identifiers	126	5	9	alex-BERT: 13/12 CNERG-BERT: 2/1	alex-BERT: 92.86/9.52 CNERG-BERT: 14.29/0.79	alex-BERT: 0 CNERG-BERT: 1
LC18: Positive statements using protected group identifiers	189	41	182	alex-BERT: 137/32 CNERG-BERT: 63/26	alex-BERT: 61.43/16.93 CNERG-BERT: 28.25/13.76	alex-BERT: 1 CNERG-BERT: 5
LC19: Denouncements of hate that quote it	173	45,616	300,704	alex-BERT: 87,197/36 CNERG-BERT: 200,285/8	alex-BERT: 25.18/20.81 CNERG-BERT: 57.83/4.62	alex-BERT: 3,379 CNERG-BERT: 5,259

Table 5: Manual study results of label consistency and linguistic capability relevancy of S^2 LCT generated test cases compared to human annotations.

Type	#Test_Cases	Label_Consistency	LC_Relevancy
Seed	100	0.83	0.90
Expanded	100	0.84	0.90

REFERENCES

- [1] 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*. Association for Computational Linguistics, Brussels, Belgium, 2890–2896. <https://doi.org/10.18653/v1/D18-1316>
- [2] Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich, Amir Kantor, George Kour, Segev Shlomov, Naama Tepper, and Naama Zwerdling. 2019. Not Enough Data? Deep Learning to the Rescue! *CoRR* abs/1911.03118 (2019). [arXiv:1911.03118](http://arxiv.org/abs/1911.03118) <http://arxiv.org/abs/1911.03118>
- [3] Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010. SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*. European Language Resources Association (ELRA), Valletta, Malta. http://www.lrec-conf.org/proceedings/lrec2010/pdf/769_Paper.pdf
- [4] Yonatan Belinkov and Yonatan Bisk. 2017. Synthetic and Natural Noise Both Break Neural Machine Translation. *CoRR* abs/1711.02173 (2017). [arXiv:1711.02173](http://arxiv.org/abs/1711.02173)
- [5] Simin Chen, Soroush Bateni, Sampath Grandhi, Xiaodi Li, Cong Liu, and Wei Yang. 2020. *DENAS: Automated Rule Generation by Knowledge Extraction from Neural Networks*. Association for Computing Machinery, New York, NY, USA, 813–825. <https://doi.org/10.1145/3368089.3409733>
- [6] Claude Coulombe. 2018. Text Data Augmentation Made Simple By Leveraging NLP Cloud APIs. *CoRR* abs/1812.04718 (2018). [arXiv:1812.04718](http://arxiv.org/abs/1812.04718) <http://arxiv.org/abs/1812.04718>
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *CoRR* abs/1810.04805 (2018). [arXiv:1810.04805](http://arxiv.org/abs/1810.04805) <http://arxiv.org/abs/1810.04805>
- [8] Mor Geva, Yoav Goldberg, and Jonathan Berant. 2019. Are we modeling the task or the annotator? an investigation of annotator bias in natural language understanding datasets. *arXiv preprint arXiv:1908.07898* (2019).
- [9] Hongyu Guo, Yongyi Mao, and Richong Zhang. 2019. Augmenting Data with Mixup for Sentence Classification: An Empirical Study. *CoRR* abs/1905.08941 (2019). [arXiv:1905.08941](http://arxiv.org/abs/1905.08941) <http://arxiv.org/abs/1905.08941>
- [10] Wenbo Guo, Dongliang Mu, Jun Xu, Purui Su, Gang Wang, and Xinyu Xing. 2018. Lemna: Explaining deep learning based security applications. In *proceedings of the 2018 ACM SIGSAC conference on computer and communications security*. 364–379.
- [11] Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R Bowman, and Noah A Smith. 2018. Annotation artifacts in natural language inference data. *arXiv preprint arXiv:1803.02324* (2018).
- [12] Pinjia He, Clara Meister, and Zhendong Su. 2019. Structure-Invariant Testing for Machine Translation. *CoRR* abs/1907.08710 (2019). [arXiv:1907.08710](http://arxiv.org/abs/1907.08710) <http://arxiv.org/abs/1907.08710>
- [13] Pinjia He, Clara Meister, and Zhendong Su. 2020. Testing Machine Translation via Referential Transparency. *CoRR* abs/2004.10361 (2020). [arXiv:2004.10361](http://arxiv.org/abs/2004.10361) <http://arxiv.org/abs/2004.10361>
- [14] Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. (2017). To appear.
- [15] Jen-tse Huang, Jianping Zhang, Wenxuan Wang, Pinjia He, Yuxin Su, and Michael R Lyu. 2022. AEON: A Method for Automatic Evaluation of NLP Test Cases. *arXiv preprint arXiv:2205.06439* (2022).
- [16] Muftaba Husnain, Malik Muhammad Saad Missen, Nadeem Akhtar, Mickaël Coustaty, Shahzad Mumtaz, and VB Prasath. 2021. A systematic study on the role of SentiWordNet in opinion mining. *Frontiers of Computer Science* 15, 4 (2021), 1–19.
- [17] Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. Adversarial Example Generation with Syntactically Controlled Paraphrase Networks. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 1875–1885. <https://doi.org/10.18653/v1/N18-1170>
- [18] Nikita Kitaev, Steven Cao, and Dan Klein. 2019. Multilingual Constituency Parsing with Self-Attention and Pre-Training. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 3499–3505. <https://doi.org/10.18653/v1/P19-1340>
- [19] Nikita Kitaev and Dan Klein. 2018. Constituency Parsing with a Self-Attentive Encoder. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Melbourne, Australia, 2676–2686. <https://doi.org/10.18653/v1/P18-1249>
- [20] Piyawat Lertvittayakumjorn, Lucia Specia, and Francesca Toni. 2020. FIND: Human-in-the-loop Debugging Deep Text Classifiers. <https://doi.org/10.48550/ARXIV.2010.04987>
- [21] Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. BERT-ATTACK: Adversarial Attack Against BERT Using BERT. *CoRR* abs/2004.09984 (2020). [arXiv:2004.09984](https://arxiv.org/abs/2004.09984) <https://arxiv.org/abs/2004.09984>
- [22] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR* abs/1907.11692 (2019). [arXiv:1907.11692](http://arxiv.org/abs/1907.11692) <http://arxiv.org/abs/1907.11692>
- [23] Lei Ma, Felix Juefei-Xu, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Chunyang Chen, Ting Su, Li Li, Yang Liu, et al. 2018. Deepgauge: Multi-granularity testing criteria for deep learning systems. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. 120–131.
- [24] Pingchuan Ma, Shuai Wang, and Jin Liu. 2020. Metamorphic Testing and Certified Mitigation of Fairness Violations in NLP Models. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, Christian Bessiere (Ed.). International Joint Conferences on Artificial Intelligence Organization, 458–465. <https://doi.org/10.24963/ijcai.2020/64> Main track.
- [25] Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a Large Annotated Corpus of English: The Penn Treebank. *Comput. Linguist.* 19, 2 (jun 1993), 313–330.
- [26] John X. Morris, Eli Lifland, Jin Yong Yoo, and Yanjun Qi. 2020. TextAttack: A Framework for Adversarial Attacks in Natural Language Processing. *CoRR* abs/2005.05909 (2020). [arXiv:2005.05909](https://arxiv.org/abs/2005.05909) <https://arxiv.org/abs/2005.05909>
- [27] Vinodkumar Prabhakaran, Ben Hutchinson, and Margaret Mitchell. 2019. Perturbation Sensitivity Analysis to Detect Unintended Model Biases. *CoRR* abs/1910.04210 (2019). [arXiv:1910.04210](http://arxiv.org/abs/1910.04210) <http://arxiv.org/abs/1910.04210>
- [28] Marco Tulio Ribeiro, Carlos Guestrin, and Sameer Singh. 2019. Are Red Roses Red? Evaluating Consistency of Question-Answering Models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 6174–6184. <https://doi.org/10.18653/v1/P19-1621>
- [29] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 1135–1144.
- [30] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. *CoRR* abs/1602.04938 (2016). [arXiv:1602.04938](http://arxiv.org/abs/1602.04938) <http://arxiv.org/abs/1602.04938>
- [31] 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 (Volume 1: Long Papers)*. Association for Computational Linguistics, Melbourne, Australia, 856–865. <https://doi.org/10.18653/v1/P18-1079>
- [32] Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond Accuracy: Behavioral Testing of NLP models with CheckList. In *Association for Computational Linguistics (ACL)*.
- [33] Paul Röttger, Bertram Vidgen, Dong Nguyen, Zeerak Waseem, Helen Margetts, and Janet B Pierrehumbert. 2020. Hatecheck: Functional tests for hate speech detection models. *arXiv preprint arXiv:2012.15606* (2020).
- [34] Barbara Rychalska, Dominika Basaj, Alicja Gosiewska, and Przemyslaw Biecek. 2019. Models in the Wild: On Corruption Robustness of Neural NLP Systems. In *Neural Information Processing*, Tom Gedeon, Kok Wai Wong, and Minh Lee (Eds.). Springer International Publishing, Cham, 235–247.
- [35] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *ArXiv* abs/1910.01108 (2019).
- [36] Sergio Segura, Gordon Fraser, Ana B. Sanchez, and Antonio Ruiz-Cortés. 2016. A Survey on Metamorphic Testing. *IEEE Transactions on Software Engineering* 42, 9 (2016), 805–824. <https://doi.org/10.1109/TSE.2016.2532875>
- [37] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Seattle, Washington, USA, 1631–1642. <https://aclanthology.org/D13-1170>
- [38] Ezekiel O. Soremekun, Sakshi Udeshi, and Sudipta Chattopadhyay. 2020. Astraea: Grammar-based Fairness Testing. *CoRR* abs/2010.02542 (2020). [arXiv:2010.02542](https://arxiv.org/abs/2010.02542) <https://arxiv.org/abs/2010.02542>
- [39] Sphinx and NLTK Theme. 2021. *NLTK Documentation*. <https://www.nltk.org/howto/corpus.html>
- [40] Sakshi Udeshi, Pryanshu Arora, and Sudipta Chattopadhyay. 2018. Automated Directed Fairness Testing. *CoRR* abs/1807.00468 (2018). [arXiv:1807.00468](http://arxiv.org/abs/1807.00468) <http://arxiv.org/abs/1807.00468>
- [41] Sakshi Udeshi and Sudipta Chattopadhyay. 2019. Grammar Based Directed Testing of Machine Learning Systems. *CoRR* abs/1902.10027 (2019). [arXiv:1902.10027](http://arxiv.org/abs/1902.10027) <http://arxiv.org/abs/1902.10027>
- [42] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems.

- [43] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2018. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. *CoRR* abs/1804.07461 (2018). arXiv:1804.07461 <http://arxiv.org/abs/1804.07461>
- [44] Jason Wei and Kai Zou. 2019. EDA: Easy Data Augmentation Techniques for Boosting Performance on Text Classification Tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 6382–6388. <https://doi.org/10.18653/v1/D19-1670>
- [45] Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2019. Errudite: Scalable, Reproducible, and Testable Error Analysis. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 747–763. <https://doi.org/10.18653/v1/P19-1073>
- [46] Qizhe Xie, Zihang Dai, Eduard H. Hovy, Minh-Thang Luong, and Quoc V. Le. 2019. Unsupervised Data Augmentation. *CoRR* abs/1904.12848 (2019). arXiv:1904.12848 <http://arxiv.org/abs/1904.12848>
- [47] Yuan Zang, Fanchao Qi, Chenghao Yang, Zhiyuan Liu, Meng Zhang, Qun Liu, and Maosong Sun. 2020. Word-level Textual Adversarial Attacking as Combinatorial Optimization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 6066–6080. <https://doi.org/10.18653/v1/2020.acl-main.540>
- [48] Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A Benchmarking Platform for Text Generation Models. *SIGIR* (2018).
- [49] Hugo Zylberajch, Piyawat Lertvittayakumjorn, and Francesca Toni. 2021. HILDIF: Interactive Debugging of NLI Models Using Influence Functions. *Proceedings of the First Workshop on Interactive Learning for Natural Language Processing* (2021).