
NATURE: Natural Auxiliary Text Utterances for Realistic Spoken Language Evaluation

David Alfonso-Hermelo
Huawei Noah’s Ark Lab
david.ah@huawei.com

Ahmad Rashid
Huawei Noah’s Ark Lab
ahmad.rashid@huawei.com

Abbas Ghaddar
Huawei Noah’s Ark Lab
abbas.ghaddar@huawei.com

Philippe Langlais
RALI/DIRO, Université de Montréal
felipe@iro.umontreal.ca

Mehdi Rezagholizadeh
Huawei Noah’s Ark Lab
mehdi.rezagholizadeh@huawei.com

Abstract

1 Slot-filling and intent detection are the backbone of conversational agents such as
2 voice assistants and they are active areas of research. Even though state-of-the-art
3 techniques on publicly available benchmarks show impressive performance,
4 their ability to generalize to realistic scenarios has yet to be improved. In
5 this work, we present NATURE, a set of simple spoken language oriented
6 transformations, applied to the evaluation set of datasets, to introduce human spoken
7 language variations while preserving the semantics of an utterance. We apply
8 NATURE to common slot-filling and intent detection benchmarks and demonstrate
9 that simple deviations from the standard test set by NATURE can deteriorate
10 model performances significantly. Additionally, we apply different strategies to
11 mitigate the effects of NATURE and report that data-augmentation leads to some
12 improvement.

13 1 Introduction

14 The growing demand for Virtual Assistant systems (Uğurlu et al. (2020), Li et al. (2021)) has led to
15 advances in conversational and spoken language oriented models, Natural Language Understanding
16 (NLU), and Spoken Language Understanding (SLU). One of the backbones of NLU and SLU is
17 the joint tasks of Intent Detection (ID, identification of the speaker’s intent) and Slot-filling (SF,
18 extraction of the semantic constituents from the utterance). In recent years, NLU models specialized
19 in ID and SF have obtained outstanding results (Qin et al. (2019), Wang et al. (2018), Yamada et al.
20 (2020)). However, these models usually lack satisfying generalization capabilities (McCoy et al.
21 (2019), Gururangan et al. (2018), Balasubramanian et al. (2020), Lin et al. (2020)).
22

23 Data Augmentation (DA) is one of the well-known solutions to this problem (Hou et al. (2020),
24 Louvan and Magnini (2020), Kale and Siddhant (2021)). However, without looking at the test set,
25 we cannot account for all the patterns which are missing in the training set. Moreover, it still does
26 not resolve the issue of a lack of generalization to out-of-distribution evaluation. This is an issue
27 in real scenarios, specially considering the paraphrase richness of spoken language. Other works
28 propose modified evaluation sets (Lin et al. (2020), Agarwal et al. (2020)). This is a valid option but
29 for some tasks (as ID and SF) the available data is scarce, rarely open-source and producing new and
30 qualitative data is labor-intensive, time-consuming, and expensive.

31 We propose a framework that focuses on transforming the existing test sets by applying simple,
32 spoken language-oriented, realistic operators that slightly modify the input sentence but without

Utterance	Task: Model Prediction Errors
play party anthems → ploy party anthems	ID: Play_Music → Search_Creative_Work
play some sixties music → plays some sixties music	SF: [sixties]:year → [sixties]:year; [plays]:album
listen to dragon ball: music collection → like listen to dragon ball: music collection	ID: Search_Creative_Work → Play_Music SF: [dragon ball: music collection]:object_name → [dragon ball]:artist; [collection]:album

Figure 1: Examples of NATURE-altered utterances with badly predicted slots and or intent. The altered utterance is preceded by a →.

33 altering the original meaning (as we shall see). **By realistic, we mean that modified utterances remain**
34 **semantically similar to the original ones.** We call this framework NATURE (*Naive Alterations of*
35 *Textual Utterances for Realistic Evaluation*). Figure 1 shows examples of altered utterances where a
36 state-of-the-art model (Qin et al. (2019)) correctly predicted the label for the original utterance but
37 failed for the altered utterance.

38 We conduct experiments that apply our framework to standard benchmarks and compare the *before*
39 *and after* performances of state-of-the-art models. The results illustrate the heuristic dependencies of
40 each model.

41

42 2 Related Work

43 2.1 Realizing model use shortcuts

44 A growing number of studies identify a tendency in NLU models to leverage the superficial features
45 and language artifacts instead of generalizing over the semantic content. **A naive way to force**
46 **generalization is to automatically add noise to the training set, however, as demonstrated by Belinkov**
47 **and Bisk (2017), models trained on synthetic noise do not necessarily perform well on natural noise,**
48 **requiring a more elaborated approach.** Given our incapacity to control what features these models
49 learn, each task requires an in-depth analysis and a data or model modification that guides it to the
50 correct answer. For the political claims detection task Padó et al. (2019) and Dayanik and Padó
51 (2020) unveil a strong bias towards the claims made by frequent actors that require masking the
52 actor and its pronouns during training to improve the performance. Other works (Gururangan et al.
53 (2018), Poliak et al. (2018), Zellers et al. (2018), McCoy et al. (2019), Naik et al. (2018)) have
54 focused on the artifact and heuristic over-fitting for the Natural Language Inference (NLI) task **or**
55 **for the Question-Answering (QA) task (Jia and Liang (2017)).** The work of Balasubramanian et al.
56 (2020) show how substituting Named-Entities (NEs) influence the robustness of BERT-based models
57 for different tasks (NLI, co-reference resolution and grammar error correction). To the best of our
58 knowledge, no work has attempted to demonstrate that the benchmarks and models for the dual tasks
59 of SF and ID rely on frequent heuristic patterns.

60

61 2.2 Alternative evaluation

62 Some researchers have proposed evaluation sets with naturally occurring adverse sentences for
63 different tasks such as HANS for MNLI (McCoy et al. (2019)) or PAWS (Zhang et al. (2019)) and
64 PAWS-X (Yang et al. (2019)) for paraphrase identification. Another strategy involves a systematic
65 alteration of the test set (Lin et al. (2020)). **This has gained popularity in recent years with a**
66 **growing interest in more challenging and adversarial evaluation frameworks.** However, a more
67 **challenging test set has to ensure high quality annotation, which is why many papers have suggested**
68 **an human-in-the-loop approach (Kaushik et al. (2019), Gardner et al. (2020), Kiela et al. (2021)).**
69 **But these approaches are costly, specially due to the number and quality of annotators necessary to**

70 produce a high-quality output. Generalization is more easily achieved when the training data is large
71 and diverse. A model can be effective, yet, if it is only fed with small and/or similar data, it will have
72 difficulties to achieve robustness. Some researchers (Louvan and Magnini (2020), Zeng et al. (2020),
73 Dai and Adel (2020), Min et al. (2020), Moosavi et al. (2020)) use DA strategies to improve the
74 training data and help boost a model’s performance.

75 Other researchers have taken a different path and suggest a whole different way of evaluating: testing
76 multiple task-agnostic requisites instead of using a test set that matches the train and validation sets
77 (Ribeiro et al. (2020), Goel et al. (2021)).
78

79 2.3 Test set alteration methods

80 There has been many proposals of spoken-language oriented alteration methods (Tsvetkov et al.
81 (2014), Simonnet et al. (2018), Li et al. (2018), Gopalakrishnan et al. (2020)) but the ones we are
82 interested in require to change the utterance form while maintaining the original semantic value of
83 each token (in the form of labels). Very few works have managed to devise methods that change the
84 form while maintaining the semantic labeling, such as the work of Yin et al. (2020) where the authors
85 suggest altering methods that emulate non-native errors or the work of Li et al. (2020) where they use
86 simple methods to produce more counterfactual versions of the original utterances.
87

88 3 Methodology

89 In this section we describe the operators used to generate new utterances out of a given one. We
90 present examples for each operator on Figure 2.
91

92 3.1 Fillers

93 Fillers are ubiquitous in everyday spoken language and often appear in human-to-human dialog
94 (transcribed to text) corpora (such as the Switchboard corpus Godfrey et al. (1992), composed of
95 approximately 1.6% fillers Shriberg (2001)). Yet they are intentionally cleaned off in SF and ID
96 benchmarks. Fillers serve as hesitation markers (e.g.: *Bring me the, like, Greek yogurt. I’ve heard*
97 *it’s really, you know, savoury.*) or as introduction/closure of a turn of speech (e.g., *Now, bring me*
98 *the Greek yogurt please and thank you. Actually, I’ve heard it’s really savoury, right?*). Fillers are
99 semantically poor and do not add essential information, and therefore, do not change the overall
100 meaning of an utterance.

101 We propose 4 different filler operators:

- 102 • **Begin-of-sentence (BOS)**: a small introductory filler phrase at the beginning of the utterance,
103 such as: *so, like, actually, okay so, so okay, so basically, now or well.*
- 104 • **End-of-sentence (EOS)**: a small conclusive filler phrase at the end of the utterance, such as:
105 *if you please, please, pretty please, please and thank you, now please, if you can, now, right*
106 *now, right away, right this minute, will you ?, would you ?, can you ?, would you mind ?.*
- 107 • **Pre-verb**: a filler word or sequence of words appearing before the utterance’s verb or verbal
108 phrase, such as: *like, basically* or *actually.*
- 109 • **Post-verb**: a filler word or sequence of words appearing after the utterance’s verb or verbal
110 phrase, such as: *basically, actually, like* or *you know.*

111 BOS and EOS operators simply add a filler at the very beginning or the end of the utterance,
112 respectively. The pre-verb and post-verb operators require us to find the part-of-speech (POS) tag
113 of the utterance tokens (we use the NLTK library to find the POS of the tokens). Then the filler is
114 placed at the correct place. We add a fail-safe rule to ensure that a filler is added if no verb is found
115 where expected. To that end, we use the overly-recurrent filler, *like*, and the first appearing Named
116 Entity as a pivot instead of the first appearing verb e.g., *let’s check like avengers).*

Test set	Example sentence
Original	add <u>tune</u> to <u>sxsw fresh</u> playlist
BOS Filler	okay so add <u>tune</u> to <u>sxsw fresh</u> playlist
Pre-V. Filler	like add <u>tune</u> to <u>sxsw fresh</u> playlist
Post-V. Filler	add <u>tune</u> actually to <u>sxsw fresh</u> playlist
EOS Filler	add <u>tune</u> to <u>sxsw fresh</u> playlist if you can
Synonym V.	play <u>tune</u> to <u>sxsw fresh</u> playlist
Synonym Adj.	add <u>tune</u> to <u>sxsw cool</u> playlist
Synonym Adv.	add prior to <u>sxsw fresh</u> playlist
Synonym Any	mix <u>tune</u> to <u>sxsw fresh</u> playlist
Synonym StopW	add <u>tune</u> the <u>sxsw fresh</u> playlist
Speako	add tua to <u>sxsw fresh</u> playlist

Figure 2: Processed variants of original utterances from the SNIPS corpus. The tokens labeled as *music_item* appear with a dotted underline and the tokens labeled as *playlist* show a dashed underline. In SNIPS, the *sxsw* token is part of a playlist name and an abbreviation of *South by Southwest*.

117 3.2 Synonymy

118 A synonym is a word that can be interchanged with another in context, without changing the meaning
 119 of the whole. To replicate this semantic operation, we select the POS corresponding to our operator
 120 (among verb, adjective, adverb, etc.). We then select a word of that type in the input utterance and
 121 make a list of potential synonym candidates (with the same POS tag) to replace it. Then we select the
 122 most probable of the candidates as our replacement. We use the pre-trained BERT-base model with a
 123 Language Modeling head on top to produce the synonym candidates instead of a human populated
 124 dictionary (such as Wiktionary) since not all dictionary entries show synonyms. We first randomly
 125 choose a POS tag and find a target token which has this tag in our utterance. Then we replace the
 126 target with a special [MASK] token. We feed this utterance into BERT and obtain a list of candidates
 127 from most to least probable.

128 In case the sentence contains no token with the target POS, we use the more common *noun* POS. We
 129 observe an example in the *Syn. Adv.* row in Table 2.

130 As we can see in Figure 3, not all BERT candidates are suitable synonyms of the target token. We
 131 remove candidates that do not have the same POS of the target token. For a better performance,
 132 we put each candidate in the context of the utterance before extracting candidate POS. We have
 133 5 different Synonymy operators based on different target POS: **verb**, **adjective**, **adverb**, **any** (at
 134 random between verb, adjective, adverb or noun), **stop-words** (grammatical and most common
 135 words).
 136

137 3.3 Speako

138 Some words sound similar to others but have a different meaning altogether (e.g., *decent* and *descent*).
 139 This operator is based on the idea that anyone can make an error, but an efficient and robust model

Token in context	Wiktionary synonyms	BERT candidates
let me <u>buy</u> it verb	purchase, accept, [...]	get, buy, present, make, purchase, offer, give, sell, [...]
is it <u>large</u> ? adj	giant, big, huge, [...]	unusual, big, dangerous, large, powerful, [...]
i said it <u>quickly</u> adv	rapidly, fast	fast, well, strong, high, good, deep, large, slow, [...]
give me <u>freedom</u> noun	liberty, license, [...]	rights, property, freedom, status, goods, liberty, [...]
i found <u>the</u> ball stopword	le	the, second, also, third, their, still, a, our, 2nd, [...]

Figure 3: Target words (underlined) of various POS and their synonyms taken from the crowd-sourced dictionary Wiktionary and candidates obtained using a pre-trained BERT language model.

140 should be able to recover a minor mistake using the context. Thus, we introduce speakos (slip of the
 141 tongue, speech-to-text misinterpretation), which are common in user-machine communication.
 142 To do so, we use a prepared dictionary of tokens appearing 1000+ times in the whole English
 143 Wikipedia¹. We convert each entry of the dictionary into its representation in International Phonetic
 144 Alphabet (IPA). We randomly select one token from the sentence, convert it to IPA, calculate the
 145 similarity between it and the dictionary’s entries (using Levenshtein distance) and replace it with
 146 the closest candidate. For instance, the sentence *let me watch (/watʃ/) a comedy video* could be
 147 transformed into *let me which (/wɪtʃ/) a comedy video*).

148 4 Experimental Setup

149 4.1 Data

150 In our work, we use 3 popular open-source benchmarks² which are summarized in Table 1:

151 **Airline Travel Information System (ATIS)**³ Hemphill et al. (1990) introduced an NLU benchmark
 152 for the SF and ID tasks with 18 different intent labels, 127 slot labels and a vocabulary of
 153 939 tokens. It contains annotated utterances corresponding to flight reservations, spoken
 154 dialogues and requests.

155 **SNIPS**⁴ Coucke et al. (2018) proposed the SNIPS voice platform, from which a dataset of queries
 156 for the SF and ID tasks with 7 intent labels, 72 slot labels and a vocabulary of 12k tokens
 157 were extracted.

158 **NLU-ED**⁵ is a dataset of 25K human annotated utterances using the Amazon Mechanical Turk
 159 service Liu et al. (2019). This NLU benchmark for the SF and ID tasks is comprised of 69
 160 intent labels, 108 slot labels and a vocabulary of 7.9k tokens.

161 Following the common practice in the field (Hakkani-Tür et al. (2016), Goo et al. (2018), Qin et al.
 162 (2019), Razumovskaia et al. (2021), Krishnan et al. (2021)), we report the performance of SF using
 163 the F1 score. Moreover, we propose an End-to-End accuracy (E2E) metric (sometimes referred in
 164 the literature as the sentence-level semantic accuracy (Qin et al. (2019))). This metric counts true
 165 positives when all the predicted labels (intent+slots) match the ground truth labels. This allows us to
 166 combine the SF and ID performance in a single more strict metric.

167

Benchmark		Train	Valid.	Test
ATIS	Sent	4 478	500	893
	Words	50 497	5 703	9 164
	Voc	867	463	448
SNIPS	Sent	13 084	700	700
	Words	117 700	6 384	6 354
	Voc	11 418	1 571	1 624
NLU-ED	Sent	20 628	2 544	2 544
	Words	145 950	18 167	17 347
	Voc	7 010	2 182	2 072

Table 1: Dataset size information of our benchmarks: ATIS, SNIPS and NLU-ED.

168 Any dialog-based dataset extracted from real user situations has the potential of containing private
 169 and security sensitive information. This is the main cause for the relative low amount of datasets for

¹We empirically observed that removing all tokens that had a co-occurrence lower than 1000 eliminated most of the nonsensical strings and extreme misspellings and conserved most functional words and very common typos.

²We did not select the SGD dataset of Rastogi et al. (2020) despite being recent and large, since it is a multi-turn dialog benchmark and cannot be used out of the box for the SF and ID tasks.

³CGNU General Public License, version 2

⁴Creative Commons Zero v1.0 Universal License

⁵Creative Commons Attribution 4.0 International License

170 SF and ID. The benchmarks we mention are well known and cautiously cleaned (as presented in
171 Section 3). Our operators purposely avoid using any type of resource that would contain personal
172 information. To the best of our knowledge, our work is not detrimental to people’s safety, privacy,
173 security, rights or to the environment in any way.
174

175 4.2 Models

176 We use two different state-of-the-art models:

177 **Stack-Prop+BERT** (Qin et al., 2019) uses BERT as a token-level encoder that feeds into two
178 different BiLSTMs, one per each task. The output of the SF BiLSTM is added to the ID
179 BiLSTM input in order to produce a token-level intent prediction which is further averaged
180 into a sentence-level prediction.

181 **Bi-RNN** (Wang et al., 2018) uses two correlated BiLSTMs that cross-impact each other by accessing
182 the other’s hidden states and come to a joint prediction for ID and SF.

183 The pre-trained version of these models were not available⁶. For ATIS and SNIPS, we trained the
184 models using the same hyperparameters proposed in the documentation by Qin et al. (2019)⁷ and
185 Wang et al. (2018)⁸, respectively. For NLU-ED, we use the hyperparameters from SNIPS, as their
186 size is comparable. Our trained models obtained comparable results to their published counterpart
187 (see in Appendix). To train the models, we used 1 NVIDIA Tesla V100 with 32Gb of internal
188 memory. It took between 3 and 71 hours to train the Stack-Prop+BERT model (Qin et al., 2019)
189 (depending on the size of the benchmark), and between 68 and 130 hours to train the Bi-RNN
190 model (Wang et al., 2018).
191

192 4.3 Modified NATURE Test Sets

193 Since the original test sets only cover a limited set of patterns, we transform them by applying our
194 NATURE patterns to obtain test sets of the same size as the original ones. As previously illustrated,
195 NATURE operators offer simple ways of altering utterances. In order to avoid rendering utterances
196 unrecognizable from their original version, we only apply one operator at a time and only once in the
197 sentence (e.g. we add 1 filler or synonymize one token or transform a token into its speako version).
198 We design 2 NATURE experimental test sets: *Random* and *Hard*. In the Random setting, for each
199 utterance, we apply one operator at random. **This random selection may cause an unbalanced
200 distribution of alterations (some operators being more used than others). To obtain a more
201 impartial score, we repeat the random operator selection 10 times and calculate the mean
202 score.**

203 For the Hard setting experiments, after applying all our operators on each utterance and gathering
204 all candidates, we use a relatively simple BERT-based model to calculate the performance of each
205 candidate. We use JointBERT⁹, which is an unofficial implementation of the SF and ID architecture
206 described in Chen et al. (2019) **to extract (for each utterance) the candidate that performs more
207 harshly. The assumption being that the candidate that performed poorly for one model will have a
208 greater chance of performing poorly on other models.**

209 The Random test set is meant to show how a random small change in the sentence can influence
210 evaluation while the Hard test set is meant to assess the lower-bound performance of how much the
211 model depends on similar pattern sentences to obtain the correct prediction.
212

⁶<https://github.com/LeepLosed/StackPropagation-SLU> and <https://github.com/ray075hl/Bi-Model-Intent-And-Slot>

⁷300 epochs, 0.001 learning rate, 0.4 dropout rate, 256 encoder hidden dimensions, 1024 attention hidden dimensions, 128 attention output dimensions, 256 word embedding dimensions for ATIS and 32 for SNIPS.

⁸500 epochs, max sentence length of 120, 0.001 learning rate, 0.2 dropout rate, 300 word embedding size, 200 LSTM hidden size

⁹<https://github.com/monologg/JointBERT>

Operator	ATIS	SNIPS	NLU-ED
BOS Filler	0.8	0.1	2.5
Pre-V. Filler	6.0	3.7	16.0
Post-V. Filler	1.9	8.6	5.1
EOS Filler	9.0	52.3	8.3
Syn. V.	25.6	5.4	16.3
Syn. Adj.	29.2	15.0	23.4
Syn. Adv.	11.8	5.6	10.2
Syn. Any	5.3	1.1	4.8
Syn. StopW	3.2	2.7	6.4
Speako	7.2	5.4	6.9

Table 2: Distribution of JointBERT-selected operators for the Hard experimental test set.

213 4.4 Augmented Training Sets

214 Even though our NATURE operators are designed for different purposes, some of these operators
 215 may look like certain DA strategies. However, in this subsection, we show to what extent our current
 216 operators are different from most famous heuristic DA techniques. In this regard, we apply standard
 217 DA strategies to the train and validation sets and illustrate their impact on the model’s generalization
 218 ability. We use common automatic DA strategies from the NLPaug library (Ma, 2019) that allow to
 219 easily relabel the augmented data using the original labels:

- 220 1. **Keyboard Augmentation:** simulates keyboard distance error.
 221 (e.g., *find a tv seriSs called armaRdvdon summer*)
- 222 2. **Spelling Augmentation:** substitutes word according to spelling mistake dictionary.
 223 (e.g., *fine a tv serie called armageddon summer*)
- 224 3. **Synonym Augmentation:** substitutes similar word according to WordNet/PPDB synonym.
 225 (e.g., *find a tv set series called armageddon summertime*)
- 226 4. **Antonym Augmentation:** substitutes opposite meaning word according to WordNet
 227 antonym.
 228 (e.g., *lose a tv series called armageddon summer*)
- 229 5. **TF-IDF Augmentation:** uses the TF-IDF measure to find out how a word should be
 230 augmented.
 231 (e.g., *find tv series called armageddon forms*)
- 232 6. **Contextual Word Embeddings Augmentation:** feeds surroundings word to BERT,
 233 DistilBERT, RoBERTa or XLNet language model to find out the most suitable word for
 234 augmentation.
 235 (e.g., *find a second series called armageddon ii*)

236 We apply the DA strategies exclusively to the train and validation sets, choosing 1 of the 6 DA
 237 functions at random and adding one output to the original dataset which will give us a training
 238 and validation data twice as large as the original training and validation sets. One might notice
 239 that some of the DA techniques implemented in this toolkit are close in nature to some of our
 240 NATURE operators, still (as we shall see) this DA toolkit does not suffice to generalize well to the
 241 transformations of NATURE.
 242

243 5 Results and Discussion

244 5.1 Qualitative Evaluation

245 Our assumption is that the operator-generated utterances share the same meaning and labeling as the
 246 original sentence. In order to measure this, we conducted a small but representative multiple-choice
 247 survey. We select 120 operator-altered utterances from the ATIS, SNIPS and NLU-ED benchmarks.
 248 We selected at random 40 utterances from each benchmark, making sure they were also evenly

249 distributed between operators (12 utterances per operator). In addition to these, we cherry-picked 12
 250 original utterances of high-quality that served as control. As we can see in the Appendix Survey
 251 Table, the control scores stayed high and therefore, there was no reason to invalidate any participant’s
 252 annotations.

253 14 participants (NLP and ML interns and colleagues, with no links to this work) volunteered to
 254 participate in this unpaid survey and consented verbally to the use of their data within the scope of
 255 this research. To avoid a decrease in annotation quality (due to fatigue), we split the participants in 2
 256 groups (of 7 members) and divided the utterances in two sets (each with 60 operator-altered + 12
 257 control utterances). We estimated the survey time to be 30-60 minutes, which was not far from the
 258 actual time (27-53 minutes).

259 For each utterance, we asked the participants to evaluate the intent and slot labels as *reasonable* or
 260 *unreasonable*.

261

	Group 1		Group 2	
	Experiment	Control	Experiment	Control
Slot	94.5	94.0	93.8	97.0
Intent	89.0	97.6	85.9	97.5

Table 3: Survey results and statistics per group. All scores appear as percentages.

262 In Table 3 we observe a sizable decrease on the experiment side for Intent, which can be partially
 263 explained by disposition of some operators to alter a words (such as verbs) that are highly associated
 264 with the intent classification. We also observe that the Slot labeling results are high and very close
 265 to the control scores. This indicates that (contrary to many DA strategies) the NATURE operators
 266 maintain a close-to-ground-truth slot labeling.

267

268 5.2 Quantitative Evaluation

269 Table 4 show the performances of the Stack-Prop+BERT and Bi-RNN models trained on the original
 270 train data of ATIS, SNIPS and NLU-ED benchmarks. Models are evaluated on the Original, Rand
 271 and Hard test sets. We also show the scores on 10 test sets, each altered with a single NATURE
 272 operator altered test sets, where one operator is applied to the whole test set. For each benchmark, we
 273 report the F1 and accuracy on the SF and ID tasks respectively, and our End-to-End (E2E) metric.
 274 Furthermore, we report the unweighted average (Avg. column) of the aforementioned scores on the
 275 three benchmarks. Altered test sets results are sorted in descending order according to the averaged
 276 E2E metric. We notice that BERT-based models outperform RNN ones not only on original, but also
 277 on all test set variants. More precisely, we observe a gap of 6.3%, 8.7% and 5.9% on the Avg. E2E
 278 metric on the Orig, Rand and Hard test sets.

279 First, we observe a noticeable lowering in the scores on Rand, and quite a radical change on Hard
 280 test set. We must consider the possibility that the hard test set incorporates more noise than the
 281 random test sets, and this could be the cause of this low score. Depending on the benchmark, the
 282 sharpest operators are not always the ones expected to be most disruptive. Yet, the decrease in score
 283 is extreme across all benchmarks and for both models.

284 Second, we notice that not all operators are equally disruptive. Models seem to handle well Filler
 285 operators (except for EOS), suggesting some syntax-level pattern independence and indicating that
 286 the models are using the position of the tokens instead of the tokens themselves achieve the correct
 287 predictions. The Synonymy operators, specially the adjective and adverb, greatly deteriorate the
 288 performances. This decrease in score shines a light on the importance of the token-level pattern,
 289 signaling that the models are using certain adjectives and adverbs to make their predictions. Since
 290 adjectives and adverbs are much less diverse than the nouns and verbs, we infer that the models are
 291 using these words as prediction clues. The Speako operator is not very disruptive either, suggesting a
 292 good capacity of the models to overcome these variants and generalize using the remaining context.
 293 Interestingly, we notice that the drop of performances is highly strong on the E2E metric. For
 294 instance, using the Stack-Prop+BERT model on the ATIS test set, altered with the EOS Filler
 295 operator, we observe a 0.3% and 6.8% drop on SF and ID respectively but a 32.1% drop on E2E.
 296 We argue that E2E is a more reliable metric compared to reporting ID accuracy and SF F1 scores
 297 separately. Specially in an industrial environment, where a Virtual Assistant can only execute the

Test Set	ATIS			SNIPS			NLU-ED			Avg.		
	Slot (F1)	Intent (Acc)	E2E (Acc)	Slot (F1)	Intent (Acc)	E2E (Acc)	Slot (F1)	Intent (Acc)	E2E (Acc)	Slot (F1)	Intent (Acc)	E2E (Acc)
Stack-Prop+BERT												
Orig	95.7	96.5	86.2	95.0	98.3	87.9	74.0	85.1	67.8	88.2	93.3	80.6
Rand	91.3	95.0	66.5	83.4	96.1	53.8	67.4	76.1	56.8	80.7	89.1	59.0
Hard	82.3	90.7	34.9	70.6	95.3	12.9	55.5	62.7	38.9	69.5	82.9	28.9
Pre-V. Filler	95.6	96.5	85.6	92.2	98.3	79.3	71.0	83.6	65.7	86.3	92.8	76.9
Syn. StopW	93.0	94.8	76.5	89.7	96.7	74.3	70.2	78.9	60.2	84.3	90.1	70.3
BOS Filler	95.6	96.2	85.8	86.5	97.1	54.9	72.5	80.8	63.9	84.9	91.4	68.2
Post-V. Filler	94.0	96.5	80.3	84.8	98.0	57.1	68.0	84.1	63.6	82.3	92.9	67.0
Syn. V.	90.1	95.3	63.6	88.4	95.1	66.7	68.5	74.2	56.5	82.3	88.2	62.3
Speako	92.9	92.7	72.5	77.9	94.6	45.3	69.5	74.2	57.6	80.1	87.2	58.5
Syn. Any	90.3	90.5	54.4	86.9	94.4	61.6	67.8	71.0	53.5	81.7	85.3	56.5
Syn. Adj.	84.7	92.7	42.4	78.2	95.4	44.4	60.2	69.7	47.2	74.4	85.9	44.7
Syn. Adv.	88.2	89.1	43.9	77.6	94.3	41.9	61.6	65.6	45.4	75.8	83.0	43.7
EOS Filler	88.9	96.3	54.1	<u>72.1</u>	<u>97.7</u>	<u>13.1</u>	63.9	<u>78.0</u>	53.6	75.0	90.7	40.3
Bi-RNN												
Orig	94.9	97.6	84.7	89.4	97.1	76.6	66.4	80.9	61.7	83.6	91.9	74.3
Rand	89.9	94.3	61.8	75.6	94.1	39.0	60.6	70.8	50.1	75.4	86.4	50.3
Hard	79.9	92.0	27.6	62.4	92.9	7.0	49.6	58.8	34.4	64.0	81.2	23.0
Pre-V. Filler	94.7	97.3	82.2	84.6	96.4	60.0	63.3	80.1	59.3	80.9	91.3	67.2
Syn. StopW	90.6	94.7	72.7	80.5	95.4	56.4	62.3	73.2	52.7	77.8	87.8	60.6
BOS Filler	<u>80.7</u>	96.7	82.6	80.9	96.7	38.4	65.8	78.8	59.6	75.8	90.7	60.2
Post-V. Filler	93.8	96.9	80.3	77.9	96.6	37.4	62.6	79.3	56.6	78.1	90.9	58.1
Syn. V.	87.6	95.9	56.6	79.5	92.1	50.6	61.3	70.5	50.7	76.1	86.2	52.6
Speako	91.8	90.3	68.1	70.1	<u>90.1</u>	33.6	61.5	69.8	51.0	74.5	83.4	50.9
Syn. Any	89.2	90.4	52.6	77.8	91.4	40.6	62.0	67.3	49.1	76.3	83.0	47.4
Syn. Adj.	81.7	94.2	<u>34.4</u>	71.7	93.9	34.9	<u>54.3</u>	65.5	42.1	69.2	84.5	37.1
Syn. Adv.	87.2	<u>85.1</u>	38.4	69.9	92.1	29.0	<u>54.7</u>	<u>61.4</u>	<u>40.3</u>	70.6	79.5	35.9
EOS Filler	88.9	96.8	52.2	<u>64.1</u>	94.1	<u>5.9</u>	56.4	65.8	42.0	69.8	85.6	33.4

Table 4: SF, ID and E2E performances of BERT and RNN based models trained on ATIS, SNIPS, and NLU-ED and evaluated on their original and altered test sets. We show results on *per-operator* as well as on Rand and Hard test sets. Furthermore, we report the unweighted average score on the 3 benchmark we considered. The lowest scores in each column appear underlined.

298 correct command if the intent and all slots are correctly predicted.
299 Additionally, to better understand the underlying processes of the state-of-the-art models, we
300 produced and analyzed the self-attention weight heat-maps. This allows us to better understand what
301 tokens the models focus on more to make their prediction. In Figure 4 we show a representative
302 excerpt heat-maps for wrongly predicted sentences (for both SF and ID). One for the unchanged
303 SNIPS test set and one for each type of operator. We observe that the self-attention often focuses
304 more heavily on verbs, nouns and certain types of stop words, such as "the". It also shows that high
305 attention is given to verbs and certain stop words at the end of the sentence. This is evident in all
306 Figures but particularly in Figure 4b, where we can see high attention on non-frequent tokens (for the
307 benchmark), such as "if" or "?".
308

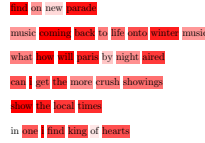
Test Set	ATIS		SNIPS		NLU-ED		Avg.	
	w/o	w Aug.	w/o	w Aug.	w/o	w Aug.	w/o	w Aug.
Orig	86.2	83.3 (-2.9)	87.9	85.3 (-2.6)	67.8	66.2 (-1.6)	80.6	78.3 (-2.3)
Rand	66.5	69.2 (+2.7)	39.0	48.2 (+9.2)	56.8	56.7 (-0.1)	54.1	58.3 (+4.2)
Hard	34.9	54.0 (+19.1)	12.9	27.1 (+15.2)	38.9	40.7 (+1.8)	28.9	40.6 (+11.7)

Table 5: End-to-End (E2E) scores of Stack-Prop+BERT models trained on ATIS, SNIPS and NLU-ED original (w/o) and augmented (w) training data. Each model is evaluated on its respective original, Rand, and Hard test set. We report the unweighted average of the 3 datasets.

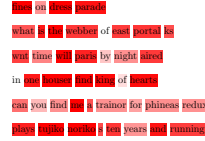


(a) Heat-map of original SNIPS utterances.

(b) Heat-map of EOS filler-altered utterances.



(c) Heat-map of Synonymy Adjective-altered utterances.



(d) Heat-map of Speako-altered utterances.

Figure 4: Heat-maps of SNIPS utterances whose SF and ID labels were wrongly predicted by the Stack-Prop+BERT model. The more intense the color, the greater the self-attention weight.

309 So far, we have shown that state-of-the-art SF and ID models do suffer when small perturbations are
 310 introduced to the test data. We now run experiments on augmented data in order to test the models'
 311 performances on larger and slightly more diverse train sets (Section 4.4). Table 5 reports E2E scores
 312 of Stack-Prop+BERT¹⁰ model when trained without (w/o) and with (w Aug) data-augmented train
 313 and validation sets. Similar to Table 4, we evaluate the model on the Original, Rand, and Hard test
 314 sets of ATIS, SNIPS and NLU-ED while also reporting the unweighted average score.

315 On one hand, we observe significant gains on the altered test sets (except on NLU-ED Rand) across
 316 all benchmarks. The largest increase in performances are obtained on the Hard sets with 19.1%
 317 and 15.2% of gain on ATIS and SNIPS respectively. The gain can be partially explained by the
 318 augmentation of training data size, forcing the model to better generalize and also to the fact that our
 319 operator shares some characteristics with the used DA toolkit (i.e., Synonymy).

320 On the other hand, the performances decrease on the 3 benchmark, by an average of 2.3%, when
 321 the model is evaluated on the Original test sets. DA is a valid strategy in NLP, specially for small
 322 sized datasets. However, even the large and more diverse NLU-ED benchmark shows only small
 323 improvement and does not solve the unobserved pattern problem exemplified by the NATURE
 324 operators. This is a strong indicator that the problem is far from solved, and that there is much room
 325 for research.

326

327 6 Conclusions

328 Neural Network models have a black-box architecture that makes it hard to discern when they
 329 correctly generalize over the input and when they resort to heuristic features that correlate to the
 330 expected output. We present the NATURE operators, apply them to test sets of standard spoken
 331 language oriented benchmarks and observe a consequential drop of the state-of-the-art model scores.
 332 The different operators in our framework help discern what surface patterns is the model misusing.
 333 We apply simple DA techniques (that are distinct from our operators) to the train and validation sets
 334 of each benchmark, allowing us to determine when and to what extent the problem is due to a small
 335 training set size. Although DA strategies tends to improve the generalization score, they do not fully
 336 recover nor catch up to their original scores.

337 In future work, we expect to improve the current operators and include more diverse and realistic
 338 speech handicap, vocabulary, syntax, and miscellaneous pattern operators.

339

340 7 Acknowledgments

341 We are grateful to the participants of the survey.

¹⁰Performances of the Bi-RNN model show very similar trends.

342 **References**

- 343 Oshin Agarwal, Yinfei Yang, Byron C Wallace, and Ani Nenkova. 2020. Entity-switched datasets:
344 An approach to auditing the in-domain robustness of named entity recognition models. [arXiv](#)
345 [preprint arXiv:2004.04123](#).
- 346 Sriram Balasubramanian, Naman Jain, Gaurav Jindal, Abhijeet Awasthi, and Sunita Sarawagi. 2020.
347 What’s in a name? are bert named entity representations just as good for any other name? [arXiv](#)
348 [preprint arXiv:2007.06897](#).
- 349 Yonatan Belinkov and Yonatan Bisk. 2017. Synthetic and natural noise both break neural machine
350 translation. [arXiv preprint arXiv:1711.02173](#).
- 351 Qian Chen, Zhu Zhuo, and Wen Wang. 2019. Bert for joint intent classification and slot filling. [arXiv](#)
352 [preprint arXiv:1902.10909](#).
- 353 Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément
354 Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, et al. 2018. Snips
355 voice platform: an embedded spoken language understanding system for private-by-design voice
356 interfaces. [arXiv preprint arXiv:1805.10190](#).
- 357 Xiang Dai and Heike Adel. 2020. An analysis of simple data augmentation for named entity
358 recognition. In [Proceedings of the 28th International Conference on Computational Linguistics](#),
359 pages 3861–3867.
- 360 Erenay Dayanik and Sebastian Padó. 2020. Masking actor information leads to fairer political claims
361 detection. In [Proceedings of the 58th Annual Meeting of the Association for Computational](#)
362 [Linguistics](#), pages 4385–4391.
- 363 Matt Gardner, Yoav Artzi, Victoria Basmova, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep
364 Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, et al. 2020. Evaluating models’ local
365 decision boundaries via contrast sets. [arXiv preprint arXiv:2004.02709](#).
- 366 John Godfrey, Edward Holliman, and Jane McDaniel. 1992. Switchboard: Telephone speech corpus
367 for research and development. In [Proceedings of the 1992 IEEE International Conference on](#)
368 [Acoustics, Speech and Signal Processing \(ICASSP\)](#), volume 1, pages 517–520. IEEE.
- 369 Karan Goel, Nazneen Rajani, Jesse Vig, Samson Tan, Jason Wu, Stephan Zheng, Caiming Xiong,
370 Mohit Bansal, and Christopher Ré. 2021. Robustness gym: Unifying the nlp evaluation landscape.
371 [arXiv preprint arXiv:2101.04840](#).
- 372 Chih-Wen Goo, Guang Gao, Yun-Kai Hsu, Chih-Li Huo, Tsung-Chieh Chen, Keng-Wei Hsu,
373 and Yun-Nung Chen. 2018. Slot-gated modeling for joint slot filling and intent prediction.
374 In [Proceedings of the 2018 Conference of the North American Chapter of the Association for](#)
375 [Computational Linguistics: Human Language Technologies, Volume 2 \(Short Papers\)](#), pages
376 753–757.
- 377 Karthik Gopalakrishnan, Behnam Hedayatnia, Longshaokan Wang, Yang Liu, and Dilek Hakkani-Tur.
378 2020. Are neural open-domain dialog systems robust to speech recognition errors in the dialog
379 history? an empirical study. [arXiv preprint arXiv:2008.07683](#).
- 380 Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R Bowman, and
381 Noah A Smith. 2018. Annotation artifacts in natural language inference data. [arXiv preprint](#)
382 [arXiv:1803.02324](#).
- 383 Dilek Hakkani-Tür, Gökhan Tür, Asli Celikyilmaz, Yun-Nung Chen, Jianfeng Gao, Li Deng, and
384 Ye-Yi Wang. 2016. Multi-domain joint semantic frame parsing using bi-directional rnn-lstm. In
385 [Interspeech](#), pages 715–719.
- 386 Charles T Hemphill, John J Godfrey, and George R Doddington. 1990. The atis spoken language
387 systems pilot corpus. In [Speech and Natural Language: Proceedings of a Workshop Held at](#)
388 [Hidden Valley, Pennsylvania, June 24-27, 1990](#).

- 389 Yutai Hou, Sanyuan Chen, Wanxiang Che, Cheng Chen, and Ting Liu. 2020. C2c-genda:
390 Cluster-to-cluster generation for data augmentation of slot filling. [arXiv preprint](#)
391 [arXiv:2012.07004](#).
- 392 Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension
393 systems. [arXiv preprint arXiv:1707.07328](#).
- 394 Mihir Kale and Aditya Siddhant. 2021. Mixout: A simple yet effective data augmentation scheme for
395 slot-filling. In [Conversational Dialogue Systems for the Next Decade](#), pages 279–288. Springer.
- 396 Divyansh Kaushik, Eduard Hovy, and Zachary C Lipton. 2019. Learning the difference that makes a
397 difference with counterfactually-augmented data. [arXiv preprint arXiv:1909.12434](#).
- 398 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie
399 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, et al. 2021. Dynabench: Rethinking
400 benchmarking in nlp. [arXiv preprint arXiv:2104.14337](#).
- 401 Jitin Krishnan, Antonios Anastasopoulos, Hemant Purohit, and Huzefa Rangwala. 2021. Multilingual
402 code-switching for zero-shot cross-lingual intent prediction and slot filling. [arXiv preprint](#)
403 [arXiv:2103.07792](#).
- 404 Chen Li, Jinha Park, Hahyeon Kim, and Dimitrios Chrysostomou. 2021. How can i help you?
405 an intelligent virtual assistant for industrial robots. In [Companion of the 2021 ACM/IEEE](#)
406 [International Conference on Human-Robot Interaction](#), pages 220–224.
- 407 Shiyang Li, Semih Yavuz, Kazuma Hashimoto, Jia Li, Tong Niu, Nazneen Rajani, Xifeng Yan,
408 Yingbo Zhou, and Caiming Xiong. 2020. Coco: Controllable counterfactuals for evaluating
409 dialogue state trackers. [arXiv preprint arXiv:2010.12850](#).
- 410 Xiang Li, Haiyang Xue, Wei Chen, Yang Liu, Yang Feng, and Qun Liu. 2018. Improving the
411 robustness of speech translation. [arXiv preprint arXiv:1811.00728](#).
- 412 Hongyu Lin, Yaojie Lu, Jialong Tang, Xianpei Han, Le Sun, Zhicheng Wei, and Nicholas Jing Yuan.
413 2020. A rigorous study on named entity recognition: Can fine-tuning pretrained model lead to the
414 promised land? [arXiv preprint arXiv:2004.12126](#).
- 415 Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. 2019. Benchmarking
416 natural language understanding services for building conversational agents. [arXiv preprint](#)
417 [arXiv:1903.05566](#).
- 418 Samuel Louvan and Bernardo Magnini. 2020. Simple is better! lightweight data augmentation for
419 low resource slot filling and intent classification. [arXiv preprint arXiv:2009.03695](#).
- 420 Edward Ma. 2019. Nlp augmentation. <https://github.com/makcedward/nlpaug>.
- 421 R Thomas McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing
422 syntactic heuristics in natural language inference. [arXiv preprint arXiv:1902.01007](#).
- 423 Junghyun Min, R Thomas McCoy, Dipanjan Das, Emily Pitler, and Tal Linzen. 2020. Syntactic
424 data augmentation increases robustness to inference heuristics. In [Proceedings of the 58th Annual](#)
425 [Meeting of the Association for Computational Linguistics](#), pages 2339–2352.
- 426 Nafise Sadat Moosavi, Marcel de Boer, Prasetya Ajie Utama, and Iryna Gurevych. 2020. Improving
427 robustness by augmenting training sentences with predicate-argument structures. [arXiv preprint](#)
428 [arXiv:2010.12510](#).
- 429 Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018.
430 Stress test evaluation for natural language inference. [arXiv preprint arXiv:1806.00692](#).
- 431 Sebastian Padó, André Blessing, Nico Blokker, Erenay Dayanik, Sebastian Haunss, and Jonas Kuhn.
432 2019. Who sides with whom? towards computational construction of discourse networks for
433 political debates. In [Proceedings of the 57th Annual Meeting of the Association for Computational](#)
434 [Linguistics](#), pages 2841–2847.

- 435 Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme.
436 2018. Hypothesis only baselines in natural language inference. [arXiv preprint arXiv:1805.01042](#).
- 437 Libo Qin, Wanxiang Che, Yangming Li, Haoyang Wen, and Ting Liu. 2019. A stack-propagation
438 framework with token-level intent detection for spoken language understanding. [arXiv preprint](#)
439 [arXiv:1909.02188](#).
- 440 Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020.
441 Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In
442 [Proceedings of the AAAI Conference on Artificial Intelligence](#), volume 34, pages 8689–8696.
- 443 Evgeniia Razumovskaia, Goran Glavaš, Olga Majewska, Anna Korhonen, and Ivan Vulić. 2021.
444 Crossing the conversational chasm: A primer on multilingual task-oriented dialogue systems.
445 [arXiv preprint arXiv:2104.08570](#).
- 446 Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy:
447 Behavioral testing of nlp models with checklist. [arXiv preprint arXiv:2005.04118](#).
- 448 Elizabeth Shriberg. 2001. To ‘errrr’ is human: Ecology and acoustics of speech disfluencies. [Journal](#)
449 [of the International Phonetic Association](#), 31:153–169.
- 450 Edwin Simonnet, Sahar Ghannay, Nathalie Camelin, and Yannick Estève. 2018. Simulating asr errors
451 for training slu systems. In [Proceedings of the Eleventh International Conference on Language](#)
452 [Resources and Evaluation \(LREC 2018\)](#).
- 453 Yulia Tsvetkov, Florian Metze, and Chris Dyer. 2014. Augmenting translation models with simulated
454 acoustic confusions for improved spoken language translation. In [Proceedings of the 14th](#)
455 [Conference of the European Chapter of the Association for Computational Linguistics](#), pages
456 616–625.
- 457 Yusuf Uğurlu, Murat Karabulut, and İslam Mayda. 2020. A smart virtual assistant answering
458 questions about covid-19. In [2020 4th International Symposium on Multidisciplinary Studies and](#)
459 [Innovative Technologies \(ISMSIT\)](#), pages 1–6. IEEE.
- 460 Yu Wang, Yilin Shen, and Hongxia Jin. 2018. A bi-model based rnn semantic frame parsing model
461 for intent detection and slot filling. [arXiv preprint arXiv:1812.10235](#).
- 462 Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. Luke:
463 deep contextualized entity representations with entity-aware self-attention. [arXiv preprint](#)
464 [arXiv:2010.01057](#).
- 465 Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. 2019. Paws-x: A cross-lingual adversarial
466 dataset for paraphrase identification. [arXiv preprint arXiv:1908.11828](#).
- 467 Fan Yin, Quanyu Long, Tao Meng, and Kai-Wei Chang. 2020. On the robustness of language
468 encoders against grammatical errors. [arXiv preprint arXiv:2005.05683](#).
- 469 Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. Swag: A large-scale adversarial
470 dataset for grounded commonsense inference. [arXiv preprint arXiv:1808.05326](#).
- 471 Xiangji Zeng, Yunliang Li, Yuchen Zhai, and Yin Zhang. 2020. Counterfactual generator: A
472 weakly-supervised method for named entity recognition. In [Proceedings of the 2020 Conference](#)
473 [on Empirical Methods in Natural Language Processing \(EMNLP\)](#), pages 7270–7280.
- 474 Yuan Zhang, Jason Baldridge, and Luheng He. 2019. Paws: Paraphrase adversaries from word
475 scrambling. [arXiv preprint arXiv:1904.01130](#).

476 **Checklist**

- 477 1. For all authors...
- 478 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
479 contributions and scope? [Yes]
- 480 (b) Did you describe the limitations of your work? [Yes]
- 481 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 482 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
483 them? [Yes]
- 484 2. If you are including theoretical results...
- 485 (a) Did you state the full set of assumptions of all theoretical results? [No]
- 486 (b) Did you include complete proofs of all theoretical results? [No]
- 487 3. If you ran experiments...
- 488 (a) Did you include the code, data, and instructions needed to reproduce the main
489 experimental results (either in the supplemental material or as a URL)? [Yes]
- 490 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
491 were chosen)? [Yes]
- 492 (c) Did you report error bars (e.g., with respect to the random seed after running
493 experiments multiple times)? [No]
- 494 (d) Did you include the total amount of compute and the type of resources used (e.g., type
495 of GPUs, internal cluster, or cloud provider)? [Yes]
- 496 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 497 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 498 (b) Did you mention the license of the assets? [Yes]
- 499 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 500 (d) Did you discuss whether and how consent was obtained from people whose data you're
501 using/curating? [Yes]
- 502 (e) Did you discuss whether the data you are using/curating contains personally identifiable
503 information or offensive content? [Yes]
- 504 5. If you used crowdsourcing or conducted research with human subjects...
- 505 (a) Did you include the full text of instructions given to participants and screenshots, if
506 applicable? [Yes]
- 507 (b) Did you describe any potential participant risks, with links to Institutional Review
508 Board (IRB) approvals, if applicable? [N/A]
- 509 (c) Did you include the estimated hourly wage paid to participants and the total amount
510 spent on participant compensation? [Yes]