

# Appendix

## A Published and reproduced models

We reproduce the Stack-Prop+BERT and Bi-RNN models. The resulting trained models obtain similar results to the published, as shown in Appendix Table 1.

Test Set	ATIS		SNIPS		NLU-ED	
	Slot	Int.	Slot	Int.		
Stack-Prop+BERT						
Published	96.1	97.5	97.0	99.0	na	na
Reproduced	95.7	96.5	95.0	98.2	74.0	85.1
Bi-RNN						
Published	94.9	97.6	89.4*	97.1*	na	na
Reproduced	95.7	96.5	95.0	98.3	65.8	78.8

Table 1: Published and reproduced SF and ID results. The numbers with \* indicate that the scores were not published in the original Wang et al. Bi-RNN paper but in the Qin et al. Stack-Prop+BERT article.

## B External perturbation-based techniques

Most DA techniques focus on modifying data to obtain a semantically valid output. The NATURE operators are designed not only to have a semantically valid output but to maintain the same token-level labels as the original data. This small distinction makes a great difference to the end result and we show in Table 2 that the DA techniques are not sufficient to cancel out NATURE’s alterations. To this end, we apply standard DA strategies to the train and validation sets, re-train the model from scratch and illustrate their impact on the model’s generalization ability. We use common automatic DA strategies from the NLPaug library <sup>1</sup> that allow to easily relabel the augmented data using the original labels. We describe these strategies in Appendix Table B.

DA strategy name	Description	Example
<b>Keyboard Augmentation</b>	Simulates keyboard distance error.	<i>find a tv <b>seri</b>Ss called <b>armaRdvdon</b> summer</i>
<b>Spelling Augmentation</b>	Substitutes word according to spelling mistake dictionary.	<i><b>fine</b> a tv <b>serie</b> called <b>armageddon</b> summer</i>
<b>Synonym Augmentation</b>	Substitutes similar word according to WordNet/PPDB synonym.	<i>find a tv <b>set</b> series called <b>armageddon</b> summertime</i>
<b>Antonym Augmentation</b>	Substitutes opposite meaning word according to WordNet antonym.	<i><b>lose</b> a tv series called <b>armageddon</b> summer</i>
<b>TF-IDF Augmentation</b>	Uses the TF-IDF measure to find out how a word should be augmented.	<i>find tv series called <b>armageddon</b> <b>forms</b></i>
<b>Contextual Word Embeddings Augmentation</b>	Feeds surroundings word to BERT, DistilBERT, RoBERTa or XLNet language model to find out the most suitable word for augmentation.	<i>find a <b>second</b> series called <b>armageddon</b> <b>ii</b></i>

We apply the DA strategies exclusively to the train and validation sets, choosing 1 of the 6 DA functions at random and adding one output to the original dataset which results in a training and validation data twice as large as the original training and validation sets.

We have shown that state-of-the-art SF and ID models do suffer when small perturbations are introduced to the test data. We now run experiments on augmented data in order to test the

<sup>1</sup><https://github.com/makcedward/nlpaug>

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 2: 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.

models’ performances on larger and slightly more diverse train sets. Table 2 reports E2E scores of Stack-Prop+BERT <sup>2</sup> model when trained without (w/o) and with (w Aug) data-augmented train and validation sets. Similar to the results table in the main article, we evaluate the model on the Original, Rand, and Hard test sets of ATIS, SNIPS and NLU-ED while also reporting the unweighted average score.

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

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

## C Qualitative Evaluation

In Appendix Tables 1a and 1b We show the instructions and an excerpt of the sentences, as presented to the surveyed participants<sup>3</sup>.

Group	1													
Participant Id	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Experiment														
Slot	95.3	96.9	95.3	91.3	94.5	96.1	92.1	86.1	98.3	98.2	95.7	90.4	97.4	90.4
Intent	83.3	93.3	87.9	83.3	90.0	91.7	93.3	76.7	90.0	88.1	93.2	87.7	84.5	81.4
Control														
Fluency	4.9	5	4.8	4.6	4.9	4.7	4.5	4.2	4.3	5	5	4.9	4.8	4.2
Slot	89.5	89.5	100	94.7	100	94.7	89.5	94.7	100	100	100	100	94.7	89.5
Intent	91.7	100	100	100	100	100	91.7	91.7	100	100	100	90.9	100	100

Table 3: Survey results and statistics per participant. The average slot score and the average intent score appear as percentages, the average sentence fluency score appears as a scale from 1 to 5.

<sup>2</sup>Performances of the Bi-RNN model show very similar trends.

<sup>3</sup>We asked the participants to rate the fluency of each utterance (from 1 to 5) in order to average it over the control utterances. Allowing us to establish the annotator capacity of our volunteer participants. We expected this metric to reflect the high quality of the cherry-picked control utterances. As expected, our participants score remained between 4.2 and 5 out of 5.

### Nature Huawei Human Evaluation

This is a survey to evaluate the quality of the commands/questions/responses a human user gives (orally) to a virtual assistant (such as Siri, Alexa, Google Home, etc.). The sentences do not contain upper-case (hi bob, i think it's 5 am), they might not contain any punctuation (where is the restaurant) and some concepts are joined with underscores (meryl\_streep won the 1982\_academy\_award).

For each sentence, you will be asked to evaluate 3 things:

- Sentence quality: Give a rating scale (from 1 to 5), 1 being 'gibberish', 3 being 'pretty much understandable' and 5 being 'natural'.
- Labeling quality: Evaluate with REASONABLE or UNREASONABLE all labeled tokens (and only the labeled ones).
- Classification quality: Evaluate with REASONABLE or UNREASONABLE the classification of the sentence to the type of order given.

Examples:

Sentence: clean the floor please  
Sentence quality: What is the quality rating of the sentence 'clean the floor please'? 5  
Labels: - - house\_place -  
Labeling quality: Is it reasonable to qualify 'floor' as a 'house\_place'? REASONABLE  
Classification quality: Is it reasonable to classify the sentence 'clean the floor please' as a 'floor\_cleaning' type of order? REASONABLE

Sentence: that's perfect, much appreciated.  
Sentence quality: What is the quality rating of the sentence 'that's perfect, much appreciated.'? 4  
Labels: - - - - -  
Labeling quality: general\_praise  
Classification quality: Is it reasonable to classify the sentence 'that's perfect, much appreciated.' as a 'general\_praise' type of order? REASONABLE

Sentence: google  
Sentence quality: What is the quality rating of the sentence 'google'? 2  
Labels: -  
Labeling quality: qa\_maths  
Classification quality: Is it reasonable to classify the sentence 'google' as a 'qa\_maths' (question-answering math) type of order? UNREASONABLE

Example-1) need to actually see mother\_joan\_of\_the\_angels in\_one\_second.

1 gibberish 2 3 understandable but ungrammatical 4 5 natural

Fluency ☐ ☐ ☐ ☐ ☐

Slot-1) need to actually see mother\_joan\_of\_the\_angels in\_one\_second.

Reasonable Unreasonable

mother\_joan\_of\_the\_angels ☐ ☐

in\_one\_second-timeframe ☐ ☐

Intent-1) need to actually see mother\_joan\_of\_the\_angels in\_one\_second.

Reasonable Unreasonable

SearchScreeningEvent ☐ ☐

Example-2) add on track to the another\_glass\_playlist.

1 gibberish 2 3 understandable but ungrammatical 4 5 natural

Fluency ☐ ☐ ☐ ☐ ☐

Slot-2) add on track to the another\_glass\_playlist.

Reasonable Unreasonable

track-music\_item ☐ ☐

another\_glass\_playlist ☐ ☐

Intent-2) add on track to the another\_glass\_playlist.

Reasonable Unreasonable

AddToPlaylist ☐ ☐

Example-77) please play me a popular little from 1984.

1 gibberish 2 3 understandable but ungrammatical 4 5 natural

Fluency ☐ ☐ ☐ ☐ ☐

Slot-77) please play me a popular little from 1984.

Reasonable Unreasonable

popular-sort ☐ ☐

little-music\_item ☐ ☐

1984-year ☐ ☐

Intent-77) please play me a popular little from 1984.

Reasonable Unreasonable

PlayMusic ☐ ☐

Example-78) whats actually needed to make pizza.

1 gibberish 2 3 understandable but ungrammatical 4 5 natural

Fluency ☐ ☐ ☐ ☐ ☐

Slot-78) whats actually needed to make pizza.

Reasonable Unreasonable

pizza-food\_type ☐ ☐

Intent-78) whats actually needed to make pizza.

Reasonable Unreasonable

cooking\_recipe ☐ ☐

(a) Print-screen of the survey instructions.

(b) Print-screen excerpts of the survey.

## 38 D Quantitative Evaluation

39 In the Appendix Figure2 we show a more concise illustration of the quantitative experiments' results  
40 than Table 7. Appendix Figure2 shows the E2E score averaged between the benchmarks (ATIS,  
41 SNIPS, NLU-ED) and between the two models (Stack-Prop+BERT and Bi-RNN).

## 42 E Complete table of NATURE operators applied to ATIS, SNIPS and 43 NLU-ED

44 In the Appendix Table 4 we present all obtained scores ran on 2 models trained on the original train  
45 and validation sets of ATIS, SNIPS and NLU-ED and evaluated on the original, random and hard  
46 altered test sets.

Test Set	ATIS			SNIPS			NLU-ED		
	Slot (F1)	Intent (Acc)	E2E (Acc)	Slot (F1)	Intent (Acc)	E2E (Acc)	Slot (F1)	Intent (Acc)	E2E (Acc)
Stack-Prop+BERT									
Original	95.7	96.5	86.2	95.0	98.3	87.9	74.0	85.1	67.8
Random	91.3	95.0	66.5	83.4	96.1	53.8	67.4	76.1	56.8
	$\pm 0.1$	$\pm 0.3$	$\pm 1.0$	$\pm 0.5$	$\pm 0.3$	$\pm 3.2$	$\pm 0.1$	$\pm 0.2$	$\pm 0.2$
Hard	82.3	90.7	34.9	70.6	95.3	12.9	55.5	62.7	38.9
Bi-RNN									
Original	94.7	97.6	84.3	88.9	97.6	77.3	65.9	82.1	61.9
Random	89.9	94.3	61.8	75.6	94.1	39.0	60.6	70.8	50.1
	$\pm 0.1$	$\pm 0.1$	$\pm 1.6$	$\pm 0.5$	$\pm 0.1$	$\pm 2.5$	$\pm 0.4$	$\pm 0.4$	$\pm 0.3$
Hard	79.9	92.0	27.6	62.4	92.9	7.0	49.6	58.8	34.5

Table 4: Stack-Prop+BERT and Bi-RNN performances for ATIS, SNIPS and NLU-ED. We report F1 slot filling, accuracy for intent detection and end-to-end accuracy overall. The reported scores of the Random altered test set are a mean of 10 random distribution of processes and is accompanied by the variance score.

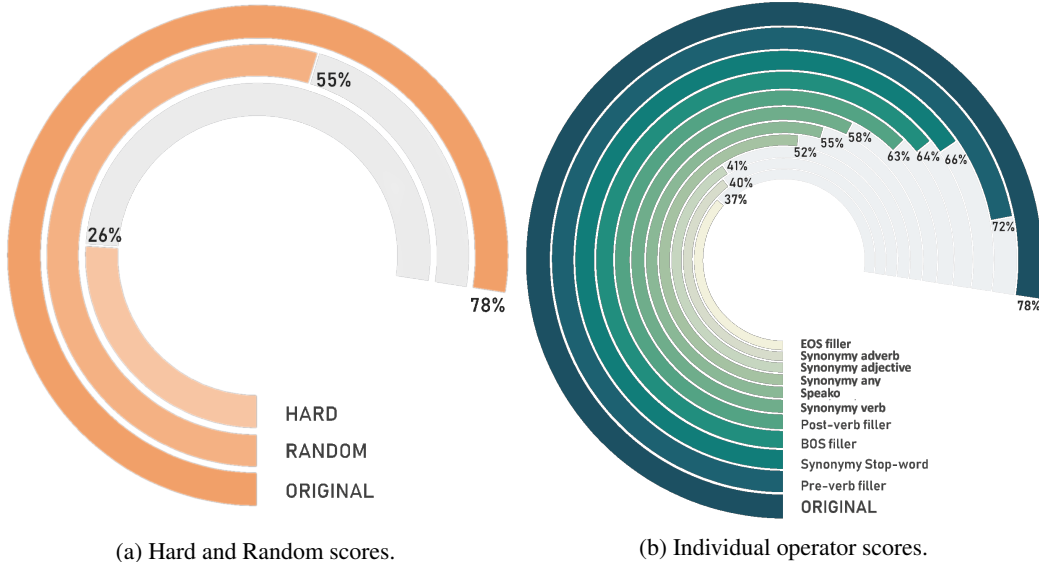


Figure 2: Unweighted average End-to-End score performances averaged between benchmarks (ATIS, SNIPS, and NLU-ED) and models (Stack-Prop+BERT and Bi-RNN). The models were trained using their original train and validation sets and evaluated on altered test sets. Figure 2a shows the scores for the Random and Hard evaluation sets while Figure 2b shows the scores on 10 evaluation sets, each perturbed with a single NATURE operator, where one operator is applied once to each utterance of the evaluation set

## 47 F Manual analysis of utterance weight

48 To better understand the underlying processes of the state-of-the-art models, we use the LIME tool<sup>4</sup>  
 49 to produce and analyze multiple self-attention weight heat-maps. This process allows us to better  
 50 understand what tokens the models focus best to make their prediction. In Figure 3 we show a  
 51 representative excerpt heat-maps for wrongly predicted sentences (for both SF and ID). One for the  
 52 unchanged SNIPS test set and one for each type of operator. At first sight, we notice that the attention  
 53 is quite evenly distributed among all tokens in the sentence. However, if we carefully examine the  
 54 small differences between them, we observe a tendency to often focus more heavily on verbs, nouns,  
 55 certain types of stop words (such as *"the"*) as well as tokens appearing at the extremities of the  
 56 utterance (although it might not be immediately evident in these small samples of small utterances).  
 57 It also shows that higher attention is given to verbs and certain stop words at the end of the sentence.  
 58 This is evident in all Figures but particularly in Figure 3b, where we can see high attention on  
 59 non-frequent tokens (for the benchmark), such as *"if"* or *"?"*. After more careful analysis, we observe  
 60 in Figure 3b that the attention is often high for the added filler. This is not the case in Figure 3c,  
 61 where the attention of the altered synonym is usually low if it doesn't replace a noun. As for the  
 62 Figure 3d, if we take for example the utterance *what time will paris by night aired*, we observe that  
 63 just as for the original utterance (and the Synonymy Adjective-altered) the self attention is just as  
 64 high in the tokens *will*, *paris* and *aired* but it also introduces a high weight on the Speako altered  
 65 token *want* → *wnt*, which doesn't appear in the original utterance.

## 66 G Complete NATURE operators applied to Data Augmented versions of 67 ATIS, SNIPS and NLU-ED

68 In the Appendix Table 6 we present all obtained scores ran on 2 models trained on a Data Augmented  
 69 version of ATIS, SNIPS and NLU-ED.

<sup>4</sup><https://github.com/marcotcr/lime>

find on new series  
what time will start by night area  
in how soon will start of series  
need to see mother hour of the night in one second  
play the new noise theology e  
want to watch supernatural the season seven of animal

(a) Heat-map of original SNIPS utterances.

find on new series  
music return back to life into winter music  
what will start by night area  
can i get the more crush showings  
can you find the series  
in the find king of hearts

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

find on new series  
can i get the butterfly crush showings you be bad  
what is there the girl cooper foundation you don't mind  
in one hour find king of hearts you don't mind  
i need a series in uruguni in 212 days when i challenge you please  
need to see mother hour of the night in one second we understand each other

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

find on new series  
what is the relation of our portal is  
what time will start by night area  
in the series find the of series  
can you find the trainor for phineas rebus  
what time series the years are running

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

Figure 3: 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 LIME weight.

Original: find a tv series called armageddon summer			
NATURE		DA	
BOS Filler	yeah so find a tv series called armageddon summer	Keyb.	find a tv seriSs called armaRdvdon summer
PreV Filler	basically find a tv series called armageddon summer	Spell.	fine a tv serie called armageddon summer
PosV Filler	find you know a tv series called armageddon summer	Syn.	find a tv set series called armageddon summertime
EOS Filler	find a tv series called armageddon summer if it pleases mi liege	Ant.	lose a tv series called armageddon summer
Syn. V.	finds a tv series called armageddon summer	TF IDF	find tv series called armageddon forms
Syn. Adj.	find a tv series called last summer	Ctxt. WE.	find a second series called armageddon ii
Syn. Adv.	find a another series called armageddon summer		
Syn. SW	find and tv series called armageddon summer		
Speako	find a tv serie called armageddon summer		

Table 5: Nature and DA candidates for the same utterance.

Test Set	ATIS			SNIPS			NLU-ED		
	Slot (F1)	Intent (Acc)	E2E (Acc)	Slot (F1)	Intent (Acc)	E2E (Acc)	Slot (F1)	Intent (Acc)	E2E (Acc)
Stack-Prop+BERT									
Original	94.7	95.7	83.3	93.8	97.7	85.3	72.4	83.8	66.2
Random	91.7	94.3	69.2	85.7	96.0	64.4	67.3	75.6	56.7
	$\pm 0.0$	$\pm 0.1$	$\pm 0.9$	$\pm 0.2$	$\pm 0.4$	$\pm 1.5$	$\pm 0.2$	$\pm 0.1$	$\pm 0.2$
Hard	87.2	91.0	54.0	72.7	95.1	27.1	55.3	64.0	40.7
Bi-RNN									
Original	93.7	96.9	81.8	86.2	97.6	69.7	66.3	82.5	61.8
Random	90.3	93.9	65.6	77.4	95.3	48.2	61.2	73.4	51.8
Random	$\pm 0.1$	$\pm 0.2$	$\pm 1.1$	$\pm 0.3$	$\pm 0.2$	$\pm 1.8$	$\pm 0.1$	$\pm 0.2$	$\pm 0.2$
Hard	83.2	92.8	43.0	65.0	94.1	19.1	62.1	50.2	38.6

Table 6: Stack-Prop+BERT and Bi-RNN performances for ATIS, SNIPS and NLU-ED using data augmentation on the train and validation sets. We report F1 slot filling, accuracy for intent detection and end-to-end accuracy overall. The reported scores of the Random altered test set are a mean of 10 random distribution of processes and is accompanied by the variance score.